# A Two-Timescale Stochastic Algorithm Framework for Bilevel Optimization: Complexity Analysis and Application to Actor-Critic

Mingyi Hong[*†]      Hoi-To Wai [‡]      Zhaoran Wang[§]      Zhuoran Yang[¶]

June 9, 2022

## Abstract

This paper analyzes a two-timescale stochastic algorithm framework for bilevel optimization. Bilevel optimization is a class of problems which exhibits a two-level structure, and its goal is to minimize an outer objective function with variables which are constrained to be the optimal solution to an (inner) optimization problem. We consider the case when the inner problem is unconstrained and strongly convex, while the outer problem is constrained and has a smooth objective function. We propose a two-timescale stochastic approximation (TTSA) algorithm for tackling such a bilevel problem. In the algorithm, a stochastic gradient update with a larger step size is used for the inner problem, while a projected stochastic gradient update with a smaller step size is used for the outer problem. We analyze the convergence rates for the TTSA algorithm under various settings: when the outer problem is strongly convex (resp. weakly convex), the TTSA algorithm finds an $\mathcal{O}(K_{\max}^{-2/3})$-optimal (resp. $\mathcal{O}(K_{\max}^{-2/5})$-stationary) solution, where $K_{\max}$ is the total iteration number. As an application, we show that a two-timescale natural actor-critic proximal policy optimization algorithm can be viewed as a special case of our TTSA framework. Importantly, the natural actor-critic algorithm is shown to converge at a rate of $\mathcal{O}(K_{\max}^{-1/4})$ in terms of the gap in expected discounted reward compared to a global optimal policy.

## 1   Introduction

Consider bilevel optimization problems of the form:

$$\min_{x \in X \subseteq \mathbb{R}^{d_1}} \ell(x) := f(x, y^\star(x)) \quad \text{subject to} \quad y^\star(x) \in \arg\min_{y \in \mathbb{R}^{d_2}} g(x, y), \tag{1}$$

where $d_1, d_2 \geq 1$ are integers; $X$ is a closed and convex subset of $\mathbb{R}^{d_1}$, $f : X \times \mathbb{R}^{d_2} \to \mathbb{R}$ and $g : X \times \mathbb{R}^{d_2} \to \mathbb{R}$ are continuously differentiable functions with respect to (w.r.t.) $x, y$. Problem (1)

---

[*]Authors listed in alphabetical order.

[†]University of Minnesota, `email:mhong@umn.edu`

[‡]The Chinese University of Hong Kong, `email:htwai@se.cuhk.edu.hk`

[§]Northwestern University, `email:zhaoranwang@gmail.com`

[¶]Princeton University, `email:zy6@princeton.edu`

Table 1: Summary of the main results. SC stands for strongly convex, WC for weakly convex, C for convex; $k$ is the iteration counter, $K_{\max}$ is the total number of iterations.

| $\ell(x)$ | Constraint | Step Size $(\alpha_k, \beta_k)$ | Rate (outer) | Rate (inner) |
|---|---|---|---|---|
| SC | $X \subseteq \mathbb{R}^{d_1}$ | $\mathcal{O}(k^{-1}), \ \mathcal{O}(k^{-2/3})$ | $\mathcal{O}(K_{\max}^{-2/3})^\dagger$ | $\mathcal{O}(K_{\max}^{-2/3})^\star$ |
| C | $X \subseteq \mathbb{R}^{d_1}$ | $\mathcal{O}(K_{\max}^{-3/4}), \ \mathcal{O}(K_{\max}^{-1/2})$ | $\mathcal{O}(K_{\max}^{-1/4})^\P$ | $\mathcal{O}(K_{\max}^{-1/2})^\star$ |
| WC | $X \subseteq \mathbb{R}^{d_1}$ | $\mathcal{O}(K_{\max}^{-3/5}), \ \mathcal{O}(K_{\max}^{-2/5})$ | $\mathcal{O}(K_{\max}^{-2/5})^\#$ | $\mathcal{O}(K_{\max}^{-2/5})^\star$ |

