

# High-probability sample complexities for policy evaluation with linear function approximation

Gen Li, *Member, IEEE*, Weichen Wu, Yuejie Chi, *Fellow, IEEE*, Cong Ma, *Member, IEEE*,  
Alessandro Rinaldo, Yuting Wei, *Member, IEEE*

**Abstract**—This paper is concerned with the problem of policy evaluation with linear function approximation in discounted infinite horizon Markov decision processes. We investigate the sample complexities required to guarantee a predefined estimation error of the best linear coefficients for two widely-used policy evaluation algorithms: the temporal difference (TD) learning algorithm and the two-timescale linear TD with gradient correction (TDC) algorithm. In both the on-policy setting, where observations are generated from the target policy, and the off-policy setting, where samples are drawn from a behavior policy potentially different from the target policy, we establish the first sample complexity bound with high-probability convergence guarantee that attains the optimal dependence on the tolerance level. We also exhibit an explicit dependence on problem-related quantities, and show in the on-policy setting that our upper bound matches the minimax lower bound on crucial problem parameters, including the choice of the feature map and the problem dimension.

**Index Terms**—policy evaluation, temporal difference learning, two-timescale stochastic approximation, minimax optimal, function approximation

## I. INTRODUCTION

POLICY evaluation plays a critical role in many scientific and engineering applications in which practitioners aim to evaluate the performance of a target strategy based on either sequentially collected or a batch of offline data samples [1], [2], [3], [4]. For example, in clinical trials [3], real-time data acquisition might be expensive and risky; it is thus of essential value if historical data can be analyzed and information can be transferred to new tasks. While in other applications, such

as mobile health [5], it is practical to implement the desired policy and collect its feedback in a timely manner.

Mathematically, Markov decision processes (MDPs) provide a general framework to design policy evaluation methods in dynamic settings; reinforcement learning (RL) is often modeled using MDPs when the exact model configuration is not available [6], [7]. In this framework, a target policy is assessed through its corresponding value function. In practice, evaluating value functions often require an overwhelming number of samples due to the large dimensionality of the underlying state space. For this reason, RL methods are normally concerned with some form of function approximation. Dating back to the seminal work of [8], there has been an extensive line of works that consider different types of function approximation, including linear function approximation [9], [10], reproducing kernel Hilbert space [11], [12], deep neural networks [13], [14] or function approximation on the model itself (see, e.g. [15], [16], [17]), with a focus on improving the sample efficiency of RL algorithms.

*a) Two settings: on-policy vs. off-policy.*: The main goal of this paper is to provide sharp statistical guarantees of policy evaluation algorithms with linear function approximation in two different settings. As the aforementioned examples already indicated, there are typically two different types of data-generating mechanisms to consider: the *on-policy* setting when we have access to the outcomes of the target policy and the *off-policy* setting, in which the only available data are generated from a behavior policy that is potentially different from the target policy.

In the on-policy setting, temporal difference (TD) learning is arguably the most popular algorithm [18] for policy evaluation in RL practice, partly because it is easy to implement and lends itself well to function approximations. As a model-free algorithm, TD learning processes data in an online manner without explicitly modeling the environment and is, therefore, memory efficient. While the asymptotic convergence of TD with linear function approximation has been known since [8], the finite-sample minimax optimality of TD has been established only recently for the tabular MDP [19]. For TD learning with linear function approximation, several recent contributions have produced new non-asymptotic analyses and insights (e.g. [9], [20], [21], [22]), which partially unveil impacts of both the tolerance level and various problem-related parameters on its sample efficiency. However, minimax-optimal dependence on the tolerance level (i.e. target level of estimation accuracy) is only established in expectation instead

The first two authors contributed equally.

The work of Gen Li was supported in part by the CUHK Direct Grant for Research. The work of Weichen Wu and Alessandro Rinaldo were supported in part by the NIH Grant R01 NS121913. The work of Yuejie Chi was supported in part by the grants ONR N00014-19-1-2404 and NSF CCF-2106778. The work of Yuting Wei was supported in part by the NSF grants CCF-2106778, DMS-2147546/2015447 and NSF CAREER award DMS-2143215 and Google Research Scholar Award. (Correspondent author: Yuting Wei.)

Gen Li is with the Department of Statistics, The Chinese University of Hong Kong, Hong Kong.

Weichen Wu is with the Department of Statistics and Data Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

Yuejie Chi is with the Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

Cong Ma is with the Department of Statistics, University of Chicago, Chicago, IL 60637, USA.

Alessandro Rinaldo is with the Department of Statistics and Data Science, University of Texas at Austin, Austin, UT 78712, USA.

Yuting Wei is with the Department of Statistics and Data Science, The Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA; email: ytwei@wharton.upenn.edu

of with high probability; furthermore, the optimal dependence on problem-related parameters, such as the size of the state space and the effective horizon, still remains unsettled, and it is unclear whether existing sample complexity bounds can be further improved. Failing to understand these questions, however, casts doubt on whether TD with linear function approximation is statistically efficient in practice, and brings difficulties to performing statistical inference based on TD estimators. In this paper, we seek to answer these questions by providing tighter characterizations of the performance of TD with linear function approximation.

In the off-policy setting, it is known that the error of TD learning with linear function approximation may diverge to infinity [23]. In order to address this issue, [24] proposed a now popular alternative with two-timescale learning rates, called the linear TD with gradient correction (TDC) algorithm, which enjoys convergence guarantees in the off-policy case. In terms of finite-sample guarantees, although a number of recent efforts (see, e.g. [25], [26], [27], [28], [29], [30]) tried to characterize the statistical performance of TDC for both *i.i.d.* and Markovian data, they remain inadequate in providing either a convergence guarantee with high-probability, an explicit dependence on salient problem parameters, or a sharp dependence on the sample size. The challenge lies in dealing with the statistical dependence between two separate iterate sequences at different timescales. To tackle this challenge, it calls for a new analysis framework for the TDC algorithm.

#### A. Our main contributions

This paper is concerned with evaluating the performance of a given target policy  $\pi$  in an infinite-horizon  $\gamma$ -discounted MDP with a finite but large number of states. The goal is to learn the best linear approximation of the value function in a pre-specified feature space given *i.i.d.* transition pairs drawn from the stationary distribution. In the on-policy setting, we focus on the TD learning algorithm; in the off-policy setting, we shift gear to the TDC learning algorithm. We summarize our main contributions as follows, with their exact statements and consequences postponed to later sections.

- Via a careful analysis of TD learning with Polyak-Ruppert averaging, we show that, in the on-policy setting, a number of samples of order

$$\tilde{O}\left(\frac{\max_s\{\phi(s)^\top \Sigma^{-1} \phi(s)\}(1 + \|\theta^*\|_{\Sigma}^2)}{(1 - \gamma)^2 \varepsilon^2}\right)$$

is sufficient to achieve an accuracy level (estimation error) of  $\varepsilon > 0$ , with high probability. Here,  $\phi(s) \in \mathbb{R}^d$  indicates the linear feature vector for the state  $s$  in the state space  $\mathcal{S}$ ,  $\theta^*$  is the best linear approximation coefficient of the value function, and  $\Sigma$  corresponds to the feature covariance matrix weighted by the stationary distribution. See Section II for the definitions of these parameters. Compared to prior work by [9] and [21], our sample complexity bound can be tighter by a factor of  $\text{cond}(\Sigma)$  which can be as large as  $|\mathcal{S}|$  (the cardinality of the state space). Our result also controls  $\varepsilon$ -convergence with high probability that matches the minimax-optimal

dependence on the tolerance level  $\varepsilon$  with lowest burn-in cost. To assess the tightness of this upper bound, we provide a minimax lower bound in Section III-C, which certifies the optimal dependence of our bound on both the tolerance level  $\varepsilon$  and problem-related parameters  $\Sigma$  and  $\theta^*$ .

- In the off-policy setting, we establish a sample complexity bound for the TDC algorithm of order

$$\tilde{O}\left(\frac{\rho_{\max}^7}{\lambda_1^4 \lambda_2^3} \frac{\|\tilde{\Sigma}\|^2}{\varepsilon^2} (1 + \|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2)\right),$$

where  $\tilde{\theta}$  corresponds to the best linear approximation coefficient of the value function in the off-policy setting,  $\tilde{\Sigma}$  is the feature covariance matrix under the behavior policy,  $\rho_{\max}$  denotes the largest importance sampling ratio measuring the discrepancy between the target policy and the behavior policy, and lastly,  $\lambda_1$  and  $\lambda_2$  denote the smallest eigenvalues of some problem-dependent matrices. Details about these constants are deferred to Section IV. To the best of our knowledge, our bound is the first one to control  $\varepsilon$ -convergence with high probability that matches the minimax-optimal dependence on the tolerance level  $\varepsilon$ . At the same time, our sample complexity bound also provides an explicit dependence on the salient parameters.

Comparisons of our results to existing bounds and relevant commentary can be found in Table I and II.

#### B. Other related works

In this section, we review several recent lines of works and provide a broader context of the current paper.

##### a) Finite-sample guarantees for policy evaluation.:

Classical analyses of policy evaluation algorithms have mainly focused on providing asymptotic guarantees given a fixed model [8], [32]. New tools developed in high-dimensional statistics and probability allow for a fine-grained understanding of these algorithms especially from a finite-sample and finite-time perspective. As argued in this paper, understanding how statistical errors depend on the effective horizon, dimension of the problem and the number of samples, is essential as it provides important insights on how these RL algorithms perform in practice. A highly incomplete list of prior art includes [33], [34], [21], [9], [20], [22], [35] with a focus on the non-asymptotic analyses for model-free algorithms, and [36], [37], [38], [39] which derive non-asymptotic bounds for model-based algorithms.

##### b) Stochastic approximation.:

The idea of stochastic approximation (SA) [40], [41] lies at the core of the TD and TDC learning algorithms considered in this paper. With the intention of solving a deterministic fixed-point equation, SA methods perform stochastic updates based on approximations of the current residual. The asymptotic theory of SA methods are relatively well-developed, where SA iterates provably track the trajectory of a limiting ordinary differential equation [42], [43] and with properly decaying step sizes, the Polyak-Ruppert averaged iterates asymptotically follow the central limit theorem. Recently, non-asymptotic results have also been

paper	algorithm	stepsize	sample complexity	error control
[9]	TD	$\eta_t \asymp t^{-1}$	$O\left(\frac{\ \Sigma^{-2}\  \ \Sigma\ }{(1-\gamma)^4 \varepsilon^2}\right)$	in expectation
[21]	TD	$\eta_t \asymp T^{-1}$	$O\left(\frac{\ \Sigma^{-2}\  \ \Sigma\ }{(1-\gamma)^4 \varepsilon^2}\right)$	in expectation
[20]	TD	$\eta_t = t^{-1}$	$O\left(\frac{1}{\varepsilon^{\max\{2, 1+\frac{1}{\lambda}\}}}\right), \lambda \in (0, \lambda_{\min}(\mathbf{A}))$	w. high-prob
[31]	Averaged TD	$\eta_t \asymp T^{-1/2}$	$O\left(\frac{\ \Sigma^{-1}\ }{(1-\gamma)^4 \varepsilon^2} \vee \frac{\ \Sigma^2\  \ \Sigma^{-4}\ }{(1-\gamma)^6}\right)$	w. high-prob
<b>This work</b>	Averaged TD	$\eta_t = \eta$	$O\left(\frac{\ \Sigma^{-1}\ }{(1-\gamma)^4 \varepsilon^2} \vee \frac{\ \Sigma^2\  \ \Sigma^{-3}\ }{(1-\gamma)^4}\right)$	w. high-prob

**TABLE I.** Comparisons with prior results (up to logarithmic terms) in finding an  $\varepsilon$ -optimal solution using TD learning. Using the Polyak-Ruppert averaging, our high-probability sample complexity bound improves upon previous works in the dependence on the tolerance level  $\varepsilon$  and problem-related parameters.

obtained for SA for different problems especially in the RL setting; see [44], [45], [22], [46] and references therein. The TDC algorithm is a special case of two-timescale linear SA, whose convergence rates have also been investigated in [28], [47], [48], [30], among others.

c) *Off-policy learning.*: Policy evaluation in the off-policy setting is closely related to offline or batch RL, which aims to learn purely based on historical data without actively exploring the environment. The main challenge here lies in the discrepancy between the behavior policy and the target or optimal policy. One natural approach is to use importance sampling (IS) in order to form an unbiased estimator of the target policy [49], and various different techniques have been applied to reduce the high variance of IS (see, e.g. [50], [51], [52], [53], [54], [55]). Non-asymptotic guarantees are also provided for off-policy evaluation using a fitted  $Q$ -iteration approach under linear function approximation in [56]. A recent line of works also considered finding the optimal policy using batch datasets [57], [58], [59], [60], [61].

### C. Notation

Throughout this paper, we denote by  $\Delta(\mathcal{S})$  (resp.  $\Delta(\mathcal{A})$ ) the probability simplex over the finite set  $\mathcal{S}$  (resp.  $\mathcal{A}$ ). For any positive integer  $n$ , we use  $[n]$  to denote the set of positive integers that are no larger than  $n$ :  $[n] = \{1, 2, \dots, n\}$ . When a function is applied to a vector, it should be understood as being applied in a component-wise fashion; for example,  $\sqrt{\mathbf{z}} := [\sqrt{z_i}]_{1 \leq i \leq n}$  and  $|\mathbf{z}| := [|z_i|]_{1 \leq i \leq n}$ . For any vectors  $\mathbf{z} = [z_i]_{1 \leq i \leq n}$  and  $\mathbf{w} = [w_i]_{1 \leq i \leq n}$ , the notation  $\mathbf{z} \geq \mathbf{w}$  (resp.  $\mathbf{z} \leq \mathbf{w}$ ) stands for  $z_i \geq w_i$  (resp.  $z_i \leq w_i$ ) for all  $1 \leq i \leq n$ . Additionally, we write  $\mathbf{1}$  for the all-one vector,  $\mathbf{I}$  for the identity matrix, and  $\mathbf{1}\{\cdot\}$  for the indicator function.

For any matrix  $\mathbf{P} = [P_{ij}]$ , we denote  $\|\mathbf{P}\|_1 := \max_i \sum_j |P_{ij}|$ . Given a symmetric positive definite matrix  $\mathbf{D}$ , define the inner product  $\langle \cdot, \cdot \rangle_{\mathbf{D}}$  as  $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{D}} = \mathbf{u}^\top \mathbf{D} \mathbf{v}$  and the associated norm  $\|\mathbf{v}\|_{\mathbf{D}} = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle_{\mathbf{D}}}$ . For any matrix  $\mathbf{M}$ , we use  $\|\mathbf{M}\|$  to denote its operator norm (i.e. the largest singular value), if not specified otherwise. Throughout this paper, we

use  $c, c_0, c_1, C, \dots$  to denote universal constants that do not depend either on the parameters of the MDP or the target levels  $(\varepsilon, \delta)$ ; their exact values may change from line to line. Given two sequences,  $\{f_t\}_{t \geq 0}$  and  $\{g_t\}_{t \geq 0}$ , we write  $f_t \lesssim g_t$  (resp.  $f_t \gtrsim g_t$ ) or  $f_t = O(g_t)$  (resp.  $g_t = O(f_t)$ ) if there exists some universal constant  $c_1 > 0$ , such that  $f_t \leq c_1 g_t$  (resp.  $f_t \geq c_1 g_t$ ). If both  $f = O(g)$  and  $g = O(f)$  hold simultaneously, we write  $f \asymp g$  or  $f = \Theta(g)$ . We adopt the notation  $f = \tilde{O}(g)$  to indicate  $f = O(g)$  up to logarithmic factors in  $g$ . For any symmetric matrix  $\mathbf{X}$ , we use  $\lambda_{\min}(\mathbf{X})$  to denote its smallest eigenvalue.

## II. PROBLEM FORMULATION

### A. Model and settings

a) *Markov decision process*: Consider an infinite-horizon MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$  with discounted rewards, where  $\mathcal{S}$  and  $\mathcal{A}$  denote respectively the (finite) state space and action space, and  $\gamma \in (0, 1)$  indicates the discount factor [7]. The probability transition kernel of the MDP is given by  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \mapsto \Delta(\mathcal{S})$ , where for each state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,  $\mathcal{P}(\cdot | s, a) \in \Delta(\mathcal{S})$  denotes the transition probability distribution from state  $s$  when action  $a$  is executed. The reward function is represented by the function  $r : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ , where  $r(s, a)$  denotes the immediate reward from state  $s$  when action  $a$  is taken; for simplicity, we assume throughout that all immediate rewards lie within  $[0, 1]$ .

A policy  $\pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$  is an action selection rule that maps a state to a distribution over the set of actions; in particular, it is said to be stationary if it is time-invariant. The value function  $V^\pi : \mathcal{S} \mapsto \mathbb{R}$  is used to measure the quality of a policy  $\pi$ , defined as

$$\forall s \in \mathcal{S} : \quad V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right], \quad (1)$$

which is the expected discounted cumulative reward received by following the policy  $\pi$  under the MDP  $\mathcal{M}$  when initialized at state  $s_0 = s$ . Here,  $a_t \sim \pi(\cdot | s_t)$  and  $s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)$

paper	algorithm	stepsize	sample complexity	error control
[28]	<b>Projected</b> TDC	$\alpha_t = t^{-\alpha}, \beta_t = t^{-\beta}$	$O\left(\frac{1}{\varepsilon^{2\alpha}}\right), \alpha < 1$	w. high-prob
[29]	TDC	$\alpha_t, \beta_t \asymp \frac{1}{T}$	$O\left(\frac{1}{\varepsilon^2}\right)$	in expectation
[26]	<b>Batched</b> TDC	$\alpha_t = \alpha, \beta_t = \beta$	$O\left(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon}\right)$	in expectation
<b>This work</b>	TDC	$\alpha_t, \beta_t \asymp \frac{1}{T}$	$O\left(\frac{1}{\varepsilon^2}\right)$	w. high-prob

**TABLE II.** Comparisons with prior results (up to logarithmic terms) in finding an  $\varepsilon$ -optimal solution using TDC learning. We omit dependence on problem-related parameters in this table. Our sample complexity bound for TDC is the first to achieve high-probability convergence guarantee with non-varying stepsizes and without using projection steps or batched updates; in the mean time, we also provide explicit dependence on problem-related parameters.

for all  $t \geq 0$ . It can be easily verified that  $0 \leq V^\pi(s) \leq \frac{1}{1-\gamma}$  for any  $\pi$ .

For a given policy  $\pi$ , we can define the reward function of every state  $s \in \mathcal{S}$  as the expected reward for  $(s, a)$  when  $a$  is chosen according to  $\pi$ :

$$r(s) = \mathbb{E}_{a \sim \pi(\cdot|s)}[r(s, a)]. \quad (2)$$

For simplicity, we introduce the vector notation for the reward function  $\mathbf{r} := [r(s)]_{1 \leq s \leq |\mathcal{S}|} \in \mathbb{R}^{|\mathcal{S}|}$ , and the value function  $\mathbf{V}^\pi = [V^\pi(s)]_{1 \leq s \leq |\mathcal{S}|} \in \mathbb{R}^{|\mathcal{S}|}$ . We can also define the transition matrix  $\mathbf{P}^\pi$  for this given policy  $\pi$ , such that its  $(i, j)$  element represents the probability that state  $i$  is transited to state  $j$  under the policy  $\pi$ ; formally,

$$P_{ij}^\pi = \sum_{a \in \mathcal{A}} \mathcal{P}(s_{t+1} = j \mid s_t = i, a_t = a) \pi(a_t = a \mid s_t = i). \quad (3)$$

We denote by  $\mu$  the stationary distribution corresponding to the Markov chain when the transition follows  $\mathbf{P}^\pi$ , which we assume to be well-defined, and introduce the vector notation  $\boldsymbol{\mu} := [\mu(s)]_{1 \leq s \leq |\mathcal{S}|} \in \mathbb{R}^{|\mathcal{S}|}$ .

*b) Linear approximation for the value function:* As discussed previously, it is often infeasible to collect a number of samples that scales with the ambient dimension  $|\mathcal{S}|$ . This motivates the search for lower dimension approximation of the value function, of which linear approximation emerges as a convenient option. Mathematically, for  $\boldsymbol{\theta} \in \mathbb{R}^d$ , define  $V_{\boldsymbol{\theta}}(s)$  as

$$\forall s \in \mathcal{S}: \quad V_{\boldsymbol{\theta}}(s) = \phi(s)^\top \boldsymbol{\theta},$$

where  $\phi(s) \in \mathbb{R}^d$  is the feature vector associated with state  $s \in \mathcal{S}$ , with  $d \leq |\mathcal{S}|$ . The vector  $\boldsymbol{\theta}$  of linear coefficients is shared across states.

Using matrix notation, we let

$$\boldsymbol{\Phi} := [\phi(1), \phi(2), \dots, \phi(|\mathcal{S}|)]^\top \in \mathbb{R}^{|\mathcal{S}| \times d}, \quad (4)$$

be the feature matrix that concatenates the feature vectors for all states and  $\mathbf{V}_{\boldsymbol{\theta}} = [V_{\boldsymbol{\theta}}(s)]_{s \in \mathcal{S}} \in \mathbb{R}^{|\mathcal{S}|}$  be the linear

approximation vector to the value function. It follows that

$$\mathbf{V}_{\boldsymbol{\theta}} = \boldsymbol{\Phi} \boldsymbol{\theta}.$$

We impose the following mild assumption on the feature vectors.

**Assumption 1.** *The columns of  $\boldsymbol{\Phi}$  are linearly independent with Euclidean norm uniformly bounded by one, i.e.  $\max_{s \in \mathcal{S}} \|\phi(s)\|_2 \leq 1$ .*

#### B. Policy evaluation with linear approximation

*a) On-policy evaluation with linear approximation:* The task of policy evaluation is to measure the value function  $V^\pi(s)$  for every  $s \in \mathcal{S}$  (see definition (1)) given a policy  $\pi$  of interest. In the **on-policy** setting, data samples are collected while the policy  $\pi$  is executed and a sequence of samples are obtained

$$\{(s_0, a_0, r_0), \dots, (s_T, a_T, r_T)\},$$

where  $a_t \sim \pi(\cdot \mid s_t)$ ,  $r_t = r(s_t, a_t)$ .

In this setting, in order to find the best linear approximation to  $\mathbf{V}^\pi$ , we find it helpful to first introduce some shorthand notation. First, given the stationary distribution  $\mu$  for  $\mathbf{P}^\pi$ , we let

$$\mathbf{D}_\mu = \text{diag}(\mu(1), \mu(2), \dots, \mu(|\mathcal{S}|)) \quad (5)$$

and denote with

$$\boldsymbol{\Sigma} := \boldsymbol{\Phi}^\top \mathbf{D}_\mu \boldsymbol{\Phi} = \mathbb{E}_{s \sim \mu} [\phi(s) \phi(s)^\top] \in \mathbb{R}^{d \times d} \quad (6)$$

the feature covariance matrix with respect to this stationary distribution.

The best linear approximation coefficients,  $\boldsymbol{\theta}^*$ , is defined as the unique solution to the following projected Bellman equation [8]

$$\boldsymbol{\Phi} \boldsymbol{\theta} = \Pi_{\mathbf{D}_\mu} \mathcal{T}^\pi (\boldsymbol{\Phi} \boldsymbol{\theta}). \quad (7)$$

Here,  $\Pi_{\mathbf{D}_\mu}$  denotes the projection operator onto the column space of  $\boldsymbol{\Phi}$  (namely, the subspace  $\{\boldsymbol{\Phi} \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d\}$ ) w.r.t. the



inner product  $\langle \cdot, \cdot \rangle_{D_\mu}$ , where for any vector  $\mathbf{v} \in \mathbb{R}^{|S|}$  one has

$$\Pi_{D_\mu}(\mathbf{v}) := \arg \min_{\mathbf{z} \in \{\Phi \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d\}} \|\mathbf{z} - \mathbf{v}\|_{D_\mu}^2.$$

The function  $\mathcal{T}^\pi : \mathbb{R}^{|S|} \mapsto \mathbb{R}^{|S|}$  is known as the *Bellman operator*, which is given by

$$\mathbf{v} \mapsto \mathcal{T}^\pi(\mathbf{v}) := \mathbf{r} + \gamma \mathbf{P}^\pi \mathbf{v}. \quad (8)$$

*b) Off-policy evaluation with linear approximation:* In contrast, in the **off-policy** setting, we observe a trajectory from a behavior policy  $\pi_b$  instead of the target policy  $\pi$ . The goal is then to learn the value function for the target policy  $\pi$  based on

$$\{(s_0, a_0, r_0), \dots, (s_T, a_0, r_T)\},$$

where  $a_t \sim \pi_b(\cdot \mid s_t)$ ,  $r_t = r(s_t, a_t)$ .

Let  $\mu_b$  be the stationary distribution over  $\mathcal{S}$  induced by the behavior  $\pi_b$ , and correspondingly let

$$\mathbf{D}_{\mu_b} := \text{diag}(\mu_b(1), \mu_b(2), \dots, \mu_b(|S|)).$$

We denote with  $\Pi_{D_{\mu_b}}$  the projection operator associated with  $\mathbf{D}_{\mu_b}$ , which is given explicitly as

$$\Pi_{D_{\mu_b}} \mathbf{v} := \arg \min_{\mathbf{z} \in \{\Phi \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d\}} \|\mathbf{z} - \mathbf{v}\|_{D_{\mu_b}}^2.$$

In the off-policy setting, instead of trying to solve the projected Bellman's equation (7), we aim at minimizing the Mean-Squared Projected Bellman Error (MSPBE):

$$\text{minimize}_{\boldsymbol{\theta}} \quad \text{MSPBE}(\boldsymbol{\theta}) := \frac{1}{2} \|\mathbf{V}_{\boldsymbol{\theta}} - \Pi_{D_{\mu_b}} \mathcal{T}^\pi \mathbf{V}_{\boldsymbol{\theta}}\|_{D_{\mu_b}}^2. \quad (9)$$

Throughout, we shall denote the minimizer of the above problem (9) as  $\tilde{\boldsymbol{\theta}}^*$ . We remark here that the norm and the projection are both induced by  $\mathbf{D}_{\mu_b}$ , while the Bellman operator is again in terms of the target policy  $\pi$ . For this reason, solving (9) is different from solving the projected Bellman's equation (7); as a result, in general,  $\boldsymbol{\theta}^* \neq \tilde{\boldsymbol{\theta}}^*$ .

### III. ON-POLICY EVALUATION WITH TD LEARNING

In this section, we study the accuracy of the estimator of  $\boldsymbol{\theta}^*$  (cf. (7)) returned by the TD learning algorithm in the on-policy setting. Specifically, we seek to determine the tightest sample complexity for this algorithm that ensures an  $\varepsilon$ -close solution. To better highlight our analysis strategy, we only consider the stylized generative model<sup>1</sup> whereby, at each time stamp  $t$ , one acquires an independent sample pair

$$(s_t, s'_t), \quad \text{where} \quad s_t \stackrel{\text{i.i.d.}}{\sim} \mu, \quad a_t \sim \pi(s_t), \quad \text{and} \quad s'_t \sim \mathcal{P}(\cdot \mid s_t, a_t). \quad (10)$$

Here recall that  $\mu$  is the stationary distribution corresponding to  $\mathbf{P}^\pi$ . Notice that in the on-policy setting, since we are focused on a fixed policy  $\pi$  and interested only in the state pairs  $\{(s_t, s'_t)\}_{t=0}^T$  and not the actions  $\{a_t\}_{t=0}^T$ , the Markov

decision process reduces to a Markov reward process (MRP). Given a sequence of sample pairs  $\{(s_t, s'_t)\}_{t=0}^T$  and a given level of tolerance  $\varepsilon > 0$ , our goal is to derive a sharp lower bound on the number of samples  $T$  that is required for TD learning to produce an estimator  $\hat{\boldsymbol{\theta}}$  such that, with high probability,

$$\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_{\Sigma} \leq \varepsilon.$$

#### A. The TD learning algorithm

To motivate TD learning, it is helpful to first consider the properties of the best linear approximation coefficients  $\boldsymbol{\theta}^*$ ; see (7). For any sample transition  $(s_t, s'_t)$  (see (10)), define the random quantities

$$\mathbf{A}_t := \phi(s_t) (\phi(s_t) - \gamma \phi(s'_t))^\top \in \mathbb{R}^{d \times d}, \quad (11a)$$

$$\mathbf{b}_t := \phi(s_t) r(s_t) \in \mathbb{R}^d, \quad (11b)$$

whose means are given respectively by

$$\mathbf{A} := \mathbb{E}_{s \sim \mu, s' \sim P^\pi(\cdot \mid s)} [\phi(s) (\phi(s) - \gamma \phi(s'))^\top] \quad (12a)$$

$$= \Phi^\top \mathbf{D}_\mu (\mathbf{I} - \gamma \mathbf{P}^\pi) \Phi \in \mathbb{R}^{d \times d}, \quad (12b)$$

$$\mathbf{b} := \mathbb{E}_{s \sim \mu} [\phi(s) r(s)] = \Phi^\top \mathbf{D}_\mu \mathbf{r} \in \mathbb{R}^d. \quad (12c)$$

It turns out that the target vector  $\boldsymbol{\theta}^*$  satisfies the equation [8]

$$\boldsymbol{\theta}^* := \mathbf{A}^{-1} \mathbf{b}. \quad (13)$$

The TD learning algorithm leverages this representation by iteratively improving the linear approximation of the value function at each time stamp through the updates

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t (\mathbf{A}_t \boldsymbol{\theta}_t - \mathbf{b}_t), \quad t = 0, 1, 2, \dots, \quad (14a)$$

where, for each  $t$ ,  $\eta_t > 0$  denotes the learning rate or stepsize. After  $T$  iterations, the TD learning algorithm returns  $\boldsymbol{\theta}_T$  as the estimator. In contrast, TD learning with Polyak-Ruppert averaging, or *averaged TD learning* in short, returns an average across all iterates

$$\bar{\boldsymbol{\theta}}_T = \frac{1}{T} \sum_{i=1}^T \boldsymbol{\theta}_i. \quad (14b)$$

While we are mainly concerned with the averaged estimator  $\bar{\boldsymbol{\theta}}_T$ , we also obtain some theoretical properties of  $\boldsymbol{\theta}_T$  as a by-product of our analysis.

