

Sharp high-probability sample complexities for policy evaluation with linear function approximation

Gen Li^{*†}
UPenn Statistics

Weichen Wu^{*‡}
CMU Statistics

Yuejie Chi[§]
CMU ECE

Cong Ma[¶]
UChicago Statistics

Alessandro Rinaldo[‡]
CMU Statistics

Yuting Wei[†]
UPenn Statistics

May 31, 2023

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.

Keywords: policy evaluation, temporal difference learning, two-timescale stochastic approximation, minimax optimal, function approximation

Contents

1	Introduction	2
1.1	Our main contributions	3
1.2	Other related works	4
1.3	Notation	5
2	Problem formulation	6
2.1	Model and settings	6
2.2	Policy evaluation with linear approximation	7
3	On-policy evaluation with TD learning	8
3.1	The TD learning algorithm	8
3.2	Sample complexity of TD learning	9
3.3	Minimax lower bounds	10

*The first two authors contributed equally.

[†]Department of Statistics and Data Science, Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA.

[‡]Department of Statistics and Data Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

[§]Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

[¶]Department of Statistics, University of Chicago, Chicago, IL 60637, USA.

4	Off-policy evaluation with TDC learning	11
4.1	The TDC algorithm	11
4.2	Sample complexity of TDC	12
5	Numerical experiments	13
5.1	On-policy evaluation: averaged TD learning	13
5.2	Off-policy evaluation: TDC learning	14
6	Proof of Theorem 1 (TD learning)	15
7	Proof of Theorem 3 (TDC learning)	19
7.1	Population analysis	19
7.2	Finite-sample analysis	21
8	Discussion	23
A	Preliminary facts	24
B	Proof of Theorem 2 (minimax lower bounds)	27
C	Proofs of auxiliary lemmas and claims	29
C.1	Proof of Lemma 1	29
C.2	Proof of the inequalities (58a) and (58b)	31
C.3	Proof of Lemma 2	33
C.4	Proof of Lemma 3	34
C.5	Proof of Lemma 6 and Lemma 7	36
C.6	Proof of Lemma 8	38
D	Comparisons with previous works	39
D.1	Comparisons with Srikant and Ying (2019)	39
D.2	Comparisons with Bhandari et al. (2021)	42

1 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 (Bojinov and Shephard, 2019; Dann et al., 2014; Murphy, 2003; Tang and Wiens, 2021). For example, in clinical trials (Tang and Wiens, 2021), 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 (Bertsimas et al., 2022), 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 (Bertsekas, 2017; Sutton and Barto, 2018). 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 Tsitsiklis and Van Roy (1997), there has been an extensive line of works that consider different types of function approximation, including linear function approximation (Bhandari et al., 2021; Fan et al., 2020), reproducing kernel Hilbert space (Duan et al., 2021; Farahmand et al., 2016), deep neural networks (Arulkumaran et al., 2017; Bertsekas and Tsitsiklis, 1995) or function approximation on the model itself (see, e.g. Jin et al. (2020); Li et al. (2021a); Wang et al. (2021a)), with a focus on improving the sample efficiency of RL algorithms.

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 (Sutton, 1988) 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 Tsitsiklis and Van Roy (1997), the finite-sample minimax optimality of TD has been established only recently for the tabular MDP (Li et al., 2023a). For TD learning with linear function approximation, several recent contributions have produced new non-asymptotic analyses and insights (e.g. Bhandari et al. (2021); Dalal et al. (2018a); Lakshminarayanan and Szepesvari (2018); Srikant and Ying (2019)), 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 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 (Baird, 1995). In order to address this issue, Sutton et al. (2009) 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. Dalal et al. (2020, 2018b); Gupta et al. (2019); Kaledin et al. (2020); Wang et al. (2021b); Xu and Liang (2021)) 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.

1.1 Our main contributions

This paper is concerned with evaluating the performance of a given target policy π in an infinite-horizon γ -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} , θ^* is the best linear approximation coefficient of the value function, and Σ corresponds to the feature covariance matrix weighted by the stationary distribution. See Section 2 for the definitions of these parameters. Compared to prior work by Bhandari et al. (2021) and Srikant and Ying (2019), 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

paper	algorithm	stepsize	sample complexity	error control
Bhandari et al. (2021)	TD	$\eta_t \asymp t^{-1}$	$O\left(\frac{\ \Sigma^{-2}\ \ \Sigma\ \ \theta^*\ _{\Sigma}^2}{(1-\gamma)^2 \varepsilon^2}\right)$	in expectation
Srikant and Ying (2019)	TD	$\eta_t \asymp T^{-1}$	$O\left(\frac{\ \Sigma^{-2}\ \ \Sigma\ \ \theta^*\ _{\Sigma}^2}{(1-\gamma)^2 \varepsilon^2}\right)$	in expectation
Dalal et al. (2018a)	TD	$\eta_t = t^{-1}$	$O\left(\frac{1}{\varepsilon^{\max\{2, 1+\frac{1}{\lambda}\}}}\right)$	w. high-prob
This work	Averaged TD	$\eta_t = \eta$	$O\left(\frac{\ \Sigma^{-1}\ \ \theta^*\ _{\Sigma}^2}{(1-\gamma)^2 \varepsilon^2}\right)$	w. high-prob

Table 1. Comparisons with prior results (up to logarithmic terms) in finding an ε -optimal solution using TD learning. Using the Polyak-Ruppert averaging, our result improves upon previous sample complexity bounds by a multiplicative factor of the condition number of Σ , and is the first sample complexity with high-probability guarantee of convergence to match the minimax-optimal dependence on the tolerance level ε .

space). Our result is also the first to control ε -convergence with high probability that matches the minimax-optimal dependence on the tolerance level ε . To assess the tightness of this upper bound, we provide a minimax lower bound in Section 3.3, which certifies the optimal dependence of our bound on both the tolerance level ε and problem-related parameters Σ and θ^* .

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

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

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, ρ_{\max} denotes the largest importance sampling ratio measuring the discrepancy between the target policy and the behavior policy, and lastly, λ_1 and λ_2 denote the smallest eigenvalues of some problem-dependent matrices. Details about these constants are deferred to Section 4. To the best of our knowledge, our bound is the first one to control ε -convergence with high probability that matches the minimax-optimal dependence on the tolerance level ε . 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 1 and 2.

1.2 Other related works

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

Finite-sample guarantees for policy evaluation. Classical analyses of policy evaluation algorithms have mainly focused on providing asymptotic guarantees given a fixed model ([Szepesvári, 1998](#); [Tsitsiklis and Van Roy, 1997](#)). 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 [Bhandari et al. \(2021\)](#); [Boyan \(1999\)](#); [Dalal et al. \(2018a\)](#); [Jin et al. \(2018\)](#); [Khamaru et al. \(2020\)](#); [Lakshminarayanan and Szepesvari \(2018\)](#); [Srikant and Ying \(2019\)](#) with a focus on the non-asymptotic analyses for model-free algorithms, and [Agarwal et al. \(2020\)](#); [Li et al. \(2023b\)](#); [Pananjady and Wainwright \(2021\)](#); [Sidford et al. \(2018\)](#) which derive non-asymptotic bounds for model-based algorithms.

paper	algorithm	stepsize	sample complexity	error control
Dalal et al. (2020)	Projected TDC	$\alpha_t = t^{-\alpha}, \beta_t = t^{-\beta}$	$O\left(\frac{1}{\varepsilon^{2\alpha}}\right), \alpha < 1$	w. high-prob
Kaledin et al. (2020)	TDC	$\alpha_t, \beta_t \asymp \frac{1}{T}$	$O\left(\frac{1}{\varepsilon^2}\right)$	in expectation
Xu and Liang (2021)	Batched TDC	$\alpha_t = \alpha, \beta_t = \beta$	$O\left(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon}\right)$	in expectation
This work	TDC	$\alpha_t, \beta_t \asymp \frac{1}{T}$	$O\left(\frac{1}{\varepsilon^2}\right)$	w. high-prob

Table 2. Comparisons with prior results (up to logarithmic terms) in finding an ε -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.

Stochastic approximation. The idea of stochastic approximation (SA) (Lai, 2003; Robbins and Monro, 1951) 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 (Borkar, 2009; Borkar and Meyn, 2000) and with properly decaying step sizes, the Polyak-Ruppert averaged iterates asymptotically follow the central limit theorem. Recently, non-asymptotic results have also been obtained for SA for different problems especially in the RL setting; see Lakshminarayanan and Szepesvari (2018); Mou et al. (2020); Moulines and Bach (2011); Nemirovski et al. (2009) and references therein. The TDC algorithm is a special case of two-timescale linear SA, whose convergence rates have also been investigated in Dalal et al. (2020); Gupta et al. (2019); Wu et al. (2020); Xu et al. (2019), among others.

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 (Precup, 2000), and various different techniques have been applied to reduce the high variance of IS (see, e.g. Jiang and Li (2016); Kallus and Uehara (2020); Ma et al. (2022); Thomas and Brunskill (2016); Xie et al. (2019); Yang et al. (2020)). Non-asymptotic guarantees are also provided for off-policy evaluation using a fitted Q -iteration approach under linear function approximation in Duan et al. (2020). A recent line of works also considered finding the optimal policy using batch datasets (Jin et al., 2021; Li et al., 2022; Rashidinejad et al., 2021; Shi et al., 2022; Xie et al., 2021).

1.3 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 $\mathbb{1}\{\cdot\}$ 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 (ε, δ) ; 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 $f_t \asymp g_t$ or $f_t = \Theta(g_t)$. 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.

2 Problem formulation

2.1 Model and settings

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 (Bertsekas, 2017). 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 π , 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 π under the MDP \mathcal{M} when initialized at state $s_0 = s$. Here, $a_t \sim \pi(\cdot | s_t)$ and $s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)$ for all $t \geq 0$. It can be easily verified that $0 \leq V^\pi(s) \leq \frac{1}{1-\gamma}$ for any π .

For a given policy π , 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 π :

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

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}^π for this given policy π , such that its (i, j) element represents the probability that state i is transited to state j under the policy π ; 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 μ the stationary distribution corresponding to the Markov chain when the transition follows \mathbf{P}^π , which we assume to be well-defined, and introduce the vector notation $\boldsymbol{\mu} := [\mu(s)]_{1 \leq s \leq |\mathcal{S}|} \in \mathbb{R}^{|\mathcal{S}|}$.

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

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

where $\boldsymbol{\phi}(s) \in \mathbb{R}^d$ 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

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

be the feature matrix that concatenates the feature vectors for all states and $\mathbf{V}_\theta = [V_\theta(s)]_{s \in \mathcal{S}} \in \mathbb{R}^{|\mathcal{S}|}$ be the linear approximation vector to the value function. It follows that

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

We impose the following mild assumption on the feature vectors.

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

2.2 Policy evaluation with linear approximation

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 π of interest. In the **on-policy** setting, data samples are collected while the policy π is executed and a sequence of samples are obtained

$$\{(s_0, a_0, r_0), \dots, (s_T, a_0, r_T)\}, \quad \text{where } a_t \sim \pi(\cdot | s_t), \quad r_t = r(s_t, a_t).$$

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

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

and denote with

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

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 (Tsitsiklis and Van Roy, 1997)

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

Here, $\Pi_{\mathbf{D}_\mu}$ denotes the projection operator onto the column space of Φ (namely, the subspace $\{\Phi \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d\}$) w.r.t. the inner product $\langle \cdot, \cdot \rangle_{\mathbf{D}_\mu}$, where for any vector $\mathbf{v} \in \mathbb{R}^{|\mathcal{S}|}$ one has

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

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

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

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

$$\{(s_0, a_0, r_0), \dots, (s_T, a_0, r_T)\}, \quad \text{where } a_t \sim \pi_b(\cdot | s_t), \quad r_t = r(s_t, a_t).$$

Let μ_b be the stationary distribution over \mathcal{S} induced by the behavior π_b , and correspondingly let

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

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

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

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):

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

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 \mathcal{D}_{μ_b} , while the Bellman operator is again in terms of the target policy π . For this reason, solving (9) is different from solving the projected Bellman’s equation (7); as a result, in general, $\boldsymbol{\theta}^* \neq \tilde{\boldsymbol{\theta}}^*$.

3 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 ε -close solution. To better highlight our analysis strategy, we only consider the stylized generative model¹ whereby, at each time stamp t , one acquires an independent sample pair

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

Here recall that μ is the stationary distribution corresponding to \mathbf{P}^{π} . Notice that in the on-policy setting, since we are focused on a fixed policy π 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,

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

3.1 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

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

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

whose means are given respectively by

$$\mathbf{A} := \mathbb{E}_{s \sim \mu, s' \sim \mathbf{P}^{\pi}(\cdot | s)} \left[\boldsymbol{\phi}(s) (\boldsymbol{\phi}(s) - \gamma \boldsymbol{\phi}(s'))^{\top} \right] = \boldsymbol{\Phi}^{\top} \mathbf{D}_{\mu} (\mathbf{I} - \gamma \mathbf{P}^{\pi}) \boldsymbol{\Phi} \in \mathbb{R}^{d \times d}, \quad (12a)$$

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

It turns out that the target vector $\boldsymbol{\theta}^*$ satisfies the equation (Tsitsiklis and Van Roy, 1997)

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

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

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

¹We believe that our framework can be potentially generalized to Markovian samples using similar techniques in Li et al. (2021b) which is beyond the scope of the current paper.

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.

3.2 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 κ the condition number of $\boldsymbol{\Sigma}$ as follows

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

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

$$\begin{aligned} \|\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*\|_{\boldsymbol{\Sigma}} \leq C_0 \left\{ \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\phi}(s) \log(\frac{d}{\delta})}{T(1-\gamma)^2}} (\|\boldsymbol{\theta}^*\|_{\boldsymbol{\Sigma}} + 1) \right. \\ \left. + \frac{\|\boldsymbol{\Sigma}^{-1}\| \left[(\|\boldsymbol{\theta}^*\|_{\boldsymbol{\Sigma}} + 1) \sqrt{\frac{\kappa \log(\frac{dT}{\delta})}{\eta(1-\gamma)^3}} + \frac{1}{\eta(1-\gamma)} \|\boldsymbol{\theta}^*\|_{\boldsymbol{\Sigma}} \right]}{T} \right\} \end{aligned} \quad (16)$$

with probability at least $1 - \delta$, provided that $\boldsymbol{\theta}_0 = \mathbf{0}$, $\eta_0 = \dots = \eta_T = \eta < \frac{c_0(1-\gamma)}{\kappa \log(Td/\delta)}$ and

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

The proof of the theorem and the other results from this section can be found in Section 6. 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 $\varepsilon \in (0, \|\boldsymbol{\theta}^*\|_{\boldsymbol{\Sigma}})$ and $\delta \in (0, 1)$, the averaged TD estimator (14b) achieves*

$$\|\mathbf{V}_{\bar{\boldsymbol{\theta}}_T} - \mathbf{V}_{\boldsymbol{\theta}^*}\|_{D_\mu} = \|\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*\|_{\boldsymbol{\Sigma}} \leq \varepsilon \quad (17)$$

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

$$T \geq \frac{c \{ \max_s \boldsymbol{\phi}(s)^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\phi}(s) \} (1 + \|\boldsymbol{\theta}^*\|_{\boldsymbol{\Sigma}}^2) \log(\frac{d}{\delta})}{(1-\gamma)^2 \varepsilon^2}. \quad (18)$$

Comparisons to prior literature. We remark that the best finite-sample results for TD learning obtained so far are given by (Bhandari et al., 2021, Theorem 2(c)) and (Srikant and Ying, 2019, 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 ε , namely

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

it suffices to take (up to some logarithmic factors)

$$T^{\text{prior}} \asymp \frac{\kappa \|\boldsymbol{\Sigma}^{-1}\| (\|\boldsymbol{\theta}^*\|_{1+\boldsymbol{\Sigma}}^2)}{(1-\gamma)^2} \frac{1}{\varepsilon^2}. \quad (19)$$

We refer readers to Appendix D.1 and D.2 for a detailed translation of their results. Comparing (18) and (19), our result improves upon previous works by a multiplicative factor of

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

the condition number of Σ ; κ can be as large as d , the dimension of the features, which can scale with $|\mathcal{S}|$.