$^\dagger$in terms of $\|x^{K_{\max}} - x^\star\|^2$, where $x^\star$ is the optimal solution; $^\star$in terms of $\|y^{K_{\max}} - y^\star(x^{K_{\max}-1})\|^2$, where $y^\star(x^{K_{\max}-1})$ is the optimal inner solution for fixed $x^{K_{\max}-1}$; $^\P$measured using $\ell(x) - \ell(x^\star)$; $^\#$measured using distance to a fixed point with the Moreau proximal map $\widehat{x}(\cdot)$; see (18).

involves two optimization problems following a two-level structure. We refer to $\min_{y \in \mathbb{R}^{d_2}} g(x, y)$ as the *inner problem* whose solution depends on $x$, and $g(x, y)$ is called the *inner objective function*; $\min_{x \in X} \ell(x)$ is referred as the *outer problem*, which represents the outer objective function that we wish to minimize and $\ell(x) \equiv f(x, y^\star(x))$ is called the *outer objective function*. Moreover, both $f, g$ can be stochastic functions whose gradient may be difficult to compute. Despite being a non-convex stochastic problem in general, (1) has a wide range of applications, e.g., reinforcement learning [33], hyperparameter optimization [22], game theory [51], etc..

Tackling (1) is challenging as it involves solving the inner and outer optimization problems *simultaneously*. Even in the simplest case when $\ell(x)$ and $g(x, y)$ are strongly convex in $x$, $y$, respectively, solving (1) is difficult. For instance, if we aim to minimize $\ell(x)$ via a gradient method, at any iterate $x^{\text{cur}} \in \mathbb{R}^{d_1}$ – applying the gradient method for (1) involves a double-loop algorithm that (a) solves the inner optimization problem $y^\star(x^{\text{cur}}) = \arg\min_{y \in \mathbb{R}^{d_2}} g(x^{\text{cur}}, y)$ and then (b) evaluates the gradient as $\nabla \ell(x^{\text{cur}})$ based on the solution $y^\star(x^{\text{cur}})$. Depending on the application, step (a) is usually accomplished by applying yet another gradient method for solving the inner problem (unless a closed form solution for $y^\star(x^{\text{cur}})$ exists). In this way, the resulting algorithm necessitates a double-loop structure.

To this end, [23] and the references therein proposed a stochastic algorithm for (1) involving a *double-loop* update. During the iterations, the inner problem $\min_{y \in \mathbb{R}^{d_2}} g(x^{\text{cur}}, y)$ is solved using a stochastic gradient (SGD) method, with the solution denoted by $\widehat{y}^\star(x^{\text{cur}})$. Then, the outer problem is optimized with an SGD update using estimates of $\nabla f(x^{\text{cur}}, \widehat{y}^\star(x^{\text{cur}}))$. Such a double-loop algorithm is proven to converge to a stationary solution, yet a practical issues lingers: *What if the (stochastic) gradients of the inner and outer problems are only revealed sequentially? For example, when these problems are required to be updated at the same time such as in a sequential game.*

To address the above issues, this paper investigates a *single-loop* stochastic algorithm for (1). Focusing on a class of the bilevel optimization problem (1) where the inner problem is unconstrained and strongly convex, and the outer objective function is smooth, our contributions are three-fold:

- We study a two-timescale stochastic approximation (TTSA) algorithm [7] for the concerned class of bilevel optimization. The TTSA algorithm updates both outer and inner solutions simultaneously, by using some cheap estimates of stochastic gradients of both inner and outer objectives. The algorithm guarantees convergence by improving the inner (resp., outer) solution with a larger (resp., smaller) step size, also known as using a faster (resp., slower) timescale.

2

- We analyze the expected convergence rates of the TTSA algorithm. Our results are summarized in Table 1. Our analysis is accomplished by building a set of *coupled inequalities* for the one-step update in TTSA. For strongly convex outer function, we show inequalities that couple between the outer and inner optimality gaps. For convex or weakly convex outer functions, we establish inequalities coupling between the difference of outer iterates, the optimality of function values, and the inner optimality gap. We also provide new and generic results for solving coupled inequalities. The distinction of timescales between step sizes of the inner and outer updates plays a crucial role in our convergence analysis.