#### B. Sample complexity of TD learning

In this section, we present a finite-sample bound for the estimation error of  $\bar{\boldsymbol{\theta}}_T$  assuming independent data, from which we derive a novel sample complexity guarantee for TD learning. Below, we denote by  $\kappa$  the condition number of  $\Sigma$  as follows

$$\kappa := \lambda_{\max}(\Sigma) / \lambda_{\min}(\Sigma) \geq 1. \quad (15)$$

**Theorem 1.** *There exist universal, positive constants  $C_0, c_0 > 0$  and  $c_1 > 0$ , such that for any given  $0 < \delta < 1$ , the averaged*

<sup>1</sup>We believe that our framework can be potentially generalized to Markovian samples using similar techniques in [62]. We will briefly discuss the techniques and difficulties in the following sections, but the full details are beyond the scope of the current paper.

TD learning estimator  $\bar{\theta}_T$  (14) after  $T$  iterations satisfies the bound

$$\|\bar{\theta}_T - \theta^*\|_{\Sigma} \quad (16)$$

$$\leq C_0 \left\{ \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log(\frac{d}{\delta})}{T(1-\gamma)^2}} (\|\theta^*\|_{\Sigma} + 1) + \frac{\|\Sigma^{-1}\| \left[ (\|\theta^*\|_{\Sigma} + 1) \sqrt{\frac{\kappa \log(\frac{dT}{\delta})}{\eta(1-\gamma)^3}} + \frac{1}{\eta(1-\gamma)} \|\theta^*\|_{\Sigma} \right]}{T} \right\} \quad (17)$$

with probability at least  $1 - \delta$ , provided that  $\theta_0 = \mathbf{0}$ ,  $\eta_0 = \dots = \eta_T = \eta < \frac{c_0(1-\gamma)}{\kappa \log(Td/\delta)}$  and

$$T \geq \frac{c_1 \kappa (\|\theta^*\|_{\Sigma} + 1)^2 \log^2 \frac{\kappa d T (\|\theta^*\|_2 + 1)}{(1-\gamma)\delta}}{\eta(1-\gamma) \lambda_{\min}(\Sigma)}.$$

a) *Proof sketch and technical novelty:* An essential step in our proof of Theorem 1 is to guarantee with high probability that the estimation error of the original TD estimator,  $\Delta_t$ , is bounded by a time-invariant value. Towards this end, we combine the matrix Freedman's inequality with an induction argument. For the technical details, we refer the readers to Steps 2 and 3 in Section 6. Controlling the norm of  $\Delta_t$  with high probability in a time-invariant manner, paves way for bounding the norm of  $\bar{\Delta}_T$  with high probability without the need of performing another projection step to restrict the norm of  $\Delta_t$  during the TD learning iterations.

b) *Generalize to Markov samples:* We give a brief introduction of how our results can be extended to Markov samples. The main difficulty towards this end lies in bounding the sequence of temporal difference errors

$$\{A_t \theta_t - b_t\}_{t \geq 0}$$

which is no longer a Martingale difference process, so the Freedman's inequality is not directly applicable anymore. In order to tackle this problem, the popular strategy is to assume that the Markov chain mixes at a geometric rate. Specifically, for arbitrarily small  $\varepsilon > 0$ , there exists a positive integer  $t_{\text{mix}}(\varepsilon) \asymp \log(\frac{1}{\varepsilon})$ , such that for any  $t \geq t_{\text{mix}}$ ,

$$\|\mathbb{E}_{t-t_{\text{mix}}}[A_t] - A\| < \varepsilon, \quad \text{and} \quad \|\mathbb{E}_{t-t_{\text{mix}}}[b_t] - b\|_2 < \varepsilon. \quad (18)$$

Here,  $\mathbb{E}_i[\cdot]$  represents the expectation conditioned on the filtration  $\mathcal{F}_i$  — the  $\sigma$ -algebra generated by  $\{(s_j, s'_j)\}_{j \leq i}$ . Under this assumption, the temporal difference error can be decomposed as

$$\begin{aligned} A_t \theta_t - b_t &= A_t (\theta_t - \theta_{t-t_{\text{mix}}}) \\ &+ [(A_t - \mathbb{E}_{t-t_{\text{mix}}}[A_t]) \theta_{t-t_{\text{mix}}} - (b_t - \mathbb{E}_{t-t_{\text{mix}}}[b_t])] \\ &+ [(\mathbb{E}_{t-t_{\text{mix}}}[A_t] - A) \theta_{t-t_{\text{mix}}} - (\mathbb{E}_{t-t_{\text{mix}}}[b_t] - b)]. \end{aligned}$$

On the right hand side, the last term can be bounded by the mixing property (18); the first term can be further expanded as

$$\theta_t - \theta_{t-t_{\text{mix}}} = \sum_{j=t-t_{\text{mix}}}^{t-1} (\theta_{j+1} - \theta_j)$$

$$= \sum_{j=t-t_{\text{mix}}}^{t-1} \eta (A_j \theta_j - b_j),$$

and bounded in terms of the step size  $\eta$  and the mixing time  $t_{\text{mix}}$ ; the second term can be controlled by separating the sequence

$$\{(A_t - \mathbb{E}_{t-t_{\text{mix}}}[A_t]) \theta_{t-t_{\text{mix}}} - (b_t - \mathbb{E}_{t-t_{\text{mix}}}[b_t])\}_{t \geq 0}$$

into  $t_{\text{mix}}$  disjoint Martingale difference processes and invoking the matrix Freedman's inequality. We leave the details to our future work.

Theorem 1 directly implies the following corollary, which gives an upper bound for the sample complexity of TD learning with independent samples.

**Corollary 1** (Sample complexity of TD learning). *There exists a universal constant  $c > 0$  such that, for any  $\varepsilon \in (0, \|\theta^*\|_{\Sigma})$  and  $\delta \in (0, 1)$ , the averaged TD estimator (14b) achieves*

$$\|V_{\bar{\theta}_T} - V_{\theta^*}\|_{D_{\mu}} = \|\bar{\theta}_T - \theta^*\|_{\Sigma} \leq \varepsilon \quad (19)$$

with probability exceeding  $1 - \delta$ , provided that

$$T \geq \frac{c \{ \max_s \phi(s)^\top \Sigma^{-1} \phi(s) \} (1 + \|\theta^*\|_{\Sigma}^2) \log(\frac{d}{\delta})}{(1-\gamma)^2 \varepsilon^2}. \quad (20)$$

c) *Comparisons to prior literature:* We remark that the best finite-sample results for TD learning obtained so far are given by [9, Theorem 2(c)] and [21, Corollary 1], with decaying stepsizes  $\eta_t \asymp t^{-1}$  and sample size-related stepsizes  $\eta_t \asymp T^{-1}$  respectively. Translated into our notation, they both prove that in order for the *expected* estimation error to be controlled by  $\varepsilon$ , namely

$$\mathbb{E} \|\theta_T - \theta^*\|_{\Sigma}^2 \leq \varepsilon^2,$$

it suffices to take (up to some logarithmic factors)

$$T^{\text{prior}} \asymp \frac{\kappa \|\Sigma^{-1}\| (\|\theta^*\|_{\Sigma}^2 + 1)}{(1-\gamma)^2} \frac{1}{\varepsilon^2}. \quad (21)$$

We refer readers to Appendix D-A and D-B for a detailed translation of their results. Comparing (20) and (21), our result improves upon previous works by a multiplicative factor of

$$\frac{T^{\text{prior}}}{T^{\text{ours}}} = \kappa,$$

the condition number of  $\Sigma$ ;  $\kappa$  can be as large as  $d$ , the dimension of the features, which can scale with  $|S|$ .

As for sample complexity with high-probability convergence guarantees, another recent work [20] shows that in order for (19) to hold with probability at least  $1 - \delta$ , it suffices to take

$$T \asymp \max \left\{ \left( \frac{1}{\varepsilon} \right)^2 \left( \log \frac{1}{\delta} \right)^3, \right. \quad (22)$$

$$\left. \left( \frac{1}{\varepsilon} \right)^{1+1/\lambda_{\min}(A)} \left( \log \frac{1}{\delta} \right)^{1+1/\lambda_{\min}(A)} \right\}. \quad (23)$$

Comparing (20) and (22), we can see that our result improves on the dependence of both the error tolerance  $\varepsilon$  and the probability tolerance  $\delta$ ; in fact, our result is the first sample

complexity for TD learning with high-probability convergence guarantee that matches the minimax-optimal dependence of  $\varepsilon$  and displays a clear dependence on the problem-related parameters, as would be shown in the following section.

After the initial post of the current paper, we are pointed to the work [31], which provides a general treatment of linear stochastic approximation with Polyak-Ruppert averaging. Their results lead to the same sample complexity as Corollary 1 with a slightly higher burn-in cost. We include the detailed comparisons of their result in Section D-D. We also point out the works of [63] and [64], which derived similar results regarding the error bound for averaged TD learning in expectation. Our result, as shown in Theorem 1, improves upon theirs in the sense that we provide high probability guarantees and offer explicit dependencies on problem-related parameters. Detailed comparisons can be found in Section D.3.

### C. Minimax lower bounds

To assess the tightness of our upper bounds in Corollary 1, in this section, we provide a minimax lower bound for the value function estimation problem with linear approximation. More specifically, the question we intend to answer is: for any target accuracy level  $\varepsilon$ , do there exist estimators that achieve an  $\varepsilon$ -approximation of  $V_{\theta^*}$  with fewer samples? As shown in the following result, the answer is, by and large, negative.

**Theorem 2** (Minimax lower bound). *Consider any  $\frac{1}{2} < \gamma < 1$ ,  $1 < d \leq |\mathcal{S}|$ , and  $0 < \varepsilon < c_1 \max\{1, \|\theta^*\|_{\Sigma}\}$  for some universal constant  $c_1 > 0$ . There exist universal constants  $c_2, c_3 > 0$  such that for any estimator  $\hat{\theta}$  based on  $T$  independent pairs  $\{(s_t, s'_t)\}_{t=1}^T$  as in (10), there exists a Markov reward process and a choice of the feature matrix  $\Phi$  such that*

$$\mathbb{P}\left\{\|\hat{\theta} - \theta^*\|_{\Sigma} > c_2 \varepsilon\right\} \geq \frac{1}{4}, \quad (24)$$

provided that the number of samples  $T$  satisfies

$$T \leq \frac{c_3 \{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)\} (1 + \|\theta^*\|_{\Sigma}^2)}{(1 - \gamma) \varepsilon^2}. \quad (25)$$

**Remark 1.** We remark that minimax lower bounds are also previously investigated in a general framework in [12] where the value function is approximated using a general reproducing kernel Hilbert space (RKHS). When it comes to linear function approximation, for completeness, we include in Section B a different but simpler construction tailored to the linear space. Compared to the results of [12], our lower bound is stated in terms of different parameters, which allows us to evaluate the tightness of Corollary 1 directly. Instantiating both lower bounds, they do agree and equal to

$$O\left(\frac{d}{\varepsilon^2(1 - \gamma)^3}\right), \quad (26)$$

as one plugs in the exact parameters from our construction.

As asserted by this theorem, no algorithm whatsoever can attain an  $\varepsilon$ -approximation of the best linear coefficient — in

a minimax sense — unless the total sample size exceeds

$$O\left(\frac{\{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)\} (1 + \|\theta^*\|_{\Sigma}^2)}{(1 - \gamma) \varepsilon^2}\right).$$

Consequently, the upper bounds developed in Corollary 1 are sharp in terms of the accuracy level  $\varepsilon$ , the dependence of the feature map  $\Phi$ , the underlying coefficient  $\theta^*$ , and the covariance matrix  $\Sigma$ . Therefore, it implies that the performances of the TD learning algorithms can not be further improved in the minimax sense other than a factor of  $\frac{1}{1 - \gamma}$  — the effective horizon.

We believe that the gap in terms of  $\frac{1}{1 - \gamma}$  mainly comes from the function approximation paradigm. In fact, as far as we know, with linear function approximation, there has been no minimax optimal results established for this problem either for the model based method (e.g. LSTD) or for the variance-reduced approach, both of which are known to be minimax optimal in the tabular setting; the latter is also proved to be instance-optimal from [33] and [64]. We conjecture the minimax optimal dependency of  $\frac{1}{1 - \gamma}$  to be the same as that of the tabular setting and TD with LFA to be minimax optimal. Establishing this result, however, requires developing completely new analysis tools, particularly in dealing with the structure of variance across different steps, which we leave as an interesting open direction.

### IV. OFF-POLICY EVALUATION WITH TDC LEARNING

In this section, we aim to estimate the optimizer  $\tilde{\theta}^*$  of the optimization problem (9) in the off-policy setting by means of the TDC algorithm. We continue to focus on the case when samples are generated in the i.i.d. fashion by the behavior policy  $\pi_b$ . At each time stamp  $t$ , one obtains

$$(s_t, a_t, s'_t), \quad \text{where} \\ s_t \stackrel{\text{i.i.d.}}{\sim} \mu_b, \quad a_t \sim \pi_b(\cdot | s_t), \quad \text{and} \quad s'_t \sim \mathcal{P}(\cdot | s_t, a_t). \quad (27)$$

Here, recall that  $\mu_b$  is the stationary distribution corresponding to the behavior policy  $\pi_b$ . We first provide some intuition behind the TDC algorithm before describing novel bounds on its sample complexity for obtaining an  $\varepsilon$ -accurate solution.

#### A. The TDC algorithm

The TDC algorithm is designed to solve the optimization problem (9) using a two-timescale linear TD with gradient correction [24]. To provide some high-level ideas behind the design of this algorithm, it is helpful to rewrite the objective function in the following form by directly expanding the terms in expression (9).

**Claim 1.** *The quantity  $\text{MSPBE}(\theta)$  can be equivalently written as*

$$\text{MSPBE}(\theta) = \frac{1}{2} \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t) \delta_t]^\top \left\{ \mathbb{E}_{\mu_b} [\phi(s_t) \phi(s_t)^\top] \right\}^{-1} \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t) \delta_t], \quad (28)$$

where  $\delta_t := r_t + \gamma \phi(s'_t)^\top \theta - \phi(s_t)^\top \theta$  is the temporal difference error.

In light of the above expression, the gradient of MSPBE( $\theta$ ) with respect to  $\theta$  equals to

$$\begin{aligned} \nabla_{\theta} \text{MSPBE}(\theta) &= \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [(\gamma\phi(s'_t) - \phi(s_t)) \phi(s_t)^\top] \\ &\quad \{ \mathbb{E}_{\mu_b} [\phi(s_t)\phi(s_t)^\top] \}^{-1} \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t)\delta_t] \\ &= -\mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t)\delta_t] + \gamma \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s'_t)\phi(s_t)^\top] \\ &\quad \{ \mathbb{E}_{\mu_b} [\phi(s_t)\phi(s_t)^\top] \}^{-1} \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t)\delta_t] \\ &= -\mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t)\delta_t] + \gamma \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s'_t)\phi(s_t)^\top] \mathbf{w}_t, \end{aligned} \quad (29)$$

where in the last step we have defined

$$\begin{aligned} \mathbf{w}_t &= \mathbf{w}(\theta_t) \\ &= \{ \mathbb{E}_{\mu_b} [\phi(s_t)\phi(s_t)^\top] \}^{-1} \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t)\delta_t]. \end{aligned} \quad (30)$$

and have used the importance weights

$$\rho_t := \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} \quad (31)$$

to replace the expectation w.r.t.  $\pi$  with the expectation w.r.t.  $\pi_b$ .

The high-level idea of TDC is to estimate the right hand side of (29) based on the sample trajectory (27), and then perform stochastic gradient updates for  $\theta_t$ . However, the challenge is that the second term in the gradient of MSPBE (29) involves the product of two expectations. Simultaneously sampling and using the sample product is inappropriate due to their correlation. In order to address this issue, [65] and [24] introduced an auxiliary parameter  $\mathbf{w}$  to estimate  $\mathbf{w}(\theta_t)$  by solving a linear stochastic approximation (SA) problem corresponding to the linear system

$$\mathbb{E}_{\mu_b} [\phi(s_t)\phi(s_t)^\top] \mathbf{w} = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t)\delta_t]. \quad (32)$$

Putting these ideas together, TDC amounts to the following two-timescale linear stochastic method

$$\begin{aligned} \tilde{\theta}_{t+1} &= \tilde{\theta}_t - \alpha_t [\gamma \rho_t \phi(s'_t)\phi(s_t)^\top \mathbf{w}_t - \rho_t \delta_t \phi(s_t)]; \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \beta_t [\phi(s_t)\phi(s_t)^\top \mathbf{w}_t - \rho_t \delta_t \phi(s_t)]. \end{aligned}$$

Here, the update of  $\tilde{\theta}_t$  corresponds to a gradient step regarding (28), the update of  $\mathbf{w}_t$  corresponds to linear SA for solving (32), and  $\delta_t := r_t + \gamma\phi(s'_t)^\top \tilde{\theta}_t - \phi(s_t)^\top \tilde{\theta}_t$  is the temporal difference error. In addition,  $\alpha_t, \beta_t$  are the corresponding stepsizes. For notational convenience, let us denote

$$\begin{aligned} \tilde{\mathbf{A}}_t &= \rho_t \phi(s_t) (\phi(s_t) - \gamma\phi(s'_t))^\top, \\ \tilde{\mathbf{b}}_t &:= \rho_t \phi(s_t) r_t, \\ \Pi_t &:= \rho_t \phi(s_t)\phi(s'_t)^\top, \\ \tilde{\Sigma}_t &:= \phi(s_t)\phi(s_t)^\top. \end{aligned} \quad (33)$$

With these definitions, the TDC iterates can be written compactly as

$$\tilde{\theta}_{t+1} = \tilde{\theta}_t - \alpha_t (\tilde{\mathbf{A}}_t \tilde{\theta}_t - \tilde{\mathbf{b}}_t + \gamma \Pi_t^\top \mathbf{w}_t); \quad (34a)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \beta_t (\tilde{\mathbf{A}}_t \mathbf{w}_t - \tilde{\mathbf{b}}_t + \tilde{\Sigma}_t \mathbf{w}_t). \quad (34b)$$

## B. Sample complexity of TDC

Our finite-sample characterization of TDC builds upon a careful analysis of the population dynamics of TDC, which we then show to be uniformly well approximated by the empirical dynamics of TDC via matrix concentration inequalities. Before stating our main result, we find it helpful to introduce some extra pieces of notation. Specifically, define the population parameters as

$$\tilde{\mathbf{A}} := \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\tilde{\mathbf{A}}_t] = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t) (\phi(s_t) - \gamma\phi(s'_t))^\top]; \quad (35a)$$

$$\tilde{\mathbf{b}} := \mathbb{E}_{\mu_b} [\tilde{\mathbf{b}}_t] = \mathbb{E}_{\mu_b, \pi_b} [\rho_t \phi(s_t) r_t]; \quad (35b)$$

$$\Pi := \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\Pi_t] = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t)\phi(s'_t)^\top]; \quad (35c)$$

$$\tilde{\Sigma} := \mathbb{E}_{\mu_b} [\tilde{\Sigma}_t] = \mathbb{E}_{\mu_b} [\phi(s_t)\phi(s_t)^\top]. \quad (35d)$$

In addition, denote the parameters

$$\begin{aligned} \lambda_1 &= \lambda_{\min}(\tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}), \\ \lambda_2 &= \lambda_{\min}(\tilde{\Sigma}), \\ \lambda_\Sigma &= \|\tilde{\Sigma}^{-1}\| = 1/\lambda_2, \\ \tilde{\kappa} &= \lambda_\Sigma \cdot \|\tilde{\Sigma}\|, \\ \rho_{\max} &= \max_{s,a} [\pi(a|s)/\pi_b(a|s)]. \end{aligned} \quad (36)$$

With these notation in place, we are ready to state our main result for TDC learning, with its proof deferred to Section VII.

**Theorem 3.** *There exist universal constants  $\tilde{C}_0, \tilde{c}_1 > 0$ , such that for any given  $0 \leq \delta \leq 1$ , the output  $\tilde{\theta}_T$  of the TDC learning iterate (34) at time  $T$  satisfies the bound*

$$\|\tilde{\theta}_T - \tilde{\theta}^*\|_{\tilde{\Sigma}} \leq \tilde{C}_0 \frac{\rho_{\max}^2 \|\tilde{\Sigma}\|}{\lambda_1} \sqrt{\frac{\beta}{\lambda_2} \log \frac{2dT}{\delta}} (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + 2), \quad (37)$$

with probability at least  $1 - \delta$ , provided that

$$\begin{aligned} \tilde{\theta}_0 &= \mathbf{0}, \\ \alpha_0 &= \dots = \alpha_T = \alpha, \\ \beta_0 &= \dots = \beta_T = \beta, \\ 0 < \alpha &< \frac{1}{\lambda_1 \lambda_\Sigma^2 \|\tilde{\Sigma}\| \log \frac{2dT}{\delta}}, \\ \frac{\alpha}{\beta} &= \frac{1}{128 \rho_{\max}^2 (1 + \lambda_\Sigma \rho_{\max})}, \\ T &\geq \tilde{c}_1 \frac{\log \|\tilde{\theta}^*\|_2}{\alpha \lambda_1} \log \max \left\{ \sqrt{\tilde{\kappa}}, \|\tilde{\theta}^*\|_{\tilde{\Sigma}} \sqrt{\frac{\alpha \lambda_1}{\log \frac{2dT}{\delta}}} \right\}. \end{aligned} \quad (38)$$

**Remark 2.** A similar result in terms of the  $\ell_2$  error (namely,  $\|\tilde{\theta}_T - \tilde{\theta}^*\|_2$ ) can be derived in the same way as in (37). In particular, under the same conditions as in (38), it can be derived with probability at least  $1 - \delta$  that

$$\|\tilde{\theta}_T - \tilde{\theta}^*\|_2 \lesssim \tilde{C}_0 \frac{\rho_{\max}^2}{\lambda_1} \sqrt{\frac{\beta}{\lambda_2} \log \frac{2dT}{\delta}} (\|\tilde{\theta}^*\|_2 + 2). \quad (39)$$

Since the proof follows in the similar fashion, we omit here for brevity.



a) *Proof sketch and technical novelty:* Our proof of Theorem 3 considers the convergence of the vector

$$\mathbf{x}_t := \left[ \varkappa \left[ \tilde{\mathbf{w}}_t + \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}(\tilde{\boldsymbol{\theta}}_t - \tilde{\boldsymbol{\theta}}^*) \right] \right],$$

where  $\varkappa \in (0, 1)$  is a constant to be specified. Firstly, we identify the conditions on  $\alpha, \beta$  and  $\varkappa$  that guarantee the exponential convergence of  $\mathbf{x}_t$  in the noise-free scenario; after this, we again combine the matrix Freedman's inequality and an induction argument to bound the norm of  $\mathbf{x}_t$  for *i.i.d.* samples with high probability by a time-invariant value in terms of the step size. And finally, with a careful choice of  $\alpha, \beta$  and  $\varkappa$ , we establish the finite-sample guarantee as is shown in Theorem 3. The main technical novelty of this proof lies in the construction of the vector  $\mathbf{x}_t$  and the choice of the parameter  $\varkappa$ .

Next, we state a direct consequence of Theorem 3 below, which gives an upper bound for the sample complexity of TDC.

**Corollary 2.** *There exists a universal constant  $\tilde{c}$  such that, for any  $\delta \in (0, 1)$  and  $\varepsilon \in (0, \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}})$ , the TDC estimator  $\tilde{\boldsymbol{\theta}}_T$  at iterate  $T$  satisfies the bound*

$$\|\mathbf{V}_{\tilde{\boldsymbol{\theta}}_T} - \mathbf{V}_{\tilde{\boldsymbol{\theta}}^*}\|_{\mathbf{D}_{\mu_b}} = \|\tilde{\boldsymbol{\theta}}_T - \tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} \leq \varepsilon \quad (40)$$

with probability exceeding  $1 - \delta$ , provided that

$$T \geq \tilde{c} \frac{\rho_{\max}^7}{\lambda_1^4 \lambda_2^3} \frac{\|\tilde{\Sigma}\|^2}{\varepsilon^2} (1 + \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}^2) \log \left( \frac{d \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}}{\delta} \right), \quad (41)$$

and the stepsize parameters  $\alpha_t$  and  $\beta_t$  are chosen as

$$\alpha_t \asymp \frac{\log \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}}{T \lambda_1}, \quad \beta_t = 128 \frac{\rho_{\max}^2 (1 + \lambda_{\Sigma} \rho_{\max})}{\lambda_1 \lambda_2} \alpha. \quad (42)$$

b) *Comparisons to other sample complexity bounds for TDC:* Let us compare our results in Theorem 3 and Corollary 2 with the state-of-the-art sample complexities for the TDC algorithm. The result that is most comparable to ours is obtained by [28], where a projected version of TDC is considered with decaying stepsizes  $\alpha_t = O(t^{-\alpha})$  and  $\beta_t = O(t^{-\beta})$  for  $0 < \beta < \alpha < 1$ . The sample complexity therein, with high-probability convergence guarantee at tolerance level  $\varepsilon$ , is of order  $O\left(\left(\frac{1}{\varepsilon}\right)^{2\alpha}\right)$  without explicit dependence on the problem-related parameters. If one chooses  $\alpha = 1 - \delta$  with  $\delta$  sufficiently small, their sample complexity bound can be improved, but it cannot achieve the rate  $\Theta\left(\frac{1}{\varepsilon^2}\right)$ . Regarding finite-sample in-expectation error control for TDC, the best result so far is developed by [29], who shows that with the choice of  $\alpha_t, \beta_t \asymp \frac{1}{T}$ , the sample complexity for TDC with tolerance level  $\varepsilon$  can be upper bounded by  $O\left(\frac{1}{\varepsilon^2}\right)$ . Our result in Corollary 2 is the first sample complexity for the original TDC algorithm that guarantees high-probability convergence and achieves the minimax-optimal rate of  $O\left(\frac{1}{\varepsilon^2}\right)$ ; it is also noteworthy that we display an explicit dependence on problem-related parameters. The key to achieving this again lies in our combination of the matrix Freedman's inequality with an induction argument; the details of the proof is postponed to Section VII. We also remark that [26] considers

a variant of TDC where  $\tilde{\boldsymbol{\theta}}_t$  is updated *not* with every sample tuple  $(s_t, a_t, s'_t)$ , but with every batch of samples, and obtains a sample complexity of order  $O\left(\frac{1}{\varepsilon^2} \log\left(\frac{1}{\varepsilon}\right)\right)$ .

## V. NUMERICAL EXPERIMENTS

In this section, we corroborate our theoretical results with illustrative numerical experiments. In what follows, we will consider the on-policy and off-policy settings respectively.

### A. On-policy evaluation: averaged TD learning

In the on-policy setting, we will investigate the empirical performance of the averaged TD learning algorithm.

a) *MDP setting:* We consider a member of the family of MDPs constructed in proof of Theorem 2, which provides a minimax lower bound. This family of MDPs is designed to be difficult to distinguish between each other, and hence, is a natural instance for evaluating the performance of TD learning. For construction details of this MDP, we refer the reader to Appendix B. In these simulations, we set  $|S| = 10$ ,  $\gamma = 0.2$ , and choose the stepsize of TD as  $\eta = 0.01$ . We examine both the original and the averaged TD iterates when the feature dimension equals to  $d = 3$  and  $d = 9$ . Under each setting, 100 independent trials for  $T = 10^5$  iterations were conducted, and we report the mean value as well as the 95% confidence band for the estimation error  $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_{\Sigma}$  for TD and  $\|\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_{\Sigma}$  for averaged TD.

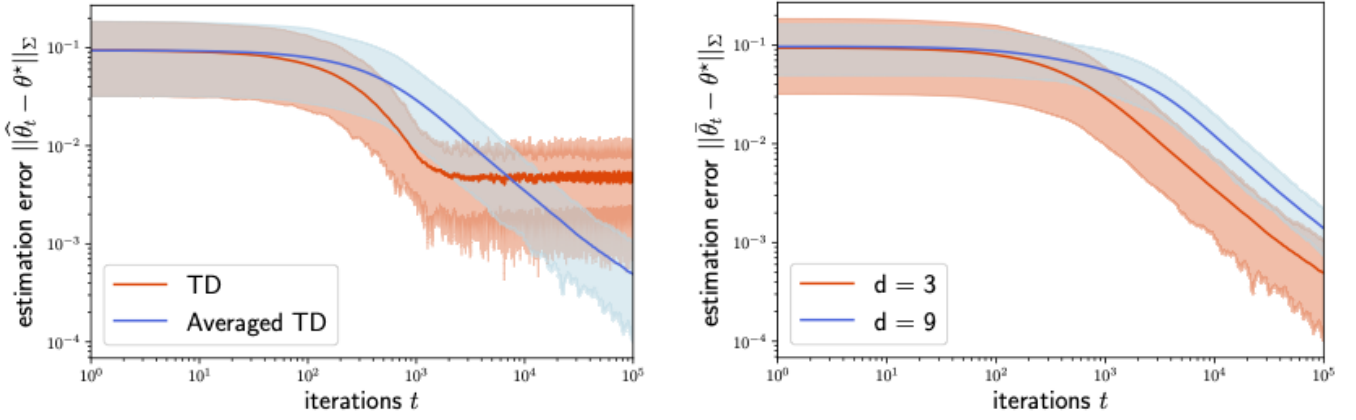
b) *Experimental results:* Figure 1(a) compares the performances of TD and averaged TD of an MDP with feature dimension  $d = 3$ . While the estimation error of TD levels off at around  $5 \times 10^{-3}$  after  $10^3$  iterations, the error of averaged TD keeps decreasing to below  $5 \times 10^{-4}$  when  $T = 10^5$ . In addition, Figure 1(b) demonstrates the estimation error of averaged TD for MDPs with feature dimension  $d = 3$  and  $d = 9$ . The slopes of these curves on the right part of this log-log plot match our theoretical prediction: the estimation error decreases in the order of  $O(t^{-1/2})$ . Moreover, the difference between the two curves indicates that the lower-dimension problem enjoys a faster convergence rate.