As for sample complexity with high-probability convergence guarantees, the best result so far is given by Dalal et al. (2018a), who shows that in order for (17) to hold with probability at least $1 - \delta$, it suffices to take

$$T \asymp \max \left\{ \left(\frac{1}{\varepsilon} \right)^2 \left(\log \frac{1}{\delta} \right)^3, \left(\frac{1}{\varepsilon} \right)^{1+1/\lambda_{\min}(\mathbf{A})} \left(\log \frac{1}{\delta} \right)^{1+1/\lambda_{\min}(\mathbf{A})} \right\}. \quad (20)$$

Comparing (18) and (20), we can see that our result improves on the dependence of both the error tolerance ε and the probability tolerance δ ; in fact, our result is the first sample complexity for TD learning with high-probability convergence guarantee that matches the minimax-optimal dependence of ε and displays a clear dependence on the problem-related parameters, as would be shown in the following section.

3.3 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 ε , do there exist estimators that achieve an ε -approximation of \mathbf{V}_{θ^*} 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 Φ such that*

$$\mathbb{P} \left\{ \|\hat{\theta} - \theta^*\|_{\Sigma} > c_2 \varepsilon \right\} \geq \frac{1}{4}, \quad (21)$$

provided that the number of samples T satisfies

$$T \leq \frac{c_3 \{ \max_s \phi(s)^\top \Sigma^{-1} \phi(s) \} (1 + \|\theta^*\|_{\Sigma}^2)}{(1 - \gamma) \varepsilon^2}. \quad (22)$$

Remark 1. We remark that minimax lower bounds are also previously investigated in a general framework in (Duan et al., 2021) 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 Duan et al. (2021), 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

$$O \left(\frac{d}{\varepsilon^2 (1 - \gamma)^3} \right), \quad (23)$$

as one plugs in the exact parameters from our construction.

As asserted by this theorem, no algorithm whatsoever can attain an ε -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 ε , the dependence of the feature map Φ , the underlying coefficient θ^* , and the covariance matrix Σ . 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.

4 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 π_b . At each time stamp t , one obtains

$$(s_t, a_t, s'_t) \quad \text{where } s_t \stackrel{\text{i.i.d.}}{\sim} \mu_b, \quad a_t \sim \pi_b(\cdot | s_t), \quad \text{and } s'_t \sim \mathcal{P}(\cdot | s_t, a_t). \quad (24)$$

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

4.1 The TDC algorithm

The TDC algorithm is designed to solve the optimization problem (9) using a two-timescale linear TD with gradient correction (Sutton et al., 2009). 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*

$$\text{MSPBE}(\theta) = \frac{1}{2} \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t) \delta_t]^\top \{ \mathbb{E}_{\mu_b} [\phi(s_t) \phi(s_t)^\top] \}^\top \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t) \delta_t], \quad (25)$$

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 $\text{MSPBE}(\theta)$ with respect to θ equals to

$$\begin{aligned} \nabla_{\theta} \text{MSPBE}(\theta) &= \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [(\gamma \phi(s'_t) - \phi_t) \phi(s_t)^\top] \{ \mathbb{E}_{\mu_b} [\phi(s_t) \phi(s_t)^\top] \}^{-1} \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t) \delta_t] \\ &= -\mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t) \delta_t] + \gamma \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s'_t) \phi(s_t)^\top] \{ \mathbb{E}_{\mu_b} [\phi(s_t) \phi(s_t)^\top] \}^{-1} \mathbb{E}_{\mu_b, \pi, \mathcal{P}} [\phi(s_t) \delta_t] \\ &= -\mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t) \delta_t] + \gamma \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s'_t) \phi(s_t)^\top] \mathbf{w}_t, \end{aligned} \quad (26)$$

where in the last step we have defined

$$\mathbf{w}_t = \mathbf{w}(\theta_t) = \{ \mathbb{E}_{\mu_b} [\phi(s_t) \phi(s_t)^\top] \}^{-1} \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t) \delta_t]. \quad (27)$$

and have used the importance weights

$$\rho_t := \frac{\pi(a_t | s_t)}{\pi_b(a_t | s_t)} \quad (28)$$

to replace the expectation w.r.t. π with the expectation w.r.t. π_b .

The high-level idea of TDC is to estimate the right hand side of (26) based on the sample trajectory (24), and then perform stochastic gradient updates for θ_t . However, the challenge is that the second term in the gradient of MSPBE (26) 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, Sutton et al. (2008) and Sutton et al. (2009) 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

$$\mathbb{E}_{\mu_b} [\phi(s_t) \phi(s_t)^\top] \mathbf{w} = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}} [\rho_t \phi(s_t) \delta_t]. \quad (29)$$

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 $\tilde{\boldsymbol{\theta}}_t$ corresponds to a gradient step regarding (25), the update of \mathbf{w}_t corresponds to linear SA for solving (29), and $\delta_t := r_t + \gamma \boldsymbol{\phi}(s'_t)^\top \tilde{\boldsymbol{\theta}}_t - \boldsymbol{\phi}(s_t)^\top \tilde{\boldsymbol{\theta}}_t$ is the temporal difference error. In addition, α_t, β_t are the corresponding stepsizes. For notational convenience, let us denote

$$\begin{aligned}\tilde{\mathbf{A}}_t &= \rho_t \boldsymbol{\phi}(s_t) (\boldsymbol{\phi}(s_t) - \gamma \boldsymbol{\phi}(s'_t))^\top, & \tilde{\mathbf{b}}_t &:= \rho_t \boldsymbol{\phi}(s_t) r_t, \\ \boldsymbol{\Pi}_t &:= \rho_t \boldsymbol{\phi}(s_t) \boldsymbol{\phi}(s'_t)^\top, & \tilde{\boldsymbol{\Sigma}}_t &:= \boldsymbol{\phi}(s_t) \boldsymbol{\phi}(s_t)^\top.\end{aligned}\quad (30)$$

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

$$\tilde{\boldsymbol{\theta}}_{t+1} = \tilde{\boldsymbol{\theta}}_t - \alpha_t (\tilde{\mathbf{A}}_t \tilde{\boldsymbol{\theta}}_t - \tilde{\mathbf{b}}_t + \gamma \boldsymbol{\Pi}_t^\top \mathbf{w}_t); \quad (31a)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \beta_t (\tilde{\mathbf{A}}_t \tilde{\boldsymbol{\theta}}_t - \tilde{\mathbf{b}}_t + \tilde{\boldsymbol{\Sigma}}_t \mathbf{w}_t). \quad (31b)$$

4.2 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

$$\tilde{\mathbf{A}} := \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\tilde{\mathbf{A}}_t] = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\rho_t \boldsymbol{\phi}(s_t) (\boldsymbol{\phi}(s_t) - \gamma \boldsymbol{\phi}(s'_t))^\top]; \quad (32a)$$

$$\tilde{\mathbf{b}} := \mathbb{E}_{\mu_b}[\tilde{\mathbf{b}}_t] = \mathbb{E}_{\mu_b, \pi_b}[\rho_t \boldsymbol{\phi}(s_t) r_t]; \quad (32b)$$

$$\boldsymbol{\Pi} := \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\boldsymbol{\Pi}_t] = \mathbb{E}_{\mu_b, \pi_b, \mathcal{P}}[\rho_t \boldsymbol{\phi}(s_t) \boldsymbol{\phi}(s'_t)^\top]; \quad (32c)$$

$$\tilde{\boldsymbol{\Sigma}} := \mathbb{E}_{\mu_b}[\tilde{\boldsymbol{\Sigma}}_t] = \mathbb{E}_{\mu_b}[\boldsymbol{\phi}(s_t) \boldsymbol{\phi}(s_t)^\top]. \quad (32d)$$

In addition, denote the parameters

$$\begin{aligned}\lambda_1 &= \lambda_{\min}(\tilde{\mathbf{A}}^\top \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}}), & \lambda_2 &= \lambda_{\min}(\tilde{\boldsymbol{\Sigma}}), & \lambda_\Sigma &= \|\tilde{\boldsymbol{\Sigma}}^{-1}\| = 1/\lambda_2, \\ \tilde{\kappa} &= \lambda_\Sigma \cdot \|\tilde{\boldsymbol{\Sigma}}\|, & \rho_{\max} &= \max_{s,a}[\pi(a|s)/\pi_b(a|s)].\end{aligned}\quad (33)$$

With these notation in place, we are ready to state our main result for TDC learning, with its proof deferred to Section 7.

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 $\tilde{\boldsymbol{\theta}}_T$ of the TDC learning iterate (31) at time T satisfies the bound*

$$\|\tilde{\boldsymbol{\theta}}_T - \tilde{\boldsymbol{\theta}}^*\|_{\tilde{\boldsymbol{\Sigma}}} \leq \tilde{C}_0 \frac{\rho_{\max}^2 \|\tilde{\boldsymbol{\Sigma}}\|}{\lambda_1} \sqrt{\frac{\beta}{\lambda_2} \log \frac{2dT}{\delta}} (\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\boldsymbol{\Sigma}}} + 2), \quad (34)$$

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

$$\begin{aligned}\tilde{\boldsymbol{\theta}}_0 &= \mathbf{0}, \\ \alpha_0 &= \dots = \alpha_T = \alpha, & \beta_0 &= \dots = \beta_T = \beta, \\ 0 < \alpha &< \frac{1}{\lambda_1 \lambda_\Sigma^2 \|\tilde{\boldsymbol{\Sigma}}\| \log \frac{2dT}{\delta}}, & \frac{\alpha}{\beta} &= \frac{1}{128} \frac{\lambda_1 \lambda_2}{\rho_{\max}^2 (1 + \lambda_\Sigma \rho_{\max})}, \\ T &\geq \tilde{c}_1 \frac{\log \|\tilde{\boldsymbol{\theta}}^*\|_2}{\alpha \lambda_1} \log \max \left\{ \sqrt{\tilde{\kappa}}, \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\boldsymbol{\Sigma}}} \sqrt{\frac{\alpha \lambda_1}{\log \frac{2dT}{\delta}}} \right\}.\end{aligned}\quad (35)$$

Remark 2. A similar result in terms of the ℓ_2 error (namely, $\|\tilde{\boldsymbol{\theta}}_T - \tilde{\boldsymbol{\theta}}^*\|_2$) can be derived in the same way as in (34). In particular, under the same conditions as in (35), it can be derived with probability at least $1 - \delta$ that

$$\|\tilde{\boldsymbol{\theta}}_T - \tilde{\boldsymbol{\theta}}^*\|_2 \lesssim \tilde{C}_0 \frac{\rho_{\max}^2}{\lambda_1} \sqrt{\frac{\beta}{\lambda_2} \log \frac{2dT}{\delta}} (\|\tilde{\boldsymbol{\theta}}^*\|_2 + 2). \quad (36)$$

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

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 $\varepsilon \in (0, \|\tilde{\theta}^*\|_{\tilde{\Sigma}})$, the TDC estimator $\tilde{\theta}_T$ at iterate T satisfies the bound*

$$\|V_{\tilde{\theta}_T} - V_{\tilde{\theta}^*}\|_{D_{\mu_b}} = \|\tilde{\theta}_T - \tilde{\theta}^*\|_{\tilde{\Sigma}} \leq \varepsilon \quad (37)$$

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

$$T \geq \tilde{c} \frac{\rho_{\max}^7}{\lambda_1^4 \lambda_2^3} \frac{\|\tilde{\Sigma}\|^2}{\varepsilon^2} (1 + \|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2) \log\left(\frac{d\|\tilde{\theta}^*\|_{\tilde{\Sigma}}}{\delta}\right), \quad (38)$$

and the stepsize parameters α_t and β_t are chosen as

$$\alpha_t \asymp \frac{\log\|\tilde{\theta}^*\|_{\tilde{\Sigma}}}{T\lambda_1}, \quad \beta_t = 128 \frac{\rho_{\max}^2(1 + \lambda_{\Sigma}\rho_{\max})}{\lambda_1\lambda_2} \alpha. \quad (39)$$

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 Dalal et al. (2020), 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 ε , 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 δ 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 Kaledin et al. (2020), who shows that with the choice of $\alpha_t, \beta_t \asymp \frac{1}{T}$, the sample complexity for TDC with tolerance level ε 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. We also remark that Xu and Liang (2021) considers a variant of TDC where θ_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 $O\left(\frac{1}{\varepsilon^2} \log\left(\frac{1}{\varepsilon}\right)\right)$.

5 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.

5.1 On-policy evaluation: averaged TD learning

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

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 $|\mathcal{S}| = 10$, $\gamma = 0.2$, 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 $\|\theta_t - \theta^*\|_{\Sigma}$ for TD and $\|\tilde{\theta}_t - \theta^*\|_{\Sigma}$ for averaged TD.

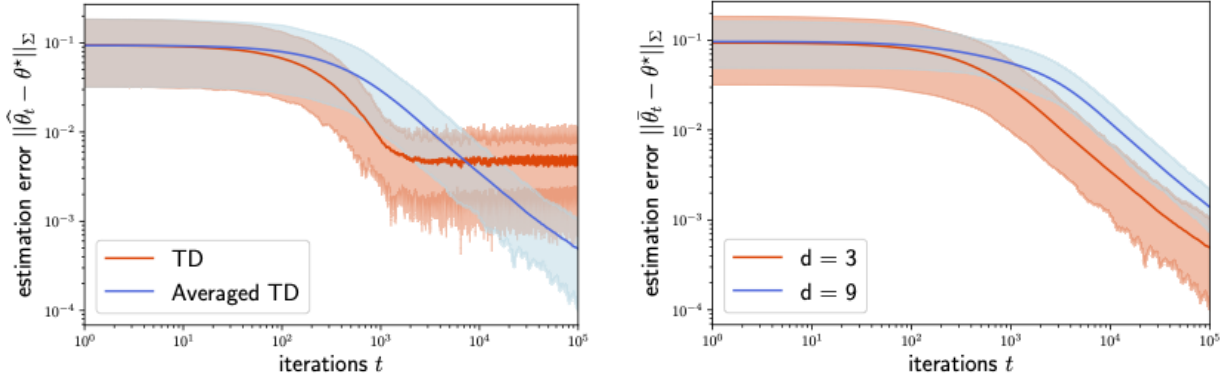


Figure 1. (a) Comparisons of the estimation error of TD and averaged TD when $d = 3$. (b) Comparisons of the estimation error for averaged TD with $d = 3$ and $d = 9$. Two curves in the middle represent their average errors, while the shaded areas represent the 95% confidence bands.

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 5×10^{-3} after 10^3 iterations, the error of averaged TD keeps decreasing to below 5×10^{-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.

5.2 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* (Baird, 1995).

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 π selects $a = 1$ at any given state s , while the behavior policy π_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)

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

In this example, it is easy to check that the stationary distribution corresponding to the behavior policy π_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 $\theta \in \mathbb{R}^8$,

$$\begin{cases} V(i) = 2\theta_i + \theta_8, & \text{for } 1 \leq i \leq 6; \\ V(7) = \theta_7 + 2\theta_8. \end{cases} \quad (40)$$

We remark that with this approximation, the feature space has a higher dimension ($d = 8$) than the state space ($|\mathcal{S}| = 7$). Consequently, the optimal estimator $\tilde{\theta}^*$ is not unique, and instead can be any $\theta \in \mathbb{R}^8$ such

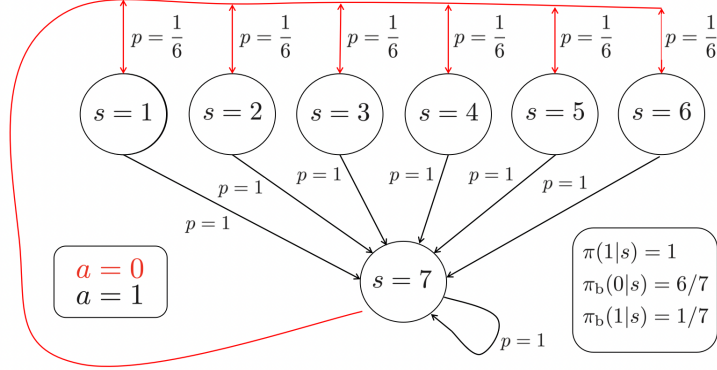


Figure 2. Baird’s counterexample. Taking action $a = 1$ always leads to state $s = 7$, while taking $a = 0$ leads to one of the other six states with equal probability. The reward is set to be always zero.