- Finally, we illustrate the application of our analysis results on a two-timescale natural actor-critic policy optimization algorithm with linear function approximation [30, 45]. The natural actor-critic algorithm converges at the rate $\mathcal{O}(K^{-1/4})$ to an optimal policy, which is comparable to the state-of-the-art results.

The rest of this paper is organized as follows. §2 formally describes the problem setting of bilevel optimization and specify the problem class of interest. In addition, the TTSA algorithm is introduced and some application examples are discussed. §3 presents the main convergence results for the generic bilevel optimization problem (1). The convergence analysis is also presented where we highlight the main proof techniques used. Lastly, §4 discusses the application to reinforcement learning. Notice that some technical details of the proof have been relegated to the online appendix [24].

## 1.1 Related Works

The study of bilevel optimization problem (1) can be traced to that of Stackelberg games [51], where the outer (resp. inner) problem optimizes the action taken by a leader (resp. the follower). In the optimization literature, bilevel optimization was introduced in [10] for resource allocation problems, and later studied in [9]. Furthermore, bilevel optimization is a special case of the broader class problem of Mathematical Programming with Equilibrium Constraints [39].

Many related algorithms have been proposed for bilevel optimization. This includes approximate descent methods [19,56], and penalty-based method [26,58]. The approximate descent methods deal with a subclass of problem where the outer problem possesses certain (local) differentiability property, while the penalty-based methods approximate the inner problems and/or the outer problems with an appropriate penalty functions. It is noted in [12] that descent based methods have relatively strong assumptions about the inner problem (such as non-degeneracy), while the penalty based methods are typically slow. Moreover, these works typically focus on asymptotic convergence analysis, without characterizing the convergence rates; see [12] for a comprehensive survey.

In [13,23,27], the authors considered bilevel problems in the (stochastic) unconstrained setting, when the outer problem is non-convex and the inner problem is strongly (or strictly) convex. These works are more related to the algorithms and results to be developed in the current paper. In this case, the (stochastic) gradient of the outer problem may be computed using the chain rule. However, to obtain an accurate estimate, one has to either use *double loop* structure where the inner loop solves the inner sub-problem to a high accuracy [13,23], or use a large batch-size (e.g., $\mathcal{O}(1/\epsilon)$) [27]. Both of these methods could be difficult to implement in practice as the batch-size selection or number of required inner loop iterations are difficult to adjust. In reference [48], the

authors analyzed a special bilevel problem where there is a single optimization variable in both outer and inner levels. The authors proposed a Sequential Averaging Method (SAM) algorithm which can provably solve a problem with strongly convex outer problem, and convex inner problems. Building upon the SAM, [35, 38] developed first-order algorithms for bilevel problem, without requiring that for each fixed outer-level variable, the inner-level solution must be a singleton.

In a different line of recent works, references [36, 50] proposed and analyzed different versions of the so-called truncated back-propagation approach for approximating the (stochastic) gradient of the outer-problem, and established convergence for the respective algorithms. The idea is to use a dynamical system to model an optimization algorithm that solves the inner problem, and then replace the optimal solution of the inner problem by unrolling a few iterations of the updates. However, computing the (hyper-)gradient of the objective function $\ell(x)$ requires using back-propagation through the optimization algorithm, and can be computationally very expensive. It is important to note that none of the methods discussed above have considered *single-loop* stochastic algorithms, in which a small batch of samples are used to approximate the inner and outer gradients at each iteration. Later we will see that the ability of being able to update using a small number of samples for both outer and inner problems is critical in a number of applications, and it is also beneficial numerically.