### B. Off-policy evaluation: TDC learning

In order to demonstrate the efficiency of TDC for off-policy evaluation, we compare its performance with that of the off-policy TD learning on Baird's counterexample [23].

a) *Baird's counterexample:* We start by introducing Baird's counterexample, which was constructed to illustrate the instability of TD learning in the off-policy regime. Consider an MDP  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$ , with the discount factor  $\gamma = 0.9$ , state space  $\mathcal{S} = [7]$ , action space  $\mathcal{A} = \{0, 1\}$  and the reward function  $r = 0$  for all states and actions. The action  $a = 1$  transitions any initial state  $s$  to  $s' = 7$ , while the action  $a = 0$  transitions any initial state  $s$  to  $s' \in [6]$  with the same probability. The target policy  $\pi$  selects  $a = 1$  at any given state  $s$ , while the behavior policy  $\pi_b$  takes  $a = 0$  with probability  $\frac{6}{7}$  and  $a = 1$  with probability  $\frac{1}{7}$ . Formally, the MDP satisfies the equations (see also Figure 2 for an illustration)

$$\mathcal{P}(s'|s, 1) = \mathbb{1}\{s' = 7\}, \quad \forall s \in [7];$$



**Fig. 1.** (a) Comparisons of the estimation error of TD and averaged TD when  $d = 3$ . (b) Comparisons of the estimation error for averaged TD with  $d = 3$  and  $d = 9$ . Two curves in the middle represent their average errors, while the shaded areas represent the 95% confidence bands.

$$\begin{aligned} \mathcal{P}(s'|s, 0) &= \frac{1}{6} \mathbb{1}\{1 \leq s' \leq 6\}, \quad \forall s \in [7]; \\ \pi(1|s) &= 1, \quad \forall s \in [7]; \\ \pi_b(0|s) &= \frac{6}{7}, \quad \forall s \in [7]; \\ \pi_b(1|s) &= \frac{1}{7}, \quad \forall s \in [7]. \end{aligned}$$

In this example, it is easy to check that the stationary distribution corresponding to the behavior policy  $\pi_b$  is the uniform distribution among all states, and that the value function is 0 for all states. We apply the following linear approximation of the value function: for  $\theta \in \mathbb{R}^8$ ,

$$\begin{cases} V(i) = 2\theta_i + \theta_8, & \text{for } 1 \leq i \leq 6; \\ V(7) = \theta_7 + 2\theta_8. \end{cases} \quad (43)$$

We remark that with this approximation, the feature space has a higher dimension ( $d = 8$ ) than the state space ( $|S| = 7$ ). Consequently, the optimal estimator  $\theta^*$  is not unique, and instead can be any  $\theta \in \mathbb{R}^8$  such that the estimated value vector is  $V_\theta = \mathbf{0}$ . Technically, this issue can be circumvented by creating several identical states as state  $s = 7$ ; we omit this detail here for simplicity, since we use  $\|\hat{\theta}_t - \theta^*\|_{\Sigma} = \|\hat{V}_{\hat{\theta}_t} - V^*\|_{D_{\mu_b}}$  to evaluate the estimation error, and our experimental results would remain the same.

*b) Experimental results:* We perform 100 independent trials for both off-policy averaged TD learning (with stepsize  $\eta = 0.02$ ) and TDC (with stepsizes  $\alpha = 0.02, \beta = 0.002$ ), starting at  $\hat{\theta}_0 = (1, 1, 1, 1, 1, 1, 10, 1)^\top$ , as suggested by [23]. In these experiments, we set  $\alpha = \eta$  to ensure that the stepsize for  $\theta$ -updates are the same between the two algorithms. Figure 3 demonstrates how the estimation error  $\|\hat{\theta}_t - \theta^*\|_{\Sigma}$  changes as two algorithms execute. As can be seen in this figure, TDC converges to an error of below 0.01 after  $T = 10^5$  iterations while the off-policy averaged TD diverges to infinity.

## VI. PROOF OF THEOREM 1 (TD LEARNING)

For the sake of convenience, let us introduce the following notation

$$\Delta_t := \theta_t - \theta^*, \quad \text{and} \quad \bar{\Delta}_t := \bar{\theta}_t - \theta^*. \quad (44)$$

*a) Step 1: a recursive relation:* To understand the convergence behavior of  $\bar{\Delta}_t$ , the idea is to first look at the following decomposition

$$\begin{aligned} \Delta_{t+1} &= \theta_{t+1} - \theta^* = \theta_t - \theta^* - \eta(A_t \theta_t - b_t) \\ &= \theta_t - \theta^* - \eta(A_t \theta_t - b_t - (A \theta^* - b)) \\ &= \theta_t - \theta^* - \eta(A(\theta_t - \theta^*) + (A_t - A)\theta_t - (b_t - b)) \\ &= (I - \eta A) \Delta_t - \eta \xi_t, \end{aligned}$$

where we define

$$\xi_t := (A_t - A)\theta_t - (b_t - b). \quad (45)$$

Here, the second line invokes the update rule (14a) and the identity  $A \theta^* = b$ , whereas the third line is obtained by properly rearranging terms. Applying the above relation recursively, one arrives at

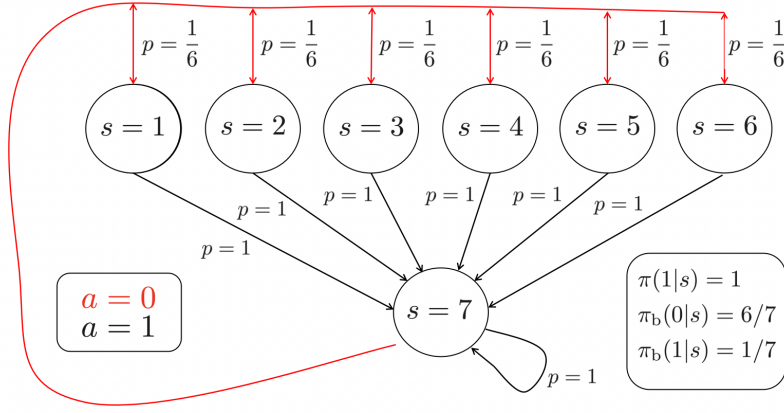
$$\begin{aligned} \Delta_t &= (I - \eta A) \Delta_{t-1} - \eta \xi_{t-1} \\ &= (I - \eta A)^t \Delta_0 - \eta \sum_{i=0}^{t-1} (I - \eta A)^{t-i-1} \xi_i. \end{aligned} \quad (46)$$

*b) Step 2: a crude bound on  $\|\Delta_t\|_{\Sigma}$ :* We aim to establish, via an induction argument, that with probability at least  $1 - \delta$ ,

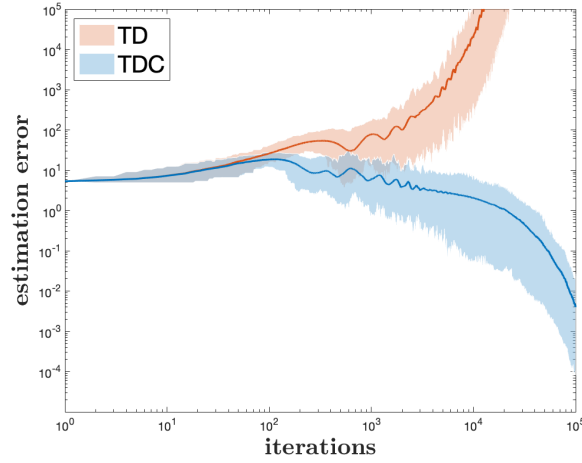
$$\|\Delta_t\|_{\Sigma} \leq 32 \sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1 - \gamma}} (1 + \|\theta^*\|_{\Sigma}) + 2\sqrt{\kappa} \|\Delta_0\|_{\Sigma} =: R_0 \quad (47)$$

simultaneously over all  $0 \leq t \leq T$ , as long as  $0 < \eta_t \leq \frac{c_3(1-\gamma)}{\kappa \log \frac{2dT}{\delta}}$  for some sufficiently small constant  $c_3 > 0$ . As a side remark, this boundedness property saves us from enforcing additional projection steps as adopted in [9].

To start with, note that the inequality (47) holds trivially for



**Fig. 2.** Baird's counterexample. Taking action  $a = 1$  always leads to state  $s = 7$ , while taking  $a = 0$  leads to one of the other six states with equal probability. The reward is set to be always zero.



**Fig. 3.** Performances of off-policy averaged TD (red,  $\eta = 0.02$ ) and TDC (blue,  $\alpha = 0.02$ ,  $\beta = 0.002$ ). Two curves in the middle represent their average errors, while the shaded areas correspond to 95% confidence bands.

the base case with  $t = 0$ , given that  $\kappa \geq 1$ . Next, suppose that the hypothesis (47) holds for  $\Delta_0, \dots, \Delta_{t-1}$ , and we intend to establish it for  $\Delta_t$  as well. Towards this end, invoking the decomposition (46) and the triangle inequality yields

$$\|\Delta_t\|_{\Sigma} \leq \|(\mathbf{I} - \eta\mathbf{A})^t \Delta_0\|_{\Sigma} + \eta \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta\mathbf{A})^{t-i-1} \tilde{\xi}_i \right\|_{\Sigma}. \quad (48)$$

As for the first term of (48), it is seen that

$$\begin{aligned} & \|(\mathbf{I} - \eta\mathbf{A})^t \Delta_0\|_{\Sigma} \\ &= \|\Sigma^{1/2} (\mathbf{I} - \eta\mathbf{A})^t \Sigma^{-1/2} \Sigma^{1/2} \Delta_0\|_2 \\ &\leq \|\Sigma^{1/2}\| \cdot \|\Sigma^{-1/2}\| \cdot \|\mathbf{I} - \eta\mathbf{A}\|^t \cdot \|\Sigma^{1/2} \Delta_0\|_2 \\ &\leq \sqrt{\kappa} \left(1 - \frac{1}{2}\eta(1 - \gamma)\lambda_{\min}(\Sigma)\right)^t \|\Delta_0\|_{\Sigma} \leq \sqrt{\kappa} \|\Delta_0\|_{\Sigma}, \end{aligned} \quad (49)$$

where the last inequality arises from the definition of  $\kappa$  and the property (88g) (with the restriction that  $\eta \leq (1 - \gamma)/(4\|\Sigma\|)$ ). When it comes to the second term of (48), the following lemma comes in handy.

**Lemma 1.** Fix any quantity  $R > 0$  and, for each  $0 \leq i \leq T - 1$ , define the auxiliary random vector

$$\tilde{\xi}_i := \xi_i \mathbb{1}\{\mathcal{H}_i\}, \quad \text{where } \mathcal{H}_i := \left\{ \|\Delta_i\|_{\Sigma} \leq R \right\}. \quad (50)$$

Then, with probability at least  $1 - \delta/T$ , simultaneously over the indices  $(l, u, t)$  such that  $0 \leq l \leq u \leq t - 1 < T$  it holds that

$$\begin{aligned} & \left\| \sum_{i=l}^u (\mathbf{I} - \eta\mathbf{A})^{t-i-1} \tilde{\xi}_i \right\|_{\Sigma} \\ &\leq 16 \left(1 - \frac{1}{2}\eta(1 - \gamma)\lambda_{\min}(\Sigma)\right)^{t-u-1} \\ &\quad \cdot (\|\theta^*\|_{\Sigma} + R + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1 - \gamma)}}, \end{aligned}$$

provided that  $0 < \eta \leq \frac{1 - \gamma}{\kappa \log \frac{2dT}{\delta}}$ .

*Proof.* See Section C-A.  $\square$

Under the induction hypothesis that  $\|\Delta_i\|_{\Sigma} \leq R_0$  for  $0 \leq i \leq t - 1$ , we can invoke Lemma 1 (with  $R = R_0, l = 0$  and

$u = t - 1$ ) to show that

$$\begin{aligned} & \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} \\ &= \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \mathbb{1}_{\{\|\Delta_i\|_{\Sigma} \leq R_0\}} \right\|_{\Sigma} \\ &\leq 16(\|\theta^*\|_{\Sigma} + R_0 + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \end{aligned} \quad (51)$$

holds with probability at least  $1 - \delta/T$ , provided that  $0 < \eta \leq \frac{1-\gamma}{\kappa \log \frac{2dT}{\delta}}$ . Combining (48), (49) and (51) together and recalling the definition (47) of  $R_0$ , we can easily verify that

$$\|\Delta_t\|_{\Sigma} \leq \sqrt{\kappa} \|\Delta_0\|_{\Sigma} + 16(\|\theta^*\|_{\Sigma} + R_0 + 1) \sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}} \leq R_0 \leq \left\| \sum_{i=t-t_{\text{seg}}+1}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} + \left\| \sum_{i=0}^{t-t_{\text{seg}}} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma}, \quad (52)$$

with the proviso that  $32\sqrt{\frac{\eta \kappa \log(2dT/\delta)}{1-\gamma}} \leq 1$ . The induction argument coupled with the union bound then establishes the claim (47).

c) *Step 3: a refined bound on  $\|\Delta_t\|_{\Sigma}$ .* It turns out that the upper bound (47) is somewhat loose due to the complete ignorance of the contraction effect of  $\mathbf{I} - \eta \mathbf{A}$ ; see (49). In what follows, we develop a strengthened bound. Define

$$t_{\text{seg}} := \frac{c_1 \log \max\{4\sqrt{\kappa}, \frac{16\kappa\|\Delta_0\|_{\Sigma}}{\|\theta^*\|_{\Sigma}+1}, \|\Delta_0\|_{\Sigma} \sqrt{\frac{1-\gamma}{\kappa \log \frac{2dT}{\delta}}}\}}{\eta(1-\gamma)\lambda_{\min}(\Sigma)} \quad (53)$$

for some sufficiently large constant  $c_1 > 0$ . For any integer  $k \geq 1$ , we aim to establish that

$$\|\Delta_t\|_{\Sigma} \leq 32\sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}} \left( \|\theta^*\|_{\Sigma} + \frac{\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}{2^{k-1}} + \frac{3}{2} \right) =: R_k \quad (54)$$

for any  $t$  obeying  $kt_{\text{seg}} \leq t \leq T$ , provided that  $0 < \eta \leq \frac{c_3(1-\gamma)}{\kappa \log \frac{2dT}{\delta}}$  for some small enough constant  $c_3 > 0$ .

Because of relation (48), we claim that it suffices to prove that

$$\begin{aligned} & \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} \\ &\leq 32 \left( \|\theta^*\|_{\Sigma} + \frac{2\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}{2^k} + 1 \right) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}}, \end{aligned} \quad (55)$$

$$\text{when } kt_{\text{seg}} \leq t \leq T. \quad (56)$$

To see this: note that the first term on the right-hand side of (48) has already been bounded in (49), which combined with the definition (53) of  $t_{\text{seg}}$  indicates that

$$\|\Sigma^{1/2}(\mathbf{I} - \eta \mathbf{A})^t \Delta_0\|_2 \quad (57)$$

$$\leq \sqrt{\kappa} \left( 1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma) \right)^{t_{\text{seg}}} \|\Delta_0\|_{\Sigma} \quad (58)$$

$$\leq \sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}} \quad (59)$$

for any  $t \geq t_{\text{seg}}$ . Clearly, combining (56) with (48) and (59) shall immediately lead to the claim (54). The remainder of this step is thus devoted to demonstrating (56) inductively.

The base case (i.e.  $k = 1$ ) follows immediately from our bounds (47) and (51) in Step 2, given that  $\sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}}$  is sufficiently small. Suppose now that the claim (56) holds for a given integer  $k \geq 1$  and any  $t$  obeying  $kt_{\text{seg}} \leq t \leq T$ , and we intend to show that (56) continues to hold for  $k+1$  and any  $t$  obeying  $(k+1)t_{\text{seg}} \leq t \leq T$ . Towards this, we first single out the following straightforward decomposition

$$\left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} = \left\| \sum_{i=t-t_{\text{seg}}+1}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} + \left\| \sum_{i=0}^{t-t_{\text{seg}}} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma},$$

which allows us to upper bound the two terms on the right-hand side above separately.

- Under the induction hypothesis that  $\|\Delta_i\|_{\Sigma} \leq R_k$  for all  $i$  obeying  $kt_{\text{seg}} \leq i \leq T$ , one can invoke Lemma 1 with  $R = R_k, l = t - t_{\text{seg}} + 1$  and  $u = t - 1$  to see that

$$\begin{aligned} & \left\| \sum_{i=t-t_{\text{seg}}+1}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} \\ &= \left\| \sum_{i=t-t_{\text{seg}}+1}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \mathbb{1}_{\{\|\Delta_i\|_{\Sigma} \leq R_k\}} \right\|_{\Sigma} \\ &\leq 16(\|\theta^*\|_{\Sigma} + R_k + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \\ &\leq 24 \left( \|\theta^*\|_{\Sigma} + \frac{2\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}{2^{k+1}} + 1 \right) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}}, \end{aligned}$$

where the last line uses the definition (54) of  $R_k$  and holds as long as  $\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}$  is sufficiently small.

- In addition, we make the observation that: for any  $t \geq t_{\text{seg}}$ ,

$$\begin{aligned} & \left\| \sum_{i=0}^{t-t_{\text{seg}}} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} \\ &= \left\| \sum_{i=0}^{t-t_{\text{seg}}} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \mathbb{1}_{\{\|\Delta_i\|_{\Sigma} \leq R_0\}} \right\|_{\Sigma} \\ &\leq 16(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma))^{t_{\text{seg}}-1} \\ &\quad \cdot (\|\theta^*\|_{\Sigma} + R_0 + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \\ &\leq 8(\|\theta^*\|_{\Sigma} + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}}. \end{aligned}$$

Here, the first equality uses the crude bound  $\|\Delta_i\|_{\Sigma} \leq R_0$  for all  $i$  (see (47)), the second to last inequality utilizes Lemma 1 with  $R = R_0, l = 0$  and  $u = t - t_{\text{seg}}$ , whereas the last inequality relies on the definition (47) of  $R_0$  and invokes



the fact that  $\sqrt{\kappa} \left(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma)\right)^{t_{\text{seg}}-1} \leq 1 - \delta$  (with their proofs deferred to Section C-B)  
 $\min \left\{ \frac{1}{4}, \frac{1}{4\sqrt{\kappa}\|\Delta_0\|_{\Sigma}} \right\}$  with our choice (53) of  $t_{\text{seg}}$ .

Combine the previous two bounds to reach

$$\begin{aligned} & \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta\mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} \\ & \leq 24 \left( \|\theta^*\|_{\Sigma} + \frac{2\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}{2^{k+1}} + 1 \right) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \\ & \quad + 8(\|\theta^*\|_{\Sigma} + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \\ & \leq 32 \left( \|\theta^*\|_{\Sigma} + \frac{2\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}{2^k} + 1 \right) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}}. \end{aligned}$$

This finishes the induction step and in turn establishes (56) (and hence (54)).

As a straightforward consequence, the bounds (47) and (54) imply that

$$\|\Delta_t\|_{\Sigma} \leq \begin{cases} R_0, & 0 \leq t < t'_{\text{seg}}, \\ 32\sqrt{\frac{\eta\kappa \log \frac{2dT}{\delta}}{1-\gamma}}(\|\theta^*\|_{\Sigma} + 2), & t'_{\text{seg}} \leq t < T, \end{cases} \quad (60)$$

where

$$t'_{\text{seg}} := c_2 t_{\text{seg}} \log(\kappa(\|\Delta_0\|_2 + 1)) \quad (61)$$

for some large enough constant  $c_2 > 0$ . To see this, note that for any  $t \geq t'_{\text{seg}}$ , it is guaranteed that the second term on the right-hand side of (54) obeys  $\frac{4\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}{2^{\lfloor t/t_{\text{seg}} \rfloor}} \leq 2$ , thus confirming the second case in (60).

d) *Step 4: controlling  $\|\bar{\Delta}_T\|_{\Sigma}$ :* Now we are positioned to control  $\bar{\Delta}_T$ . The key is to write  $\bar{\Delta}_T$  as a linear combination of  $\{\xi_i\}_{0 \leq i \leq T-1}$  as follows, which is a direct consequence of the relation (46):

$$\begin{aligned} \bar{\Delta}_T &= \frac{1}{T} \sum_{j=1}^T \Delta_j \\ &= \frac{1}{T} \sum_{j=1}^T (\mathbf{I} - \eta\mathbf{A})^j \Delta_0 - \frac{1}{T} \sum_{j=1}^T \eta \sum_{i=0}^{j-1} (\mathbf{I} - \eta\mathbf{A})^{j-i-1} \xi_i \\ &= \frac{1}{T} \sum_{j=1}^T (\mathbf{I} - \eta\mathbf{A})^j \Delta_0 - \frac{1}{T} \sum_{i=0}^{T-1} \eta \sum_{j=i+1}^T (\mathbf{I} - \eta\mathbf{A})^{j-i-1} \xi_i \\ &= \frac{1}{T\eta} \mathbf{A}_0^{(T+1)} \Delta_0 - \frac{1}{T} \Delta_0 - \frac{1}{T} \sum_{i=0}^{T-1} \mathbf{A}_i^{(T)} \xi_i, \end{aligned} \quad (62)$$

where the middle line follows from swapping the summation over  $i$  and  $j$ , and in the last line we define

$$\mathbf{A}_i^{(t)} := \eta \sum_{j=i+1}^t (\mathbf{I} - \eta\mathbf{A})^{j-i-1} = \mathbf{A}^{-1}(\mathbf{I} - (\mathbf{I} - \eta\mathbf{A})^{t-i}). \quad (63)$$

We claim that the following two inequalities hold, the first deterministically and the second with probability of at least

$$\|\mathbf{A}_0^{(T+1)} \Delta_0\|_{\Sigma} \leq \frac{2\|\Sigma^{-1}\|}{1-\gamma} \|\Delta_0\|_{\Sigma}; \quad (64a)$$

$$\left\| \sum_{i=0}^{T-1} \mathbf{A}_i^{(T)} \xi_i \right\|_{\Sigma} \lesssim \left\{ \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T(1-\gamma)^2}} \right. \quad (64b)$$

$$\left. + \frac{\|\Sigma^{-1}\|}{T} \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)^3}} \right\} (\|\theta^*\|_{\Sigma} + 1). \quad (64c)$$

Putting the above two inequalities together with (62), we arrive at

$$\begin{aligned} \|\bar{\Delta}_T\|_{\Sigma} &\leq \left\| \frac{1}{T\eta} \mathbf{A}_0^{(T+1)} \Delta_0 \right\|_{\Sigma} + \left\| \frac{1}{T} \Delta_0 \right\|_{\Sigma} + \left\| \frac{1}{T} \sum_{i=0}^{T-1} \mathbf{A}_i^{(T)} \xi_i \right\|_{\Sigma} \\ &\lesssim \frac{1}{\eta T} \frac{\|\Sigma^{-1}\|}{1-\gamma} \|\Delta_0\|_{\Sigma} + \frac{1}{T} \|\Delta_0\|_{\Sigma} \\ &\quad + \left\{ \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T(1-\gamma)^2}} \right. \\ &\quad \left. + \frac{\|\Sigma^{-1}\|}{T} \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)^3}} \right\} (\|\theta^*\|_{\Sigma} + 1) \\ &\asymp \frac{1}{\eta T} \frac{\|\Sigma^{-1}\|}{1-\gamma} \|\Delta_0\|_{\Sigma} + \left\{ \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T(1-\gamma)^2}} \right. \\ &\quad \left. + \frac{\|\Sigma^{-1}\|}{T} \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)^3}} \right\} (\|\theta^*\|_{\Sigma} + 1), \end{aligned}$$

where the last line follows since  $\|\Sigma^{-1}\| \geq 1$  (see (88h)) and  $\eta < 1$ . This finishes the proof of Theorem 1.

## VII. PROOF OF THEOREM 3 (TDC LEARNING)

Firstly, let us analyze the population dynamics of TDC. It turns out that the convergence of this dynamics can be described via a contractive linear mapping. Given this nice property of population TDC, we shall decompose the empirical TDC into two parts: the first part can be controlled via the aforementioned population dynamics, and the rest is treated as a stochastic component, which is controlled via matrix martingale concentration.

### A. Population analysis

First recall that the population parameters are defined as

$$\begin{aligned} \tilde{\mathbf{A}} &:= \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\tilde{\mathbf{A}}_t] = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\rho_t \phi(s_t) (\phi(s_t) - \gamma \phi(s'_t))^{\top}]; \\ \tilde{\mathbf{b}} &:= \mathbb{E}_{\mu_b, \pi_b}[\tilde{\mathbf{b}}_t] = \mathbb{E}_{\mu_b, \pi_b}[\rho_t \phi(s_t) r_t]; \\ \Pi &:= \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\Pi_t] = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\rho_t \phi(s_t) \phi(s'_t)^{\top}]; \\ \tilde{\Sigma} &:= \mathbb{E}_{\mu_b}[\Sigma_t] = \mathbb{E}_{\mu_b}[\phi(s_t) \phi(s_t)^{\top}]. \end{aligned}$$

Corresponding to the empirical version of TDC as given in (34), we can define its population analogue of TDC as

$$\begin{aligned} \check{\theta}_{t+1} &= \check{\theta}_t - \alpha(\tilde{\mathbf{A}}\check{\theta}_t - \tilde{\mathbf{b}} + \gamma\Pi^{\top}\check{w}_t), \\ \check{w}_{t+1} &= \check{w}_t - \beta(\tilde{\mathbf{A}}\check{\theta}_t - \tilde{\mathbf{b}} + \tilde{\Sigma}\check{w}_t), \end{aligned} \quad (65)$$

where sampled parameters are replaced by their corresponding expectations. In this section, we analyze the population dynamics of TDC as given above; in order to control the finite-sample dynamics, we bound the difference of these two in the section to follow.

Since  $\phi(s_t)$  is independent of the transition, the expectation of  $\tilde{\Sigma}_t$  is independent of which policy is being adopted. Hence,  $\tilde{\Sigma}$  can also be presented as

$$\begin{aligned}\tilde{\Sigma} &= \sum_{s_t \in \mathcal{S}} \mu_b(s_t) \phi(s_t) \phi(s_t)^\top \\ &= \sum_{s_t \in \mathcal{S}} \mu_b(s_t) \left( \sum_{a_t \in \mathcal{A}} \pi(a_t|s_t) \right) \phi(s_t) \phi(s_t)^\top \\ &= \sum_{s_t \in \mathcal{S}} \sum_{a_t \in \mathcal{A}} \mu_b(s_t) \pi_b(a_t|s_t) \left( \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} \right) \phi(s_t) \phi(s_t)^\top \\ &= \mathbb{E}_{\mu_b, \pi_b} [\rho_t \phi(s_t) \phi(s_t)^\top].\end{aligned}\quad (66)$$

$$(67)$$

In view of this relation,  $\tilde{\mathbf{A}}$  admits another characterization, namely

$$\tilde{\mathbf{A}} = \tilde{\Sigma} - \gamma \Pi. \quad (68)$$

Consequently, the fixed point  $(\tilde{\theta}^*, w^*)$  of the population dynamics obeys

$$\begin{cases} \tilde{\mathbf{A}}\tilde{\theta}^* - \tilde{\mathbf{b}} + \gamma \Pi^\top w^* = \mathbf{0}, \\ \tilde{\mathbf{A}}\tilde{\theta}^* - \tilde{\mathbf{b}} + \tilde{\Sigma} w^* = \mathbf{0}. \end{cases}$$

As long as  $\tilde{\mathbf{A}}$  is invertible, this set of conditions is equivalent to

$$\tilde{\mathbf{A}}\tilde{\theta}^* = \tilde{\mathbf{b}}, \quad \text{and} \quad w^* = \mathbf{0}.$$

In order to study the population dynamics, it is useful to consider two auxiliary parameters

$$\begin{aligned}\check{\Delta}_t &:= \tilde{\theta}_t - \tilde{\theta}^*, \\ \check{z}_t &:= \tilde{w}_t + \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} \check{\Delta}_t;\end{aligned}$$

here  $\check{\Delta}_t$  tracks the convergence of  $\tilde{\theta}_t$  to  $\tilde{\theta}^*$ , and  $\check{z}_t$  tracks the size of the residual  $\tilde{\mathbf{A}}\tilde{\theta}_t - \tilde{\mathbf{b}} + \tilde{\Sigma}\tilde{w}_t$ . With these two parameters in place, the population dynamics satisfy

$$\begin{aligned}\begin{bmatrix} \check{\Delta}_t \\ \check{z}_t \end{bmatrix} &= \\ \begin{bmatrix} \mathbf{I} - \alpha \tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} & -\alpha \gamma \Pi^\top \\ -\alpha (\mathbf{I} - \gamma \Sigma^{-1} \Pi) \tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} & \mathbf{I} - \beta \tilde{\Sigma} - \alpha \gamma (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \Pi) \Pi^\top \end{bmatrix} \begin{bmatrix} \check{\Delta}_{t-1} \\ \check{z}_{t-1} \end{bmatrix}.\end{aligned}\quad (69)$$

To analyze this optimization dynamics, for every positive constant  $\varkappa \in (0, 1)$ , consider

$$\check{x}_t := \begin{bmatrix} \check{\Delta}_t \\ \varkappa \check{z}_t \end{bmatrix}$$

then  $\check{x}_t$  yields

$$\check{x}_t = \Psi \check{x}_{t-1}, \quad (70)$$

where  $\Psi$  represents the matrix

$$\begin{bmatrix} \mathbf{I} - \alpha \tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} & -\frac{1}{\varkappa} \alpha \gamma \Pi^\top \\ -\varkappa \alpha (\mathbf{I} - \gamma \Sigma^{-1} \Pi) \tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} & \mathbf{I} - \beta \tilde{\Sigma} - \alpha \gamma (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \Pi) \Pi^\top \end{bmatrix}. \quad (71)$$

It is known that how fast  $\check{x}_t$  converges to  $\mathbf{0}$  is determined by the spectral norm of  $\Psi$ , which is characterized in the lemma below.