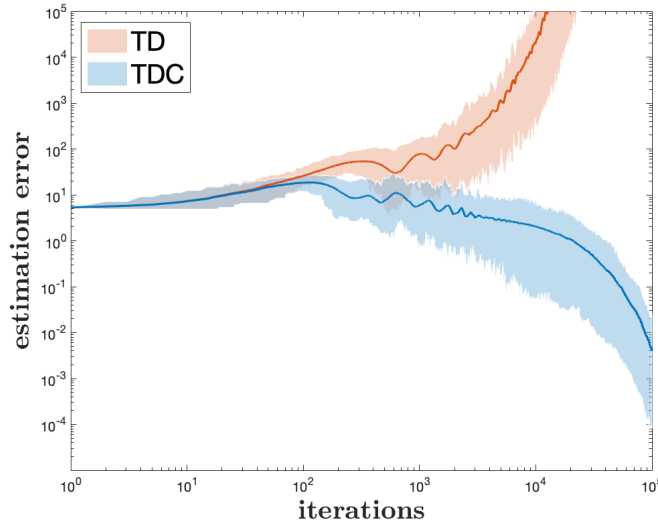


Figure 3. Performances of off-policy averaged TD (red, $\eta = 0.02$) and TDC (blue, $\alpha = 0.02$, $\beta = 0.002$). Two curves in the middle represent their average errors, while the shaded areas correspond to 95% confidence bands.

that the estimated value vector is $\mathbf{V}_{\hat{\theta}} = \mathbf{0}$. 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 $\|\hat{\theta}_t - \tilde{\theta}^*\|_{\tilde{\Sigma}} = \|\mathbf{V}_{\hat{\theta}_t} - \mathbf{V}^*\|_{D_{\mu_b}}$ to evaluate the estimation error, and our experimental results would remain the same.

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 Baird (1995). In these experiments, we set $\alpha = \eta$ to ensure that the stepsize for θ -updates are the same between the two algorithms. Figure 3 demonstrates how the estimation error $\|\hat{\theta}_t - \tilde{\theta}^*\|_{\tilde{\Sigma}}$ 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.

6 Proof of Theorem 1 (TD learning)

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

$$\Delta_t := \theta_t - \theta^*, \quad \text{and} \quad \bar{\Delta}_t := \bar{\theta}_t - \theta^*. \quad (41)$$

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

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

where we define

$$\xi_t := (\mathbf{A}_t - \mathbf{A})\theta_t - (\mathbf{b}_t - \mathbf{b}). \quad (42)$$

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

$$\Delta_t = (\mathbf{I} - \eta \mathbf{A})\Delta_{t-1} - \eta \xi_{t-1} = (\mathbf{I} - \eta \mathbf{A})^t \Delta_0 - \eta \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i. \quad (43)$$

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

$$\|\Delta_t\|_{\Sigma} \leq 32 \sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1 - \gamma}} (1 + \|\theta^*\|_{\Sigma}) + 2\sqrt{\kappa} \|\Delta_0\|_{\Sigma} =: R_0 \quad (44)$$

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 Bhandari et al. (2021).

To start with, note that the inequality (44) holds trivially for the base case with $t = 0$, given that $\kappa \geq 1$. Next, suppose that the hypothesis (44) holds for $\Delta_0, \dots, \Delta_{t-1}$, and we intend to establish it for Δ_t as well. Towards this end, invoking the decomposition (43) and the triangle inequality yields

$$\|\Delta_t\|_{\Sigma} \leq \|(\mathbf{I} - \eta \mathbf{A})^t \Delta_0\|_{\Sigma} + \eta \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma}. \quad (45)$$

As for the first term of (45), it is seen that

$$\begin{aligned}
\|(\mathbf{I} - \eta \mathbf{A})^t \Delta_0\|_{\Sigma} &= \|\Sigma^{1/2} (\mathbf{I} - \eta \mathbf{A})^t \Sigma^{-1/2} \Sigma^{1/2} \Delta_0\|_2 \leq \|\Sigma^{1/2}\| \cdot \|\Sigma^{-1/2}\| \cdot \|\mathbf{I} - \eta \mathbf{A}\|^t \cdot \|\Sigma^{1/2} \Delta_0\|_2 \\
&\leq \sqrt{\kappa} \left(1 - \frac{1}{2} \eta (1 - \gamma) \lambda_{\min}(\Sigma)\right)^t \|\Delta_0\|_{\Sigma} \leq \sqrt{\kappa} \|\Delta_0\|_{\Sigma},
\end{aligned} \quad (46)$$

where the last inequality arises from the definition of κ and the property (79g) (with the restriction that $\eta \leq (1 - \gamma)/(4\|\Sigma\|)$). When it comes to the second term of (45), 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

$$\tilde{\xi}_i := \xi_i \mathbb{1}\{\mathcal{H}_i\}, \quad \text{where } \mathcal{H}_i := \left\{ \|\Delta_i\|_{\Sigma} \leq R \right\}. \quad (47)$$

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

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

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

Proof. See Section C.1. □

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

$$\begin{aligned} \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} &= \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \mathbf{1}_{\{\|\Delta_i\|_{\Sigma} \leq R_0\}} \right\|_{\Sigma} \\ &\leq 16(\|\theta^*\|_{\Sigma} + R_0 + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \end{aligned} \quad (48)$$

holds with probability at least $1 - \delta/T$, provided that $0 < \eta \leq \frac{1-\gamma}{\kappa \log \frac{2dT}{\delta}}$. Combining (45), (46) and (48) together and recalling the definition (44) of R_0 , we can easily verify that

$$\|\Delta_t\|_{\Sigma} \leq \sqrt{\kappa} \|\Delta_0\|_{\Sigma} + 16(\|\theta^*\|_{\Sigma} + R_0 + 1) \sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}} \leq R_0, \quad (49)$$

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 (44).

Step 3: a refined bound on $\|\Delta_t\|_{\Sigma}$. It turns out that the upper bound (44) is somewhat loose due to the complete ignorance of the contraction effect of $\mathbf{I} - \eta \mathbf{A}$; see (46). In what follows, we develop a strengthened bound. Define

$$t_{\text{seg}} := \frac{c_1 \log \max\{4\sqrt{\kappa}, \frac{16\kappa\|\Delta_0\|_{\Sigma}}{\|\theta^*\|_{\Sigma}+1}, \|\Delta_0\|_{\Sigma} \sqrt{\frac{1-\gamma}{\eta \kappa \log \frac{2dT}{\delta}}}\}}{\eta(1-\gamma)\lambda_{\min}(\Sigma)} \quad (50)$$

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

$$\|\Delta_t\|_{\Sigma} \leq 32 \sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}} \left(\|\theta^*\|_{\Sigma} + \frac{\sqrt{\kappa} \|\Delta_0\|_{\Sigma}}{2^{k-1}} + \frac{3}{2} \right) =: R_k \quad (51)$$

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 (45), we claim that it suffices to prove that

$$\left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \xi_i \right\|_{\Sigma} \leq 32 \left(\|\theta^*\|_{\Sigma} + \frac{2\sqrt{\kappa} \|\Delta_0\|_{\Sigma}}{2^k} + 1 \right) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}}, \quad kt_{\text{seg}} \leq t \leq T. \quad (52)$$

To see this: note that the first term on the right-hand side of (45) has already been bounded in (46), which combined with the definition (50) of t_{seg} indicates that

$$\|\Sigma^{1/2} (\mathbf{I} - \eta \mathbf{A})^t \Delta_0\|_2 \leq \sqrt{\kappa} \left(1 - \frac{1}{2} \eta(1-\gamma) \lambda_{\min}(\Sigma) \right)^{t_{\text{seg}}} \|\Delta_0\|_{\Sigma} \leq \sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}} \quad (53)$$

for any $t \geq t_{\text{seg}}$. Clearly, combining (52) with (45) and (53) shall immediately lead to the claim (51). The remainder of this step is thus devoted to demonstrating (52) inductively.

The base case (i.e. $k = 1$) follows immediately from our bounds (44) and (48) in Step 2, given that $\sqrt{\frac{\eta \kappa \log \frac{2dT}{\delta}}{1-\gamma}}$ is sufficiently small. Suppose now that the claim (52) 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 (52) 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

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

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 (51) 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 $t \geq t_{\text{seg}}$,

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

Here, the first equality uses the crude bound $\|\Delta_i\|_{\Sigma} \leq R_0$ for all i (see (44)), 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 (44) of R_0 and invokes the fact that $\sqrt{\kappa} \left(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma)\right)^{t_{\text{seg}}-1} \leq \min\left\{\frac{1}{4}, \frac{1}{4\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}\right\}$ with our choice (50) of t_{seg} .

Combine the previous two bounds to reach

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

This finishes the induction step and in turn establishes (52) (and hence (51)).

As a straightforward consequence, the bounds (44) and (51) imply that

$$\|\Delta_t\|_{\Sigma} \leq \begin{cases} R_0, & 0 \leq t < t'_{\text{seg}}, \\ 32\sqrt{\frac{\eta\kappa \log \frac{2dT}{\delta}}{1-\gamma}} (\|\theta^*\|_{\Sigma} + 2), & t'_{\text{seg}} \leq t < T, \end{cases} \quad (54)$$

where

$$t'_{\text{seg}} := c_2 t_{\text{seg}} \log(\kappa(\|\Delta_0\|_2 + 1)) \quad (55)$$

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 (51) obeys $\frac{4\sqrt{\kappa}\|\Delta_0\|_{\Sigma}}{2^{\lfloor t/t_{\text{seg}} \rfloor}} \leq 2$, thus confirming the second case in (54).

Step 4: controlling $\|\overline{\Delta}_T\|_{\Sigma}$. Now we are positioned to control $\overline{\Delta}_T$. The key is to write $\overline{\Delta}_T$ as a linear combination of $\{\xi_i\}_{0 \leq i \leq T-1}$ as follows, which is a direct consequence of the relation (43):

$$\overline{\Delta}_T = \frac{1}{T} \sum_{j=1}^T \Delta_j = \frac{1}{T} \sum_{j=1}^T (\mathbf{I} - \eta\mathbf{A})^j \Delta_0 - \frac{1}{T} \sum_{j=1}^T \eta \sum_{i=0}^{j-1} (\mathbf{I} - \eta\mathbf{A})^{j-i-1} \xi_i$$

$$\begin{aligned}
&= \frac{1}{T} \sum_{j=1}^T (\mathbf{I} - \eta \mathbf{A})^j \Delta_0 - \frac{1}{T} \sum_{i=0}^{T-1} \eta \sum_{j=i+1}^T (\mathbf{I} - \eta \mathbf{A})^{j-i-1} \xi_i \\
&= \frac{1}{T\eta} \mathbf{A}_0^{(T+1)} \Delta_0 - \frac{1}{T} \Delta_0 - \frac{1}{T} \sum_{i=0}^{T-1} \mathbf{A}_i^{(T)} \xi_i,
\end{aligned} \tag{56}$$

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

$$\mathbf{A}_i^{(t)} := \eta \sum_{j=i+1}^t (\mathbf{I} - \eta \mathbf{A})^{j-i-1} = \mathbf{A}^{-1} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{A})^{t-i}). \tag{57}$$

We claim that the following two inequalities hold, the first deterministically and the second with probability of at least $1 - \delta$ (with their proofs deferred to Section C.2)

$$\|\mathbf{A}_0^{(T+1)} \Delta_0\|_{\Sigma} \leq \frac{2\|\Sigma^{-1}\|}{1-\gamma} \|\Delta_0\|_{\Sigma}; \tag{58a}$$

$$\left\| \sum_{i=0}^{T-1} \mathbf{A}_i^{(T)} \xi_i \right\|_{\Sigma} \lesssim \left\{ \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log \frac{2d}{\delta}}{T(1-\gamma)^2}} + \frac{\|\Sigma^{-1}\|}{T} \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)^3}} \right\} (\|\theta^*\|_{\Sigma} + 1). \tag{58b}$$

Putting the above two inequalities together with (56), we arrive at

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

where the last line follows since $\|\Sigma^{-1}\| \geq 1$ (see (79h)) and $\eta < 1$. This finishes the proof of Theorem 1.

7 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.

7.1 Population analysis

First recall that the population parameters are defined as

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

Corresponding to the empirical version of TDC as given in (31), we can define its population analogue of TDC as

$$\check{\theta}_{t+1} = \check{\theta}_t - \alpha (\tilde{\mathbf{A}} \check{\theta}_t - \tilde{\mathbf{b}} + \gamma \tilde{\mathbf{\Pi}}^\top \check{w}_t),$$

$$\check{\mathbf{w}}_{t+1} = \check{\mathbf{w}}_t - \beta(\tilde{\mathbf{A}}\check{\boldsymbol{\theta}}_t - \tilde{\mathbf{b}} + \tilde{\boldsymbol{\Sigma}}\check{\mathbf{w}}_t), \quad (59)$$

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 $\tilde{\boldsymbol{\Sigma}}_t$ is independent of which policy is being adopted. Hence, $\tilde{\boldsymbol{\Sigma}}$ can also be presented as

$$\begin{aligned} \tilde{\boldsymbol{\Sigma}} &= \sum_{s_t \in \mathcal{S}} \mu_b(s_t) \phi(s_t) \phi(s_t)^\top \\ &= \sum_{s_t \in \mathcal{S}} \mu_b(s_t) \left(\sum_{a_t \in \mathcal{A}} \pi(a_t | s_t) \right) \phi(s_t) \phi(s_t)^\top \\ &= \sum_{s_t \in \mathcal{S}} \sum_{a_t \in \mathcal{A}} \mu_b(s_t) \pi_b(a_t | s_t) \left(\frac{\pi(a_t | s_t)}{\pi_b(a_t | s_t)} \right) \phi(s_t) \phi(s_t)^\top = \mathbb{E}_{\mu_b, \pi_b} [\rho_t \phi(s_t) \phi(s_t)^\top]. \end{aligned} \quad (60)$$

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

$$\tilde{\mathbf{A}} = \tilde{\boldsymbol{\Sigma}} - \gamma \boldsymbol{\Pi}. \quad (61)$$

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

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

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

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

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

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

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

$$\begin{bmatrix} \check{\boldsymbol{\Delta}}_t \\ \check{\mathbf{z}}_t \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \alpha \tilde{\mathbf{A}}^\top \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} & -\alpha \gamma \boldsymbol{\Pi}^\top \\ -\alpha (\mathbf{I} - \gamma \boldsymbol{\Sigma}^{-1} \boldsymbol{\Pi}) \tilde{\mathbf{A}}^\top \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} & \mathbf{I} - \beta \tilde{\boldsymbol{\Sigma}} - \alpha \gamma (\mathbf{I} - \gamma \tilde{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Pi}) \boldsymbol{\Pi}^\top \end{bmatrix} \begin{bmatrix} \check{\boldsymbol{\Delta}}_{t-1} \\ \check{\mathbf{z}}_{t-1} \end{bmatrix}. \quad (62)$$

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

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

then $\check{\mathbf{x}}_t$ yields

$$\check{\mathbf{x}}_t = \underbrace{\begin{bmatrix} \mathbf{I} - \alpha \tilde{\mathbf{A}}^\top \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} & -\frac{1}{\varkappa} \alpha \gamma \boldsymbol{\Pi}^\top \\ -\varkappa \alpha (\mathbf{I} - \gamma \boldsymbol{\Sigma}^{-1} \boldsymbol{\Pi}) \tilde{\mathbf{A}}^\top \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} & \mathbf{I} - \beta \tilde{\boldsymbol{\Sigma}} - \alpha \gamma (\mathbf{I} - \gamma \tilde{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Pi}) \boldsymbol{\Pi}^\top \end{bmatrix}}_{=:\boldsymbol{\Psi}} \check{\mathbf{x}}_{t-1}. \quad (63)$$

It is known that how fast $\check{\mathbf{x}}_t$ converges to $\mathbf{0}$ is determined by the spectral norm of $\boldsymbol{\Psi}$, 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:

$$\beta \gtrsim \lambda_\Sigma \rho_{\max} \alpha, \tag{64a}$$

$$\varkappa \beta \gtrsim \alpha, \tag{64b}$$

$$\alpha \gamma (\rho_{\max} + \gamma \lambda_\Sigma \rho_{\max}^2) \ll \beta \lambda_w, \tag{64c}$$

$$\frac{\alpha \gamma \rho_{\max}}{\varkappa} + \varkappa \alpha (1 + \gamma \lambda_\Sigma \rho_{\max}) \lambda_\Sigma (2\rho_{\max})^2 \ll \sqrt{\alpha \lambda_1 \beta \lambda_w} \tag{64d}$$

it holds true that

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

Therefore, the mapping Ψ is contractive, thus ensuring the linear convergence of \mathbf{x}_t , with the proviso that $\alpha \lambda_1 < 2$.