In contrast to the above mentioned works, this paper considers a TTSA algorithm for stochastic bilevel optimization, which is a single-loop algorithm employing cheap stochastic estimates of the gradient. Notice that TTSA [7] is a class of algorithms designed to solve coupled system of (nonlinear) equations. While its asymptotic convergence property has been well understood, e.g., [7, 8, 32], the convergence rate analysis have been focused on *linear* cases, e.g., [14, 31, 34]. In general, the bilevel optimization problem (1) requires a nonlinear TTSA algorithm. For this case, an asymptotic convergence rate is analyzed in [42] under a restricted form of nonlinearity. For convergence rate analysis, [48] considered a single-loop algorithm for deterministic bilevel optimization with only one variable, and [17] studied the convergence rate when the expected updates are strongly monotone.

Finally, it is worthwhile mentioning that various forms of TTSA have been applied to tackle compositional stochastic optimization [57], policy evaluation methods [6, 54], and actor-critic methods [5, 33, 40]. Notice that some of these optimization problems can be cast as a bilevel optimization, as we will demonstrate next.

**Notations** Unless otherwise noted, $\|\cdot\|$ is the Euclidean norm on finite dimensional Euclidean space. For a twice differentiable function $f : X \times Y \to \mathbb{R}$, $\nabla_x f(x,y)$ (resp. $\nabla_y f(x,y)$) denotes its partial gradient taken w.r.t. $x$ (resp. $y$), and $\nabla^2_{yx} f(x,y)$ (resp. $\nabla^2_{xy} f(x,y)$) denotes the Jacobian of $\nabla_y f(x,y)$ at $y$ (resp. $\nabla_x f(x)$ at $x$). A function $\ell(\cdot)$ is said to be weakly convex with modulus $\mu_\ell \in \mathbb{R}$ if

$$\ell(w) \geq \ell(v) + \langle \nabla \ell(v), w - v \rangle + \mu_\ell \|w - v\|^2, \ \forall \ w, v \in X. \tag{2}$$

Notice that if $\mu_\ell \geq 0$ (resp. $\mu_\ell > 0$), then $\ell(\cdot)$ is convex (resp. strongly convex).

## 2 Two-Timescale Stochastic Approximation Algorithm for (1)

To formally define the problem class of interest, we state the following conditions on the bilevel optimization problem (1).

**Assumption 1.** *The outer functions $f(x,y)$ and $\ell(x) := f(x, y^\star(x))$ satisfy:*

1. *For any $x \in \mathbb{R}^{d_1}$, $\nabla_x f(x, \cdot)$ and $\nabla_y f(x, \cdot)$ are Lipschitz continuous with respect to (w.r.t.) $y \in \mathbb{R}^{d_2}$, and with constants $L_{fx}$ and $L_{fy}$, respectively.*

2. *For any $y \in \mathbb{R}^{d_2}$, $\nabla_y f(\cdot, y)$ is Lipschitz continuous w.r.t. $x \in X$, and with constant $\bar{L}_{fy}$.*

3. *For any $x \in X, y \in \mathbb{R}^{d_2}$, we have $\|\nabla_y f(x,y)\| \leq C_{f_y}$, for some $C_{f_y} > 0$.*

**Assumption 2.** *The inner function $g(x,y)$ satisfies:*

1. *For any $x \in X$ and $y \in \mathbb{R}^{d_2}$, $g(x,y)$ is twice continuously differentiable in $(x,y)$;*

2. *For any $x \in X$, $\nabla_y g(x, \cdot)$ is Lipschitz continuous w.r.t. $y \in \mathbb{R}^{d_2}$, and with constant $L_g$.*

3. *For any $x \in X$, $g(x, \cdot)$ is strongly convex in $y$, and with modulus $\mu_g > 0$.*

4. *For any $x \in X$, $\nabla^2_{xy} g(x, \cdot)$ and $\nabla^2_{yy} g(x, \cdot)$ are Lipschitz continuous w.r.t. $y \in \mathbb{R}^{d_2}$, and with constants $L_{gxy} > 0$ and $L_{gyy} > 0$, respectively.*