**Lemma 2.** Suppose that

$$\lambda_1 = \lambda_{\min}(\tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}), \quad \lambda_2 = \lambda_{\min}(\tilde{\Sigma}), \quad \lambda_\Sigma = \|\tilde{\Sigma}^{-1}\| = 1/\lambda_2.$$

Then as long as the following conditions hold:

$$\beta \gtrsim \lambda_\Sigma \rho_{\max} \alpha, \quad (72a)$$

$$\varkappa \beta \gtrsim \alpha, \quad (72b)$$

$$\alpha \gamma (\rho_{\max} + \gamma \lambda_\Sigma \rho_{\max}^2) \ll \beta \lambda_w, \quad (72c)$$

$$\frac{\alpha \gamma \rho_{\max}}{\varkappa} + \varkappa \alpha (1 + \gamma \lambda_\Sigma \rho_{\max}) \lambda_\Sigma (2 \rho_{\max})^2 \ll \sqrt{\alpha \lambda_1 \beta \lambda_w} \quad (72d)$$

it holds true that

$$\|\Psi\| \leq 1 - \frac{1}{2} \alpha \lambda_1.$$

Therefore, the mapping  $\Psi$  is contractive, thus ensuring the linear convergence of  $x_t$ , with the proviso that  $\alpha \lambda_1 < 2$ .

## B. Finite-sample analysis

Armed with the population analysis, the proof for Theorem 3 is completed if we can make a connection of the finite-sample performances to that of the population ones.

a) *Step 1: a recursive relation:* Firstly, we define two noise variables

$$\begin{aligned}\nu_t &:= (\tilde{\mathbf{A}}_t - \tilde{\mathbf{A}}) \tilde{\theta}_t - (\tilde{\mathbf{b}}_t - \tilde{\mathbf{b}}) + \gamma (\Pi_t - \Pi)^\top w_t, \\ \eta_t &:= (\tilde{\mathbf{A}}_t - \tilde{\mathbf{A}}) \tilde{\theta}_t - (\tilde{\mathbf{b}}_t - \tilde{\mathbf{b}}) + (\tilde{\Sigma}_t - \tilde{\Sigma}) w_t.\end{aligned}$$

As a result, TDC can be rewritten as

$$\begin{aligned}\tilde{\theta}_{t+1} &= \tilde{\theta}_t - \alpha (\tilde{\mathbf{A}} \tilde{\theta}_t - \tilde{\mathbf{b}} + \gamma \Pi^\top w_t) - \alpha \nu_t; \\ w_{t+1} &= w_t - \beta (\tilde{\mathbf{A}} \tilde{\theta}_t - \tilde{\mathbf{b}} + \tilde{\Sigma} w_t) - \beta \eta_t.\end{aligned}$$

Using the same notations as in Section VII-A, we observe that the following iteration holds true for finite-sample TDC:

$$x_{t+1} = \Psi x_t - \zeta_t,$$

$$\zeta_t = \begin{bmatrix} \alpha \nu_t \\ \varkappa (\alpha (1 - \gamma \tilde{\Sigma}^{-1} \Pi) \nu_t + \beta \eta_t) \end{bmatrix}. \quad (73)$$

Hence,

$$x_t = \Psi^t x_0 - \sum_{i=0}^{t-1} \Psi^{t-i-1} \zeta_i, \quad (74)$$

where  $x_0 = [\Delta_0^\top, \varkappa z_0^\top]^\top$ . Since the norm of  $\Psi$  has been bounded by Lemma 2, bounding the norm of  $x_t$  boils down to bounding the second term of (74). In the following, with

a slight abuse of notation, for any  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^{2d}$  with  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$ , we will define  $\|\mathbf{x}\|_{\tilde{\Sigma}}^2$  as

$$\|\mathbf{x}\|_{\tilde{\Sigma}}^2 = \|\mathbf{x}_1\|_{\tilde{\Sigma}}^2 + \|\mathbf{x}_2\|_{\tilde{\Sigma}}^2.$$

with this definition, it is easy to see that

$$\|\mathbf{x}_t\|_{\tilde{\Sigma}}^2 = \|\tilde{\Delta}_t\|_{\tilde{\Sigma}}^2 + \kappa^2 \|\mathbf{w}_t + \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} \tilde{\Delta}_t\|_{\tilde{\Sigma}}^2.$$

Hence, the norms of  $\tilde{\Delta}_t$ ,  $\mathbf{w}_t$  and  $\mathbf{x}_t$  can be related by the inequalities

$$\begin{cases} \|\tilde{\Delta}_t\|_{\tilde{\Sigma}} \leq \|\mathbf{x}_t\|_{\tilde{\Sigma}}; \\ \|\mathbf{w}_t\|_{\tilde{\Sigma}} \lesssim \frac{1}{\kappa} \|\mathbf{x}_t\|_{\tilde{\Sigma}}; \\ \|\mathbf{x}_t\|_{\tilde{\Sigma}} \lesssim \|\tilde{\Delta}_t\|_{\tilde{\Sigma}} + \|\mathbf{w}_t\|_{\tilde{\Sigma}}. \end{cases} \quad (75)$$

*b) Step 2: crude bound for  $\|\mathbf{x}_t\|_{\tilde{\Sigma}}$ :* We first aim to establish, via an induction argument, that with probability at least  $1 - \delta$ ,

$$\begin{aligned} \|\mathbf{x}_t\|_{\tilde{\Sigma}} &\leq 2\|\tilde{\Delta}_0\|_{\tilde{\Sigma}} + 80\kappa\beta\rho_{\max} \sqrt{\frac{1}{\alpha\lambda_1} \log \frac{2dT}{\delta}} (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + 1) \\ &=: \tilde{R}_0 \end{aligned} \quad (76)$$

holds simulatanesouly for all  $0 \leq t \leq T$ . To start with, note that the inequality (76) holds trivially for the base case with  $t = 0$ . Next, suppose that the hypothesis (76) holds for  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}$ , and we intend to establish it for  $\mathbf{x}_t$  as well. Towards this end, involking the decomposition (74) and the triangle inequality yields

$$\|\mathbf{x}_t\|_{\tilde{\Sigma}} \leq \|\Psi^t \mathbf{x}_0\|_{\tilde{\Sigma}} + \left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \zeta_i \right\|_{\tilde{\Sigma}}. \quad (77)$$

As for the first term of (77), it is seen that

$$\|\Psi^t \mathbf{x}_0\|_{\tilde{\Sigma}} \leq \|\mathbf{x}_0\|_{\tilde{\Sigma}} = \|\tilde{\Delta}_0\|_{\tilde{\Sigma}}. \quad (78)$$

When it comes to the second term of (77), the following lemma comes in handy.

**Lemma 3.** Fix any quantity  $\tilde{R} > 0$  and, for each  $0 \leq i \leq T - 1$ , define the random vector

$$\tilde{\zeta}_i := \zeta_i \mathbf{1}\{\tilde{\mathcal{H}}_i\}, \quad \text{where } \tilde{\mathcal{H}}_i := \left\{ \|\mathbf{x}_i\|_{\tilde{\Sigma}} \leq \tilde{R} \right\}. \quad (79)$$

Then, with probability at least  $1 - \delta/T$ ,

$$\left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \tilde{\zeta}_i \right\|_{\tilde{\Sigma}} \lesssim \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \kappa\beta\rho_{\max} (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{1}{\kappa} \tilde{R} + 1) \quad (80)$$

provided that the stepsizes  $\alpha, \beta$  satisfy the conditions (72) and that  $0 < \alpha < \frac{1}{\lambda_1 \lambda_{\tilde{\Sigma}}^2 \|\tilde{\Sigma}\| \log \frac{2dT}{\delta}}$ .

*Proof.* See Section C-D.  $\square$

Putting relations (77) and (80) together, we find

$$\begin{aligned} \|\mathbf{x}_t\|_{\tilde{\Sigma}} &= \|\tilde{\Delta}_0\|_{\tilde{\Sigma}} \\ &+ C \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \kappa\beta\rho_{\max} (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{1}{\kappa} \tilde{R}_0 + 1) \end{aligned}$$

$$\leq \tilde{R}_0$$

by definition of  $\tilde{R}_0$  in (76), provided that  $\sqrt{\frac{1}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \kappa\beta\rho_{\max} \leq c$  for some constant  $c > 0$  small enough. Therefore, by induction assumption, one has

$$\mathbb{P} \left\{ \max_{0 \leq i \leq t} \|\mathbf{x}_i\|_{\tilde{\Sigma}} > \tilde{R}_0 \right\} \quad (81)$$

$$\begin{aligned} &\leq \mathbb{P} \left\{ \max_{0 \leq i < t-1} \|\mathbf{x}_i\|_{\tilde{\Sigma}} > \tilde{R}_0 \right\} \\ &+ \mathbb{P} \left\{ \max_{0 \leq i < t-1} \|\mathbf{x}_i\|_{\tilde{\Sigma}} \leq \tilde{R}_0, \|\mathbf{x}_t\|_{\tilde{\Sigma}} > \tilde{R}_0 \right\} \\ &\leq \frac{(t-1)\delta}{T} + \mathbb{P} \left\{ \left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \tilde{\zeta}_i \right\|_{\tilde{\Sigma}} \right. \\ &\quad \left. \gtrsim \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \kappa\beta\rho_{\max} (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \tilde{R}_0 + 1) \right\} \\ &\leq \frac{(t-1)\delta}{T} + \frac{\delta}{T} = \frac{t\delta}{T}. \end{aligned} \quad (82)$$

This completes our claim at this step.

*c) Step 3: refined bound for  $\|\mathbf{x}_t\|_{\tilde{\Sigma}}$ :* It turns out that the upper bound (76) can be tightened by taking into account the contraction effect of  $\Psi$ . In what follows, we develop a strengthened bound. Define

$$\tilde{t}_{\text{seg}} := \frac{\tilde{c}_1 \log \max \left\{ \sqrt{\tilde{\kappa}}, \frac{\sqrt{\tilde{\kappa}} \|\tilde{\Delta}_0\|_{\tilde{\Sigma}}}{\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + 1}, \|\tilde{\Delta}_0\|_{\tilde{\Sigma}} \sqrt{\frac{\alpha\lambda_1}{\|\tilde{\Sigma}\| \log \frac{2dT}{\delta}}} \frac{1}{\kappa\beta\rho_{\max}} \right\}}{\alpha\lambda_1} \quad (83)$$

for some sufficiently large constant  $\tilde{c}_1 > 0$ , where  $\tilde{\kappa}$  is the condition number of  $\tilde{\Sigma}$ . For any integer  $k > 1$ , we claim that with probability at least  $1 - \delta$ ,

$$\begin{aligned} \|\mathbf{x}_t\|_{\tilde{\Sigma}} &\lesssim \kappa\beta\rho_{\max} \sqrt{\frac{1}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{\|\tilde{\Delta}_0\|_{\tilde{\Sigma}}}{2^{k-1}} + \frac{3}{2} \right) \\ &=: \tilde{R}_k \end{aligned} \quad (84)$$

for any  $t$  obeying  $k\tilde{t}_{\text{seg}} \leq t \leq T$ , provided that  $\sqrt{\frac{1}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \kappa\beta\rho_{\max} \leq c$  for some constant  $c$  small enough. The proof of this claim is essentially the same as that of Step 3 for proving Theorem 1, and we will omit it here. Therefore, by defining

$$\begin{aligned} \tilde{t}'_{\text{seg}} &:= \left( 2 + \frac{1}{\log 2} \log \|\tilde{\theta}^*\|_{\tilde{\Sigma}} \right) \tilde{t}_{\text{seg}}, \end{aligned} \quad (85)$$

we can conclude that with probability at least  $1 - \delta$ , for all  $t \geq \tilde{t}'_{\text{seg}}$ ,

$$\|\mathbf{x}_t\|_{\tilde{\Sigma}} \lesssim \kappa\beta\rho_{\max} \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha\lambda_1} \log \frac{2dT}{\delta}} (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + 2). \quad (86)$$

Recall that this bound holds for any  $\kappa \in (0, 1)$  satisfying the conditions (72). Hence, Theorem 3 follows by taking  $\kappa =$

$$8\rho_{\max}\sqrt{\frac{\alpha}{\lambda_1\beta\lambda_2}} \text{ and}$$

$$\frac{\alpha}{\beta} = \frac{1}{128} \frac{\lambda_1\lambda_2}{\rho_{\max}^2(1+\lambda_{\Sigma}\rho_{\max})}.$$

### VIII. DISCUSSION

Our primary contribution in this paper is obtaining high-probability sample complexity bounds for both the TD and TDC algorithms for policy evaluation in the  $\gamma$ -discounted infinite-horizon MDPs. For TD learning with Polyak-Ruppert averaging, we improve upon existing results in terms of both the accuracy level  $\varepsilon$  and other problem-related parameters like the effective horizon  $\frac{1}{1-\gamma}$ , the weighted feature covariance  $\Sigma$  and the optimal linear estimator  $\theta^*$ . We have also established a minimax lower bound and showed that our upper bound is near-minimax optimal by a factor of  $\frac{1}{1-\gamma}$ . For TDC with linear function approximation, we provide the first sample complexity bound that achieves the optimal dependence on the error tolerance  $\varepsilon$ , and characterize the exact dependence on problem-related constants at the same time.

Our analysis leaves open several directions for future investigation; we close by sampling a few of them. Regarding TD learning, a natural direction of future work is to close the  $\frac{1}{1-\gamma}$  gap between our upper bound and the minimax lower bound. Notably, this gap also appears in the bounds of [12] for least-square TD in general when no restriction of the variance for the temporal difference residual is imposed. In terms of TDC, while our result provides a tight control of the same size  $T$ , the dependence on problem-related constants can be potentially improved. Moreover, it is noteworthy that the analysis in this work is based on the assumption of *i.i.d.* transition pairs drawn from the stationary distribution; it is of natural interest to generalize these results to other scenarios such as Markovian trajectories. Moving beyond linear function approximation, understanding the sample complexities for policy evaluation with other function classes is also an interesting direction.

### APPENDIX A PRELIMINARY FACTS

The following two lemmas consider the basic properties of important matrices and vectors that would be useful in the proof of the main theorems in the paper.

**Lemma 4.** Recall the definitions of  $\Phi$ ,  $D_\mu$  and  $\Sigma$  in (4), (5) and (6), respectively. Then one has

$$\|D_\mu^{\frac{1}{2}}\Phi\Sigma^{-\frac{1}{2}}\| = 1, \quad \text{and} \quad \|D_\mu^{\frac{1}{2}}P^\pi D_\mu^{-\frac{1}{2}}\| = 1. \quad (87)$$

*Proof.* For notational convenience, let  $\tilde{\Phi} := D_\mu^{\frac{1}{2}}\Phi\Sigma^{-\frac{1}{2}}$  and  $P_{D_\mu} := D_\mu^{\frac{1}{2}}P^\pi D_\mu^{-\frac{1}{2}}$ . First of all, it is seen that

$$\begin{aligned} \|\tilde{\Phi}\| &= \sqrt{\|\tilde{\Phi}^\top \tilde{\Phi}\|} \\ &= \sqrt{\|\Sigma^{-\frac{1}{2}}\Phi^\top D_\mu^{\frac{1}{2}} D_\mu^{\frac{1}{2}} \Phi \Sigma^{-\frac{1}{2}}\|} \\ &= \sqrt{\|\Sigma^{-\frac{1}{2}}\Sigma\Sigma^{-\frac{1}{2}}\|} = 1. \end{aligned}$$

When it comes to  $\|P_{D_\mu}\|$ , we make the observation that

$$\begin{aligned} \|P_{D_\mu}\| &= \sqrt{\|P_{D_\mu} P_{D_\mu}^\top\|} \\ &= \sqrt{\|D_\mu^{\frac{1}{2}} P D_\mu^{-1} P^\top D_\mu^{\frac{1}{2}}\|} \\ &= \sqrt{\|D_\mu^{\frac{1}{2}} (P D_\mu^{-1} P^\top D_\mu) D_\mu^{-\frac{1}{2}}\|} = 1. \end{aligned}$$

To see why the last identity holds, observe that  $P D_\mu^{-1} P^\top D_\mu$  is a stochastic matrix, that is  $P D_\mu^{-1} P^\top D_\mu$  contains nonnegative elements, and

$$P D_\mu^{-1} P^\top D_\mu \mathbf{1} = \mathbf{1}.$$

In addition,  $D_\mu^{\frac{1}{2}} (P D_\mu^{-1} P^\top D_\mu) D_\mu^{-\frac{1}{2}}$  is similar to  $P D_\mu^{-1} P^\top D_\mu$ . As a result, by the Perron-Frobenius theorem,

$$\begin{aligned} &\|D_\mu^{\frac{1}{2}} (P D_\mu^{-1} P^\top D_\mu) D_\mu^{-\frac{1}{2}}\| \\ &= \max_i |\lambda_i(D_\mu^{\frac{1}{2}} (P D_\mu^{-1} P^\top D_\mu) D_\mu^{-\frac{1}{2}})| \\ &= \max_i |\lambda_i(P D_\mu^{-1} P^\top D_\mu)| = 1, \end{aligned}$$

where  $\lambda_i(B)$  denotes the  $i$ -th eigenvalue of the matrix  $B$ .  $\square$

**Lemma 5.** Suppose that  $\|r\|_\infty \leq 1$ . For any  $0 \leq \gamma < 1$ , the matrix  $\Sigma$  defined in (6) and the vector  $b$  defined in (12c) obey

$$\Sigma^{-\frac{1}{2}} A^\top \Sigma^{-1} A \Sigma^{-\frac{1}{2}} \succeq (1-\gamma)^2 \mathbf{I}, \quad (88a)$$

$$\Sigma^{-\frac{1}{2}} A \Sigma^{-1} A^\top \Sigma^{-\frac{1}{2}} \succeq (1-\gamma)^2 \mathbf{I}, \quad (88b)$$

$$\|\Sigma^{\frac{1}{2}} (A^\top)^{-1} \Sigma A^{-1} \Sigma^{\frac{1}{2}}\| \leq (1-\gamma)^{-2}, \quad (88c)$$

$$\|\Sigma^{\frac{1}{2}} A^{-1} \Sigma (A^\top)^{-1} \Sigma^{\frac{1}{2}}\| \leq (1-\gamma)^{-2}, \quad (88d)$$

$$\|\Sigma^{\frac{1}{2}} A^{-1} \Sigma^{\frac{1}{2}}\| \leq (1-\gamma)^{-1}, \quad (88e)$$

$$\|\Sigma^{-1/2} \Phi^\top D_\mu\| \leq \max_{s \in \mathcal{S}} \phi(s)^\top \Sigma^{-1} \phi(s), \quad (88f)$$

$$\|I - \eta A\| \leq 1 - \frac{1}{2} \eta (1-\gamma) \lambda_{\min}(\Sigma), \quad \forall 0 < \eta < \frac{1-\gamma}{4\|\Sigma\|}, \quad (88g)$$

$$\|\Sigma\| \leq 1, \quad \|\Sigma^{-1}\| \geq 1, \quad (88h)$$

$$\|\Sigma^{-\frac{1}{2}} b\|_2 \leq 1. \quad (88i)$$

*Proof.* We shall establish each of these claims separately as follows.

a) *Proof of Eqn. (88a) and (88b):* We start with the lower bound on  $\Sigma^{-\frac{1}{2}} A^\top \Sigma^{-1} A \Sigma^{-\frac{1}{2}}$ . To begin with, observe that

$$\begin{aligned} \Sigma^{-\frac{1}{2}} A \Sigma^{-\frac{1}{2}} &= \Sigma^{-\frac{1}{2}} \Phi^\top D_\mu (I - \gamma P) \Phi \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} \Phi^\top D_\mu \Phi \Sigma^{-\frac{1}{2}} \\ &\quad - \gamma \Sigma^{-\frac{1}{2}} \Phi^\top D_\mu^{\frac{1}{2}} \left( D_\mu^{\frac{1}{2}} P D_\mu^{-\frac{1}{2}} \right) D_\mu^{\frac{1}{2}} \Phi \Sigma^{-\frac{1}{2}} \\ &= I - \gamma \tilde{\Phi}^\top P_{D_\mu} \tilde{\Phi}, \end{aligned}$$

where

$$\tilde{\Phi} := D_\mu^{\frac{1}{2}} \Phi \Sigma^{-\frac{1}{2}} \quad \text{and} \quad P_{D_\mu} := D_\mu^{\frac{1}{2}} P D_\mu^{-\frac{1}{2}}. \quad (89)$$



Therefore, any unit vector  $\mathbf{x}$  (i.e.  $\|\mathbf{x}\|_2 = 1$ ) necessarily satisfies

$$\begin{aligned} \mathbf{x}^\top \Sigma^{-\frac{1}{2}} \mathbf{A}^\top \Sigma^{-1} \mathbf{A} \Sigma^{-\frac{1}{2}} \mathbf{x} &= \|\Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-\frac{1}{2}} \mathbf{x}\|_2^2 \\ &\geq (\mathbf{x}^\top \Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-\frac{1}{2}} \mathbf{x})^2 \\ &= (1 - \gamma \mathbf{x}^\top \tilde{\Phi}^\top P_{D_\mu} \tilde{\Phi} \mathbf{x})^2. \end{aligned}$$

Further, Lemma 4 tells us that

$$|\mathbf{x}^\top \tilde{\Phi}^\top P_{D_\mu} \tilde{\Phi} \mathbf{x}| \leq \|\tilde{\Phi}^\top P_{D_\mu} \tilde{\Phi}\| \leq \|\tilde{\Phi}\|^2 \|P_{D_\mu}\| = 1. \quad (90)$$

Putting the preceding two bounds together, we demonstrate that

$$\mathbf{x}^\top \Sigma^{-\frac{1}{2}} \mathbf{A}^\top \Sigma^{-1} \mathbf{A} \Sigma^{-\frac{1}{2}} \mathbf{x} \geq (1 - \gamma)^2$$

for any unit vector  $\mathbf{x}$ , thus concluding the proof of (88a). The proof for (88b) follows from an identical argument and is omitted for brevity.

b) *Proof of Eqn. (88c), (88d) and (88e):* With the above bounds in place, we can further obtain

$$\begin{aligned} &\|\Sigma^{\frac{1}{2}} (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \\ &= \|(\Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-1} \mathbf{A}^\top \Sigma^{-\frac{1}{2}})^{-1}\| \\ &\leq \frac{1}{\lambda_{\min}(\Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-1} \mathbf{A}^\top \Sigma^{-\frac{1}{2}})} \leq \frac{1}{(1 - \gamma)^2}, \end{aligned}$$

where  $\lambda_{\min}(\mathbf{B})$  denotes the smallest eigenvalue of  $\mathbf{B}$ , and the last inequality comes from (88b). This establishes (88c). The inequality (88d) follows from a similar argument. This also implies that

$$\|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| = \sqrt{\|\Sigma^{\frac{1}{2}} (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\|} \leq \frac{1}{1 - \gamma},$$

as claimed in (88e).

c) *Proof of Eqn. (88g):* Recalling that  $\Sigma = \Phi^\top D_\mu \Phi$ , we can arrange terms to derive

$$\begin{aligned} \mathbf{A} + \mathbf{A}^\top &= \Phi^\top D_\mu (\mathbf{I} - \gamma P) \Phi + \Phi^\top (\mathbf{I} - \gamma P^\top) D_\mu \Phi \\ &= 2\Sigma - \gamma \Sigma^{\frac{1}{2}} \left\{ \Sigma^{-\frac{1}{2}} \Phi^\top D_\mu P \Phi \Sigma^{-\frac{1}{2}} \right. \\ &\quad \left. + \Sigma^{-\frac{1}{2}} \Phi^\top P^\top D_\mu \Phi \Sigma^{-\frac{1}{2}} \right\} \Sigma^{\frac{1}{2}} \\ &= \Sigma^{\frac{1}{2}} \left\{ 2\mathbf{I} - \gamma (\tilde{\Phi}^\top P_{D_\mu} \tilde{\Phi} + \tilde{\Phi}^\top P_{D_\mu}^\top \tilde{\Phi}) \right\} \Sigma^{\frac{1}{2}} \\ &\succeq \Sigma^{\frac{1}{2}} \left\{ 2\mathbf{I} - 2\gamma \|\tilde{\Phi}^\top P_{D_\mu} \tilde{\Phi}\| \mathbf{I} \right\} \Sigma^{\frac{1}{2}} \\ &\succeq 2(1 - \gamma) \Sigma, \end{aligned}$$

where  $\tilde{\Phi}$  and  $P_{D_\mu}$  are defined in (89). Here, the last line follows since  $\|\tilde{\Phi}^\top P_{D_\mu} \tilde{\Phi}\| \leq 1$  — a fact that has already been shown in (90). In addition, the following identity

$$\mathbf{A} \mathbf{A}^\top = \Sigma^{\frac{1}{2}} \tilde{\Phi}^\top (\mathbf{I} - \gamma P_{D_\mu}) \tilde{\Phi} \Sigma \tilde{\Phi}^\top (\mathbf{I} - \gamma P_{D_\mu}^\top) \tilde{\Phi} \Sigma^{\frac{1}{2}}$$

allows us to bound

$$\begin{aligned} \|\Sigma^{-\frac{1}{2}} \mathbf{A} \mathbf{A}^\top \Sigma^{-\frac{1}{2}}\| &= \|\tilde{\Phi}^\top (\mathbf{I} - \gamma P_{D_\mu}) \tilde{\Phi} \Sigma \tilde{\Phi}^\top (\mathbf{I} - \gamma P_{D_\mu}^\top) \tilde{\Phi}\| \\ &\leq \|\mathbf{I} - \gamma P_{D_\mu}\|^2 \|\tilde{\Phi}\|^4 \|\Sigma\| \end{aligned}$$

$$\begin{aligned} &= \|\mathbf{I} - \gamma P_{D_\mu}\|^2 \|\Sigma\| \\ &\leq (1 + \gamma \|P_{D_\mu}\|)^2 \|\Sigma\| \leq 4\|\Sigma\|, \end{aligned}$$

where the last line makes use of Lemma 4. This essentially tells us that

$$\begin{aligned} \mathbf{0} &\preceq \Sigma^{-\frac{1}{2}} \mathbf{A} \mathbf{A}^\top \Sigma^{-\frac{1}{2}} \preceq 4\|\Sigma\| \mathbf{I} \\ \implies \mathbf{A} \mathbf{A}^\top &\preceq 4\|\Sigma\| \Sigma. \end{aligned}$$

Putting the preceding bounds together implies that: for any  $0 < \eta < \frac{1-\gamma}{4\|\Sigma\|}$  one has

$$\begin{aligned} \mathbf{0} &\preceq (\mathbf{I} - \eta \mathbf{A}) (\mathbf{I} - \eta \mathbf{A}^\top) = \mathbf{I} - \eta(\mathbf{A} + \mathbf{A}^\top) + \eta^2 \mathbf{A} \mathbf{A}^\top \\ &\preceq \mathbf{I} - 2\eta(1 - \gamma) \Sigma + 4\eta^2 \|\Sigma\| \Sigma \\ &= \mathbf{I} - \{2\eta(1 - \gamma) - 4\eta^2 \|\Sigma\|\} \Sigma \\ &\preceq \mathbf{I} - \eta(1 - \gamma) \Sigma \\ &\preceq (1 - \eta(1 - \gamma) \lambda_{\min}(\Sigma)) \mathbf{I}, \end{aligned}$$

thus indicating that

$$\begin{aligned} \|\mathbf{I} - \eta \mathbf{A}\| &\leq \sqrt{\|(\mathbf{I} - \eta \mathbf{A}) (\mathbf{I} - \eta \mathbf{A}^\top)\|} \\ &\leq \sqrt{1 - \eta(1 - \gamma) \lambda_{\min}(\Sigma)} \\ &\leq 1 - \frac{1}{2} \eta(1 - \gamma) \lambda_{\min}(\Sigma). \end{aligned}$$

d) *Proof of Eqn. (88h):* For any unit vector  $\mathbf{u}$ , the assumption  $\max_s \|\phi(s)\|_2 \leq 1$  guarantees that

$$\|\Phi \mathbf{u}\|_\infty \leq \max_s |\phi(s)^\top \mathbf{u}| \leq \max_s \|\phi(s)\|_2 \|\mathbf{u}\|_2 \leq 1,$$

where in the last inequality we have used Cauchy-Schwartz inequality. Consequently, for any unit vector  $\mathbf{u}$ , by Hölder's inequality,

$$\mathbf{u}^\top \Phi^\top D_\mu \Phi \mathbf{u} \leq \|\Phi \mathbf{u}\|_\infty \cdot \mathbf{1}^\top D_\mu \mathbf{1} \leq 1,$$

thus proving that  $\|\Sigma\| \leq 1$ . This immediately implies that  $\|\Sigma^{-1}\| \geq 1/\|\Sigma\| \geq 1$ .

e) *Proof of Eqn. (88i):* Finally, we observe that

$$\begin{aligned} \|\Sigma^{-\frac{1}{2}} \mathbf{b}\|_2 &= \|\Sigma^{-\frac{1}{2}} \Phi^\top D_\mu^{\frac{1}{2}} D_\mu^{\frac{1}{2}} \mathbf{r}\|_2 \\ &\leq \|\Sigma^{-\frac{1}{2}} \Phi^\top D_\mu^{\frac{1}{2}}\| \cdot \|D_\mu^{\frac{1}{2}} \mathbf{r}\|_2 \\ &\stackrel{(i)}{\leq} \|D_\mu^{\frac{1}{2}} \mathbf{r}\|_2 \leq 1 \end{aligned}$$

as claimed. Here, (i) follows from Lemma 4 and (ii) holds true since  $\|D_\mu^{\frac{1}{2}} \mathbf{r}\|_2 = \sqrt{\sum_s \mu(s) (r(s))^2} \leq \sqrt{\sum_s \mu(s)} = 1$ .  $\square$

The following lemmas, about the concentration of  $\hat{\mathbf{A}}$ , will be useful in our analysis.