7.2 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.

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

$$\begin{aligned} \boldsymbol{\nu}_t &:= (\tilde{\mathbf{A}}_t - \tilde{\mathbf{A}}) \tilde{\boldsymbol{\theta}}_t - (\tilde{\mathbf{b}}_t - \tilde{\mathbf{b}}) + \gamma (\boldsymbol{\Pi}_t - \boldsymbol{\Pi})^\top \mathbf{w}_t, \\ \boldsymbol{\eta}_t &:= (\tilde{\mathbf{A}}_t - \tilde{\mathbf{A}}) \tilde{\boldsymbol{\theta}}_t - (\tilde{\mathbf{b}}_t - \tilde{\mathbf{b}}) + (\tilde{\Sigma}_t - \tilde{\Sigma}) \mathbf{w}_t. \end{aligned}$$

As a result, TDC can be rewritten as

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_{t+1} &= \tilde{\boldsymbol{\theta}}_t - \alpha (\tilde{\mathbf{A}} \tilde{\boldsymbol{\theta}}_t - \tilde{\mathbf{b}} + \gamma \boldsymbol{\Pi}^\top \mathbf{w}_t) - \alpha \boldsymbol{\nu}_t; \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \beta (\tilde{\mathbf{A}} \tilde{\boldsymbol{\theta}}_t - \tilde{\mathbf{b}} + \tilde{\Sigma} \mathbf{w}_t) - \beta \boldsymbol{\eta}_t. \end{aligned}$$

Using the same notations as in Section 7.1, we observe that the following iteration holds true for finite-sample TDC:

$$\mathbf{x}_{t+1} = \Psi \mathbf{x}_t - \zeta_t,$$

in which

$$\zeta_t = \begin{bmatrix} \alpha \boldsymbol{\nu}_t \\ \varkappa (\alpha (1 - \gamma \tilde{\Sigma}^{-1} \boldsymbol{\Pi}) \boldsymbol{\nu}_t + \beta \boldsymbol{\eta}_t) \end{bmatrix}. \tag{65}$$

Hence,

$$\mathbf{x}_t = \Psi^t \mathbf{x}_0 - \sum_{i=0}^{t-1} \Psi^{t-i-1} \zeta_i, \tag{66}$$

where $\mathbf{x}_0 = [\boldsymbol{\Delta}_0^\top, \varkappa \mathbf{z}_0^\top]^\top$. Since the norm of Ψ has been bounded by Lemma 2, bounding the norm of \mathbf{x}_t boils down to bounding the second term of (66). 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

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

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

$$\begin{cases} \|\tilde{\Delta}_t\|_{\tilde{\Sigma}} \leq \|\mathbf{x}_t\|_{\tilde{\Sigma}}; \\ \|\mathbf{w}_t\|_{\tilde{\Sigma}} \lesssim \frac{1}{\varkappa} \|\mathbf{x}_t\|_{\tilde{\Sigma}}; \\ \|\mathbf{x}_t\|_{\tilde{\Sigma}} \lesssim \|\tilde{\Delta}_t\|_{\tilde{\Sigma}} + \|\mathbf{w}_t\|_{\tilde{\Sigma}}. \end{cases} \quad (67)$$

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

$$\|\mathbf{x}_t\|_{\tilde{\Sigma}} \leq 2\|\tilde{\Delta}_0\|_{\tilde{\Sigma}} + 80\varkappa\beta\rho_{\max} \sqrt{\frac{1}{\alpha\lambda_1} \log \frac{2dT}{\delta}} (\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + 1) =: \tilde{R}_0 \quad (68)$$

holds simulatanesouly for all $0 \leq t \leq T$. To start with, note that the inequality (68) holds trivially for the base case with $t = 0$. Next, suppose that the hypothesis (68) 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 (66) and the triangle inequality yields

$$\|\mathbf{x}_t\|_{\tilde{\Sigma}} = \|\Psi^t \mathbf{x}_0\|_{\tilde{\Sigma}} + \left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \zeta_i \right\|_{\tilde{\Sigma}}. \quad (69)$$

As for the first term of (69), it is seen that

$$\|\Psi^t \mathbf{x}_0\|_{\tilde{\Sigma}} \leq \|\mathbf{x}_0\|_{\tilde{\Sigma}} = \|\tilde{\Delta}_0\|_{\tilde{\Sigma}}. \quad (70)$$

When it comes to the second term of (69), 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

$$\tilde{\zeta}_i := \zeta_i \mathbb{1}\{\tilde{\mathcal{H}}_i\}, \quad \text{where } \tilde{\mathcal{H}}_i := \left\{ \|\mathbf{x}_i\|_{\tilde{\Sigma}} \leq \tilde{R} \right\}. \quad (71)$$

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

$$\left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \tilde{\zeta}_i \right\|_{\tilde{\Sigma}} \lesssim \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \varkappa\beta\rho_{\max} (\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \frac{1}{\varkappa} \tilde{R} + 1), \quad (72)$$

provided that the stepsizes α, β satisfy the conditions (64) and that $0 < \alpha < \frac{1}{\lambda_1 \lambda_{\tilde{\Sigma}}^2 \|\tilde{\Sigma}\| \log \frac{2dT}{\delta}}$.

Proof. See Section C.4. □

Putting relations (69) and (72) together, we find

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

by definition of \tilde{R}_0 in (68), provided that $\sqrt{\frac{1}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \beta\rho_{\max} \leq c$ for some constant $c > 0$ small enough. Therefore, by induction assumption, one has

$$\begin{aligned} \mathbb{P} \left\{ \max_{0 \leq i \leq t} \|\mathbf{x}_i\|_{\tilde{\Sigma}} > \tilde{R}_0 \right\} &\leq \mathbb{P} \left\{ \max_{0 \leq i < t-1} \|\mathbf{x}_i\|_{\tilde{\Sigma}} > \tilde{R}_0 \right\} + \mathbb{P} \left\{ \max_{0 \leq i < t-1} \|\mathbf{x}_i\|_{\tilde{\Sigma}} \leq \tilde{R}_0, \|\mathbf{x}_t\|_{\tilde{\Sigma}} > \tilde{R}_0 \right\} \\ &\leq \frac{(t-1)\delta}{T} + \mathbb{P} \left\{ \left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \tilde{\zeta}_i \right\|_{\tilde{\Sigma}} \gtrsim \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha\lambda_1} \log \frac{2dT}{\delta}} \varkappa\beta\rho_{\max} (\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \tilde{R}_0 + 1) \right\} \\ &\leq \frac{(t-1)\delta}{T} + \frac{\delta}{T} = \frac{t\delta}{T}. \end{aligned} \quad (73)$$

This completes our claim at this step.

Step 3: refined bound for $\|\mathbf{x}_t\|_{\tilde{\Sigma}}$. It turns out that the upper bound (68) can be tightened by taking into account the contraction effect of Ψ . In what follows, we develop a strengthened bound. Define

$$\tilde{t}_{\text{seg}} := \frac{\tilde{c}_1 \log \max \left\{ \sqrt{\tilde{\kappa}}, \frac{\sqrt{\tilde{\kappa}} \|\tilde{\Delta}_0\|_{\tilde{\Sigma}}}{\|\tilde{\theta}^*\|_{\tilde{\Sigma}+1}}, \|\tilde{\Delta}_0\|_{\tilde{\Sigma}} \sqrt{\frac{\alpha \lambda_1}{\|\tilde{\Sigma}\| \log \frac{2dT}{\delta}}} \frac{1}{\varkappa \beta \rho_{\max}} \right\}}{\alpha \lambda_1} \quad (74)$$

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 $1 - \delta$,

$$\|\mathbf{x}_t\|_{\tilde{\Sigma}} \lesssim \varkappa \beta \rho_{\max} \sqrt{\frac{1}{\alpha \lambda_1} \log \frac{2dT}{\delta}} \left(\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{\|\tilde{\Delta}_0\|_{\tilde{\Sigma}}}{2^{k-1}} + \frac{3}{2} \right) =: \tilde{R}_k \quad (75)$$

for any t obeying $k\tilde{t}_{\text{seg}} \leq t \leq T$, provided that $\sqrt{\frac{1}{\alpha \lambda_1} \log \frac{2dT}{\delta}} \varkappa \beta \rho_{\max} \leq c$ 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

$$\tilde{t}'_{\text{seg}} := \left(2 + \frac{1}{\log 2} \log \|\tilde{\theta}^*\|_{\tilde{\Sigma}} \right) \tilde{t}_{\text{seg}}, \quad (76)$$

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

$$\|\mathbf{x}_t\|_{\tilde{\Sigma}} \lesssim \varkappa \beta \rho_{\max} \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha \lambda_1} \log \frac{2dT}{\delta}} \left(\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + 2 \right). \quad (77)$$

Recall that this bound holds for any $\varkappa \in (0, 1)$ satisfying the conditions (64). Hence, Theorem 3 follows by taking $\varkappa = 8\rho_{\max} \sqrt{\frac{\alpha}{\lambda_1 \beta \lambda_2}}$ and

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

8 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 γ -discounted infinite-horizon MDPs. For TD learning with Polyak-Ruppert averaging, we improve upon existing results in terms of both the accuracy level ε and other problem-related parameters like the effective horizon $\frac{1}{1-\gamma}$, the weighted feature covariance Σ and the optimal linear estimator θ^* . 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 ε , 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 Duan et al. (2021) 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.

Acknowledgements

W. Wu and A. Rinaldo are supported in part by the NIH Grant R01 NS121913. Y. Chi is supported in part by the grants ONR N00014-19-1-2404 and NSF CCF-2106778. Y. Wei is supported in part by the NSF grants CCF-2106778, DMS-2147546/2015447 and NSF CAREER award DMS-2143215 and Google Research Scholar Award.

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 Φ , D_μ and Σ in (4), (5) and (6), respectively. Then one has*

$$\|D_\mu^{\frac{1}{2}}\Phi\Sigma^{-\frac{1}{2}}\| = 1, \quad \text{and} \quad \|D_\mu^{\frac{1}{2}}P^\pi D_\mu^{-\frac{1}{2}}\| = 1. \quad (78)$$

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

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

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

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

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

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

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

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

where $\lambda_i(\mathbf{B})$ denotes the i -th eigenvalue of the matrix \mathbf{B} . □

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

$$\Sigma^{-\frac{1}{2}}\mathbf{A}^\top\Sigma^{-1}\mathbf{A}\Sigma^{-\frac{1}{2}} \succeq (1-\gamma)^2\mathbf{I}, \quad (79a)$$

$$\Sigma^{-\frac{1}{2}}\mathbf{A}\Sigma^{-1}\mathbf{A}^\top\Sigma^{-\frac{1}{2}} \succeq (1-\gamma)^2\mathbf{I}, \quad (79b)$$

$$\|\Sigma^{\frac{1}{2}}(\mathbf{A}^\top)^{-1}\Sigma\mathbf{A}^{-1}\Sigma^{\frac{1}{2}}\| \leq (1-\gamma)^{-2}, \quad (79c)$$

$$\|\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}\Sigma(\mathbf{A}^\top)^{-1}\Sigma^{\frac{1}{2}}\| \leq (1-\gamma)^{-2}, \quad (79d)$$

$$\|\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}\Sigma^{\frac{1}{2}}\| \leq (1-\gamma)^{-1}, \quad (79e)$$

$$\|\Sigma^{-1/2}\Phi^\top D_\mu\| \leq \max_{s \in \mathcal{S}} \phi(s)^\top \Sigma^{-1} \phi(s), \quad (79f)$$

$$\|\mathbf{I} - \eta\mathbf{A}\| \leq 1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma), \quad \forall 0 < \eta < \frac{1-\gamma}{4\|\Sigma\|}, \quad (79g)$$

$$\|\Sigma\| \leq 1, \quad \|\Sigma^{-1}\| \geq 1, \quad (79h)$$

$$\|\Sigma^{-\frac{1}{2}}\mathbf{b}\|_2 \leq 1. \quad (79i)$$

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

Proof of Eqn. (79a) and (79b). We start with the lower bound on $\Sigma^{-\frac{1}{2}}\mathbf{A}^\top\Sigma^{-1}\mathbf{A}\Sigma^{-\frac{1}{2}}$. To begin with, observe that

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

where

$$\tilde{\Phi} := D_\mu^{\frac{1}{2}}\Phi\Sigma^{-\frac{1}{2}} \quad \text{and} \quad P_{D_\mu} := D_\mu^{\frac{1}{2}}PD_\mu^{-\frac{1}{2}}. \quad (80)$$

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

$$\left|\mathbf{x}^\top\tilde{\Phi}^\top P_{D_\mu}\tilde{\Phi}\mathbf{x}\right| \leq \|\tilde{\Phi}^\top P_{D_\mu}\tilde{\Phi}\| \leq \|\tilde{\Phi}\|^2\|P_{D_\mu}\| = 1. \quad (81)$$

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 (79a). The proof for (79b) follows from an identical argument and is omitted for brevity.

Proof of Eqn. (79c), (79d) and (79e). With the above bounds in place, we can further obtain

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

where $\lambda_{\min}(\mathbf{B})$ denotes the smallest eigenvalue of \mathbf{B} , and the last inequality comes from (79b). This establishes (79c). The inequality (79d) 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 (79e).

Proof of Eqn. (79g). Recalling that $\Sigma = \Phi^\top D_\mu\Phi$, we can arrange terms to derive

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

where $\tilde{\Phi}$ and P_{D_μ} are defined in (80). Here, the last line follows since $\|\tilde{\Phi}^\top P_{D_\mu}\tilde{\Phi}\| \leq 1$ — a fact that has already been shown in (81). 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

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

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

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

Proof of Eqn. (79h). 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 $\|\Sigma^{-1}\| \geq 1/\|\Sigma\| \geq 1$.

Proof of Eqn. (79i). Finally, we observe that

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

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 $T \gtrsim \log\left(\frac{d}{\delta}\right)$. Then the vector \mathbf{b} defined in (12b) obeys that, with probability exceeding $1 - \delta$,

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

Proof. See Section C.5. \square

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

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

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

Proof. See Section C.5. \square

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.