5. *For any $y \in \mathbb{R}^m$, $\nabla^2_{xy} g(\cdot, y)$ and $\nabla^2_{yy} g(\cdot, y)$ are Lipschitz continuous w.r.t. $x \in X$, and with constants $\bar{L}_{gxy} > 0$ and $\bar{L}_{gyy} > 0$, respectively.*

6. *For any $x \in X$ and $y \in \mathbb{R}^{d_2}$, we have $\|\nabla^2_{xy} g(x,y)\| \leq C_{gxy}$ for some $C_{gxy} > 0$.*

Basically, Assumption 1, 2 require that the inner and outer functions $f, g$ are well-behaved. In particular, $\nabla_x f$, $\nabla_y f$, $\nabla^2_{xy} g$, and $\nabla^2_{yy} g$ are Lipschitz continuous w.r.t. $x$ when $y$ is fixed, and Lipschitz continuous w.r.t. $y$ when $x$ is fixed. These assumptions are satisfied by common problems in machine learning and optimization, e.g., the application examples discussed in Sec. 2.1.

Our first endeavor is to develop a single-loop stochastic algorithm for tackling (1). Focusing on solutions which satisfy the first-order stationary condition of (1), we aim at finding a pair of solution $(x^\star, y^\star)$ such that

$$\nabla_y g(x^\star, y^\star) = 0, \quad \langle \nabla \ell(x^\star), x - x^\star \rangle \geq 0, \ \forall \ x \in X. \tag{3}$$

Given $x^\star$, a solution $y^\star$ satisfying the first condition in (3) may be found by a cheap stochastic gradient recursion such as $y \leftarrow y - \beta h_g$ with $\mathbb{E}[h_g] = \nabla_y g(x^\star, y)$. On the other hand, given $y^\star(x)$ and suppose that we can obtain a cheap stochastic gradient estimate $h_f$ with $\mathbb{E}[h_f] = \nabla \ell(x) = \overline{\nabla}_x f(x, y^\star(x))$, where $\overline{\nabla}_x f(x,y)$ is a surrogate for $\nabla \ell(x)$ (to be described later), then the second condition can be satisfied by a simple projected stochastic gradient recursion as $x \leftarrow P_X(x - \alpha h_f)$, where $P_X(\cdot)$ denotes the Euclidean projection onto $X$.

A challenge in designing a *single-loop* algorithm for satisfying (3) is to ensure that the outer function's gradient $\overline{\nabla}_x f(x,y)$ is evaluated at an inner solution $y$ that is close to $y^\star(x)$. This led us to develop a *two-timescale stochastic approximation* (TTSA) [7] framework, as summarized in Algorithm 1. An important feature is that the algorithm utilizes two step sizes $\alpha_k$, $\beta_k$ for the outer $(x^k)$, inner $(y^k)$ solution, respectively, designed with different timescales as $\alpha_k/\beta_k \to 0$. As a larger step size is taken to optimize $y^k$, the latter shall stay close to $y^\star(x^k)$. Using this strategy, it is expected that $y^k$ will converge to $y^\star(x^k)$ asymptotically.

**Algorithm 1. Two-Timescale Stochastic Approximation (TTSA)**

**S0)** Initialize the variable $(x^0, y^0) \in X \times \mathbb{R}^{d_2}$ and the step size sequence $\{\alpha_k, \beta_k\}_{k \geq 0}$;
**S1)** For iteration $k = 0, ..., K$,

$$y^{k+1} = y^k - \beta_k \cdot h_g^k, \tag{4a}$$

$$x^{k+1} = P_X(x^k - \alpha_k \cdot h_f^k), \tag{4b}$$

where $h_g^k$, $h_f^k$ are stochastic estimates of $\nabla_y g(x^k, y^k)$, $\overline{\nabla}_x f(x^k, y^{k+1})$ [cf. (6)], respectively, satisfying Assumption 3 given below. Moreover, $P_X(\cdot)$ is the Euclidean projection operator onto the convex set $X$.