**Lemma 6.** Consider any  $0 < \delta < 1$ , and suppose that  $T \gtrsim \log(\frac{d}{\delta})$ . Then the vector  $\mathbf{b}$  defined in (12c) obeys that, with probability exceeding  $1 - \delta$ ,

$$\begin{aligned} &\|\mathbf{A}^{-1}(\hat{\mathbf{b}} - \mathbf{b})\|_\Sigma \\ &\lesssim \sqrt{\frac{\max_{s \in \mathcal{S}} \phi(s)^\top \Sigma^{-1} \phi(s)}{T(1 - \gamma)^2} \log\left(\frac{d}{\delta}\right)}. \end{aligned}$$

*Proof.* See Section C-E.  $\square$

**Lemma 7.** For any  $0 < \delta < 1$ , it follows that  $\hat{\mathbf{A}}$  is invertible and that

$$\begin{aligned} & \|\Sigma^{1/2} \mathbf{A}^{-1} (\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-1/2}\| \\ & \lesssim \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{T(1-\gamma)^2} \log\left(\frac{d}{\delta}\right)} \end{aligned}$$

with probability at least  $1 - \delta$ , as long as  $T \geq c_2 \max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log(\frac{d}{\delta})$  for some sufficiently large constant  $c_2 > 0$ .

*Proof.* See Section C-E.  $\square$

## APPENDIX B

### PROOF OF THEOREM 2 (MINIMAX LOWER BOUNDS)

This theorem is proved by constructing a set of MDP instances that are hard to distinguish among each other. Based on this construction, the estimation error can be lower bounded via Fano's inequality, which reduces to control the KL-divergence between marginal likelihood functions. We start by constructing a sequence of hard MDP instances.

*a) Construction of MDP instances and their properties:* Given the state space  $\mathcal{S}$ , define a sequence of MDP  $\{\mathcal{M}_{\mathbf{q}}\}$  indexed by  $\mathbf{q} \in \mathcal{Q} \subset \{q_+, q_-\}^{d-1}$  where for each  $\mathbf{q}$ , the transition kernel equals to

$$\begin{aligned} P_{\mathbf{q}}(s' | s) &= \begin{cases} q_s \mathbb{1}(s' = s) + \frac{1-q_s}{|\mathcal{S}|-d+1} \mathbb{1}(s' \geq d) & \text{for } s < d; \\ \frac{\gamma}{|\mathcal{S}|-d+1} \mathbb{1}(s' \geq d) + \frac{1-q_{s'}}{d-1} \mathbb{1}(s' < d) & \text{for } s \geq d. \end{cases} \end{aligned} \quad (91)$$

and the reward function equals to  $r(s) = \mathbb{1}(s \geq d)$ .

Here, for each  $i \in [d-1]$ ,  $q_i$  is taken to be either  $q_+$  or  $q_-$  where

$$q_+ := \gamma + (1-\gamma)^2 \varepsilon, \quad \text{and} \quad q_- := \gamma - (1-\gamma)^2 \varepsilon.$$

We further impose the constraint that the number of  $q_+$ 's and  $q_-$ 's in  $\mathbf{q}$  are the same, namely,

$$\sum_{s=1}^{d-1} \mathbb{1}(q_s = q_+) = \sum_{s=1}^{d-1} \mathbb{1}(q_s = q_-) = (d-1)/2. \quad (92)$$

Here without loss of generality, assume  $d$  is an odd number. With these definitions in place, it can be easily verified that the stationary distribution for  $\mathbf{P}$  obeys

$$\mu(s) = \begin{cases} \frac{1}{2(d-1)} & \text{for } s < d; \\ \frac{1}{2(|\mathcal{S}|-d+1)} & \text{for } s \geq d. \end{cases} \quad (93)$$

Moreover, suppose the feature map is taken to be

$$\phi(s) = \mathbf{e}_{s \wedge d} \in \mathbb{R}^d,$$

then one can further verify that

$$\theta^*(d) = V^*(s) = \frac{1}{1-\gamma^2 - \sum_{i=1}^{d-1} \frac{\gamma^2(1-q_i)^2}{(d-1)(1-\gamma q_i)}}, \quad (94)$$

$$\theta^*(i) = V^*(i) = \frac{\gamma(1-q_i)}{1-\gamma q_i} V^*(s), \text{ for } s \geq d \text{ and } i < d. \quad (95)$$

From the expressions above, we remark that, the values of  $q$  and  $V^*(s)$  with  $s \geq d$  are fixed for all  $\mathbf{q} \in \mathcal{Q}$  which is ensured by the construction (92).

*b) Calculations of several key quantities:* Based on the above constructions, let us compute several key quantities. To begin with, some direct algebra leads to

$$\Sigma = \Phi^\top D_\mu \Phi = \sum_{s=1}^{d-1} \frac{1}{2(d-1)} \mathbf{e}_s \mathbf{e}_s^\top + \frac{1}{2} \mathbf{e}_d \mathbf{e}_d^\top,$$

as well as

$$\phi(s)^\top \Sigma^{-1} \phi(s) = \begin{cases} 2(d-1) & \text{for } s < d; \\ 2 & \text{for } s \geq d. \end{cases}$$

As a consequence, one has

$$\max_s \{\phi(s)^\top \Sigma^{-1} \phi(s)\} \asymp d. \quad (96)$$

Next, we move on to compute  $\|\theta^*\|_\Sigma$ . First notice that for  $\varepsilon \leq \frac{c_1 \gamma}{1-\gamma}$  with constant  $c_1$  small enough,  $(1-\gamma)^2 \varepsilon \leq c_1 \gamma (1-\gamma)$  and hence,  $1-\gamma q_+, 1-\gamma q_- \asymp 1-\gamma$ , which guarantees that  $V^*(s) \asymp \frac{1}{1-\gamma}$ . In view of these calculations, it satisfies that

$$\|\theta^*\|_\Sigma^2 = \sum_{i=1}^{d-1} \frac{1}{2(d-1)} \theta^{*2}(i) + \frac{1}{2} \theta^{*2}(d) \quad (97)$$

$$\begin{aligned} &= \sum_{i=1}^{d-1} \frac{1}{2(d-1)} \left[ \frac{\gamma(1-q_i)}{1-\gamma q_i} V^*(s) \right]^2 + \frac{1}{2} [V^*(s)]^2 \\ &\asymp \sum_{i=1}^{d-1} \frac{1}{2(d-1)} \left[ \frac{\gamma(1-\gamma)}{1-\gamma} \frac{1}{1-\gamma} \right]^2 + \frac{1}{2} \left[ \frac{1}{1-\gamma} \right]^2 \\ &\asymp \frac{1}{(1-\gamma)^2}. \end{aligned} \quad (98)$$

*c) Application of Fano's inequality:* Armed with the properties derived above, we are ready to establish the desired lower bound. First notice that for  $\mathbf{q}, \mathbf{q}' \in \mathcal{Q}$ , if at some  $i \in [d-1]$ ,  $q_i \neq q'_i$ , then

$$\begin{aligned} |\theta^*(i) - \theta'^*(i)| &= \gamma V^*(s) \left| \frac{1-q_i}{1-\gamma q_i} - \frac{1-q'_i}{1-\gamma q'_i} \right| \\ &= \gamma V^*(s) \frac{2\varepsilon(1-\gamma)^3}{(1-\gamma q_i)(1-\gamma q'_i)} \\ &\gtrsim (2\gamma) \frac{1}{1-\gamma} \frac{\varepsilon(1-\gamma)^3}{(1-\gamma)^2} \gtrsim \varepsilon, \end{aligned}$$

where the penultimate inequality follows from  $V^*(s) \asymp \frac{1}{1-\gamma}$ . Consequently, we can bound  $\|\theta^* - \theta'^*\|_\Sigma^2$  as

$$\begin{aligned} \|\theta^* - \theta'^*\|_\Sigma^2 &\geq \sum_{s=1}^{d-1} |\theta^*(s) - \theta'^*(s)|^2 \frac{1}{2(d-1)} \\ &\gtrsim \varepsilon^2 \frac{1}{d-1} \sum_{s=1}^{d-1} \mathbb{1}(q_s \neq q'_s). \end{aligned}$$

This relation guarantees that if  $\sum_{s=1}^{d-1} \mathbb{1}(q_s \neq q'_s) \geq (d-1)/16$ , one has

$$\|\theta^* - \theta'^*\|_\Sigma \gtrsim \varepsilon. \quad (99)$$

In other words, if we want each  $\theta^*$  to be  $\varepsilon$  apart from each other, it is sufficient to construct a set  $\mathcal{Q}$  where every  $q$  and  $q'$  are  $(d-1)/16$  apart in Hamming distance. By virtue of the Gilbert-Varshamov lemma [66], there exists a set  $\mathcal{Q}$  such that

$$M := |\mathcal{Q}| \geq e^{d/16} \quad \text{and} \quad \sum_{s=1}^{d-1} \mathbb{1}(q_s \neq q'_s) \geq \frac{d}{16} \quad (100)$$

for any  $q, q' \in \mathcal{Q}$  obeying  $q \neq q'$ .

The Fano method transforms the problem of estimating  $\theta^*$  into an  $M$ -ary testing problem among the above MDPs  $\{\mathbb{P}_{q^1}, \mathbb{P}_{q^2}, \dots, \mathbb{P}_{q^M}\}$ . More specifically, in view of Fano's inequality ([67]), the probability of interest thus satisfies

$$\begin{aligned} & \mathbb{P}(\|\hat{\theta} - \theta^*\|_{\Sigma} \gtrsim \varepsilon) \\ & \geq 1 - \frac{1}{\log M} \left( \frac{1}{M^2} \sum_{j,k=1}^M \text{KL}(\mathbb{P}_{q^j}^T \parallel \mathbb{P}_{q^k}^T) + \log 2 \right), \quad (101) \end{aligned}$$

given  $T$  independent sample pairs  $\{(s_t, s'_t)\}_{t=1}^T$ . To control the right hand side, we proceed by computing the KL-divergence between every  $\mathbb{P}_q$  and  $\mathbb{P}_{q'}$ . Here  $\mathbb{P}_q$  denotes the joint distribution of  $(s, s')$  when the transition is made according to  $P_q(s'|s)$  (cf. (91)). More specifically, given  $s \sim \mu_q$  and  $s'|s \sim P_q(s'|s)$ , one has

$$\begin{aligned} \mathbb{P}_q(s, s') &= \mu(s)P(s'|s) \\ &= \begin{cases} \frac{1}{2(d-1)} q_s \mathbb{1}(s' = s), & \text{for } s < d, s' < d; \\ \frac{1-q_s}{2(d-1)(S-d+1)}, & \text{for } s < d, s' > d; \\ \frac{1-q_{s'}}{2(d-1)(S-d+1)}, & \text{for } s > d, s' < d; \\ \frac{\gamma}{2(S-d+1)^2}, & \text{for } s > d, s' > d. \end{cases} \end{aligned}$$

Recognizing the relation between the KL divergence and the  $\chi^2$  divergence,  $\text{KL}(\mathbb{P}_q \parallel \mathbb{P}_{q'})$  satisfies

$$\begin{aligned} & \text{KL}(\mathbb{P}_q \parallel \mathbb{P}_{q'}) \\ & \leq \chi^2(\mathbb{P}_{q'} \parallel \mathbb{P}_q) \\ & = \sum_{s, s'} \frac{(\mathbb{P}_q(s, s') - \mathbb{P}_{q'}(s, s'))^2}{\mathbb{P}_q(s, s')} \\ & = \sum_{s < d, s' < d} \frac{(\mathbb{P}_q(s, s') - \mathbb{P}_{q'}(s, s'))^2}{\mathbb{P}_q(s, s')} \\ & + \sum_{s < d, s' \geq d} \frac{(\mathbb{P}_q(s, s') - \mathbb{P}_{q'}(s, s'))^2}{\mathbb{P}_q(s, s')} \\ & + \sum_{s \geq d, s' < d} \frac{(\mathbb{P}_q(s, s') - \mathbb{P}_{q'}(s, s'))^2}{\mathbb{P}_q(s, s')} \\ & + \sum_{s \geq d, s' \geq d} \frac{(\mathbb{P}_q(s, s') - \mathbb{P}_{q'}(s, s'))^2}{\mathbb{P}_q(s, s')} \\ & = \sum_{s=1}^{d-1} \frac{1}{2(d-1)} \frac{(q_s - q'_s)^2}{q_s} \\ & + \sum_{s < d, s' \geq d} \frac{1}{2(d-1)(S-d+1)} \frac{[(1-q_s) - (1-q'_s)]^2}{1-q_s} \\ & + \sum_{s \geq d, s' < d} \frac{1}{2(d-1)(S-d+1)} \frac{[(1-q_{s'}) - (1-q'_{s'})]^2}{1-q_{s'}} + 0 \end{aligned}$$

$$\begin{aligned} & \lesssim \sum_{s=1}^{d-1} \frac{1}{2(d-1)} \frac{[2\varepsilon(1-\gamma)^2]^2}{1-\gamma} \\ & + \sum_{s < d} \frac{1}{2(d-1)} \frac{[2\varepsilon(1-\gamma)^2]^2}{1-\gamma} \\ & + \sum_{s' < d} \frac{1}{2(d-1)} \frac{[2\varepsilon(1-\gamma)^2]^2}{1-\gamma} \\ & \asymp \varepsilon^2(1-\gamma)^3. \end{aligned}$$

As a result, we have

$$\text{KL}(\mathbb{P}_q^T \parallel \mathbb{P}_{q'}^T) \lesssim \varepsilon^2(1-\gamma)^3 T. \quad (102)$$

Substituting the above relation into (101) gives

$$\mathbb{P}(\|\hat{\theta} - \theta^*\|_{\Sigma} \gtrsim \varepsilon) \geq 1 - \frac{1}{d/16} (c\varepsilon^2(1-\gamma)^3 T + \log 2).$$

To prove Theorem 2, it is enough to take the above together with relations (96) and (98).

## APPENDIX C

### PROOFS OF AUXILIARY LEMMAS AND CLAIMS

#### A. Proof of Lemma 1

Here and throughout, we denote by  $\mathbb{E}_i[\cdot]$  the expectation conditioned on the probability space generated by the samples  $\{(s_j, s'_j)\}_{j \leq i}$  (more formally,  $\mathbb{E}_i[\cdot]$  represents the expectation conditioned on the filtration  $\mathcal{F}_i$  — the  $\sigma$ -algebra generated by  $\{(s_j, s'_j)\}_{j \leq i}$ ). It is then easy to check that  $\{(\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i\}$  forms a martingale difference sequence, which motivates us to apply matrix Freedman's inequality.

To this end, one needs to upper bound the following two quantities

$$\begin{aligned} W &:= \sum_{i=l}^u \mathbb{E}_{i-1} \left[ \|\Sigma^{1/2}(\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\} \right], \quad \text{and} \\ B &:= \max_{i: l \leq i \leq u} \|\Sigma^{1/2}(\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i\|_2 \mathbb{1}\{\mathcal{H}_i\}, \quad (103) \end{aligned}$$

which we accomplish in the sequel. For notational convenience, we set

$$\alpha := \left(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma)\right)^{t-u-1}. \quad (104)$$

a) *Control of  $W$* : Direct calculations yield

$$\begin{aligned} W &= \sum_{i=l}^u \mathbb{E}_{i-1} \left[ \xi_i^\top (\mathbf{I} - \eta \mathbf{A}^\top)^{t-i-1} \Sigma (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \mathbb{1}\{\mathcal{H}_i\} \right] \\ &\leq \sum_{i=l}^u \left\| (\mathbf{I} - \eta \mathbf{A}^\top)^{t-i-1} \Sigma (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \right\| \\ &\quad \cdot \mathbb{E}_{i-1} \left[ \|\xi_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\} \right] \\ &\stackrel{(i)}{\leq} \sum_{i=l}^u \|\Sigma\| \left( 1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma) \right)^{2t-2i-2} \\ &\quad \cdot 2 \max_{i: l \leq i \leq u} \left\{ \mathbb{E}_{i-1} \left[ \|(\mathbf{A}_i - \mathbf{A})\theta_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\} \right] \right. \\ &\quad \left. + \mathbb{E}_{i-1} \left[ \|\mathbf{b}_i - \mathbf{b}\|_2^2 \right] \right\} \\ &\stackrel{(ii)}{\leq} \frac{4\|\Sigma\|\alpha^2}{\eta(1-\gamma)\lambda_{\min}(\Sigma)} \cdot \max_{i: l \leq i \leq u} \left\{ \mathbb{E}_{i-1} \left[ \|(\mathbf{A}_i - \mathbf{A})\theta_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\} \right] \right\} \end{aligned}$$

$$+ \mathbb{E}_{i-1} [\|\mathbf{b}_i - \mathbf{b}\|_2^2], \quad (105)$$

where (i) follows from the property (88g) (together with the assumption  $\eta < (1 - \gamma)/(4\|\Sigma\|)$ ) and the elementary inequality  $\|\mathbf{a} + \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$ , and (ii) uses the elementary upper bound for the sum of geometric series as well as the definition (104) of  $\alpha$ .

We then turn attention to  $\mathbb{E}_{i-1} [\|(\mathbf{A}_i - \mathbf{A})\boldsymbol{\theta}_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\}]$  and  $\mathbb{E}_{i-1} [\|\mathbf{b}_i - \mathbf{b}\|_2^2]$ . First, given that  $\mathbb{E}_{i-1}[\mathbf{A}_i\boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}] = \mathbf{A}\boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}$ , one can derive

$$\begin{aligned} & \mathbb{E}_{i-1} [\|(\mathbf{A}_i - \mathbf{A})\boldsymbol{\theta}_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\}] \\ & \leq \mathbb{E}_{i-1} [\|\mathbf{A}_i\boldsymbol{\theta}_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\}] \\ & = \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top (\phi(s_i) - \gamma\phi(s'_i)) \phi(s_i)^\top \\ & \quad \phi(s_i) (\phi(s_i) - \gamma\phi(s'_i))^\top \boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}] \\ & \leq \max_s \|\phi(s)\|_2^2 \cdot \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top (\phi(s_i) - \gamma\phi(s'_i)) \\ & \quad (\phi(s_i) - \gamma\phi(s'_i))^\top \boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}] \\ & \stackrel{(i)}{\leq} 2 \max_s \|\phi(s)\|_2^2 \left( \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \phi(s_i) \phi(s_i)^\top \boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}] \right. \\ & \quad \left. + \gamma^2 \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \phi(s'_i) \phi(s'_i)^\top \boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}] \right) \\ & \stackrel{(ii)}{=} 2 \max_s \|\phi(s)\|_2^2 \left( \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \Sigma \boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}] \right. \\ & \quad \left. + \gamma^2 \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \Sigma \boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}] \right) \\ & \stackrel{(iii)}{\leq} 4 \|\boldsymbol{\theta}_i\|_2^2 \mathbb{1}\{\mathcal{H}_i\} \leq 4(\|\boldsymbol{\theta}^*\|_\Sigma + \|\boldsymbol{\Delta}_i\|_\Sigma)^2 \mathbb{1}\{\mathcal{H}_i\} \\ & \leq 4(\|\boldsymbol{\theta}^*\|_\Sigma + R)^2, \end{aligned} \quad (106)$$

where (i) relies on the elementary inequality  $(\mathbf{a} + \mathbf{b})(\mathbf{a} + \mathbf{b})^\top \preceq 2\mathbf{a}\mathbf{a}^\top + 2\mathbf{b}\mathbf{b}^\top$ , (ii) follows from the definition (6) of  $\Sigma$  and the fact that  $s_i, s'_i \sim \mu$  in this case, (iii) holds due to the assumption  $\max_s \|\phi(s)\|_2 \leq 1$ , and the last inequality results from the definition (50) of the event  $\mathcal{H}_i$ . Similarly, one can derive

$$\mathbb{E}_{i-1} [\|\mathbf{b}_i - \mathbf{b}\|_2^2] \leq \mathbb{E}_{i-1} [\|\mathbf{b}_i\|_2^2] = \mathbb{E}_{i-1} [\|\phi(s_i)r(s_i)\|_2^2] \leq 1, \quad (107)$$

where the last inequality holds since  $\max_s \|\phi(s)\|_2 \leq 1$  and  $\max_s |r(s)| \leq 1$ . Substitution into (105) yields

$$W \leq \frac{4\kappa}{\eta(1-\gamma)} \alpha^2 \left\{ 4(\|\boldsymbol{\theta}^*\|_\Sigma + R)^2 + 1 \right\} =: W_{\max}. \quad (108)$$

b) *Control of B*: By definition of  $B$ , one can write

$$\begin{aligned} B &= \max_{i:l \leq i \leq u} \left\| \Sigma^{\frac{1}{2}} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \boldsymbol{\xi}_i \right\|_2 \mathbb{1}\{\mathcal{H}_i\} \\ &= \max_{i:l \leq i \leq u} \left\| \Sigma^{\frac{1}{2}} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \Sigma^{\frac{1}{2}} \Sigma^{-\frac{1}{2}} \boldsymbol{\xi}_i \right\|_2 \mathbb{1}\{\mathcal{H}_i\} \\ &\leq \|\Sigma\| \max_{i:l \leq i \leq u} \left\| \mathbf{I} - \eta \mathbf{A} \right\|^{t-i-1} \cdot \max_{i:l \leq i \leq u} \left\| \Sigma^{-\frac{1}{2}} \boldsymbol{\xi}_i \right\|_2 \mathbb{1}\{\mathcal{H}_i\} \\ &\leq \alpha \|\Sigma\| \max_{i:l \leq i \leq u} \left\{ \left\| \Sigma^{-\frac{1}{2}} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 \mathbb{1}\{\mathcal{H}_i\} \right. \\ & \quad \left. + \left\| \Sigma^{-\frac{1}{2}} (\mathbf{b}_i - \mathbf{b}) \right\|_2 \right\}, \end{aligned} \quad (109)$$

where the last step results from (88g) (with the restriction that  $\eta < (1 - \gamma)/(4\|\Sigma\|)$ ) and the definition (104) of  $\alpha$ . It then suffices to control the two terms on the right-hand side of

(109). To begin with, we have

$$\begin{aligned} & \left\| \Sigma^{-\frac{1}{2}} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 \\ & \leq \left\| \Sigma^{-\frac{1}{2}} (\mathbf{A}_i - \mathbf{A}) \Sigma^{-\frac{1}{2}} \right\| \left\| \boldsymbol{\theta}_i \right\|_\Sigma \\ & \leq \left( \left\| \Sigma^{-\frac{1}{2}} \mathbf{A}_i \Sigma^{-\frac{1}{2}} \right\| + \left\| \Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-\frac{1}{2}} \right\| \right) (\|\boldsymbol{\theta}^*\|_\Sigma + \|\boldsymbol{\Delta}_i\|_\Sigma). \end{aligned}$$

Recall from (146) that  $\left\| \Sigma^{-\frac{1}{2}} \mathbf{A}_i \Sigma^{-\frac{1}{2}} \right\| \leq 2 \max_s \left\| \Sigma^{-1/2} \phi(s) \right\|_2^2$ , and similarly  $\left\| \Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-\frac{1}{2}} \right\| \leq 2 \max_s \left\| \Sigma^{-1/2} \phi(s) \right\|_2^2$ . We then have

$$\begin{aligned} & \left\| \Sigma^{-\frac{1}{2}} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 \\ & \leq 4 \max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\} (\|\boldsymbol{\theta}^*\|_\Sigma + \|\boldsymbol{\Delta}_i\|_\Sigma). \end{aligned} \quad (110)$$

Regarding the second term of (109), direct calculations give

$$\begin{aligned} & \left\| \Sigma^{-\frac{1}{2}} (\mathbf{b}_i - \mathbf{b}) \right\|_2^2 \\ & \leq 2 \left\| \Sigma^{-\frac{1}{2}} \mathbf{b}_i \right\|_2^2 + 2 \left\| \Sigma^{-\frac{1}{2}} \mathbf{b} \right\|_2^2 \\ & = 2 \left\| \Sigma^{-\frac{1}{2}} \phi(s_i) r(s_i) \right\|_2^2 + 2 \left\| \Sigma^{-\frac{1}{2}} \mathbb{E}_{s \sim \mu} [\phi(s) r(s)] \right\|_2^2 \\ & \leq 4 \max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\} \max_s |r(s)|^2 \\ & \leq 4 \max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\}. \end{aligned} \quad (111)$$

Substituting the preceding two bounds into (109), we arrive at

$$\begin{aligned} B &\leq 4\alpha \|\Sigma\| \left( \max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\} \max_{i:i < t} (\|\boldsymbol{\theta}^*\|_\Sigma + \|\boldsymbol{\Delta}_i\|_\Sigma) \right. \\ & \quad \left. \mathbb{1}\{\mathcal{H}_i\} + \sqrt{\max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\}} \right) \\ &\leq 4\alpha \|\Sigma\| \left( \max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\} (\|\boldsymbol{\theta}^*\|_\Sigma + R) \right. \\ & \quad \left. + \sqrt{\max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\}} \right) \\ &\leq 4\alpha \|\Sigma\| \max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\} (\|\boldsymbol{\theta}^*\|_\Sigma + R + 1) \\ &\leq 4\alpha \|\Sigma\| \left\| \Sigma^{-1} \right\| (\|\boldsymbol{\theta}^*\|_\Sigma + R + 1) \\ &= 4\kappa\alpha (\|\boldsymbol{\theta}^*\|_\Sigma + R + 1) =: B_{\max}. \end{aligned} \quad (112)$$

Here, the last line follows from the assumption  $\max \|\phi(s)\|_2 \leq 1$ , while the second to last inequality holds since  $\max_s \left\{ \phi(s)^\top \Sigma^{-1} \phi(s) \right\} \geq 1$  (cf. (148)).

c) *Invoking matrix Freedman's inequality*: Equipped with the above bounds (108) and (112), we are ready to apply Freedman's inequality [68, Corollary 1.3] (or a version in [19, Section A]), which asserts that

$$\begin{aligned} & \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \tilde{\boldsymbol{\xi}}_i \right\|_\Sigma \\ & \leq 2 \sqrt{W_{\max} \log \frac{2dT}{\delta}} + \frac{4}{3} B_{\max} \log \frac{2dT}{\delta} \\ & = \alpha \cdot \left\{ 2 \sqrt{\frac{4\kappa}{\eta(1-\gamma)}} \left\{ 4(\|\boldsymbol{\theta}^*\|_\Sigma + R)^2 + 1 \right\} \log \frac{2dT}{\delta} \right. \\ & \quad \left. + \frac{16\kappa}{3} (\|\boldsymbol{\theta}^*\|_\Sigma + R + 1) \log \frac{2dT}{\delta} \right\} \\ & \leq 16(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma))^{t-u-1} \\ & \quad (\|\boldsymbol{\theta}^*\|_\Sigma + R + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \end{aligned} \quad (113)$$



holds with probability at least  $1 - \delta/T$ , provided that  $0 < \eta \leq \frac{1}{\kappa(1-\gamma)\log \frac{2dT}{\delta}}$ . Here in the last line, we identify  $\alpha$  with  $(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma))^{t-u-1}$ . The proof is completed by observing that any  $0 < \eta \leq \frac{1-\gamma}{\kappa \log \frac{2dT}{\delta}}$  satisfies the two requirements  $0 < \eta \leq \frac{1}{\kappa(1-\gamma)\log \frac{2dT}{\delta}}$  and  $\eta < (1-\gamma)/(4\|\Sigma\|)$  (given that  $\|\Sigma\| \leq 1$  according to (88h)).