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

$$P_{\mathbf{q}}(s' | s) = \begin{cases} q_s \mathbb{1}(s' = s) + \frac{1-q_s}{|\mathcal{S}|-d+1} \mathbb{1}(s' \geq d) & \text{for } s < d; \\ \frac{\gamma}{|\mathcal{S}|-d+1} \mathbb{1}(s' \geq d) + \frac{1-q_{s'}}{d-1} \mathbb{1}(s' < d) & \text{for } s \geq d. \end{cases} \quad (82)$$

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,

$$\sum_{s=1}^{d-1} \mathbb{1}(q_s = q_+) = \sum_{s=1}^{d-1} \mathbb{1}(q_s = q_-) = (d-1)/2. \quad (83)$$

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

$$\mu(s) = \begin{cases} \frac{1}{2(d-1)} & \text{for } s < d; \\ \frac{1}{2(|\mathcal{S}|-d+1)} & \text{for } s \geq d. \end{cases} \quad (84)$$

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

$$\theta^*(d) = V^*(s) = \frac{1}{1-\gamma^2 - \sum_{i=1}^{d-1} \frac{\gamma^2(1-q_i)^2}{(d-1)(1-\gamma q_i)}}, \quad (85)$$

$$\theta^*(i) = V^*(i) = \frac{\gamma(1-q_i)}{1-\gamma q_i} V^*(s), \text{ for } s \geq d \text{ and } i < d. \quad (86)$$

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 (83).

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

$$\max_s \{\phi(s)^\top \Sigma^{-1} \phi(s)\} \asymp d. \quad (87)$$

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

$$\begin{aligned} \|\theta^*\|_\Sigma^2 &= \sum_{i=1}^{d-1} \frac{1}{2(d-1)} \theta^{*2}(i) + \frac{1}{2} \theta^{*2}(d) = \sum_{i=1}^{d-1} \frac{1}{2(d-1)} \left[\frac{\gamma(1-q_i)}{1-\gamma q_i} V^*(s) \right]^2 + \frac{1}{2} [V^*(s)]^2 \\ &\asymp \sum_{i=1}^{d-1} \frac{1}{2(d-1)} \left[\frac{\gamma(1-\gamma)}{1-\gamma} \frac{1}{1-\gamma} \right]^2 + \frac{1}{2} \left[\frac{1}{1-\gamma} \right]^2 \\ &\asymp \frac{1}{(1-\gamma)^2}. \end{aligned} \quad (88)$$

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

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

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

$$\|\theta^* - \theta'^*\|_\Sigma \gtrsim \varepsilon. \quad (89)$$

In other words, if we want each θ^* to be ε apart from each other, it is sufficient to construct a set \mathcal{Q} where every \mathbf{q} and \mathbf{q}' are $(d-1)/16$ apart in Hamming distance. By virtue of the Gilbert-Varshamov lemma (Gilbert, 1952), there exists a set \mathcal{Q} such that

$$M := |\mathcal{Q}| \geq e^{d/16} \quad \text{and} \quad \sum_{s=1}^{d-1} \mathbb{1}(q_s \neq q'_s) \geq \frac{d}{16} \quad \text{for any } q, q' \in \mathcal{Q} \text{ obeying } q \neq q'. \quad (90)$$

The Fano method transforms the problem of estimating θ^* into an M -ary testing problem among the above MDPs $\{\mathbb{P}_{\mathbf{q}^1}, \mathbb{P}_{\mathbf{q}^2}, \dots, \mathbb{P}_{\mathbf{q}^M}\}$. More specifically, in view of Fano's inequality (Tsybakov (2009)), the probability of interest thus satisfies

$$\mathbb{P}\left(\|\hat{\theta} - \theta^*\|_\Sigma \gtrsim \varepsilon\right) \geq 1 - \frac{1}{\log M} \left(\frac{1}{M^2} \sum_{j,k=1}^M \text{KL}(\mathbb{P}_{\mathbf{q}^j}^T \parallel \mathbb{P}_{\mathbf{q}^k}^T) + \log 2 \right), \quad (91)$$

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 $\mathbb{P}_{\mathbf{q}}$ and $\mathbb{P}_{\mathbf{q}'}$. Here $\mathbb{P}_{\mathbf{q}}$ denotes the joint distribution of (s, s') when the transition is made according to $P_{\mathbf{q}}(s' | s)$ (cf. (82)). More specifically, given $s \sim \mu_{\mathbf{q}}$ and $s' | s \sim P_{\mathbf{q}}(s' | s)$, one has

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

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

$$\begin{aligned}
\text{KL}(\mathbb{P}_{\mathbf{q}} \parallel \mathbb{P}_{\mathbf{q}'}) &\leq \chi^2(\mathbb{P}_{\mathbf{q}'} \parallel \mathbb{P}_{\mathbf{q}}) \\
&= \sum_{s,s'} \frac{(\mathbb{P}_{\mathbf{q}}(s,s') - \mathbb{P}_{\mathbf{q}'}(s,s'))^2}{\mathbb{P}_{\mathbf{q}}(s,s')} \\
&= \left(\sum_{s < d, s' < d} + \sum_{s < d, s' \geq d} + \sum_{s \geq d, s' < d} + \sum_{s \geq d, s' \geq d} \right) \frac{(\mathbb{P}_{\mathbf{q}}(s,s') - \mathbb{P}_{\mathbf{q}'}(s,s'))^2}{\mathbb{P}_{\mathbf{q}}(s,s')} \\
&= \sum_{s=1}^{d-1} \frac{1}{2(d-1)} \frac{(q_s - q'_s)^2}{q_s} + \sum_{s < d, s' \geq d} \frac{1}{2(d-1)(S-d+1)} \frac{[(1-q_s) - (1-q'_s)]^2}{1-q_s} \\
&\quad + \sum_{s \geq d, s' < d} \frac{1}{2(d-1)(S-d+1)} \frac{[(1-q_{s'}) - (1-q'_{s'})]^2}{1-q_{s'}} + 0 \\
&\lesssim \sum_{s=1}^{d-1} \frac{1}{2(d-1)} \frac{[2\varepsilon(1-\gamma)^2]^2}{1-\gamma} + \sum_{s < d} \frac{1}{2(d-1)} \frac{[2\varepsilon(1-\gamma)^2]^2}{1-\gamma} + \sum_{s' < d} \frac{1}{2(d-1)} \frac{[2\varepsilon(1-\gamma)^2]^2}{1-\gamma} \\
&\asymp \varepsilon^2(1-\gamma)^3.
\end{aligned}$$

As a result, we have

$$\text{KL}(\mathbb{P}_{\mathbf{q}}^T \parallel \mathbb{P}_{\mathbf{q}'}^T) \lesssim \varepsilon^2(1-\gamma)^3 T. \quad (92)$$

Substituting the above relation into (91) gives

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

To prove Theorem 2, it is enough to take the above together with relations (87) and (88).

C Proofs of auxiliary lemmas and claims

C.1 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 σ -algebra generated by $\{(s_j, s'_j)\}_{j \leq i}$). It is then easy to check that $\{(\mathbf{I} - \eta \mathbf{A})^{t-i-1} \tilde{\boldsymbol{\xi}}_i\}$ 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

$$W := \sum_{i=l}^u \mathbb{E}_{i-1} \left[\left\| \boldsymbol{\Sigma}^{1/2} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \boldsymbol{\xi}_i \right\|_2^2 \mathbf{1}\{\mathcal{H}_i\} \right], \quad \text{and} \quad B := \max_{i:l \leq i \leq u} \left\| \boldsymbol{\Sigma}^{1/2} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \boldsymbol{\xi}_i \right\|_2, \quad (93)$$

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

$$\alpha := \left(1 - \frac{1}{2} \eta (1-\gamma) \lambda_{\min}(\boldsymbol{\Sigma})\right)^{t-u-1}. \quad (94)$$

Control of W . Direct calculations yield

$$\begin{aligned}
W &= \sum_{i=l}^u \mathbb{E}_{i-1} \left[\boldsymbol{\xi}_i^\top (\mathbf{I} - \eta \mathbf{A}^\top)^{t-i-1} \boldsymbol{\Sigma} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \boldsymbol{\xi}_i \mathbf{1}\{\mathcal{H}_i\} \right] \\
&\leq \sum_{i=l}^u \left\| (\mathbf{I} - \eta \mathbf{A}^\top)^{t-i-1} \boldsymbol{\Sigma} (\mathbf{I} - \eta \mathbf{A})^{t-i-1} \right\| \cdot \mathbb{E}_{i-1} \left[\left\| \boldsymbol{\xi}_i \right\|_2^2 \mathbf{1}\{\mathcal{H}_i\} \right]
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(i)}{\leq} \sum_{i=l}^u \|\Sigma\| \left(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma)\right)^{2t-2i-2} 2 \max_{i:l \leq i \leq u} \left\{ \mathbb{E}_{i-1} [\|(\mathbf{A}_i - \mathbf{A})\boldsymbol{\theta}_i\|_2^2 \mathbf{1}\{\mathcal{H}_i\}] + \mathbb{E}_{i-1} [\|\mathbf{b}_i - \mathbf{b}\|_2^2] \right\} \\
&\stackrel{(ii)}{\leq} \frac{4\|\Sigma\|\alpha^2}{\eta(1-\gamma)\lambda_{\min}(\Sigma)} \max_{i:l \leq i \leq u} \left\{ \mathbb{E}_{i-1} [\|(\mathbf{A}_i - \mathbf{A})\boldsymbol{\theta}_i\|_2^2 \mathbf{1}\{\mathcal{H}_i\}] + \mathbb{E}_{i-1} [\|\mathbf{b}_i - \mathbf{b}\|_2^2] \right\}, \tag{95}
\end{aligned}$$

where (i) follows from the property (79g) (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 (94) of α .

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

$$\begin{aligned}
&\mathbb{E}_{i-1} [\|(\mathbf{A}_i - \mathbf{A})\boldsymbol{\theta}_i\|_2^2 \mathbf{1}\{\mathcal{H}_i\}] \leq \mathbb{E}_{i-1} [\|\mathbf{A}_i \boldsymbol{\theta}_i\|_2^2 \mathbf{1}\{\mathcal{H}_i\}] \\
&= \mathbb{E}_{i-1} \left[\boldsymbol{\theta}_i^\top (\boldsymbol{\phi}(s_i) - \gamma\boldsymbol{\phi}(s'_i)) \boldsymbol{\phi}(s_i)^\top \boldsymbol{\phi}(s_i) (\boldsymbol{\phi}(s_i) - \gamma\boldsymbol{\phi}(s'_i))^\top \boldsymbol{\theta}_i \mathbf{1}\{\mathcal{H}_i\} \right] \\
&\leq \max_s \|\boldsymbol{\phi}(s)\|_2^2 \cdot \mathbb{E}_{i-1} \left[\boldsymbol{\theta}_i^\top (\boldsymbol{\phi}(s_i) - \gamma\boldsymbol{\phi}(s'_i)) (\boldsymbol{\phi}(s_i) - \gamma\boldsymbol{\phi}(s'_i))^\top \boldsymbol{\theta}_i \mathbf{1}\{\mathcal{H}_i\} \right] \\
&\stackrel{(i)}{\leq} 2 \max_s \|\boldsymbol{\phi}(s)\|_2^2 \left(\mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \boldsymbol{\phi}(s_i) \boldsymbol{\phi}(s_i)^\top \boldsymbol{\theta}_i \mathbf{1}\{\mathcal{H}_i\}] + \gamma^2 \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \boldsymbol{\phi}(s'_i) \boldsymbol{\phi}(s'_i)^\top \boldsymbol{\theta}_i \mathbf{1}\{\mathcal{H}_i\}] \right) \\
&\stackrel{(ii)}{=} 2 \max_s \|\boldsymbol{\phi}(s)\|_2^2 \left(\mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \Sigma \boldsymbol{\theta}_i \mathbf{1}\{\mathcal{H}_i\}] + \gamma^2 \mathbb{E}_{i-1} [\boldsymbol{\theta}_i^\top \Sigma \boldsymbol{\theta}_i \mathbf{1}\{\mathcal{H}_i\}] \right) \\
&\stackrel{(iii)}{\leq} 4 \|\boldsymbol{\theta}_i\|_\Sigma^2 \mathbf{1}\{\mathcal{H}_i\} \leq 4(\|\boldsymbol{\theta}^*\|_\Sigma + \|\boldsymbol{\Delta}_i\|_\Sigma)^2 \mathbf{1}\{\mathcal{H}_i\} \\
&\leq 4(\|\boldsymbol{\theta}^*\|_\Sigma + R)^2, \tag{96}
\end{aligned}$$

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 Σ and the fact that $s_i, s'_i \sim \mu$ in this case, (iii) holds due to the assumption $\max_s \|\boldsymbol{\phi}(s)\|_2 \leq 1$, and the last inequality results from the definition (47) of the event \mathcal{H}_i . Similarly, one can derive

$$\mathbb{E}_{i-1} [\|\mathbf{b}_i - \mathbf{b}\|_2^2] \leq \mathbb{E}_{i-1} [\|\mathbf{b}_i\|_2^2] = \mathbb{E}_{i-1} [\|\boldsymbol{\phi}(s_i)r(s_i)\|_2^2] \leq 1, \tag{97}$$

where the last inequality holds since $\max_s \|\boldsymbol{\phi}(s)\|_2 \leq 1$ and $\max_s |r(s)| \leq 1$. Substitution into (95) yields

$$W \leq \frac{4\kappa}{\eta(1-\gamma)} \alpha^2 \left\{ 4(\|\boldsymbol{\theta}^*\|_\Sigma + R)^2 + 1 \right\} =: W_{\max}. \tag{98}$$

Control of B . By definition of B , one can write

$$\begin{aligned}
B &= \max_{i:l \leq i \leq u} \|\Sigma^{\frac{1}{2}}(\mathbf{I} - \eta\mathbf{A})^{t-i-1} \boldsymbol{\xi}_i\|_2 \mathbf{1}\{\mathcal{H}_i\} = \max_{i:l \leq i \leq u} \|\Sigma^{\frac{1}{2}}(\mathbf{I} - \eta\mathbf{A})^{t-i-1} \Sigma^{\frac{1}{2}} \Sigma^{-\frac{1}{2}} \boldsymbol{\xi}_i\|_2 \mathbf{1}\{\mathcal{H}_i\} \\
&\leq \|\Sigma\| \max_{i:l \leq i \leq u} \|\mathbf{I} - \eta\mathbf{A}\|^{t-i-1} \cdot \max_{i:l \leq i \leq u} \|\Sigma^{-\frac{1}{2}} \boldsymbol{\xi}_i\|_2 \mathbf{1}\{\mathcal{H}_i\} \\
&\leq \alpha \|\Sigma\| \max_{i:l \leq i \leq u} \left\{ \|\Sigma^{-\frac{1}{2}}(\mathbf{A}_i - \mathbf{A})\boldsymbol{\theta}_i\|_2 \mathbf{1}\{\mathcal{H}_i\} + \|\Sigma^{-\frac{1}{2}}(\mathbf{b}_i - \mathbf{b})\|_2 \right\}, \tag{99}
\end{aligned}$$

where the last step results from (79g) (with the restriction that $\eta < (1-\gamma)/(4\|\Sigma\|)$) and the definition (94) of α . It then suffices to control the two terms on the right-hand side of (99). To begin with, we have

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

Recall from (135) that $\|\Sigma^{-\frac{1}{2}}\mathbf{A}_i\Sigma^{-\frac{1}{2}}\| \leq 2 \max_s \|\Sigma^{-1/2}\boldsymbol{\phi}(s)\|_2^2$, and similarly $\|\Sigma^{-\frac{1}{2}}\mathbf{A}\Sigma^{-\frac{1}{2}}\| \leq 2 \max_s \|\Sigma^{-1/2}\boldsymbol{\phi}(s)\|_2^2$. We then have

$$\|\Sigma^{-\frac{1}{2}}(\mathbf{A}_i - \mathbf{A})\boldsymbol{\theta}_i\|_2 \leq 4 \max_s \left\{ \boldsymbol{\phi}(s)^\top \Sigma^{-1} \boldsymbol{\phi}(s) \right\} (\|\boldsymbol{\theta}^*\|_\Sigma + \|\boldsymbol{\Delta}_i\|_\Sigma). \tag{100}$$

Regarding the second term of (99), direct calculations give

$$\begin{aligned} \|\Sigma^{-\frac{1}{2}}(\mathbf{b}_i - \mathbf{b})\|_2^2 &\leq 2\|\Sigma^{-\frac{1}{2}}\mathbf{b}_i\|_2^2 + 2\|\Sigma^{-\frac{1}{2}}\mathbf{b}\|_2^2 = 2\|\Sigma^{-\frac{1}{2}}\phi(s_i)r(s_i)\|_2^2 + 2\|\Sigma^{-\frac{1}{2}}\mathbb{E}_{s\sim\mu}[\phi(s)r(s)]\|_2^2 \\ &\leq 4\max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\} \max_s |r(s)|^2 \leq 4\max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\}. \end{aligned} \quad (101)$$

Substituting the preceding two bounds into (99), we arrive at

$$\begin{aligned} B &\leq 4\alpha\|\Sigma\| \left(\max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\} \max_{i:i < t} (\|\theta^*\|_\Sigma + \|\Delta_i\|_\Sigma) \mathbb{1}\{\mathcal{H}_i\} + \sqrt{\max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\}} \right) \\ &\leq 4\alpha\|\Sigma\| \left(\max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\} (\|\theta^*\|_\Sigma + R) + \sqrt{\max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\}} \right) \\ &\leq 4\alpha\|\Sigma\| \max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\} (\|\theta^*\|_\Sigma + R + 1) \\ &\leq 4\alpha\|\Sigma\| \|\Sigma^{-1}\| (\|\theta^*\|_\Sigma + R + 1) = 4\kappa\alpha (\|\theta^*\|_\Sigma + R + 1) =: B_{\max}. \end{aligned} \quad (102)$$

Here, the last line follows from the assumption $\max\|\phi(s)\|_2 \leq 1$, while the second to last inequality holds since $\max_s \{\phi(s)^\top \Sigma^{-1}\phi(s)\} \geq 1$ (cf. (137)).