### B. Proof of the inequalities (64a) and (64c)

a) *Proof of the inequality (64a):* Combining the triangle inequality with the definition (63) ensures that

$$\begin{aligned} & \|\mathbf{A}_0^{(T+1)} \Delta_0\|_{\Sigma} \\ & \leq \|\mathbf{A}^{-1} \Delta_0\|_{\Sigma} + \|\mathbf{A}^{-1}(\mathbf{I} - \eta\mathbf{A})^{T+1} \Delta_0\|_{\Sigma} \\ & = \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \Sigma^{-1} \Sigma^{\frac{1}{2}} \Delta_0\|_2 \\ & + \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \Sigma^{-\frac{1}{2}} (\mathbf{I} - \eta\mathbf{A})^{T+1} \Sigma^{-\frac{1}{2}} \Sigma^{\frac{1}{2}} \Delta_0\|_2 \\ & \leq \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \cdot \|\Sigma^{-1}\| \cdot \|\Delta_0\|_{\Sigma} \\ & + \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \cdot \|\Sigma^{-\frac{1}{2}}\|^2 \cdot \|\mathbf{I} - \eta\mathbf{A}\|^{T+1} \cdot \|\Delta_0\|_{\Sigma} \\ & \leq \frac{\|\Sigma^{-1}\|}{1-\gamma} \left\{ 1 + \left( 1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma) \right)^{T+1} \right\} \|\Delta_0\|_{\Sigma} \\ & \leq \frac{2\|\Sigma^{-1}\|}{1-\gamma} \|\Delta_0\|_{\Sigma} \end{aligned} \quad (114)$$

as claimed. Here, the second to last step follows from (88e) and (88g), provided that  $\eta \leq (1-\gamma)/(4\|\Sigma\|)$ .

b) *Proof of the inequality (64c):* Again, the triangle inequality together with the definition (63) yields

$$\begin{aligned} & \left\| \sum_{i=0}^{T-1} \mathbf{A}_i^{(T)} \xi_i \right\|_{\Sigma} \\ & \leq \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} \xi_i \right\|_{\Sigma} + \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_{\Sigma} \\ & \leq \left\| \mathbf{A}^{-1} \sum_{i=0}^{T-1} (\mathbf{A}_i - \mathbf{A}) \theta_i \right\|_{\Sigma} + \left\| \mathbf{A}^{-1} \sum_{i=0}^{T-1} (\mathbf{b}_i - \mathbf{b}) \right\|_{\Sigma} \\ & + \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_{\Sigma}, \end{aligned} \quad (115)$$

leaving us with three terms to handle. Here in the second line, we substitute in the definition of  $\xi_i$  (45).

- The second term of (115) can be bounded by Lemma 6, which asserts that

$$\begin{aligned} & \frac{1}{T} \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{b}_i - \mathbf{b}) \right\|_{\Sigma} \\ & \lesssim \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s)}{T(1-\gamma)^2} \log \left( \frac{d}{\delta} \right)} \end{aligned} \quad (116)$$

holds with probability at least  $1 - \delta$ , as long as  $T \gtrsim \log \frac{d}{\delta}$ .

- For the third term of (115), invoking the property (88e) again yields

$$\left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_{\Sigma}$$

$$\begin{aligned} & = \left\| \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \Sigma^{-1} \sum_{i=0}^{T-1} \Sigma^{\frac{1}{2}} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_2 \\ & \leq \left\| \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \right\| \cdot \|\Sigma^{-1}\| \cdot \left\| \sum_{i=0}^{T-1} \Sigma^{\frac{1}{2}} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_2 \\ & \leq \frac{\|\Sigma^{-1}\|}{1-\gamma} \left\| \sum_{i=0}^{T-1} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_{\Sigma}. \end{aligned} \quad (117)$$

Repeating the same analysis as in Step 3 to see that

$$\left\| \sum_{i=0}^{T-1} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_{\Sigma} \leq 16(2\|\theta^*\|_{\Sigma} + 3) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \quad (118)$$

with probability at least  $1 - \delta$ . Substitution into (117) leads to

$$\begin{aligned} & \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{I} - \eta\mathbf{A})^{T-i} \xi_i \right\|_{\Sigma} \\ & \leq 16(2\|\theta^*\|_{\Sigma} + 3) \|\Sigma^{-1}\| \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)^3}}. \end{aligned} \quad (119)$$

- It then boils down to bounding the first term of (115). In light of (60), we decompose it as follows

$$\begin{aligned} & \left\| \frac{1}{T} \sum_{i=0}^{T-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta_i \right\|_2 \\ & \leq \left\| \frac{1}{T} \sum_{i=0}^{\tilde{t}_{\text{seg}}-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta_i \right\|_2 \\ & + \left\| \frac{1}{T} \sum_{i=\tilde{t}_{\text{seg}}}^{T-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta_i \right\|_2. \end{aligned} \quad (120)$$

Bounding these terms requires the following lemma, whose proof is deferred to Section C-F.

**Lemma 8.** Fix any  $R > 0$  and define a collection of auxiliary random vectors for  $0 \leq i \leq T-1$

$$\theta'_i := \theta_i \mathbb{1}\{\mathcal{H}_i\}, \quad \mathcal{H}_i := \{\|\Delta_i\|_{\Sigma} \leq R\}, \quad (121)$$

Then for any indices  $(l, u, t)$  obeying  $0 \leq l \leq u \leq T-1$ , one has with probability at least  $1 - \delta$  that

$$\begin{aligned} & \left\| \frac{1}{u-l+1} \sum_{i=l}^u \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta'_i \right\|_2 \\ & \leq \frac{16(\|\theta^*\|_{\Sigma} + R)}{1-\gamma} \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{u-l+1}} \end{aligned} \quad (122)$$

provided that

$$u-l+1 \geq \frac{4 \max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{9}.$$

Apply Lemma 8 with  $R = R_0$ ,  $l = 0$  and  $u = t'_{\text{seg}} - 1$  to obtain with probability of at least  $1 - \delta$  that

$$\left\| \frac{1}{t'_{\text{seg}}} \sum_{i=0}^{t'_{\text{seg}}-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta_i \right\|_2$$

$$= \left\| \frac{1}{t'_{\text{seg}}} \sum_{i=0}^{t'_{\text{seg}}-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \mathbb{1}\{\|\Delta_i\| \leq R_0\} \right\|_2$$

$$\leq \frac{16(\|\boldsymbol{\theta}^*\|_{\Sigma} + R_0)}{1 - \gamma} \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{t'_{\text{seg}}}},$$

as long as  $\frac{4\|\Sigma^{-1}\| \log \frac{2d}{\delta}}{9} \geq \frac{4\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{9}$ . Here, the identity holds since  $\|\Delta_i\|_{\Sigma} \leq R_0$  for  $i \leq t'_{\text{seg}} - 1$  with probability of at least  $1 - \delta$ . Similarly, invoke Lemma 8 with  $R = 32\sqrt{\frac{\eta\kappa \log \frac{2dT}{\delta}}{1-\gamma}}(\|\boldsymbol{\theta}^*\|_{\Sigma} + 2)$ ,  $l = t'_{\text{seg}}$  and  $u = T - 1$  to obtain with probability of at least  $1 - \delta$  that

$$\left\| \frac{1}{T - t'_{\text{seg}}} \sum_{i=t'_{\text{seg}}}^{T-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2$$

$$\leq \frac{16(\|\boldsymbol{\theta}^*\|_{\Sigma} + 32\sqrt{\frac{\eta\kappa \log \frac{2dT}{\delta}}{1-\gamma}}(\|\boldsymbol{\theta}^*\|_{\Sigma} + 2))}{1 - \gamma}$$

$$\sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T - t'_{\text{seg}}}}$$

$$\leq \frac{16(1.5\|\boldsymbol{\theta}^*\|_{\Sigma} + 2)}{1 - \gamma} \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T - t'_{\text{seg}}}}$$

provided that  $T - t'_{\text{seg}} \geq \frac{4\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{9}$ . Here, the last inequality arises from the relation

$$32\sqrt{\frac{\eta\kappa \log \frac{2dT}{\delta}}{1-\gamma}}(\|\boldsymbol{\theta}^*\|_{\Sigma} + 2) \leq 0.5\|\boldsymbol{\theta}^*\|_{\Sigma} + 2,$$

which is an immediate consequence of the assumption that  $\frac{\eta\kappa \log \frac{2dT}{\delta}}{1-\gamma}$  is sufficiently small. Therefore,

$$\left\| \frac{1}{T} \sum_{i=0}^{T-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2$$

$$\leq \frac{32(\|\boldsymbol{\theta}^*\|_{\Sigma} + \sqrt{\frac{t'_{\text{seg}}}{T}} R_0 + 1)}{1 - \gamma} \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T}}. \quad (123)$$

Combining the preceding bounds (116), (119) and (123) with (115), we reach the conclusion that with probability of at least  $1 - \delta$ ,

$$\left\| \sum_{i=0}^{T-1} \mathbf{A}_i^{(t)} \boldsymbol{\xi}_i \right\|_{\Sigma} \asymp \left\{ \sqrt{\frac{\max_s \phi(s)^{\top} \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T(1-\gamma)^2}} \right.$$

$$\left. + \frac{\|\Sigma^{-1}\|}{T} \sqrt{\frac{\kappa \log \frac{dT}{\delta}}{\eta(1-\gamma)^3}} \right\} (\|\boldsymbol{\theta}^*\|_{\Sigma} + 1),$$

as long as  $T \geq t'_{\text{seg}} \kappa \|\Delta_0\|_{\Sigma}^2$ , where we use the definition (47) of  $R_0$ . It thus establishes the inequality (64c).

### C. Proof of Lemma 2

We first decompose  $\Psi$  into

$$\Psi = \begin{bmatrix} \mathbf{I} - \alpha \tilde{\mathbf{A}}^{\top} \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \beta \tilde{\Sigma} \end{bmatrix}$$

$$+ \begin{bmatrix} \mathbf{0} & -\frac{1}{\kappa} \alpha \gamma \mathbf{\Pi}^{\top} \\ -\kappa \alpha (1 - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \tilde{\mathbf{A}}^{\top} \tilde{\Sigma}^{-1} \tilde{\mathbf{A}} & -\alpha \gamma (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \mathbf{\Pi}^{\top} \end{bmatrix}.$$

Then the triangle inequality together with the properties of the operator norm tells us that

$$\|\Psi\| \leq \max\{\|\mathbf{I} - \alpha \tilde{\mathbf{A}}^{\top} \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\|, \|\mathbf{I} - \beta \tilde{\Sigma}\|\} + \frac{1}{\kappa} \alpha \gamma \mathbf{\Pi}^{\top} \mathbf{\Pi}$$

$$+ \|\kappa \alpha (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \tilde{\mathbf{A}}^{\top} \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\| + \|\alpha \gamma (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \mathbf{\Pi}^{\top}\|.$$

Note that by definition of  $\lambda_w$  and  $\lambda_w$ , we find

$$\|\mathbf{I} - \alpha \tilde{\mathbf{A}}^{\top} \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\| \leq 1 - \alpha \lambda_{\theta},$$

$$\|\mathbf{I} - \beta \tilde{\Sigma}\| \leq 1 - \beta \lambda_w.$$

In addition, some direct algebra suggests

$$\left\| \frac{1}{\kappa} \alpha \gamma \mathbf{\Pi}^{\top} \right\| \leq \frac{\alpha \gamma \rho_{\max}}{\kappa},$$

$$\|\kappa \alpha (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \tilde{\mathbf{A}}^{\top} \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\|$$

$$\leq \kappa \alpha (1 + \gamma \lambda_{\Sigma} \rho_{\max}) \lambda_{\Sigma} (2\rho_{\max})^2, \quad \text{and}$$

$$\|\alpha \gamma (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \mathbf{\Pi}^{\top}\| \leq \alpha \gamma (\rho_{\max} + \gamma \lambda_{\Sigma} \rho_{\max}^2).$$

In summary, as long as

$$\alpha \gamma (\rho_{\max} + \gamma \lambda_{\Sigma} \rho_{\max}^2) \ll \beta \lambda_w,$$

$$\frac{\alpha \gamma \rho_{\max}}{\kappa} + \kappa \alpha (1 + \gamma \lambda_{\Sigma} \rho_{\max}) \lambda_{\Sigma} (2\rho_{\max})^2 \ll \sqrt{\alpha \lambda_{\theta} \beta \lambda_w},$$

one has

$$\|\Psi\| \leq 1 - \frac{1}{2} \alpha \lambda_{\theta}.$$

### D. Proof of Lemma 3

Using the same notation of  $\mathbb{E}_{i-1}$  as in Section C-A, we observe that  $\{\Psi^{t-i-1} \tilde{\zeta}_i\}$  forms a martingale difference sequence. Furthermore, define

$$\tilde{W} := \sum_{i=0}^{t-1} \mathbb{E}_{i-1} \left[ \|\Psi^{t-i-1} \tilde{\zeta}_i\|_{\tilde{\Sigma}}^2 \mathbb{1}\{\tilde{\mathcal{H}}_i\} \right], \quad \text{and}$$

$$\tilde{B} := \max_{i: 0 \leq i \leq t-1} \left\| \Psi^{t-i-1} \tilde{\zeta}_i \mathbb{1}\{\tilde{\mathcal{H}}_i\} \right\|_{\tilde{\Sigma}}. \quad (124)$$

In order to bound  $\tilde{W}$  and  $\tilde{B}$ , we will firstly need to bound the norm of  $\tilde{\zeta}_i$ , as is shown in the following paragraph.

a) *Controlling the norm of  $\tilde{\zeta}_i$ :* We firstly observe that since  $\|\phi(s)\|_2 \leq 1$  and  $r(s) \leq 1$  for all  $s \in \mathcal{S}$ , with similar logic as (106), (107), (110) and (111), the following bounds hold true:

- For any  $\mathcal{F}_{i-1}$ -measurable  $\tilde{\boldsymbol{\theta}}_i \in \mathbb{R}^d$ , the norm of  $(\tilde{\mathbf{A}}_i - \tilde{\mathbf{A}}) \tilde{\boldsymbol{\theta}}_i$  is bounded by

$$\mathbb{E}_{i-1} \left\| (\tilde{\mathbf{A}}_i - \tilde{\mathbf{A}}) \tilde{\boldsymbol{\theta}}_i \right\|_2^2 \leq 4\rho_{\max}^2 \left( \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right), \quad \text{and} \quad (125)$$

$$\left\| \tilde{\Sigma}^{-1/2} (\tilde{\mathbf{A}}_i - \tilde{\mathbf{A}}) \tilde{\boldsymbol{\theta}}_i \right\|_2$$

$$\leq 4\rho_{\max} \max_s \left\{ \phi(s)^{\top} \tilde{\Sigma}^{-1} \phi(s) \right\} \left( \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right); \quad (126)$$

- For any  $\mathcal{F}_{i-1}$ -measurable  $\mathbf{z}_i \in \mathbb{R}^d$ , the norm of  $(\mathbf{\Pi}_i - \mathbf{\Pi})^\top \mathbf{z}_i$  is bounded by

$$\mathbb{E}_{i-1} \left\| (\mathbf{\Pi}_i - \mathbf{\Pi})^\top \mathbf{z}_i \right\|_2^2 \leq \rho_{\max}^2 \|\mathbf{z}_i\|_{\tilde{\Sigma}}^2, \quad \text{and} \quad (127)$$

$$\begin{aligned} & \left\| \tilde{\Sigma}^{-1/2} (\mathbf{\Pi}_i - \mathbf{\Pi})^\top \mathbf{z}_i \right\|_2 \\ & \leq 2\rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \|\mathbf{z}_i\|_{\tilde{\Sigma}}; \end{aligned} \quad (128)$$

- For any  $\mathcal{F}_{i-1}$ -measurable  $\mathbf{z}_i \in \mathbb{R}^d$ , the norm of  $(\tilde{\Sigma}_i - \tilde{\Sigma}) \mathbf{z}_i$  is bounded by

$$\mathbb{E}_{i-1} \left\| (\tilde{\Sigma}_i - \tilde{\Sigma}) \mathbf{z}_i \right\|_2^2 \leq \|\mathbf{z}_i\|_{\tilde{\Sigma}}^2, \quad \text{and} \quad (129)$$

$$\left\| \tilde{\Sigma}^{-1/2} (\tilde{\Sigma}_i - \tilde{\Sigma}) \mathbf{z}_i \right\|_2 \leq 2 \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \|\mathbf{z}_i\|_{\tilde{\Sigma}}; \quad (130)$$

- The norm of  $\tilde{\mathbf{b}}_i - \tilde{\mathbf{b}}$  is bounded by

$$\mathbb{E}_{i-1} \left\| \tilde{\mathbf{b}}_i - \tilde{\mathbf{b}} \right\|_2^2 \leq \rho_{\max}^2, \quad \text{and} \quad (131)$$

$$\left\| \tilde{\Sigma}^{-1/2} (\tilde{\mathbf{b}}_i - \tilde{\mathbf{b}}) \right\|_2^2 \leq 4\rho_{\max}^2 \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\}. \quad (132)$$

Therefore, by triangle inequality, the norm of  $\nu_i$  can be bounded by

$$\mathbb{E}_{i-1} \|\nu_i\|_2^2 \lesssim \rho_{\max}^2 \left[ \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right) + 1 + \gamma^2 \|\mathbf{w}_i\|_{\tilde{\Sigma}}^2 \right], \quad (133)$$

and

$$\begin{aligned} & \left\| \tilde{\Sigma}^{-1/2} \nu_i \right\|_2 \\ & \lesssim \rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left[ \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + \gamma \|\mathbf{w}_i\|_{\tilde{\Sigma}} + 1 \right]; \end{aligned} \quad (134)$$

similarly, the norm of  $\eta_i$  can be bounded by

$$\mathbb{E}_{i-1} \|\eta_i\|_2^2 \lesssim \rho_{\max}^2 \left[ \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right) + 1 \right] + \|\mathbf{w}_i\|_{\tilde{\Sigma}}^2, \quad (135)$$

and

$$\begin{aligned} & \left\| \tilde{\Sigma}^{-1/2} \eta_i \right\|_2 \lesssim \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \\ & \left\{ \rho_{\max} \left[ \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + 1 \right] + \|\mathbf{w}_i\|_{\tilde{\Sigma}} \right\}. \end{aligned} \quad (136)$$

By combining (133) and (135) with the definition of  $\zeta_i$  (73), we obtain the following bound:

$$\begin{aligned} & \mathbb{E}_{i-1} \|\zeta_i\|_2^2 \\ & \lesssim \alpha^2 \mathbb{E}_{i-1} \|\nu_i\|_2^2 + \kappa^2 \alpha^2 \|\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}\|^2 \mathbb{E}_{i-1} \|\nu_i\|_2^2 \\ & + \kappa^2 \beta^2 \mathbb{E}_{i-1} \|\eta_i\|_2^2 \\ & \lesssim \alpha^2 \left( 1 + \kappa^2 (1 + \gamma \lambda_{\Sigma} \rho_{\max})^2 \right) \cdot \rho_{\max}^2 \\ & \left[ 4 \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right) + 1 + \gamma^2 \|\mathbf{w}_i\|_{\tilde{\Sigma}}^2 \right] \\ & + \kappa^2 \beta^2 \cdot \left\{ \rho_{\max}^2 \left[ \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right) + 1 \right] + \|\mathbf{w}_i\|_{\tilde{\Sigma}}^2 \right\} \end{aligned}$$

$$\lesssim \kappa^2 \beta^2 \rho_{\max}^2 \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2 + \frac{1}{\kappa^2} \|\mathbf{x}_i\|_{\tilde{\Sigma}}^2 + 1 \right), \quad (138)$$

and

$$\begin{aligned} & \left\| \tilde{\Sigma}^{-1/2} \zeta_i \right\|_2 \\ & \lesssim \alpha \|\tilde{\Sigma}^{-1/2} \nu_i\|_2 + \alpha \kappa \|\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}\| \|\tilde{\Sigma}^{-1/2} \nu_i\|_2 \\ & + \kappa \beta \left\| \tilde{\Sigma}^{-1/2} \eta_i \right\|_2 \\ & \lesssim \alpha \left( 1 + \kappa (1 + \gamma \lambda_{\Sigma} \rho_{\max}) \right) \cdot \rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \\ & \cdot \left[ 2 \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + \gamma \|\mathbf{w}_i\|_{\tilde{\Sigma}} + 1 \right] \\ & + \kappa \beta \cdot \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \\ & \cdot \left\{ \rho_{\max} \left[ \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + 1 \right] + \|\mathbf{w}_i\|_{\tilde{\Sigma}} \right\} \\ & \lesssim \kappa \beta \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left( \|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{2}{\kappa} \|\mathbf{x}_i\|_{\tilde{\Sigma}} + 1 \right) \end{aligned} \quad (139)$$

b) *Control of  $\tilde{W}$  and  $\tilde{B}$ :* With the norm of  $\zeta_i$  bounded, we can apply similar techniques as in equations (105), (108), (109) and (112) of Section C-A to construct the following bound for  $\tilde{W}$ :

$$\begin{aligned} \tilde{W} & \leq \|\tilde{\Sigma}\| \sum_{i=0}^{t-1} \|\Psi^{t-i-1}\|^2 \cdot \mathbb{E}_{i-1} \left[ \|\zeta_i\|_2^2 \mathbf{1} \left\{ \mathcal{H}_i \right\} \right] \\ & \lesssim \|\tilde{\Sigma}\| \sum_{i=0}^{t-1} \left( 1 - \frac{1}{2} \alpha \lambda_{\theta} \right)^{2t-2i-2} \\ & \cdot \kappa^2 \beta^2 \rho_{\max}^2 (2\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + 2\tilde{R} + 1)^2 \\ & \lesssim \frac{\|\tilde{\Sigma}\|}{\alpha \lambda_{\theta}} \kappa^2 \beta^2 \rho_{\max}^2 (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{1}{\kappa} \tilde{R} + 1)^2, \end{aligned} \quad (140)$$

and the following bound for  $\tilde{B}$ :

$$\begin{aligned} \tilde{B} & \leq \|\tilde{\Sigma}\| \max_{i:0 \leq i \leq t-1} \left\| \tilde{\Sigma}^{-1/2} \zeta_i \mathbf{1} \left\{ \mathcal{H}_i \right\} \right\|_2 \\ & \lesssim \|\tilde{\Sigma}\| \kappa \beta \rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} (\|\tilde{\theta}^*\|_2 + \tilde{R} + 1) \\ & =: \tilde{B}_{\max}. \end{aligned} \quad (141)$$

c) *Invoking the matrix Freedman's inequality:* With  $\tilde{W}$  and  $\tilde{B}$  bounded, we again invoke the matrix Freedman's inequality [68, Corollary 1.3] to assert that

$$\begin{aligned} & \left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \tilde{\zeta}_i \right\|_2 \\ & \leq 2\sqrt{\tilde{W}_{\max} \log \frac{2dT}{\delta}} + \frac{4}{3} \tilde{B}_{\max} \log \frac{2dT}{\delta} \\ & \lesssim \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha \lambda_{\theta}} \log \frac{2dT}{\delta}} \kappa \beta \rho_{\max} (\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{1}{\kappa} \tilde{R} + 1) \end{aligned} \quad (142)$$

holds with probability at least  $1 - \delta/T$ , provided that  $0 < \alpha < \frac{1}{\lambda_{\theta} \lambda_{\Sigma}^2 \|\tilde{\Sigma}\| \log \frac{2dT}{\delta}}$ .

#### E. Proof of Lemma 6 and Lemma 7

a) *Proof of Lemma 7: controlling  $\|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-\frac{1}{2}}\|$ :* We intend to invoke the matrix Bernstein inequality

to establish the advertised bound [69]. Note that

$$\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-\frac{1}{2}} = \frac{1}{T} \sum_{t=0}^{T-1} \underbrace{\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A} - \mathbf{A}_t) \Sigma^{-\frac{1}{2}}}_{=: \mathbf{Z}_t}. \quad (143)$$

In order to control it, we need to first control the following two quantities:

$$v := \max_t \left\{ \max \left\{ \left\| \mathbb{E} [\mathbf{Z}_t \mathbf{Z}_t^\top] \right\|, \left\| \mathbb{E} [\mathbf{Z}_t^\top \mathbf{Z}_t] \right\| \right\} \right\} \quad \text{and} \\ B := \max_t \|\mathbf{Z}_t\|.$$

*Step 1: controlling  $\left\| \mathbb{E} [\mathbf{Z}_t \mathbf{Z}_t^\top] \right\|$ .* Towards this, we first make the observation that

$$\begin{aligned} & \mathbb{E} [\mathbf{Z}_t \mathbf{Z}_t^\top] \\ &= \mathbb{E} \left[ \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A} - \mathbf{A}_t) \Sigma^{-1} (\mathbf{A} - \mathbf{A}_t)^\top (\mathbf{A}^\top)^{-1} \Sigma^{\frac{1}{2}} \right] \\ &\preceq \mathbb{E} \left[ \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbf{A}_t \Sigma^{-1} \mathbf{A}_t^\top (\mathbf{A}^\top)^{-1} \Sigma^{\frac{1}{2}} \right] \\ &= \mathbb{E}_{s \sim \mu, s' \sim P(\cdot|s)} \left[ \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \phi(s) (\phi(s) - \gamma \phi(s'))^\top \Sigma^{-1} \right. \\ &\quad \left. (\phi(s) - \gamma \phi(s')) \phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma^{\frac{1}{2}} \right] \\ &\preceq \max_{s, s'} \left\{ (\phi(s) - \gamma \phi(s'))^\top \Sigma^{-1} (\phi(s) - \gamma \phi(s')) \right\} \\ &\quad \cdot \mathbb{E}_{s \sim \mu} \left[ \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \phi(s) \phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma^{\frac{1}{2}} \right] \\ &\preceq \max_{s, s'} \left\{ 2\phi(s)^\top \Sigma^{-1} \phi(s) + 2\gamma^2 \phi(s')^\top \Sigma^{-1} \phi(s') \right\} \\ &\quad \cdot \left\{ \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma (\mathbf{A}^\top)^{-1} \Sigma^{\frac{1}{2}} \right\} \\ &\preceq \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{(1 - \gamma)^2} \mathbf{I}, \end{aligned} \quad (144)$$

where the second line holds since  $\mathbb{E}[(\mathbf{M} - \mathbb{E}[\mathbf{M}])(\mathbf{M} - \mathbb{E}[\mathbf{M}])^\top] \preceq \mathbb{E}[\mathbf{M} \mathbf{M}^\top]$  for any random matrix  $\mathbf{M}$ , the second to last inequality holds since  $(\mathbf{a} - \mathbf{b})^\top \Sigma^{-1} (\mathbf{a} - \mathbf{b}) \leq 2\mathbf{a}^\top \Sigma^{-1} \mathbf{a} + 2\mathbf{b}^\top \Sigma^{-1} \mathbf{b}$ , and the last inequality comes from the assumption  $\gamma < 1$  and Lemma 5.