Invoking matrix Freedman's inequality. Equipped with the above bounds (98) and (102), we are ready to apply Freedman's inequality (Tropp, 2011, Corollary 1.3) (or a version in (Li et al., 2023a, Section A)), which asserts that

$$\begin{aligned} \left\| \sum_{i=0}^{t-1} (\mathbf{I} - \eta\mathbf{A})^{t-i-1} \tilde{\xi}_i \right\|_\Sigma &\leq 2\sqrt{W_{\max} \log \frac{2dT}{\delta}} + \frac{4}{3} B_{\max} \log \frac{2dT}{\delta} \\ &= \alpha \cdot \left\{ 2\sqrt{\frac{4\kappa}{\eta(1-\gamma)} \left\{ 4(\|\theta^*\|_\Sigma + R)^2 + 1 \right\} \log \frac{2dT}{\delta}} + \frac{16\kappa}{3} (\|\theta^*\|_\Sigma + R + 1) \log \frac{2dT}{\delta} \right\} \\ &\leq 16(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma))^{t-u-1} (\|\theta^*\|_\Sigma + R + 1) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \end{aligned} \quad (103)$$

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 α 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 (79h)).

C.2 Proof of the inequalities (58a) and (58b)

Proof of the inequality (58a). Combining the triangle inequality with the definition (57) ensures that

$$\begin{aligned} \|\mathbf{A}_0^{(T+1)} \Delta_0\|_\Sigma &\leq \|\mathbf{A}^{-1} \Delta_0\|_\Sigma + \|\mathbf{A}^{-1} (\mathbf{I} - \eta\mathbf{A})^{T+1} \Delta_0\|_\Sigma \\ &= \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \Sigma^{-1} \Sigma^{\frac{1}{2}} \Delta_0\|_2 + \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}} \Sigma^{-\frac{1}{2}} (\mathbf{I} - \eta\mathbf{A})^{T+1} \Sigma^{-\frac{1}{2}} \Sigma^{\frac{1}{2}} \Delta_0\|_2 \\ &\leq \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \cdot \|\Sigma^{-1}\| \cdot \|\Delta_0\|_\Sigma + \|\Sigma^{\frac{1}{2}} \mathbf{A}^{-1} \Sigma^{\frac{1}{2}}\| \cdot \|\Sigma^{-\frac{1}{2}}\|^2 \cdot \|\mathbf{I} - \eta\mathbf{A}\|^{T+1} \cdot \|\Delta_0\|_\Sigma \\ &\leq \frac{\|\Sigma^{-1}\|}{1-\gamma} \left\{ 1 + \left(1 - \frac{1}{2}\eta(1-\gamma)\lambda_{\min}(\Sigma) \right)^{T+1} \right\} \|\Delta_0\|_\Sigma \\ &\leq \frac{2\|\Sigma^{-1}\|}{1-\gamma} \|\Delta_0\|_\Sigma \end{aligned} \quad (104)$$

as claimed. Here, the second to last step follows from (79e) and (79g), provided that $\eta \leq (1-\gamma)/(4\|\Sigma\|)$.

Proof of the inequality (58b). Again, the triangle inequality together with the definition (57) yields

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

$$\leq \left\| \mathbf{A}^{-1} \sum_{i=0}^{T-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_{\Sigma} + \left\| \mathbf{A}^{-1} \sum_{i=0}^{T-1} (\mathbf{b}_i - \mathbf{b}) \right\|_{\Sigma} + \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{I} - \eta \mathbf{A})^{T-i} \boldsymbol{\xi}_i \right\|_{\Sigma}, \quad (105)$$

leaving us with three terms to handle. Here in the second line, we substitute in the definition of $\boldsymbol{\xi}_i$ (42).

- The second term of (105) can be bounded by Lemma 6, which asserts that

$$\frac{1}{T} \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{b}_i - \mathbf{b}) \right\|_{\Sigma} \lesssim \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\phi}(s)}{T(1-\gamma)^2} \log\left(\frac{d}{\delta}\right)} \quad (106)$$

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

- For the third term of (105), invoking the property (79e) again yields

$$\begin{aligned} \left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{I} - \eta \mathbf{A})^{T-i} \boldsymbol{\xi}_i \right\|_{\Sigma} &= \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\Sigma}^{-1} \sum_{i=0}^{T-1} \boldsymbol{\Sigma}^{\frac{1}{2}} (\mathbf{I} - \eta \mathbf{A})^{T-i} \boldsymbol{\xi}_i \right\|_2 \\ &\leq \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \right\| \cdot \left\| \boldsymbol{\Sigma}^{-1} \right\| \cdot \left\| \sum_{i=0}^{T-1} \boldsymbol{\Sigma}^{\frac{1}{2}} (\mathbf{I} - \eta \mathbf{A})^{T-i} \boldsymbol{\xi}_i \right\|_2 \\ &\leq \frac{\left\| \boldsymbol{\Sigma}^{-1} \right\|}{1-\gamma} \left\| \sum_{i=0}^{T-1} (\mathbf{I} - \eta \mathbf{A})^{T-i} \boldsymbol{\xi}_i \right\|_{\Sigma}. \end{aligned} \quad (107)$$

Repeating the same analysis as in Step 3 to see that

$$\left\| \sum_{i=0}^{T-1} (\mathbf{I} - \eta \mathbf{A})^{T-i} \boldsymbol{\xi}_i \right\|_{\Sigma} \leq 16(2\|\boldsymbol{\theta}^*\|_{\Sigma} + 3) \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)}} \quad (108)$$

with probability at least $1 - \delta$. Substitution into (107) leads to

$$\left\| \sum_{i=0}^{T-1} \mathbf{A}^{-1} (\mathbf{I} - \eta \mathbf{A})^{T-i} \boldsymbol{\xi}_i \right\|_{\Sigma} \leq 16(2\|\boldsymbol{\theta}^*\|_{\Sigma} + 3) \left\| \boldsymbol{\Sigma}^{-1} \right\| \sqrt{\frac{\kappa \log \frac{2dT}{\delta}}{\eta(1-\gamma)^3}}. \quad (109)$$

- It then boils down to bounding the first term of (105). In light of (54), we decompose it as follows

$$\left\| \frac{1}{T} \sum_{i=0}^{T-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 \leq \left\| \frac{1}{T} \sum_{i=0}^{\tilde{t}_{\text{seg}}-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 + \left\| \frac{1}{T} \sum_{i=\tilde{t}_{\text{seg}}}^{T-1} \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2. \quad (110)$$

Bounding these terms requires the following lemma, whose proof is deferred to Section C.6.

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

$$\boldsymbol{\theta}'_i := \boldsymbol{\theta}_i \mathbb{1}\{\mathcal{H}_i\}, \quad \mathcal{H}_i := \{\|\boldsymbol{\Delta}_i\|_{\Sigma} \leq R\}, \quad (111)$$

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

$$\left\| \frac{1}{u-l+1} \sum_{i=l}^u \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}'_i \right\|_2 \leq \frac{16(\|\boldsymbol{\theta}^*\|_{\Sigma} + R)}{1-\gamma} \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\phi}(s) \log \frac{2d}{\delta}}{u-l+1}} \quad (112)$$

provided that

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

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

$$\begin{aligned} \left\| \frac{1}{t'_{\text{seg}}} \sum_{i=0}^{t'_{\text{seg}}-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 &= \left\| \frac{1}{t'_{\text{seg}}} \sum_{i=0}^{t'_{\text{seg}}-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \mathbf{1}_{\{\|\Delta_i\| \leq R_0\}} \right\|_2 \\ &\leq \frac{16(\|\boldsymbol{\theta}^*\|_{\Sigma} + R_0)}{1 - \gamma} \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \Sigma^{-1} \boldsymbol{\phi}(s) \log \frac{2d}{\delta}}{t'_{\text{seg}}}}, \end{aligned}$$

as long as $t'_{\text{seg}} \geq \frac{4\|\Sigma^{-1}\| \log \frac{2d}{\delta}}{9} \geq \frac{4\max_s \boldsymbol{\phi}(s)^\top \Sigma^{-1} \boldsymbol{\phi}(s) \log \frac{2d}{\delta}}{9}$. 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

$$\begin{aligned} \left\| \frac{1}{T - t'_{\text{seg}}} \sum_{i=t'_{\text{seg}}}^{T-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 &\leq \frac{16(\|\boldsymbol{\theta}^*\|_{\Sigma} + 32\sqrt{\frac{\eta\kappa \log \frac{2dT}{\delta}}{1-\gamma}}(\|\boldsymbol{\theta}^*\|_{\Sigma} + 2))}{1 - \gamma} \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \Sigma^{-1} \boldsymbol{\phi}(s) \log \frac{2d}{\delta}}{T - t'_{\text{seg}}}} \\ &\leq \frac{16(1.5\|\boldsymbol{\theta}^*\|_{\Sigma} + 2)}{1 - \gamma} \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \Sigma^{-1} \boldsymbol{\phi}(s) \log \frac{2d}{\delta}}{T - t'_{\text{seg}}}} \end{aligned}$$

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

$$\left\| \frac{1}{T} \sum_{i=0}^{T-1} \Sigma^{\frac{1}{2}} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}_i \right\|_2 \leq \frac{32(\|\boldsymbol{\theta}^*\|_{\Sigma} + \sqrt{\frac{t'_{\text{seg}}}{T}} R_0 + 1)}{1 - \gamma} \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \Sigma^{-1} \boldsymbol{\phi}(s) \log \frac{2d}{\delta}}{T}}. \quad (113)$$

Combining the preceding bounds (106), (109) and (113) with (105), we reach the conclusion that with probability of at least $1 - \delta$,

$$\left\| \sum_{i=0}^{T-1} \mathbf{A}_i^{(t)} \boldsymbol{\xi}_i \right\|_{\Sigma} \asymp \left\{ \sqrt{\frac{\max_s \boldsymbol{\phi}(s)^\top \Sigma^{-1} \boldsymbol{\phi}(s) \log \frac{2d}{\delta}}{T(1-\gamma)^2}} + \frac{\|\Sigma^{-1}\|}{T} \sqrt{\frac{\kappa \log \frac{dT}{\delta}}{\eta(1-\gamma)^3}} \right\} (\|\boldsymbol{\theta}^*\|_{\Sigma} + 1),$$

as long as $T \geq t'_{\text{seg}} \kappa \|\Delta_0\|_{\Sigma}^2$, where we use the definition (44) of R_0 . It thus establishes the inequality (58b).

C.3 Proof of Lemma 2

We first decompose Ψ 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

$$\begin{aligned} \|\Psi\| &\leq \max\{\|\mathbf{I} - \alpha \tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\|, \|\mathbf{I} - \beta \tilde{\Sigma}\|\} + \left\| \frac{1}{\kappa} \alpha \gamma \mathbf{\Pi}^\top \right\| \\ &\quad + \|\kappa \alpha (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\| + \|\alpha \gamma (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \mathbf{\Pi}^\top\|. \end{aligned}$$

Note that by definition of λ_w and λ_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

$$\begin{aligned} \left\| \frac{1}{\kappa} \alpha \gamma \mathbf{\Pi}^\top \right\| &\leq \frac{\alpha \gamma \rho_{\max}}{\kappa}, \\ \|\kappa \alpha (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \tilde{\mathbf{A}}^\top \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\| &\leq \kappa \alpha (1 + \gamma \lambda_\Sigma \rho_{\max}) \lambda_\Sigma (2\rho_{\max})^2, \\ \|\alpha \gamma (\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}) \mathbf{\Pi}^\top\| &\leq \alpha \gamma (\rho_{\max} + \gamma \lambda_\Sigma \rho_{\max}^2). \end{aligned}$$

In summary, as long as

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

one has

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

C.4 Proof of Lemma 3

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

$$\tilde{W} := \sum_{i=0}^{t-1} \mathbb{E}_{i-1} \left[\|\Psi^{t-i-1} \tilde{\zeta}_i\|_{\tilde{\Sigma}}^2 \mathbf{1} \left\{ \tilde{\mathcal{H}}_i \right\} \right], \quad \text{and} \quad \tilde{B} := \max_{i:0 \leq i \leq t-1} \left\| \Psi^{t-i-1} \tilde{\zeta}_i \mathbf{1} \left\{ \tilde{\mathcal{H}}_i \right\} \right\|_{\tilde{\Sigma}}. \quad (114)$$

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.

Controlling the norm of ζ_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 (96), (97), (100) and (101), the following bounds hold true:

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

$$\mathbb{E}_{i-1} \left\| (\tilde{\mathbf{A}}_i - \tilde{\mathbf{A}}) \tilde{\theta}_i \right\|_2^2 \leq 4\rho_{\max}^2 \left(\|\tilde{\theta}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right), \quad \text{and} \quad (115)$$

$$\left\| \tilde{\Sigma}^{-1/2} (\tilde{\mathbf{A}}_i - \tilde{\mathbf{A}}) \tilde{\theta}_i \right\|_2 \leq 4\rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left(\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right); \quad (116)$$

- 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

$$\mathbb{E}_{i-1} \left\| (\mathbf{\Pi}_i - \mathbf{\Pi})^\top \mathbf{z}_i \right\|_2^2 \leq \rho_{\max}^2 \|\mathbf{z}_i\|_{\tilde{\Sigma}}^2, \quad \text{and} \quad (117)$$

$$\left\| \tilde{\Sigma}^{-1/2} (\mathbf{\Pi}_i - \mathbf{\Pi})^\top \mathbf{z}_i \right\|_2 \leq 2\rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \|\mathbf{z}_i\|_{\tilde{\Sigma}}; \quad (118)$$

- 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

$$\mathbb{E}_{i-1} \left\| (\tilde{\Sigma}_i - \tilde{\Sigma}) \mathbf{z}_i \right\|_2^2 \leq \|\mathbf{z}_i\|_{\tilde{\Sigma}}^2, \quad \text{and} \quad (119)$$

$$\left\| \tilde{\Sigma}^{-1/2} (\tilde{\Sigma}_i - \tilde{\Sigma}) \mathbf{z}_i \right\|_2 \leq 2 \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \|\mathbf{z}_i\|_{\tilde{\Sigma}}; \quad (120)$$

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

$$\mathbb{E}_{i-1} \left\| \tilde{\mathbf{b}}_i - \tilde{\mathbf{b}} \right\|_2^2 \leq \rho_{\max}^2, \quad \text{and} \quad (121)$$

$$\left\| \tilde{\Sigma}^{-1/2} (\tilde{\mathbf{b}}_i - \tilde{\mathbf{b}}) \right\|_2^2 \leq 4\rho_{\max}^2 \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\}. \quad (122)$$

Therefore, by triangle inequality, the norm of $\boldsymbol{\nu}_i$ can be bounded by

$$\mathbb{E}_{i-1} \|\boldsymbol{\nu}_i\|_2^2 \lesssim \rho_{\max}^2 \left[\left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right) + 1 + \gamma^2 \|\mathbf{w}_i\|_{\tilde{\Sigma}}^2 \right], \quad \text{and} \quad (123)$$

$$\left\| \tilde{\Sigma}^{-1/2} \boldsymbol{\nu}_i \right\|_2 \lesssim \rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left[\left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + \gamma \|\mathbf{w}_i\|_{\tilde{\Sigma}} + 1 \right]; \quad (124)$$

similarly, the norm of $\boldsymbol{\eta}_i$ can be bounded by

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

$$\left\| \tilde{\Sigma}^{-1/2} \boldsymbol{\eta}_i \right\|_2 \lesssim \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left\{ \rho_{\max} \left[\left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + 1 \right] + \|\mathbf{w}_i\|_{\tilde{\Sigma}} \right\}. \quad (126)$$

By combining (123) and (125) with the definition of ζ_i (65), we obtain the following bound:

$$\begin{aligned} \mathbb{E}_{i-1} \|\zeta_i\|_2^2 &\lesssim \alpha^2 \mathbb{E}_{i-1} \|\boldsymbol{\nu}_i\|_2^2 + \varkappa^2 \alpha^2 \|\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}\|^2 \mathbb{E}_{i-1} \|\boldsymbol{\nu}_i\|_2^2 + \varkappa^2 \beta^2 \mathbb{E}_{i-1} \|\boldsymbol{\eta}_i\|_2^2 \\ &\lesssim \alpha^2 \left(1 + \varkappa^2 (1 + \gamma \lambda_{\Sigma} \rho_{\max})^2 \right) \cdot \rho_{\max}^2 \left[4 \left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right) + 1 + \gamma^2 \|\mathbf{w}_i\|_{\tilde{\Sigma}}^2 \right] \\ &\quad + \varkappa^2 \beta^2 \cdot \left\{ \rho_{\max}^2 \left[\left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}^2 + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}}^2 \right) + 1 \right] + \|\mathbf{w}_i\|_{\tilde{\Sigma}}^2 \right\} \\ &\lesssim \varkappa^2 \beta^2 \rho_{\max}^2 \left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}}^2 + \frac{1}{\varkappa^2} \|\mathbf{x}_i\|_{\tilde{\Sigma}}^2 + 1 \right), \end{aligned} \quad (127)$$

and

$$\begin{aligned} \left\| \tilde{\Sigma}^{-1/2} \zeta_i \right\|_2 &\lesssim \alpha \|\tilde{\Sigma}^{-1/2} \boldsymbol{\nu}_i\|_2 + \alpha \varkappa \|\mathbf{I} - \gamma \tilde{\Sigma}^{-1} \mathbf{\Pi}\| \|\tilde{\Sigma}^{-1/2} \boldsymbol{\nu}_i\|_2 + \varkappa \beta \left\| \tilde{\Sigma}^{-1/2} \boldsymbol{\eta}_i \right\|_2 \\ &\lesssim \alpha \left(1 + \varkappa (1 + \gamma \lambda_{\Sigma} \rho_{\max}) \right) \cdot \rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \\ &\quad \left[2 \left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + \gamma \|\mathbf{w}_i\|_{\tilde{\Sigma}} + 1 \right] + \\ &\quad \varkappa \beta \cdot \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left\{ \rho_{\max} \left[\left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \|\tilde{\Delta}_i\|_{\tilde{\Sigma}} \right) + 1 \right] + \|\mathbf{w}_i\|_{\tilde{\Sigma}} \right\} \\ &\lesssim \varkappa \beta \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \frac{2}{\varkappa} \|\mathbf{x}_i\|_{\tilde{\Sigma}} + 1 \right) \end{aligned} \quad (128)$$

Control of \tilde{W} and \tilde{B} . With the norm of $\tilde{\zeta}_i$ bounded, we can apply similar techniques as in equations (95), (98), (99) and (102) of Section C.1 to construct the following bound for \tilde{W} :

$$\begin{aligned} \tilde{W} &\leq \|\tilde{\Sigma}\| \sum_{i=0}^{t-1} \|\Psi^{t-i-1}\|^2 \cdot \mathbb{E}_{i-1} \left[\|\zeta_i\|_2^2 \mathbf{1} \left\{ \tilde{\mathcal{H}}_i \right\} \right] \\ &\lesssim \|\tilde{\Sigma}\| \sum_{i=0}^{t-1} \left(1 - \frac{1}{2} \alpha \lambda_{\theta} \right)^{2t-2i-2} \varkappa^2 \beta^2 \rho_{\max}^2 \left(2 \|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + 2\tilde{R} + 1 \right)^2 \\ &\lesssim \frac{\|\tilde{\Sigma}\|}{\alpha \lambda_{\theta}} \varkappa^2 \beta^2 \rho_{\max}^2 \left(\|\tilde{\boldsymbol{\theta}}^*\|_{\tilde{\Sigma}} + \frac{1}{\varkappa} \tilde{R} + 1 \right)^2, \end{aligned} \quad (129)$$

and the following bound for \tilde{B} :

$$\begin{aligned} \tilde{B} &\leq \|\tilde{\Sigma}\| \max_{i:0 \leq i \leq t-1} \left\| \tilde{\Sigma}^{-1/2} \zeta_i \mathbf{1} \left\{ \mathcal{H}_i \right\} \right\|_2 \\ &\lesssim \|\tilde{\Sigma}\| \varkappa \beta \rho_{\max} \max_s \left\{ \phi(s) \tilde{\Sigma}^{-1} \phi(s) \right\} \left(\|\tilde{\boldsymbol{\theta}}^*\|_2 + \tilde{R} + 1 \right) =: \tilde{B}_{\max}. \end{aligned} \quad (130)$$

Invoking the matrix Freedman's inequality. With \tilde{W} and \tilde{B} bounded, we again invoke the matrix Freedman's inequality (Tropp, 2011, Corollary 1.3) to assert that

$$\left\| \sum_{i=0}^{t-1} \Psi^{t-i-1} \tilde{\zeta}_i \right\|_2 \leq 2 \sqrt{\tilde{W}_{\max} \log \frac{2dT}{\delta}} + \frac{4}{3} \tilde{B}_{\max} \log \frac{2dT}{\delta}$$

$$\lesssim \sqrt{\frac{\|\tilde{\Sigma}\|}{\alpha\lambda_\theta} \log \frac{2dT}{\delta}} \varkappa\beta\rho_{\max}(\|\tilde{\theta}^*\|_{\tilde{\Sigma}} + \frac{1}{\varkappa}\tilde{R} + 1) \quad (131)$$

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

C.5 Proof of Lemma 6 and Lemma 7

Proof of Lemma 7: controlling $\|\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}(\mathbf{A} - \widehat{\mathbf{A}})\Sigma^{-\frac{1}{2}}\|$. We intend to invoke the matrix Bernstein inequality to establish the advertised bound [Tropp \(2015\)](#). Note that

$$\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}(\mathbf{A} - \widehat{\mathbf{A}})\Sigma^{-\frac{1}{2}} = \frac{1}{T} \sum_{t=0}^{T-1} \underbrace{\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}(\mathbf{A} - \mathbf{A}_t)\Sigma^{-\frac{1}{2}}}_{=: \mathbf{Z}_t}. \quad (132)$$

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

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

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

$$\begin{aligned} \mathbb{E}[\mathbf{Z}_t\mathbf{Z}_t^\top] &= \mathbb{E}\left[\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}(\mathbf{A} - \mathbf{A}_t)\Sigma^{-1}(\mathbf{A} - \mathbf{A}_t)^\top(\mathbf{A}^\top)^{-1}\Sigma^{\frac{1}{2}}\right] \\ &\preceq \mathbb{E}\left[\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}\mathbf{A}_t\Sigma^{-1}\mathbf{A}_t^\top(\mathbf{A}^\top)^{-1}\Sigma^{\frac{1}{2}}\right] \\ &= \mathbb{E}_{s \sim \mu, s' \sim P(\cdot|s)}\left[\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}\phi(s)(\phi(s) - \gamma\phi(s'))^\top\Sigma^{-1}(\phi(s) - \gamma\phi(s'))\phi(s)^\top(\mathbf{A}^\top)^{-1}\Sigma^{\frac{1}{2}}\right] \\ &\preceq \max_{s, s'} \left\{ (\phi(s) - \gamma\phi(s'))^\top\Sigma^{-1}(\phi(s) - \gamma\phi(s')) \right\} \mathbb{E}_{s \sim \mu} \left[\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}\phi(s)\phi(s)^\top(\mathbf{A}^\top)^{-1}\Sigma^{\frac{1}{2}} \right] \\ &\preceq \max_{s, s'} \left\{ 2\phi(s)^\top\Sigma^{-1}\phi(s) + 2\gamma^2\phi(s')^\top\Sigma^{-1}\phi(s') \right\} \cdot \left\{ \Sigma^{\frac{1}{2}}\mathbf{A}^{-1}\Sigma(\mathbf{A}^\top)^{-1}\Sigma^{\frac{1}{2}} \right\} \\ &\preceq \frac{4 \max_s \phi(s)^\top\Sigma^{-1}\phi(s)}{(1 - \gamma)^2} \mathbf{I}, \end{aligned} \quad (133)$$

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 $\|\mathbb{E}[\mathbf{Z}_t^\top\mathbf{Z}_t]\|$. Similarly, one can obtain

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

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 Σ . 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

$$\phi(s)^\top(\mathbf{A}^\top)^{-1}\Sigma\mathbf{A}^{-1}\phi(s) = \phi(s)^\top\Sigma^{-\frac{1}{2}} \left\{ \Sigma^{\frac{1}{2}}(\mathbf{A}^\top)^{-1}\Sigma\mathbf{A}^{-1}\Sigma^{\frac{1}{2}} \right\} \Sigma^{-\frac{1}{2}}\phi(s)$$

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

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

$$\mathbb{E}[\mathbf{Z}_t^\top \mathbf{Z}_t] \preceq \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{(1-\gamma)^2} \mathbf{I}. \quad (134)$$

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

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

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

$$\|\Sigma^{-\frac{1}{2}}\mathbf{A}_t\Sigma^{-\frac{1}{2}}\| \leq \max_s \|\Sigma^{-\frac{1}{2}}\phi(s)\phi(s)^\top \Sigma^{-\frac{1}{2}}\| + \gamma \max_{s,s'} \|\Sigma^{-\frac{1}{2}}\phi(s')\phi(s)^\top \Sigma^{-\frac{1}{2}}\| \leq 2 \max_s \|\Sigma^{-\frac{1}{2}}\phi(s)\|_2^2. \quad (135)$$

This combined with the preceding bounds yields

$$\|\mathbf{Z}_t\| \leq \frac{2 \max_s \|\Sigma^{-\frac{1}{2}}\phi(s)\|_2^2}{1-\gamma} + 1 \leq \frac{4 \max_s \|\Sigma^{-\frac{1}{2}}\phi(s)\|_2^2}{1-\gamma} = \frac{4 \max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{1-\gamma}. \quad (136)$$

Here, the inequality follows since

$$\max_s \|\Sigma^{-\frac{1}{2}}\phi(s)\|_2^2 \geq \mathbb{E}_{s \sim \mu} [\phi(s)^\top \Sigma^{-1} \phi(s)] = \mathbb{E}_{s \sim \mu} [\text{tr}(\Sigma^{-1} \phi(s)\phi(s)^\top)] = \text{tr}(\mathbf{I}_d) = d \geq 1. \quad (137)$$

Step 4: invoking the matrix Bernstein inequality. With the above bounds in mind, we are ready to apply the matrix Bernstein inequality (Tropp, 2015) to obtain that: with probability at least $1 - \delta$ one has

$$\begin{aligned}
\|\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}(\mathbf{A} - \widehat{\mathbf{A}})\Sigma^{-\frac{1}{2}}\| &\lesssim \sqrt{\frac{1}{T^2} \sum_{t=0}^{T-1} \max\{\|\mathbb{E}[\mathbf{Z}_t \mathbf{Z}_t^\top]\|, \|\mathbb{E}[\mathbf{Z}_t^\top \mathbf{Z}_t]\|\} \log\left(\frac{d}{\delta}\right)} + \frac{\max_t \|\mathbf{Z}_t\| \log\left(\frac{d}{\delta}\right)}{T} \\
&\stackrel{(i)}{\lesssim} \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{T(1-\gamma)^2} \log\left(\frac{d}{\delta}\right)} + \frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s) \log\left(\frac{d}{\delta}\right)}{T(1-\gamma)} \\
&\stackrel{(ii)}{\lesssim} \sqrt{\frac{\max_s \phi(s)^\top \Sigma^{-1} \phi(s)}{T(1-\gamma)^2} \log\left(\frac{d}{\delta}\right)}. \quad (138)
\end{aligned}$$

Here, (i) results from the bounds (133), (134) and (136), 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 $\|\Sigma^{\frac{1}{2}}\mathbf{A}^{-1}(\mathbf{A} - \widehat{\mathbf{A}})\Sigma^{-\frac{1}{2}}\| < 1$. Suppose that $\widehat{\mathbf{A}}$ is not invertible. Given that \mathbf{A} and Σ are both invertible, this means that one can find a unit vectors \mathbf{u} obeying $\mathbf{A}^{-1}\widehat{\mathbf{A}}\Sigma^{-\frac{1}{2}}\mathbf{u} = \mathbf{0}$, which in turn implies

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

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

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

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

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

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

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,

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

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 [Tropp \(2015\)](#) yields

$$\begin{aligned} \|\mathbf{A}^{-1}(\widehat{\mathbf{b}} - \mathbf{b})\|_{\Sigma} &= \left\| \frac{1}{T} \sum_{t=1}^T \mathbf{z}_t \right\|_2 \lesssim \sqrt{\frac{1}{T^2} \sum_{t=0}^{T-1} \mathbb{E}[\mathbf{z}_t^{\top} \mathbf{z}_t] \log\left(\frac{d}{\delta}\right)} + \frac{1}{T} \max_t \|\mathbf{z}_t\|_2 \log\left(\frac{d}{\delta}\right) \\ &\lesssim \frac{\max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2}{1-\gamma} \sqrt{\frac{1}{T} \log\left(\frac{d}{\delta}\right)} + \frac{\max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2}{1-\gamma} \cdot \frac{1}{T} \log\left(\frac{d}{\delta}\right) \\ &\asymp \frac{\max_{s \in \mathcal{S}} \|\Sigma^{-\frac{1}{2}} \phi(s)\|_2}{1-\gamma} \sqrt{\frac{1}{T} \log\left(\frac{d}{\delta}\right)} \end{aligned} \tag{139}$$

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

C.6 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):

$$W := \sum_{i=l}^u \mathbb{E}_{i-1} \left[\left\| \Sigma^{1/2} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}'_i \right\|_2^2 \right] \quad \text{and} \quad B := \max_{i:l \leq i \leq u} \left\| \Sigma^{1/2} \mathbf{A}^{-1} (\mathbf{A}_i - \mathbf{A}) \boldsymbol{\theta}'_i \right\|_2. \quad (140)$$

Control of B . To begin with, observe that

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

where the second to last inequality comes from (136) and the triangle inequality, and the last line is due to the definition of \mathcal{H}_i .