*Step 2: controlling  $\left\| \mathbb{E} [\mathbf{Z}_t^\top \mathbf{Z}_t] \right\|$ .* Similarly, one can obtain

$$\begin{aligned} & \mathbb{E} [\mathbf{Z}_t^\top \mathbf{Z}_t] \\ &= \mathbb{E} \left[ \Sigma^{-\frac{1}{2}} (\mathbf{A} - \mathbf{A}_t)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} (\mathbf{A} - \mathbf{A}_t) \Sigma^{-\frac{1}{2}} \right] \\ &\preceq \mathbb{E} \left[ \Sigma^{-\frac{1}{2}} \mathbf{A}_t^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \mathbf{A}_t \Sigma^{-\frac{1}{2}} \right] \\ &= \mathbb{E} \left[ \Sigma^{-\frac{1}{2}} (\phi(s_t) - \gamma \phi(s'_t)) \phi(s_t)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \right. \\ &\quad \left. \phi(s_t) (\phi(s_t) - \gamma \phi(s'_t))^\top \Sigma^{-\frac{1}{2}} \right] \\ &\preceq \max_s \left\{ \phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \phi(s) \right\} \\ &\quad \cdot \mathbb{E} \left[ \Sigma^{-\frac{1}{2}} (\phi(s_t) - \gamma \phi(s'_t)) (\phi(s_t) - \gamma \phi(s'_t))^\top \Sigma^{-\frac{1}{2}} \right] \\ &\preceq \max_s \left\{ \phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \phi(s) \right\} \\ &\quad \cdot 2\mathbb{E} \left[ \Sigma^{-\frac{1}{2}} (\phi(s_t) \phi(s_t)^\top + \phi(s'_t) \phi(s'_t)^\top) \Sigma^{-\frac{1}{2}} \right] \\ &\preceq 4 \max_s \left\{ \phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \phi(s) \right\} \mathbf{I}. \end{aligned}$$

Here, the second to last bound follows from the elementary inequality  $(\mathbf{a} - \mathbf{b})(\mathbf{a} - \mathbf{b})^\top \preceq 2\mathbf{a}\mathbf{a}^\top + 2\mathbf{b}\mathbf{b}^\top$  and the assumption  $\gamma < 1$ , whereas the last line makes use of the facts  $s_t \sim \mu$ ,  $s'_t \sim \mu$  and the definition (6) of  $\Sigma$ . It then boils down to upper bounding  $\max_s \left\{ \phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \phi(s) \right\}$ , which can be accomplished as follows

$$\begin{aligned} & \phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \phi(s) \\ &= \phi(s)^\top \Sigma^{-\frac{1}{2}} \left\{ \Sigma^{\frac{1}{2}} (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \right\} \Sigma^{-\frac{1}{2}} \phi(s) \\ &\leq \left\| \Sigma^{-\frac{1}{2}} \phi(s) \right\|_2^2 \cdot \left\| \Sigma^{\frac{1}{2}} (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \right\| \\ &\leq \frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{(1 - \gamma)^2}. \end{aligned}$$

Here, the last line arises from Lemma 5. Putting the above bounds together yields

$$\mathbb{E} [\mathbf{Z}_t^\top \mathbf{Z}_t] \preceq \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{(1 - \gamma)^2} \mathbf{I}. \quad (145)$$

*Step 3: controlling  $\|\mathbf{Z}_t\|$ .* Our starting point is the following triangle inequality

$$\begin{aligned} \|\mathbf{Z}_t\| &= \left\| \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A} - \mathbf{A}_t) \Sigma^{-\frac{1}{2}} \right\| \\ &\leq \left\| \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbf{A}_t \Sigma^{-\frac{1}{2}} \right\| + \left\| \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbf{A} \Sigma^{-\frac{1}{2}} \right\| \\ &\leq \left\| \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \right\| \cdot \left\| \Sigma^{-\frac{1}{2}} \mathbf{A}_t \Sigma^{-\frac{1}{2}} \right\| + 1 \\ &\leq \frac{1}{1 - \gamma} \left\| \Sigma^{-\frac{1}{2}} \mathbf{A}_t \Sigma^{-\frac{1}{2}} \right\| + 1, \end{aligned}$$

where the last inequality follows from Lemma 5. In addition, we see that

$$\begin{aligned} \left\| \Sigma^{-\frac{1}{2}} \mathbf{A}_t \Sigma^{-\frac{1}{2}} \right\| &\leq \max_s \left\| \Sigma^{-\frac{1}{2}} \phi(s) \phi(s)^\top \Sigma^{-\frac{1}{2}} \right\| \\ &\quad + \gamma \max_{s, s'} \left\| \Sigma^{-\frac{1}{2}} \phi(s') \phi(s)^\top \Sigma^{-\frac{1}{2}} \right\| \\ &\leq 2 \max_s \left\| \Sigma^{-\frac{1}{2}} \phi(s) \right\|_2^2. \end{aligned} \quad (146)$$

This combined with the preceding bounds yields

$$\begin{aligned} \|\mathbf{Z}_t\| &\leq \frac{2 \max_s \left\| \Sigma^{-\frac{1}{2}} \phi(s) \right\|_2^2}{1 - \gamma} + 1 \\ &\leq \frac{4 \max_s \left\| \Sigma^{-\frac{1}{2}} \phi(s) \right\|_2^2}{1 - \gamma} \\ &= \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{1 - \gamma}. \end{aligned} \quad (147)$$

Here, the inequality follows since

$$\begin{aligned} \max_s \left\| \Sigma^{-\frac{1}{2}} \phi(s) \right\|_2^2 &\geq \mathbb{E}_{s \sim \mu} [\phi(s)^\top \Sigma^{-1} \phi(s)] \\ &= \mathbb{E}_{s \sim \mu} [\text{tr}(\Sigma^{-1} \phi(s) \phi(s)^\top)] \\ &= \text{tr}(\mathbf{I}_d) = d \geq 1. \end{aligned} \quad (148)$$

*Step 4: invoking the matrix Bernstein inequality.* With the above bounds in mind, we are ready to apply the matrix Bernstein inequality [69] to obtain that: with probability at least  $1 - \delta$  one has

$$\left\| \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-\frac{1}{2}} \right\|$$



$$\begin{aligned}
& \lesssim \sqrt{\frac{1}{T^2} \sum_{t=0}^{T-1} \max \{ \|\mathbb{E}[\mathbf{Z}_t \mathbf{Z}_t^\top]\|, \|\mathbb{E}[\mathbf{Z}_t^\top \mathbf{Z}_t]\| \} \log\left(\frac{d}{\delta}\right)} \\
& + \frac{\max_t \|\mathbf{Z}_t\| \log\left(\frac{d}{\delta}\right)}{T} \\
& \stackrel{(i)}{\lesssim} \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{T(1-\gamma)^2} \log\left(\frac{d}{\delta}\right)} \\
& + \frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log\left(\frac{d}{\delta}\right)}{T(1-\gamma)} \\
& \stackrel{(ii)}{\lesssim} \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{T(1-\gamma)^2} \log\left(\frac{d}{\delta}\right)}. \tag{149}
\end{aligned}$$

Here, (i) results from the bounds (144), (145) and (147), while (ii) holds as long as  $T \gtrsim \max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log\left(\frac{d}{\delta}\right)$ .

In addition, if  $T \geq \frac{c_2 \max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log\left(\frac{d}{\delta}\right)}{(1-\gamma)^2}$  for some constant  $c_2$  large enough, then one has  $\|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-\frac{1}{2}}\| < 1$ . Suppose that  $\hat{\mathbf{A}}$  is not invertible. Given that  $\mathbf{A}$  and  $\Sigma$  are both invertible, this means that one can find a unit vectors  $\mathbf{u}$  obeying  $\mathbf{A}^{-1} \hat{\mathbf{A}} \Sigma^{-\frac{1}{2}} \mathbf{u} = \mathbf{0}$ , which in turn implies

$$\begin{aligned}
& \mathbf{u}^\top \Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-\frac{1}{2}} \mathbf{u} \\
& = \mathbf{u}^\top \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbf{A} \Sigma^{-\frac{1}{2}} \mathbf{u} - \mathbf{u}^\top \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \hat{\mathbf{A}} \Sigma^{-\frac{1}{2}} \mathbf{u} \\
& = 1 - 0 = 1
\end{aligned}$$

and hence contradicts the condition  $\|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-\frac{1}{2}}\| < 1$ . As a result, we conclude that  $\hat{\mathbf{A}}$  is invertible as long as  $\|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{A} - \hat{\mathbf{A}}) \Sigma^{-\frac{1}{2}}\| < 1$ .

*b) Proof of Lemma 6: controlling  $\|\mathbf{A}^{-1}(\hat{\mathbf{b}} - \mathbf{b})\|_\Sigma$ .* First of all, it is seen that

$$\|\mathbf{A}^{-1}(\hat{\mathbf{b}} - \mathbf{b})\|_\Sigma = \left\| \frac{1}{T} \sum_{t=0}^{T-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{b}_t - \mathbf{b}) \right\|_2 = \left\| \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{z}_t \right\|_2,$$

where we define the vector  $\mathbf{z}_t := \Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{b}_t - \mathbf{b})$ . Therefore, we need to look at the properties of  $\mathbf{z}_t$ . Towards this end, we observe that

$$\begin{aligned}
\mathbb{E}[\mathbf{z}_t^\top \mathbf{z}_t] & = \mathbb{E}[(\mathbf{b}_t - \mathbf{b})^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1}(\mathbf{b}_t - \mathbf{b})] \\
& \leq \mathbb{E}[\mathbf{b}_t^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \mathbf{b}_t] \\
& \stackrel{(i)}{\leq} \left\{ \max_{s \in \mathcal{S}} |r(s)|^2 \right\} \mathbb{E}[\phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \phi(s)] \\
& \stackrel{(ii)}{\leq} \mathbb{E}[\phi(s)^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \phi(s)] \\
& = \mathbb{E}[\phi(s)^\top \Sigma^{-\frac{1}{2}} \Sigma^{\frac{1}{2}} (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \Sigma^{-\frac{1}{2}} \phi(s)] \\
& \leq \left\{ \max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2^2 \right\} \cdot \|\Sigma^{\frac{1}{2}} (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \\
& \stackrel{(iii)}{\leq} \frac{1}{(1-\gamma)^2} \max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2^2,
\end{aligned}$$

where (i) holds since  $\mathbf{b}_t = \phi(s_t)r(s_t)$ , (ii) follows from the assumption  $\max_s |r(s)| \leq 1$ , and (iii) arises from Lemma 5. Additionally,

$$\max_t \|\mathbf{z}_t\|_2 \leq \max_t \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbf{b}_t\|_2 + \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbf{b}\|_2$$

$$\begin{aligned}
& \stackrel{(iv)}{\leq} 2 \max_s \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \phi(s) r(s)\|_2 \\
& \stackrel{(v)}{\leq} 2 \max_{s \in \mathcal{S}} \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \phi(s)\|_2 \\
& \leq 2 \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \cdot \max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2 \\
& \leq \frac{2}{1-\gamma} \max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2,
\end{aligned}$$

where (iv) holds since  $\mathbf{b}_t = \phi(s_t)r(s_t)$  and  $\mathbf{b} = \mathbb{E}_{s \sim \mu}[\phi(s)r(s)]$ , (v) comes from the assumption  $\max_s |r(s)| \leq 1$ , and the last line is due to Lemma 5. Consequently, the matrix Bernstein inequality [69] yields

$$\begin{aligned}
& \|\mathbf{A}^{-1}(\hat{\mathbf{b}} - \mathbf{b})\|_\Sigma \\
& = \left\| \frac{1}{T} \sum_{t=1}^T \mathbf{z}_t \right\|_2 \\
& \lesssim \sqrt{\frac{1}{T^2} \sum_{t=0}^{T-1} \mathbb{E}[\mathbf{z}_t^\top \mathbf{z}_t] \log\left(\frac{d}{\delta}\right)} + \frac{1}{T} \max_t \|\mathbf{z}_t\|_2 \log\left(\frac{d}{\delta}\right) \\
& \lesssim \frac{\max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2}{1-\gamma} \sqrt{\frac{1}{T} \log\left(\frac{d}{\delta}\right)} \\
& + \frac{\max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2}{1-\gamma} \cdot \frac{1}{T} \log\left(\frac{d}{\delta}\right) \\
& \asymp \frac{\max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2}{1-\gamma} \sqrt{\frac{1}{T} \log\left(\frac{d}{\delta}\right)} \tag{150}
\end{aligned}$$

with probability at least  $1 - \delta$ , as long as  $T \gtrsim \log\left(\frac{d}{\delta}\right)$ .

#### F. Proof of Lemma 8

Recall from the proof of Lemma 1 that  $\mathbb{E}_i[\cdot]$  represents the expectation conditioned on the probability space generated by the samples  $\{(s_j, s'_j)\}_{j \leq i}$ . It is easy to check that  $\{\Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{A}_i - \mathbf{A})\theta'_i\}$  forms a martingale difference sequence, and we seek to bound  $\left\| \frac{1}{u-l+1} \sum_{i=l}^u \Sigma^{\frac{1}{2}} \mathbf{A}^{-1}(\mathbf{A}_i - \mathbf{A})\theta'_i \right\|_2$  via matrix Freedman's inequality. The key is to control the following quantities (here, we abuse notation whenever it is clear from context):

$$\begin{aligned}
W & := \sum_{i=l}^u \mathbb{E}_{i-1} \left[ \|\Sigma^{1/2} \mathbf{A}^{-1}(\mathbf{A}_i - \mathbf{A})\theta'_i\|_2^2 \right] \quad \text{and} \\
B & := \max_{i:l \leq i \leq u} \|\Sigma^{1/2} \mathbf{A}^{-1}(\mathbf{A}_i - \mathbf{A})\theta'_i\|_2. \tag{151}
\end{aligned}$$

*a) Control of B:* To begin with, observe that

$$\begin{aligned}
B & = \max_{i:l \leq i \leq u} \|\Sigma^{1/2} \mathbf{A}^{-1}(\mathbf{A}_i - \mathbf{A})\Sigma^{-1/2}\| \cdot \|\theta'_i\|_\Sigma \\
& \leq \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{1-\gamma} \max_{i:l \leq i \leq u} \{\|\theta^*\|_\Sigma + \|\Delta_i\|_\Sigma\} \mathbb{1}\{\mathcal{H}_i\} \\
& \leq \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{1-\gamma} (\|\theta^*\|_\Sigma + R) =: B_{\max},
\end{aligned}$$

where the second to last inequality comes from (147) and the triangle inequality, and the last line is due to the definition of  $\mathcal{H}_i$ .

b) *Control of  $W$* : Moreover, one can derive

$$\begin{aligned} W &:= \sum_{i=l}^u \mathbb{E}_{i-1} \left[ \|\Sigma^{1/2} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \Sigma^{-1/2} \Sigma^{1/2} \boldsymbol{\theta}'_i\|_2^2 \right] \\ &= \sum_{i=l}^u \boldsymbol{\theta}'_i{}^\top \Sigma^{1/2} \mathbb{E}_{i-1} \left[ \Sigma^{-1/2} (\mathbf{A}_i - \mathbf{A})^\top (\mathbf{A}^\top)^{-1} \Sigma \mathbf{A}^{-1} \right. \\ &\quad \left. (\mathbf{A}_i - \mathbf{A}) \Sigma^{-1/2} \right] \Sigma^{1/2} \boldsymbol{\theta}'_i \\ &\leq \sum_{i=l}^u \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{(1-\gamma)^2} \|\Sigma^{1/2} \boldsymbol{\theta}'_i\|_2^2 \\ &\leq \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{(1-\gamma)^2} \sum_{i=l}^u (\|\boldsymbol{\theta}^*\|_\Sigma + \|\Delta_i\|_\Sigma)^2 \mathbb{1}\{\mathcal{H}_i\} \\ &\leq \frac{4(u-l+1) \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{(1-\gamma)^2} (\|\boldsymbol{\theta}^*\|_\Sigma + R)^2 \\ &=: W_{\max}, \end{aligned}$$

where the first inequality arises from (145), and the last inequality makes use of the definition of  $\mathcal{H}_i$ .

With the above bounds in place, we can apply Freedman's inequality [68, Corollary 1.3] for matrix martingales to demonstrate that

$$\begin{aligned} &\left\| \frac{1}{u-l+1} \sum_{i=l}^u \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}'_i \right\|_2 \\ &\leq \frac{2}{u-l+1} \sqrt{W_{\max} \log \frac{2d}{\delta}} + \frac{4}{3u-l+1} B_{\max} \log \frac{2d}{\delta} \\ &\leq \frac{8(\|\boldsymbol{\theta}^*\|_\Sigma + R)}{1-\gamma} \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{u-l+1}} \\ &\quad + \frac{16 \max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{3(1-\gamma)(u-l+1)} (\|\boldsymbol{\theta}^*\|_\Sigma + R) \\ &\leq \frac{16(\|\boldsymbol{\theta}^*\|_\Sigma + R)}{1-\gamma} \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{u-l+1}} \end{aligned}$$

with probability at least  $1 - \delta$ , as long as  $u - l + 1 \geq \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{9}$ .

## APPENDIX D

### COMPARISONS WITH PREVIOUS WORKS

#### A. Comparisons with [21]

[21] bounded the expectation of TD estimation error  $\mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_2^2$  with Markov samples by an iterative relation. For fair comparisons, we apply their ideas to bounding the error in  $\Sigma$ -norm with independent samples.

a) *Iterative relation on  $\mathbb{E}\|\Delta_t\|_\Sigma^2$* : Recall from the TD update rule (14) that

$$\begin{aligned} \Delta_{t+1} &= \Delta_t - \eta_t (\mathbf{A}_t \boldsymbol{\theta}_t - \mathbf{b}_t) \\ &= (\mathbf{I} - \eta_t \mathbf{A}_t) \Delta_t - \eta_t (\mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t). \end{aligned}$$

Therefore, the  $\Sigma$ -norm of  $\Delta_{t+1}$  can be expressed as

$$\begin{aligned} \|\Delta_{t+1}\|_\Sigma^2 &= \|\Delta_t\|_\Sigma^2 - 2\eta_t \langle \Delta_t, \mathbf{A}_t \Delta_t \rangle_\Sigma + \eta_t^2 \|\mathbf{A}_t \Delta_t\|_\Sigma^2 \\ &\quad - 2\eta_t \langle \Delta_t, \mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t \rangle_\Sigma + 2\eta_t^2 \langle \mathbf{A}_t \Delta_t, \mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t \rangle_\Sigma \\ &\quad + \eta_t^2 \|\mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t\|_\Sigma^2. \end{aligned}$$

Notice that by definition,

$$\mathbb{E}_t \langle \Delta_t, \mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t \rangle_\Sigma = \langle \Delta_t, \mathbf{A} \boldsymbol{\theta}^* - \mathbf{b} \rangle = 0,$$

and that a basic property of inner product yields

$$2 \langle \mathbf{A}_t \Delta_t, \mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t \rangle_\Sigma \leq \|\mathbf{A}_t \Delta_t\|_\Sigma^2 + \|\mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t\|_\Sigma^2.$$

Therefore, we can apply the law of total expectations to obtain the following iterative relation:

$$\begin{aligned} \mathbb{E} \|\Delta_{t+1}\|_\Sigma^2 &= \mathbb{E} \|\Delta_t\|_\Sigma^2 - \underbrace{2\eta_t \mathbb{E}[\Delta_t^\top (\mathbf{A}^\top \Sigma + \Sigma \mathbf{A}) \Delta_t]}_{I_1} \\ &\quad + \underbrace{2\eta_t^2 \mathbb{E} \|\mathbf{A}_t \Delta_t\|_\Sigma^2}_{I_2} + \underbrace{2\eta_t^2 \mathbb{E} \|\mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t\|_\Sigma^2}_{I_3}. \end{aligned} \quad (152)$$

We now turn to bounding  $I_1$ ,  $I_2$  and  $I_3$  in order.

b) *Bounding  $I_1$* : In order to lower bound  $I_1$  as a function of  $\|\Delta_t\|_\Sigma^2$ , we firstly express it as

$$\begin{aligned} &\Delta_t^\top (\mathbf{A}^\top \Sigma + \Sigma \mathbf{A}) \Delta_t \\ &= \Delta_t^\top \Sigma^{1/2} \Sigma^{-1/2} (\mathbf{A}^\top \Sigma + \Sigma \mathbf{A}) \Sigma^{-1} \Sigma^{1/2} \Delta_t \\ &\geq \|\Sigma^{1/2} \Delta_t\|_2^2 \lambda_{\min} \left( \Sigma^{-1/2} \mathbf{A}^\top \Sigma^{1/2} + \Sigma^{1/2} \mathbf{A} \Sigma^{-1/2} \right) \\ &= \|\Delta_t\|_\Sigma^2 \lambda_{\min} \left( \Sigma^{-1/2} \mathbf{A}^\top \Sigma^{1/2} + \Sigma^{1/2} \mathbf{A} \Sigma^{-1/2} \right). \end{aligned}$$

Recall from (88e) that

$$\|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \leq (1-\gamma)^{-1},$$

so the minimal eigenvalue of  $\Sigma^{-1/2} \mathbf{A}^\top \Sigma^{1/2} + \Sigma^{1/2} \mathbf{A} \Sigma^{-1/2}$  is lower bounded by

$$\begin{aligned} &\lambda_{\min} \left( \Sigma^{-1/2} \mathbf{A}^\top \Sigma^{1/2} + \Sigma^{1/2} \mathbf{A} \Sigma^{-1/2} \right) \\ &\geq \lambda_{\min}(\Sigma) \cdot \left[ \gamma_{\min} \left( \Sigma^{-\frac{1}{2}} \mathbf{A}^\top \Sigma^{-\frac{1}{2}} \right) + \gamma_{\min} \left( \Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-\frac{1}{2}} \right) \right] \\ &\geq \frac{2\lambda_{\min}(\Sigma)}{\|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\|} \geq 2\lambda_{\min}(\Sigma)(1-\gamma). \end{aligned}$$

This directly implies that  $I_1$  is lower bounded by

$$I_1 \geq 2\eta_t(1-\gamma)\lambda_{\min}(\Sigma)\mathbb{E}\|\Delta_t\|_\Sigma^2. \quad (153)$$

c) *Bounding  $I_2$* : We aim to upper bound  $I_2$  as a function of  $\eta_t^2$  and  $\|\Delta_t\|_\Sigma^2$ , so that when  $\eta_t$  is sufficiently small,  $I_2$  is negligible compared to  $I_1$ . Specifically, for any  $\mathbf{A}_t$  generated by (11a) and any  $\Delta_t \in \mathbb{R}^d$ , we observe

$$\begin{aligned} \|\mathbf{A}_t \Delta_t\|_\Sigma^2 &= \Delta_t^\top \mathbf{A}_t \Sigma \mathbf{A}_t \Delta_t \\ &\leq \|\Delta_t\|_2^2 \|\mathbf{A}\|^2 \|\Sigma\| \leq 4\|\Sigma\| \|\Delta_t\|_2^2 \\ &\leq 4\|\Sigma\| \|\Sigma^{-1}\| \|\Sigma^{\frac{1}{2}} \Delta_t\|_2^2 = 4\kappa \|\Delta_t\|_\Sigma^2, \end{aligned}$$

where we recall  $\kappa$  as the condition number of  $\Sigma$ . Therefore, as long as

$$\eta_t \leq \frac{(1-\gamma)\lambda_{\min}(\Sigma)}{4\kappa},$$

it can be guaranteed that  $I_2 \leq \frac{1}{2}I_1$ .

d) *Bounding  $I_3$* : In order to compare with our result (Theorem 1 and Corollary 1), we aim to bound  $I_3$  as a function of  $\|\boldsymbol{\theta}^*\|_\Sigma$ . Towards this end, we firstly notice that

$$\mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t = \phi(s_t) \phi(s_t)^\top \boldsymbol{\theta}^* - \gamma \phi(s_t) \phi(s'_t)^\top \boldsymbol{\theta}^* - r(s_t) \phi(s_t).$$

Therefore, we can upper bound  $\mathbb{E}\|\mathbf{A}_t\boldsymbol{\theta}^* - \mathbf{b}_t\|_{\Sigma}^2$  by

$$\begin{aligned}\mathbb{E}\|\mathbf{A}_t\boldsymbol{\theta}^* - \mathbf{b}_t\|_{\Sigma}^2 &\leq 3 \mathbb{E}_{s \sim \mu} \|\phi(s)\phi(s)^\top \boldsymbol{\theta}^*\|_{\Sigma}^2 \\ &\quad + 3 \mathbb{E}_{s \sim \mu, s' \sim \mathcal{P}(\cdot|s)} \|\phi(s)\phi(s')^\top \boldsymbol{\theta}^*\|_{\Sigma}^2 \\ &\quad + 3 \mathbb{E}_{s \sim \mu} \|r(s)\phi(s)\|_{\Sigma}^2,\end{aligned}$$

where the three terms on the right-hand-side can be bounded respectively by

$$\begin{aligned}&\mathbb{E}_{s \sim \mu} \|\phi(s)\phi(s)^\top \boldsymbol{\theta}^*\|_{\Sigma}^2 \\ &= \mathbb{E}_{s \sim \mu} [\boldsymbol{\theta}^{*\top} \phi(s) (\phi(s)^\top \Sigma \phi(s)) \phi(s)^\top \boldsymbol{\theta}^*] \\ &\leq \mathbb{E}_{s \sim \mu} [\boldsymbol{\theta}^{*\top} \phi(s) \|\Sigma\| \phi(s)^\top \boldsymbol{\theta}^*] \\ &= \|\Sigma\| \boldsymbol{\theta}^{*\top} \mathbb{E}_{s \sim \mu} [\phi(s)\phi(s)^\top] \boldsymbol{\theta}^* \\ &= \|\Sigma\| \boldsymbol{\theta}^{*\top} \Sigma \boldsymbol{\theta}^* = \|\Sigma\| \|\boldsymbol{\theta}^*\|_{\Sigma}^2; \\ &\mathbb{E}_{s \sim \mu, s' \sim \mathcal{P}(\cdot|s)} \|\phi(s)\phi(s')^\top \boldsymbol{\theta}^*\|_{\Sigma}^2 \\ &= \mathbb{E}_{s \sim \mu, s' \sim \mathcal{P}(\cdot|s)} [\boldsymbol{\theta}^{*\top} \phi(s') (\phi(s)^\top \Sigma \phi(s)) \phi(s')^\top \boldsymbol{\theta}^*] \\ &\leq \mathbb{E}_{s \sim \mu, s' \sim \mathcal{P}(\cdot|s)} [\boldsymbol{\theta}^{*\top} \phi(s') \|\Sigma\| \phi(s)^\top \boldsymbol{\theta}^*] \\ &= \|\Sigma\| \boldsymbol{\theta}^{*\top} \mathbb{E}_{s' \sim \mu} [\phi(s')\phi(s')^\top] \boldsymbol{\theta}^* \\ &= \|\Sigma\| \boldsymbol{\theta}^{*\top} \Sigma \boldsymbol{\theta}^* = \|\Sigma\| \|\boldsymbol{\theta}^*\|_{\Sigma}^2,\end{aligned}$$

$$\text{and } \mathbb{E}_{s \sim \mu} \|r(s)\phi(s)\|_{\Sigma}^2 \leq \max_{s \in \mathcal{S}} r^2(s) \|\phi(s)\|_2^2 \|\Sigma\| \leq \|\Sigma\|.$$

Consequently,  $I_3$  can be upper bounded by

$$I_3 \leq 6\eta_t^2 \|\Sigma\| (2\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1). \quad (154)$$

*e) Bounding  $\mathbb{E}\|\Delta_T\|_{\Sigma}^2$ :* By combining (152), (153) and (154) and recalling that  $I_2 \leq \frac{1}{2}I_1$  when  $\eta_t$  is sufficiently small, we obtain

$$\begin{aligned}\mathbb{E}\|\Delta_{t+1}\|_{\Sigma}^2 &\leq (1 - (1 - \gamma)\lambda_{\min}(\Sigma)\eta_t) \mathbb{E}\|\Delta_t\|_{\Sigma}^2 \\ &\quad + 6\eta_t^2 \|\Sigma\| (2\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1).\end{aligned} \quad (155)$$

Therefore, for constant stepsizes  $\eta_0 = \eta_1 = \dots = \eta_T = \eta$ , it is easy to verify by induction that

$$\begin{aligned}\mathbb{E}\|\Delta_T\|_{\Sigma}^2 &\leq (1 - (1 - \gamma)\lambda_{\min}(\Sigma)\eta)^T \|\Delta_0\|_{\Sigma}^2 \\ &\quad + \frac{6\eta \|\Sigma\| (2\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1)}{(1 - \gamma)\lambda_{\min}(\Sigma)}.\end{aligned}$$

Hence, in order to guarantee  $\mathbb{E}\|\Delta_T\|_{\Sigma}^2 \leq \varepsilon^2$ , it suffices to take

$$\begin{aligned}\frac{\eta \|\Sigma\| (\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1)}{(1 - \gamma)\lambda_{\min}(\Sigma)} &\lesssim \varepsilon^2; \quad \text{and} \\ \exp(-(1 - \gamma)\lambda_{\min}(\Sigma)\eta T) \|\Delta_0\|_{\Sigma}^2 &\lesssim \varepsilon^2.\end{aligned}$$

This implies the following upper bound for the sample complexity:

$$T \asymp \frac{\kappa \|\Sigma^{-1}\| (\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1)}{(1 - \gamma)^2} \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon}, \quad (156)$$

with the proviso that we take the stepsize  $\eta \asymp \frac{\|\Sigma^{-1}\|}{1 - \gamma} \frac{1}{T}$  and that  $T \gtrsim \|\Sigma^{-2}\| (1 - \gamma)^{-2}$ .

### B. Comparisons with [9]

Theorem 2(c) in [9] shows that with decaying stepsizes  $\eta_t = \frac{\beta}{\lambda + t}$  where

$$\beta = \frac{2\|\Sigma^{-1}\|}{(1 - \gamma)}, \quad \lambda = \frac{16\|\Sigma^{-1}\|}{(1 - \gamma)^2}, \quad (157)$$

the expected  $\ell_2$  norm of TD estimation error is bounded by

$$\mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_2^2 \leq \frac{\nu}{\lambda + T}, \quad (158)$$

where

$$\nu = \max \left\{ \frac{8\sigma^2 \|\Sigma^{-2}\|}{(1 - \gamma)^2}, \frac{16\|\boldsymbol{\theta}^*\|_2^2 \|\Sigma^{-1}\|}{(1 - \gamma)^2} \right\}. \quad (159)$$

- Suppose the maximum is attained at the second term for  $\nu$  and  $T$  is sufficiently large, (158) is simplified as

$$\mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_2^2 \lesssim \frac{16\|\boldsymbol{\theta}^*\|_2^2 \|\Sigma^{-1}\|}{(1 - \gamma)^2 T}.$$

In order for  $\mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_{\Sigma}^2 \leq \varepsilon^2$ , it suffices to take

$$\frac{\varepsilon^2}{\|\Sigma\|} \geq \frac{16\|\boldsymbol{\theta}^*\|_2^2 \|\Sigma^{-1}\|}{(1 - \gamma)^2 T} \geq \mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_2^2,$$

which implies the following sample complexity:

$$T \asymp \frac{\|\Sigma^{-1}\| \|\Sigma\| \|\boldsymbol{\theta}^*\|_2^2}{(1 - \gamma)^2 \varepsilon^2}$$

- Suppose that the first term on the right hand side of expression (159) is larger, (158) can be simplified as

$$\mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_2^2 \lesssim \frac{\sigma^2 \|\Sigma^{-2}\|}{(1 - \gamma)^2 T}.$$

Then similarly, the sample complexity is

$$T \asymp \frac{\|\Sigma^{-2}\| \|\Sigma\| \sigma^2}{(1 - \gamma)^2 \varepsilon^2},$$

where  $\sigma^2 = \mathbb{E}\|\mathbf{A}_t\boldsymbol{\theta}^* - \mathbf{b}_t\|_2^2$ .