Control of W . Moreover, one can derive

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

where the first inequality arises from (134), 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 (Tropp, 2011, Corollary 1.3) for matrix martingales to demonstrate that

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

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

D Comparisons with previous works

D.1 Comparisons with Srikant and Ying (2019)

Srikant and Ying (2019) 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 Σ -norm with independent samples.

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 Σ -norm of Δ_{t+1} can be expressed as

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

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:

$$\mathbb{E}\|\Delta_{t+1}\|_{\Sigma}^2 = \mathbb{E}\|\Delta_t\|_{\Sigma}^2 - \underbrace{2\eta_t\mathbb{E}[\Delta_t^{\top}(\mathbf{A}^{\top}\Sigma + \Sigma\mathbf{A})\Delta_t]}_{I_1} + \underbrace{2\eta_t^2\mathbb{E}\|\mathbf{A}_t\Delta_t\|_{\Sigma}^2}_{I_2} + \underbrace{2\eta_t^2\mathbb{E}\|\mathbf{A}_t\boldsymbol{\theta}^* - \mathbf{b}_t\|_{\Sigma}^2}_{I_3}. \quad (141)$$

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

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 (79e) 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

$$I_1 \geq 2\eta_t(1 - \gamma)\lambda_{\min}(\Sigma)\mathbb{E}\|\Delta_t\|_{\Sigma}^2. \quad (142)$$

Bounding I_2 . We aim to upper bound I_2 as a function of η_t^2 and $\|\Delta_t\|_{\Sigma}^2$, so that when η_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 κ as the condition number of Σ . Therefore, as long as

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

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

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

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

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

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

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

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

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

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

Consequently, I_3 can be upper bounded by

$$I_3 \leq 6\eta_t^2 \|\boldsymbol{\Sigma}\| (2\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1). \quad (143)$$

Bounding $\mathbb{E} \|\boldsymbol{\Delta}_T\|_{\Sigma}^2$. By combining (141), (142) and (143) and recalling that $I_2 \leq \frac{1}{2} I_1$ when η_t is sufficiently small, we obtain

$$\mathbb{E} \|\boldsymbol{\Delta}_{t+1}\|_{\Sigma}^2 \leq (1 - (1 - \gamma) \lambda_{\min}(\boldsymbol{\Sigma}) \eta_t) \mathbb{E} \|\boldsymbol{\Delta}_t\|_{\Sigma}^2 + 6\eta_t^2 \|\boldsymbol{\Sigma}\| (2\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1). \quad (144)$$

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

$$\mathbb{E} \|\boldsymbol{\Delta}_T\|_{\Sigma}^2 \leq (1 - (1 - \gamma) \lambda_{\min}(\boldsymbol{\Sigma}) \eta)^T \|\boldsymbol{\Delta}_0\|_{\Sigma}^2 + \frac{6\eta \|\boldsymbol{\Sigma}\| (2\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1)}{(1 - \gamma) \lambda_{\min}(\boldsymbol{\Sigma})}.$$

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

$$\frac{\eta \|\boldsymbol{\Sigma}\| (\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1)}{(1 - \gamma) \lambda_{\min}(\boldsymbol{\Sigma})} \lesssim \varepsilon^2; \quad \text{and} \quad \exp(-(1 - \gamma) \lambda_{\min}(\boldsymbol{\Sigma}) \eta T) \|\boldsymbol{\Delta}_0\|_{\Sigma}^2 \lesssim \varepsilon^2.$$

This implies the following upper bound for the sample complexity:

$$T \asymp \frac{\kappa \|\boldsymbol{\Sigma}^{-1}\| (\|\boldsymbol{\theta}^*\|_{\Sigma}^2 + 1)}{(1 - \gamma)^2} \frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon}, \quad (145)$$

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

D.2 Comparisons with Bhandari et al. (2021)

Theorem 2(c) in (Bhandari et al., 2021) shows that with decaying stepsizes $\eta_t = \frac{\beta}{\lambda+t}$ where

$$\beta = \frac{2\|\Sigma^{-1}\|}{(1-\gamma)}, \quad \lambda = \frac{16\|\Sigma^{-1}\|}{(1-\gamma)^2}, \quad (146)$$

the expected ℓ_2 norm of TD estimation error is bounded by

$$\mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_2^2 \leq \frac{\nu}{\lambda + T}, \quad (147)$$

where

$$\nu = \max \left\{ \frac{8\sigma^2\|\Sigma^{-2}\|}{(1-\gamma)^2}, \frac{16\|\boldsymbol{\theta}^*\|_2^2\|\Sigma^{-1}\|}{(1-\gamma)^2} \right\}. \quad (148)$$

- Suppose the maximum is attained at the second term for ν and T is sufficiently large, (147) 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 $\mathbb{E}\|\boldsymbol{\theta}_T - \boldsymbol{\theta}^*\|_2^2 \leq \varepsilon^2$, 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 (148) is larger, (147) 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 $\sigma^2 \asymp \|\boldsymbol{\theta}^*\|_2^2 + 1$. Therefore, the sample complexity implied by Theorem 2(c) of Bhandari et al. (2021) scales as

$$T \asymp \frac{\kappa\|\Sigma^{-1}\|(\|\boldsymbol{\theta}^*\|_2^2 + 1)}{(1-\gamma)^2} \frac{1}{\varepsilon^2}. \quad (149)$$

References

- Agarwal, A., Kakade, S., and Yang, L. F. (2020). Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR.
- Arulkumaran, K., Deisenroth, M. P., Brundage, M., and Bharath, A. A. (2017). A brief survey of deep reinforcement learning. *arXiv preprint arXiv:1708.05866*.

- Baird, L. (1995). Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pages 30–37. Elsevier.
- Bertsekas, D. P. (2017). *Dynamic programming and optimal control (4th edition)*. Athena Scientific.
- Bertsekas, D. P. and Tsitsiklis, J. N. (1995). Neuro-dynamic programming: an overview. In *Proceedings of 1995 34th IEEE conference on decision and control*, volume 1, pages 560–564. IEEE.
- Bertsimas, D., Klasnja, P., Murphy, S., and Na, L. (2022). Data-driven interpretable policy construction for personalized mobile health. In *2022 IEEE International Conference on Digital Health (ICDH)*, pages 13–22. IEEE.
- Bhandari, J., Russo, D., and Singal, R. (2021). A finite time analysis of temporal difference learning with linear function approximation. *Operations Research*, 69(3):950–973.
- Bojinov, I. and Shephard, N. (2019). Time series experiments and causal estimands: exact randomization tests and trading. *Journal of the American Statistical Association*, 114(528):1665–1682.
- Borkar, V. S. (2009). *Stochastic approximation: a dynamical systems viewpoint*, volume 48. Springer.
- Borkar, V. S. and Meyn, S. P. (2000). The ODE method for convergence of stochastic approximation and reinforcement learning. *SIAM Journal on Control and Optimization*, 38(2):447–469.
- Boyan, J. A. (1999). Least-squares temporal difference learning. In *ICML*, pages 49–56.
- Dalal, G., Szorenyi, B., and Thoppe, G. (2020). A tale of two-timescale reinforcement learning with the tightest finite-time bound. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3701–3708.
- Dalal, G., Szörényi, B., Thoppe, G., and Mannor, S. (2018a). Finite sample analyses for TD(0) with function approximation. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Dalal, G., Thoppe, G., Szörényi, B., and Mannor, S. (2018b). Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning. In *Conference On Learning Theory*, pages 1199–1233.
- Dann, C., Neumann, G., Peters, J., et al. (2014). Policy evaluation with temporal differences: A survey and comparison. *Journal of Machine Learning Research*, 15:809–883.
- Duan, Y., Jia, Z., and Wang, M. (2020). Minimax-optimal off-policy evaluation with linear function approximation. In *International Conference on Machine Learning*, pages 2701–2709. PMLR.
- Duan, Y., Wang, M., and Wainwright, M. J. (2021). Optimal policy evaluation using kernel-based temporal difference methods. *arXiv preprint arXiv:2109.12002*.
- Fan, J., Wang, Z., Xie, Y., and Yang, Z. (2020). A theoretical analysis of deep q-learning. In *Learning for Dynamics and Control*, pages 486–489. PMLR.
- Farahmand, A.-m., Ghavamzadeh, M., Szepesvári, C., and Mannor, S. (2016). Regularized policy iteration with nonparametric function spaces. *The Journal of Machine Learning Research*, 17(1):4809–4874.
- Gilbert, E. N. (1952). A comparison of signalling alphabets. *The Bell system technical journal*, 31(3):504–522.
- Gupta, H., Srikant, R., and Ying, L. (2019). Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning. *Advances in Neural Information Processing Systems*, 32.
- Jiang, N. and Li, L. (2016). Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 652–661. PMLR.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. (2018). Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873.

- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR.
- Jin, Y., Yang, Z., and Wang, Z. (2021). Is pessimism provably efficient for offline RL? In *International Conference on Machine Learning*, pages 5084–5096. PMLR.
- Kaledin, M., Moulines, E., Naumov, A., Tadic, V., and Wai, H.-T. (2020). Finite time analysis of linear two-timescale stochastic approximation with Markovian noise. In *Conference on Learning Theory*, pages 2144–2203. PMLR.
- Kallus, N. and Uehara, M. (2020). Double reinforcement learning for efficient off-policy evaluation in markov decision processes. *The Journal of Machine Learning Research*, 21(1):6742–6804.
- Khamaru, K., Pananjady, A., Ruan, F., Wainwright, M. J., and Jordan, M. I. (2020). Is temporal difference learning optimal? an instance-dependent analysis. *arXiv preprint arXiv:2003.07337*.
- Lai, T. L. (2003). Stochastic approximation. *The Annals of Statistics*, 31(2):391–406.
- Lakshminarayanan, C. and Szepesvari, C. (2018). Linear stochastic approximation: How far does constant step-size and iterate averaging go? In *International Conference on Artificial Intelligence and Statistics*, pages 1347–1355. PMLR.
- Li, G., Cai, C., Chen, Y., Wei, Y., and Chi, Y. (2023a). Is Q-learning minimax optimal? a tight sample complexity analysis. *Operations Research*.
- Li, G., Chen, Y., Chi, Y., Gu, Y., and Wei, Y. (2021a). Sample-efficient reinforcement learning is feasible for linearly realizable mdps with limited revisiting. *Advances in Neural Information Processing Systems*, 34:16671–16685.
- Li, G., Shi, L., Chen, Y., Chi, Y., and Wei, Y. (2022). Settling the sample complexity of model-based offline reinforcement learning. *arXiv preprint arXiv:2204.05275*.
- Li, G., Wei, Y., Chi, Y., and Chen, Y. (2023b). Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Operations Research*.
- Li, G., Wei, Y., Chi, Y., Gu, Y., and Chen, Y. (2021b). Sample complexity of asynchronous Q-learning: Sharper analysis and variance reduction. *IEEE Transactions on Information Theory*, 68(1):448–473.
- Ma, C., Zhu, B., Jiao, J., and Wainwright, M. J. (2022). Minimax off-policy evaluation for multi-armed bandits. *IEEE Transactions on Information Theory*, 68(8):5314–5339.
- Mou, W., Li, C. J., Wainwright, M. J., Bartlett, P. L., and Jordan, M. I. (2020). On linear stochastic approximation: Fine-grained Polyak-Ruppert and non-asymptotic concentration. *arXiv preprint arXiv:2004.04719*.
- Moulines, E. and Bach, F. (2011). Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Advances in neural information processing systems*, 24.
- Murphy, S. A. (2003). Optimal dynamic treatment regimes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(2):331–355.
- Nemirovski, A., Juditsky, A., Lan, G., and Shapiro, A. (2009). Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609.
- Pananjady, A. and Wainwright, M. J. (2021). Instance-dependent ℓ_∞ -bounds for policy evaluation in tabular reinforcement learning. *IEEE Transactions on Information Theory*, 67(1):566–585.
- Precup, D. (2000). Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series*, page 80.

- Rashidinejad, P., Zhu, B., Ma, C., Jiao, J., and Russell, S. (2021). Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information Processing Systems*, 34:11702–11716.
- Robbins, H. and Monro, S. (1951). A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407.
- Shi, L., Li, G., Wei, Y., Chen, Y., and Chi, Y. (2022). Pessimistic q-learning for offline reinforcement learning: Towards optimal sample complexity. In *International Conference on Machine Learning*, pages 19967–20025. PMLR.
- Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. (2018). Near-optimal time and sample complexities for solving Markov decision processes with a generative model. In *Advances in Neural Information Processing Systems*, pages 5186–5196.
- Srikant, R. and Ying, L. (2019). Finite-time error bounds for linear stochastic approximation and TD learning. In *Conference on Learning Theory*, pages 2803–2830.
- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44.
- Sutton, R. S. and Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Sutton, R. S., Maei, H. R., Precup, D., Bhatnagar, S., Silver, D., Szepesvári, C., and Wiewiora, E. (2009). Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proceedings of the 26th annual international conference on machine learning*, pages 993–1000.
- Sutton, R. S., Szepesvári, C., and Maei, H. R. (2008). A convergent $O(n)$ algorithm for off-policy temporal-difference learning with linear function approximation. *Advances in neural information processing systems*, 21(21):1609–1616.
- Szepesvári, C. (1998). The asymptotic convergence-rate of Q-learning. In *Advances in Neural Information Processing Systems*, pages 1064–1070.
- Tang, S. and Wiens, J. (2021). Model selection for offline reinforcement learning: Practical considerations for healthcare settings. In *Machine Learning for Healthcare Conference*, pages 2–35. PMLR.
- Thomas, P. and Brunskill, E. (2016). Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 2139–2148. PMLR.
- Tropp, J. (2011). Freedman’s inequality for matrix martingales. *Electronic Communications in Probability*, 16:262–270.
- Tropp, J. A. (2015). An introduction to matrix concentration inequalities. *Found. Trends Mach. Learn.*, 8(1-2):1–230.
- Tsitsiklis, J. and Van Roy, B. (1997). An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5):674–690.
- Tsybakov, A. B. (2009). *Introduction to nonparametric estimation*, volume 11. Springer.
- Wang, B., Yan, Y., and Fan, J. (2021a). Sample-efficient reinforcement learning for linearly-parameterized mdps with a generative model. *Advances in neural information processing systems*, 34:23009–23022.
- Wang, Y., Zou, S., and Zhou, Y. (2021b). Non-asymptotic analysis for two time-scale tdc with general smooth function approximation. *Advances in Neural Information Processing Systems*, 34:9747–9758.
- Wu, Y. F., Zhang, W., Xu, P., and Gu, Q. (2020). A finite-time analysis of two time-scale actor-critic methods. *Advances in Neural Information Processing Systems*, 33:17617–17628.

- Xie, T., Jiang, N., Wang, H., Xiong, C., and Bai, Y. (2021). Policy finetuning: Bridging sample-efficient offline and online reinforcement learning. *Advances in neural information processing systems*, 34:27395–27407.
- Xie, T., Ma, Y., and Wang, Y.-X. (2019). Towards optimal off-policy evaluation for reinforcement learning with marginalized importance sampling. *Advances in Neural Information Processing Systems*, 32.
- Xu, T. and Liang, Y. (2021). Sample complexity bounds for two timescale value-based reinforcement learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pages 811–819. PMLR.
- Xu, T., Zou, S., and Liang, Y. (2019). Two time-scale off-policy TD learning: Non-asymptotic analysis over Markovian samples. In *Advances in Neural Information Processing Systems*, pages 10633–10643.
- Yang, M., Nachum, O., Dai, B., Li, L., and Schuurmans, D. (2020). Off-policy evaluation via the regularized Lagrangian. *Advances in Neural Information Processing Systems*, 33:6551–6561.