In the worst-case scenario, it satisfies  $\sigma^2 \asymp \|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1$ . Therefore, the sample complexity implied by Theorem 2(c) of [9] scales as

$$T \asymp \frac{\kappa \|\Sigma^{-1}\| (\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1)}{(1 - \gamma)^2} \frac{1}{\varepsilon^2}. \quad (160)$$

### C. Comparison with [63] and [64]

[63] studied a more general problem of linear approximation for fixed point equations in Hilbert spaces, and considered its application to TD learning with linear function approximation and *i.i.d.* samples. A similar result was reached by [64] in their Theorem 2 and Theorem 3. While these works explored both the *approximation error*, which measures the difference between  $\Phi\boldsymbol{\theta}^*$  and  $\mathbf{V}^*$  under our notation, and the *statistical error*, which measures the convergence of  $\boldsymbol{\theta}_T$  to  $\boldsymbol{\theta}^*$ , it is the latter that is directly comparable to our results. Therefore, we hereby provide a comparison between the statistical error term

in their Corollary 5 and the sample complexity result of ours as shown in Theorem 1. Translated to our notation, [63] proved that with a sufficiently large sample size  $T$  and a stepsize of  $\eta \asymp \frac{1}{\sqrt{T}}$ , the estimation error of the averaged TD learning algorithm satisfies

$$\begin{aligned} & \mathbb{E}_{s \sim \mu} [V_{\bar{\theta}_T}(s) - V_{\theta^*}(s)]^2 \\ &= \|\bar{\theta}_T - \theta^*\|_{\Sigma}^2 \\ &\lesssim \frac{1}{T} \text{Tr} [(I - M)^{-1} (\Sigma_L + \Sigma_b) (I - M)^{-\top}], \end{aligned} \quad (161)$$

in which  $M$ ,  $\Sigma_L$  and  $\Sigma_b$  are defined as

$$\begin{aligned} M &= \gamma \Sigma^{-\frac{1}{2}} \mathbb{E}_{s \sim \mu, s' \sim P(\cdot|s)} [\phi(s) \phi(s')^\top] \Sigma^{-\frac{1}{2}}, \\ \Sigma_L &= \text{Cov}_{s_t \sim \mu, s'_t \sim P(\cdot|s_t)} [\Sigma^{-\frac{1}{2}} \mathbf{A}_t \theta^*], \quad \text{and} \\ \Sigma_b &= \text{Cov}_{s_t \sim \mu} [\Sigma^{-\frac{1}{2}} \mathbf{b}_t]. \end{aligned}$$

*a) Translation into our notation:* We firstly translate the upper bound (161) into our notation. By definition,

$$\mathbb{E}_{s \sim \mu, s' \sim P(\cdot|s)} [\phi(s) \phi(s')^\top] = \Phi^\top D_\mu P \Phi.$$

Therefore, the term  $I - M$  can be expressed as

$$\begin{aligned} I - M &= \Sigma^{-\frac{1}{2}} \Sigma \Sigma^{-\frac{1}{2}} - \gamma \Sigma^{-\frac{1}{2}} \Phi^\top D_\mu P \Phi \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} [\Phi^\top D_\mu \Phi - \gamma \Phi^\top D_\mu P \Phi] \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} \Phi^\top D_\mu (I - \gamma P) \Phi \Sigma^{-\frac{1}{2}} = \Sigma^{-\frac{1}{2}} \mathbf{A} \Sigma^{-\frac{1}{2}}. \end{aligned}$$

Furthermore, the terms  $\Sigma_L$  and  $\Sigma_b$  can be expressed in our notation as

$$\begin{aligned} \Sigma_L &= \Sigma^{-\frac{1}{2}} \text{Cov}_{s_t \sim \mu, s'_t \sim P(\cdot|s_t)} [\mathbf{A}_t \theta^*] \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} \mathbb{E} [[(\mathbf{A}_t - \mathbf{A}) \theta^*][(\mathbf{A}_t - \mathbf{A}) \theta^*]^\top] \Sigma^{-\frac{1}{2}}, \end{aligned}$$

and

$$\begin{aligned} \Sigma_b &= \Sigma^{-\frac{1}{2}} \text{Cov}_{s_t \sim \mu} [\mathbf{b}_t] \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} \mathbb{E}_{s_t \sim \mu} [(\mathbf{b}_t - \mathbf{b})(\mathbf{b}_t - \mathbf{b})^\top] \Sigma^{-\frac{1}{2}}. \end{aligned}$$

For simplicity, we will omit the subscript  $s_t \sim \mu, s'_t \sim P(\cdot|s_t)$  in the following. Combining these terms, the upper bound in (161) can be expressed as

$$\begin{aligned} & \frac{1}{T} \text{Tr} [(I - M)^{-1} (\Sigma_L + \Sigma_b) (I - M)^{-\top}] \\ &= \frac{1}{T} \text{Tr} \left[ \left( \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \right) \Sigma^{-\frac{1}{2}} \mathbb{E} [[(\mathbf{A}_t - \mathbf{A}) \theta^*][(\mathbf{A}_t - \mathbf{A}) \theta^*]^\top] \right. \\ &\quad \left. + (\mathbf{b}_t - \mathbf{b})(\mathbf{b}_t - \mathbf{b})^\top \right] \Sigma^{-\frac{1}{2}} \left( \Sigma^{\frac{1}{2}} \mathbf{A}^{-\top} \Sigma^{\frac{1}{2}} \right) \\ &= \frac{1}{T} \text{Tr} \left[ \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbb{E} [[(\mathbf{A}_t - \mathbf{A}) \theta^*][(\mathbf{A}_t - \mathbf{A}) \theta^*]^\top] \mathbf{A}^{-\top} \Sigma^{\frac{1}{2}} \right] \\ &\quad + \frac{1}{T} \text{Tr} \left[ \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \mathbb{E} [(\mathbf{b}_t - \mathbf{b})(\mathbf{b}_t - \mathbf{b})^\top] \mathbf{A}^{-\top} \Sigma^{\frac{1}{2}} \right] \\ &= \frac{1}{T} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_t - \mathbf{A}) \theta^*\|_{\Sigma}^2 + \frac{1}{T} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{b}_t - \mathbf{b})\|_{\Sigma}^2 \end{aligned}$$

So in summary, [63] bounds the estimation error by

$$\begin{aligned} & \|\bar{\theta}_T - \theta^*\|_{\Sigma} \\ &\lesssim \frac{1}{T} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_t - \mathbf{A}) \theta^*\|_{\Sigma}^2 + \frac{1}{T} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{b}_t - \mathbf{b})\|_{\Sigma}^2. \end{aligned} \quad (162)$$

*b) Comparison to our results:* In the following, we show that the upper bound (162) can be directly deducted from our proof of Theorem 1. Specifically, our analysis in (62), (115) and (117) reveals that  $\|\bar{\theta}_T - \theta^*\|_{\Sigma}^2$  is bounded by

$$\begin{aligned} \|\bar{\theta}_T - \theta^*\|_{\Sigma}^2 &\lesssim \frac{1}{T^2} \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta_i \right\|_{\Sigma}^2 \\ &\quad + \frac{1}{T^2} \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{b}_i - \mathbf{b}) \right\|_{\Sigma}^2 + o\left(\frac{1}{T}\right). \end{aligned}$$

Taking expectations on both sides and applying the martingale property, we obtain

$$\begin{aligned} & \mathbb{E} \|\bar{\theta}_T - \theta^*\|_{\Sigma}^2 \\ &\lesssim \frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta_i\|_{\Sigma}^2 + \frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{b}_i - \mathbf{b})\|_{\Sigma}^2 \\ &\lesssim \frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \theta^*\|_{\Sigma}^2 + \frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{b}_i - \mathbf{b})\|_{\Sigma}^2 \\ &\quad + \frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \Delta_i\|_{\Sigma}^2 \\ &= \frac{1}{T} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_t - \mathbf{A}) \theta^*\|_{\Sigma}^2 + \frac{1}{T} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{b}_t - \mathbf{b})\|_{\Sigma}^2 \\ &\quad + \frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \Delta_i\|_{\Sigma}^2. \end{aligned}$$

Notice here that the first two terms are exactly the same as the right-hand-side of (162). Hence, it now boils down to showing that

$$\frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \Delta_i\|_{\Sigma}^2 = o\left(\frac{1}{T}\right). \quad (163)$$

Towards this end, we firstly observe that

$$\begin{aligned} & \frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \Delta_i\|_{\Sigma}^2 \\ &\lesssim \frac{\|\Sigma\|^2 \|\Sigma^{-1}\|^2}{(1 - \gamma)^2 T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\Delta_i\|_{\Sigma}^2. \end{aligned}$$

For the expectation of  $\|\Delta_i\|_{\Sigma}^2$ , we again apply the iterative relation deducted in (155) and obtain

$$\begin{aligned} \mathbb{E} \|\Delta_i\|_{\Sigma}^2 &\leq (1 - (1 - \gamma) \lambda_{\min}(\Sigma) \eta)^i \|\Delta_0\|_{\Sigma}^2 \\ &\quad + \frac{6\eta \|\Sigma\| (2\|\theta^*\|_{\Sigma}^2 + 1)}{(1 - \gamma) \lambda_{\min}(\Sigma)}. \end{aligned}$$

Summing from  $i = 0$  through  $i = T - 1$  yields

$$\begin{aligned} \sum_{i=0}^{T-1} \mathbb{E} \|\Delta_i\|_{\Sigma}^2 &\leq \frac{1}{(1 - \gamma) \lambda_{\min}(\Sigma) \eta} \|\Delta_0\|_{\Sigma}^2 \\ &\quad + \frac{6\eta T \|\Sigma\| (2\|\theta^*\|_{\Sigma}^2 + 1)}{(1 - \gamma) \lambda_{\min}(\Sigma)}. \end{aligned}$$

By setting  $\eta \asymp T^{-1/2}$  as suggested by [63], this immediately implies

$$\frac{1}{T^2} \sum_{i=0}^{T-1} \mathbb{E} \|\mathbf{A}^{-1}(\mathbf{A}_i - \mathbf{A})\Delta_i\|_{\Sigma}^2 \lesssim \frac{1}{T^2} \cdot T^{1/2} = o\left(\frac{1}{T}\right).$$

In summary, we have shown that the upper bound proposed by [63] follows directly from our analysis. Our result, as is shown in Theorem 1, improves upon theirs in the sense that we use a stepsize  $\eta$  that only depends on the logarithm of  $T$ , provide a bound with high probability instead of in expectation, and reveal a clearer dependence on the problem-related parameters.

#### D. Comparison with [31]

It is difficult to place the corresponding instance dependent results in comparison, so, we focus our attention on the minimax results. In the following, we make use of the relations that  $\|\mathbf{A}(\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*)\|_2 \geq \|\mathbf{A}\Sigma^{-1/2}\| \|\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*\|_{\Sigma} \gtrsim (1 - \gamma)\sqrt{\lambda_{\min}(\Sigma)} \|\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*\|_{\Sigma}$ , and  $\mathbb{E}\|\mathbf{A}_t\boldsymbol{\theta}^* - \mathbf{b}_t\|_2^2 \lesssim \frac{1}{(1-\gamma)^2}$ ,  $\sup \|\mathbf{A}_t\boldsymbol{\theta}^* - \mathbf{b}_t\|_2 \lesssim \frac{1}{1-\gamma}$ . We also consider the situations when  $\|\boldsymbol{\theta}^*\|_{\Sigma} \lesssim \frac{1}{1-\gamma}$ , and  $\phi(s)^\top \Sigma^{-1} \phi(s) \lesssim \lambda_{\min}(\Sigma)^{-1}$ . Notice that there exists an MDP instance such that equality can be attained in all these bounds simultaneously. For ease of presentation, let us first rephrase the result [31, Corollary 1] in terms of our notation<sup>2</sup>. It is shown therein that for

$$\eta \lesssim \frac{(1-\gamma)^3 \lambda_{\min}(\Sigma)}{\kappa \sqrt{T}},$$

with probability at least  $1 - \delta$ , the averaged TD estimation error is bounded by

$$\|\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*\|_{\Sigma} \lesssim \sqrt{\frac{1}{\lambda_{\min}(\Sigma)(1-\gamma)^4 T}}, \quad (164)$$

when  $T \gtrsim \frac{1}{c_A^2} \gtrsim \frac{\kappa^2}{(1-\gamma)^6 \lambda_{\min}(\Sigma)^2}$ . Here, we omit the dependency of log factors. In comparison, our result delivers the same bound as long as  $T \gtrsim \frac{\kappa^2}{(1-\gamma)^4 \lambda_{\min}(\Sigma)}$ . We incur a lower born-in cost for the relation (164) to hold.

#### REFERENCES

- [1] S. A. Murphy, "Optimal dynamic treatment regimes," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 65, no. 2, pp. 331–355, 2003.
- [2] I. Bojinov and N. Shephard, "Time series experiments and causal estimands: exact randomization tests and trading," *Journal of the American Statistical Association*, vol. 114, no. 528, pp. 1665–1682, 2019.
- [3] S. Tang and J. Wiens, "Model selection for offline reinforcement learning: Practical considerations for healthcare settings," in *Machine Learning for Healthcare Conference*. PMLR, 2021, pp. 2–35.
- [4] C. Dann, G. Neumann, J. Peters *et al.*, "Policy evaluation with temporal differences: A survey and comparison," *Journal of Machine Learning Research*, vol. 15, pp. 809–883, 2014.
- [5] D. Bertsimas, P. Klasnja, S. Murphy, and L. Na, "Data-driven interpretable policy construction for personalized mobile health," in *2022 IEEE International Conference on Digital Health (ICDH)*. IEEE, 2022, pp. 13–22.
- [6] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [7] D. P. Bertsekas, *Dynamic programming and optimal control (4th edition)*. Athena Scientific, 2017.
- [8] J. Tsitsiklis and B. Van Roy, "An analysis of temporal-difference learning with function approximation," *IEEE Transactions on Automatic Control*, vol. 42, no. 5, pp. 674–690, 1997.
- [9] J. Bhandari, D. Russo, and R. Singal, "A finite time analysis of temporal difference learning with linear function approximation," *Operations Research*, vol. 69, no. 3, pp. 950–973, 2021.
- [10] J. Fan, Z. Wang, Y. Xie, and Z. Yang, "A theoretical analysis of deep q-learning," in *Learning for Dynamics and Control*. PMLR, 2020, pp. 486–489.
- [11] A.-m. Farahmand, M. Ghavamzadeh, C. Szepesvári, and S. Mannor, "Regularized policy iteration with nonparametric function spaces," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 4809–4874, 2016.
- [12] Y. Duan, M. Wang, and M. J. Wainwright, "Optimal policy evaluation using kernel-based temporal difference methods," *arXiv preprint arXiv:2109.12002*, 2021.
- [13] D. P. Bertsekas and J. N. Tsitsiklis, "Neuro-dynamic programming: an overview," in *Proceedings of 1995 34th IEEE conference on decision and control*, vol. 1. IEEE, 1995, pp. 560–564.
- [14] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "A brief survey of deep reinforcement learning," *arXiv preprint arXiv:1708.05866*, 2017.
- [15] C. Jin, Z. Yang, Z. Wang, and M. I. Jordan, "Provably efficient reinforcement learning with linear function approximation," in *Conference on Learning Theory*. PMLR, 2020, pp. 2137–2143.
- [16] B. Wang, Y. Yan, and J. Fan, "Sample-efficient reinforcement learning for linearly-parameterized mdps with a generative model," *Advances in neural information processing systems*, vol. 34, pp. 23 009–23 022, 2021.
- [17] G. Li, Y. Chen, Y. Chi, Y. Gu, and Y. Wei, "Sample-efficient reinforcement learning is feasible for linearly realizable mdps with limited revisiting," *Advances in Neural Information Processing Systems*, vol. 34, pp. 16 671–16 685, 2021.
- [18] R. S. Sutton, "Learning to predict by the methods of temporal differences," *Machine learning*, vol. 3, no. 1, pp. 9–44, 1988.
- [19] G. Li, C. Cai, Y. Chen, Y. Wei, and Y. Chi, "Is Q-learning minimax optimal? a tight sample complexity analysis," *Operations Research*, 2023.
- [20] G. Dalal, B. Szörényi, G. Thoppe, and S. Mannor, "Finite sample analyses for TD(0) with function approximation," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [21] R. Srikant and L. Ying, "Finite-time error bounds for linear stochastic approximation and TD learning," in *Conference on Learning Theory*, 2019, pp. 2803–2830.
- [22] 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*. PMLR, 2018, pp. 1347–1355.
- [23] L. Baird, "Residual algorithms: Reinforcement learning with function approximation," in *Machine Learning Proceedings 1995*. Elsevier, 1995, pp. 30–37.
- [24] R. S. Sutton, H. R. Maei, D. Precup, S. Bhatnagar, D. Silver, C. Szepesvári, and E. Wiewiora, "Fast gradient-descent methods for temporal-difference learning with linear function approximation," in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 993–1000.
- [25] G. Dalal, G. Thoppe, B. Szörényi, and S. Mannor, "Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning," in *Conference On Learning Theory*, 2018, pp. 1199–1233.
- [26] T. Xu and Y. Liang, "Sample complexity bounds for two timescale value-based reinforcement learning algorithms," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 811–819.
- [27] Y. Wang, S. Zou, and Y. Zhou, "Non-asymptotic analysis for two timescale tdc with general smooth function approximation," *Advances in Neural Information Processing Systems*, vol. 34, pp. 9747–9758, 2021.
- [28] G. Dalal, B. Szorenyi, and G. Thoppe, "A tale of two-timescale reinforcement learning with the tightest finite-time bound," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 3701–3708.
- [29] M. Kaledin, E. Moulines, A. Naumov, V. Tadic, and H.-T. Wai, "Finite time analysis of linear two-timescale stochastic approximation with Markovian noise," in *Conference on Learning Theory*. PMLR, 2020, pp. 2144–2203.

<sup>2</sup>We take  $C_A \lesssim (1-\gamma)^{-1}$ ,  $a \asymp \|\mathbf{Q}\|^{-1} \lesssim (1-\gamma)\lambda_{\min}(\Sigma)$  and then  $c_A \lesssim \kappa^{-1}(1-\gamma)^3 \lambda_{\min}(\Sigma)$  (see the definitions of these parameters in [31]).



- [30] H. Gupta, R. Srikant, and L. Ying, "Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [31] A. Durmus, E. Moulines, A. Naumov, and S. Samsonov, "Finite-time high-probability bounds for polyak-ruppert averaged iterates of linear stochastic approximation," *arXiv preprint arXiv:2207.04475*, 2022.
- [32] C. Szepesvári, "The asymptotic convergence-rate of Q-learning," in *Advances in Neural Information Processing Systems*, 1998, pp. 1064–1070.
- [33] K. Khamaru, A. Pananjady, F. Ruan, M. J. Wainwright, and M. I. Jordan, "Is temporal difference learning optimal? an instance-dependent analysis," *arXiv preprint arXiv:2003.07337*, 2020.
- [34] C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan, "Is Q-learning provably efficient?" in *Advances in Neural Information Processing Systems*, 2018, pp. 4863–4873.
- [35] J. A. Boyan, "Least-squares temporal difference learning," in *ICML*, 1999, pp. 49–56.
- [36] A. Sidford, M. Wang, X. Wu, L. Yang, and Y. Ye, "Near-optimal time and sample complexities for solving Markov decision processes with a generative model," in *Advances in Neural Information Processing Systems*, 2018, pp. 5186–5196.
- [37] A. Agarwal, S. Kakade, and L. F. Yang, "Model-based reinforcement learning with a generative model is minimax optimal," in *Conference on Learning Theory*. PMLR, 2020, pp. 67–83.
- [38] A. Pananjady and M. J. Wainwright, "Instance-dependent  $\ell_\infty$ -bounds for policy evaluation in tabular reinforcement learning," *IEEE Transactions on Information Theory*, vol. 67, no. 1, pp. 566–585, 2021.
- [39] G. Li, Y. Wei, Y. Chi, and Y. Chen, "Breaking the sample size barrier in model-based reinforcement learning with a generative model," *Operations Research*, 2023.
- [40] T. L. Lai, "Stochastic approximation," *The Annals of Statistics*, vol. 31, no. 2, pp. 391–406, 2003.
- [41] H. Robbins and S. Monro, "A stochastic approximation method," *The Annals of Mathematical Statistics*, pp. 400–407, 1951.
- [42] V. S. Borkar, *Stochastic approximation: a dynamical systems viewpoint*. Springer, 2009, vol. 48.
- [43] V. S. Borkar and S. P. Meyn, "The ODE method for convergence of stochastic approximation and reinforcement learning," *SIAM Journal on Control and Optimization*, vol. 38, no. 2, pp. 447–469, 2000.
- [44] 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," *arXiv preprint arXiv:2004.04719*, 2020.
- [45] E. Moulines and F. Bach, "Non-asymptotic analysis of stochastic approximation algorithms for machine learning," *Advances in neural information processing systems*, vol. 24, 2011.
- [46] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, "Robust stochastic approximation approach to stochastic programming," *SIAM Journal on optimization*, vol. 19, no. 4, pp. 1574–1609, 2009.
- [47] T. Xu, S. Zou, and Y. Liang, "Two time-scale off-policy TD learning: Non-asymptotic analysis over Markovian samples," in *Advances in Neural Information Processing Systems*, 2019, pp. 10 633–10 643.
- [48] Y. F. Wu, W. Zhang, P. Xu, and Q. Gu, "A finite-time analysis of two time-scale actor-critic methods," *Advances in Neural Information Processing Systems*, vol. 33, pp. 17 617–17 628, 2020.
- [49] D. Precup, "Eligibility traces for off-policy policy evaluation," *Computer Science Department Faculty Publication Series*, p. 80, 2000.
- [50] N. Jiang and L. Li, "Doubly robust off-policy value evaluation for reinforcement learning," in *International Conference on Machine Learning*. PMLR, 2016, pp. 652–661.
- [51] C. Ma, B. Zhu, J. Jiao, and M. J. Wainwright, "Minimax off-policy evaluation for multi-armed bandits," *IEEE Transactions on Information Theory*, vol. 68, no. 8, pp. 5314–5339, 2022.
- [52] P. Thomas and E. Brunskill, "Data-efficient off-policy policy evaluation for reinforcement learning," in *International Conference on Machine Learning*. PMLR, 2016, pp. 2139–2148.
- [53] T. Xie, Y. Ma, and Y.-X. Wang, "Towards optimal off-policy evaluation for reinforcement learning with marginalized importance sampling," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [54] N. Kallus and M. Uehara, "Double reinforcement learning for efficient off-policy evaluation in markov decision processes," *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 6742–6804, 2020.
- [55] M. Yang, O. Nachum, B. Dai, L. Li, and D. Schuurmans, "Off-policy evaluation via the regularized Lagrangian," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6551–6561, 2020.
- [56] Y. Duan, Z. Jia, and M. Wang, "Minimax-optimal off-policy evaluation with linear function approximation," in *International Conference on Machine Learning*. PMLR, 2020, pp. 2701–2709.
- [57] Y. Jin, Z. Yang, and Z. Wang, "Is pessimism provably efficient for offline RL?" in *International Conference on Machine Learning*. PMLR, 2021, pp. 5084–5096.
- [58] T. Xie, N. Jiang, H. Wang, C. Xiong, and Y. Bai, "Policy finetuning: Bridging sample-efficient offline and online reinforcement learning," *Advances in neural information processing systems*, vol. 34, pp. 27 395–27 407, 2021.
- [59] G. Li, L. Shi, Y. Chen, Y. Chi, and Y. Wei, "Settling the sample complexity of model-based offline reinforcement learning," *arXiv preprint arXiv:2204.05275*, 2022.
- [60] L. Shi, G. Li, Y. Wei, Y. Chen, and Y. Chi, "Pessimistic q-learning for offline reinforcement learning: Towards optimal sample complexity," in *International Conference on Machine Learning*. PMLR, 2022, pp. 19 967–20 025.
- [61] P. Rashidinejad, B. Zhu, C. Ma, J. Jiao, and S. Russell, "Bridging offline reinforcement learning and imitation learning: A tale of pessimism," *Advances in Neural Information Processing Systems*, vol. 34, pp. 11 702–11 716, 2021.
- [62] G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen, "Sample complexity of asynchronous Q-learning: Sharper analysis and variance reduction," *IEEE Transactions on Information Theory*, vol. 68, no. 1, pp. 448–473, 2021.
- [63] W. Mou, A. Pananjady, and M. J. Wainwright, "Optimal oracle inequalities for solving projected fixed-point equations," *arXiv preprint arXiv:2012.05299*, 2020.
- [64] T. Li, G. Lan, and A. Pananjady, "Accelerated and instance-optimal policy evaluation with linear function approximation," *arXiv preprint arXiv:2112.13109*, 2021.
- [65] R. S. Sutton, C. Szepesvári, and H. R. Maei, "A convergent  $O(n)$  algorithm for off-policy temporal-difference learning with linear function approximation," *Advances in neural information processing systems*, vol. 21, no. 21, pp. 1609–1616, 2008.
- [66] E. N. Gilbert, "A comparison of signalling alphabets," *The Bell system technical journal*, vol. 31, no. 3, pp. 504–522, 1952.
- [67] A. B. Tsybakov, *Introduction to nonparametric estimation*. Springer, 2009, vol. 11.
- [68] J. Tropp, "Freedman's inequality for matrix martingales," *Electronic Communications in Probability*, vol. 16, pp. 262–270, 2011.
- [69] J. A. Tropp, "An introduction to matrix concentration inequalities," *Found. Trends Mach. Learn.*, vol. 8, no. 1-2, pp. 1–230, May 2015. [Online]. Available: <http://dx.doi.org/10.1561/22000000048>

**Gen Li** (Member, IEEE) received his Ph.D. degree from the Department of Electronic Engineering at Tsinghua University, advised by Professor Yuantao Gu, and bachelor's degree from the Department of Electronic Engineering and Department of Mathematics at Tsinghua University in 2016. He is currently an Assistant Professor at the Department of Statistic, the Chinese University of Hong Kong. Prior to that, he was a post-doctoral researcher at the Department of Statistics and Data Science, Wharton School, University of Pennsylvania. His research interests include diffusion based generative model, reinforcement learning, high-dimensional statistics, machine learning, signal processing, and mathematical optimization.

**Weichen Wu** received the B.S. degree in Data Science and Big Data Technology from Peking University, Beijing, China, in 2015, and is currently a Ph.D. candidate from the Department of Statistics and Data Science at Carnegie Mellon University, co-advised by Alessandro Rinaldo and Yuting Wei. His research interests lie in Reinforcement Learning and Topological Data Analysis.

**Yuejie Chi** (Fellow, IEEE; S'09–M'12–SM'17–F'23) received Ph.D. and M.A. in Electrical Engineering from Princeton University in 2012 and 2009, and B.E. (Hon.) in Electrical Engineering from Tsinghua University, Beijing, China, in 2007. She is currently the Sense of Wonder Group Endowed Professor of Electrical and Computer Engineering in AI Systems at Carnegie Mellon University, with courtesy appointments in the Machine Learning Department and CyLab. Her research interests lie in the theoretical and algorithmic foundations of data science, signal processing, machine learning and inverse problems, with applications in sensing, imaging, decision making, and AI systems. Among others, she is a recipient of Presidential Early Career Award for Scientists and Engineers (PECASE), SIAM Activity Group on Imaging Science Best Paper Prize, IEEE Signal Processing Society Young Author Best Paper Award, and the inaugural IEEE Signal Processing Society Early Career Technical Achievement Award for contributions to high-dimensional structured signal processing. She was named a Goldsmith Lecturer by IEEE Information Theory Society, a Distinguished Lecturer by IEEE Signal Processing Society, and a Distinguished Speaker by ACM. She currently serves or served as an Associate Editor for IEEE Trans. on Information Theory, IEEE Trans. on Signal Processing, IEEE Trans. on Pattern Recognition and Machine Intelligence, Information and Inference: A Journal of the IMA, and SIAM Journal on Mathematics of Data Science.

**Cong Ma** (Member, IEEE) received the B.Eng. degree from Tsinghua University in 2015 and the Ph.D. degree from Princeton University in 2020. He was a Post-Doctoral Researcher with the Department of Electrical Engineering and Computer Sciences, UC Berkeley. He is currently an Assistant Professor with the Department of Statistics, The University of Chicago. His research interests include mathematics of data science, machine learning, high-dimensional statistics, convex, and nonconvex optimization. He has received the Student Paper Award from the International Chinese Statistical Association in 2017, the School of Engineering and Applied Science Award for Excellence from Princeton University in 2019, the AI Labs Fellowship from Hudson River Trading in 2019, the Hannan Graduate Student Travel Award from IMS in 2020, as well as the Best Paper Prize from SIAM Activity Group on Imaging Science in 2024.

**Alessandro Rinaldo** is a Professor of Statistics and Data Science at the University of Texas at Austin. He received his Ph.D. in Statistics at Carnegie Mellon University (2005), and remained there as a faculty member until 2023. His research interests revolve around the theoretical properties of statistical and machine learning models for high-dimensional data under various structural assumptions, such as sparsity or intrinsic low dimensionality. An IMS fellow, he has served on the editorial boards of various journals, including the Annals of Statistics, the Journal of the American Statistics Association, and the Electronic Journal of Statistics.

**Yuting Wei** (Member, IEEE) received the B.S. degree (Hons.) in statistics from Peking University, Beijing, China, in 2013, and the Ph.D. degree in statistics from the University of California, Berkeley in 2018, advised by Martin Wainwright and Aditya Guntuboyina. She is currently an Assistant Professor at the Department of Statistics and Data Science, Wharton School, University of Pennsylvania. Prior to that, She was with Carnegie Mellon University from 2019 to 2021 as an Assistant Professor of Statistics and Data Science, and with Stanford University as a Stein Fellow from 2018 to 2019. Her research interests include high-dimensional and non-parametric statistics, statistical machine learning, and reinforcement learning. She was a recipient of the NSF CAREER award, the Google Research Award, the ICASA Junior Research Award, honorable mention for Bernoulli Society's New Researcher Award, and the Erich L. Lehmann Citation from Berkeley Statistics (awarded to the best dissertation in theoretical statistics).