Inspired by [23], we provide a method for computing a surrogate of $\nabla \ell(x)$ given $y$ with general objective functions satisfying Assumption 1, 2. Given $y^\star(x)$, we observe that using chain rule, the gradient of $\ell(x)$ can be derived as

$$\nabla \ell(x) = \nabla_x f\big(x, y^\star(x)\big) - \nabla_{xy}^2 g\big(x, y^\star(x)\big)\big[\nabla_{yy}^2 g\big(x, y^\star(x)\big)\big]^{-1} \nabla_y f\big(x, y^\star(x)\big). \tag{5}$$

We note that the computation of the above gradient critically depends on the fact that the inner problem is strongly convex and unconstrained, so that the inverse function theorem can be applied when computing $\nabla y^*(x)$.

We may now define $\overline{\nabla}_x f(x, y)$ as a surrogate of $\nabla \ell(x)$ by replacing $y^\star(x)$ with $y \in \mathbb{R}^{d_2}$:

$$\overline{\nabla}_x f(x, y) := \nabla_x f(x, y) - \nabla_{xy}^2 g(x, y)[\nabla_{yy}^2 g(x, y)]^{-1} \nabla_y f(x, y). \tag{6}$$

Notice that $\nabla \ell(x) = \overline{\nabla}_x f(x, y^\star(x))$. Eq. (6) is a surrogate for $\nabla \ell(x)$ that may be used in TTSA. We emphasize that (6) is not the only construction and the TTSA can accommodate other forms of gradient surrogate. For example, see (41) in the application of our results to actor-critic.

Let $\mathcal{F}_k := \sigma\{y^0, x^0, ..., y^k, x^k\}$, $\mathcal{F}_k' := \sigma\{y^0, x^0 ..., y^k, x^k, y^{k+1}\}$ be the filtration of the random variables up to iteration $k$, where $\sigma\{\cdot\}$ denotes the $\sigma$-algebra generated by the random variables. We consider the following assumption regarding $h_f^k, h_g^k$:

**Assumption 3.** *For any $k \geq 0$, there exist constants $\sigma_g, \sigma_f$, and a nonincreasing sequence $\{b_k\}_{k \geq 0}$ such that:*

$$\mathbb{E}[h_g^k|\mathcal{F}_k] = \nabla_y g(x^k, y^k), \quad \mathbb{E}[h_f^k|\mathcal{F}_k'] = \overline{\nabla}_x f(x^k, y^{k+1}) + B_k, \quad \|B_k\| \leq b_k, \tag{7a}$$

$$\mathbb{E}[\|h_g^k - \nabla_y g(x^k, y^k)\|^2|\mathcal{F}_k] \leq \sigma_g^2 \cdot \{1 + \|\nabla_y g(x^k, y^k)\|^2\}, \tag{7b}$$

$$\mathbb{E}[\|h_f^k - B_k - \overline{\nabla}_x f(x^k, y^{k+1})\|^2|\mathcal{F}_k'] \leq \sigma_f^2. \tag{7c}$$

Notice that the conditions on $h_g^k$ are standard when the latter is taken as a stochastic gradient of $g(x^k, y^k)$, while $h_f^k$ is a potentially biased estimate of $\overline{\nabla}_x f(x^k, y^{k+1})$. As we will see in our convergence analysis, the bias shall decay polynomially to zero.

In light of (6) and as inspired by [23], we suggest to construct a stochastic estimate of $\overline{\nabla}_x f(x^k, y^{k+1})$ as follows. Let $\mathsf{t}_{\max}(k) \geq 1$ be an integer, $\mathsf{c}_h \in (0, 1]$ be a scalar parameter, and denote $x \equiv x^k, y \equiv y^{k+1}$ for brevity. Consider:

1. Select $\mathsf{p} \in \{0, \ldots, \mathsf{t}_{\max}(k) - 1\}$ uniformly at random and draw $2 + \mathsf{p}$ independent samples as

$$\xi^{(1)} \sim \mu^{(1)}, \ \xi_0^{(2)}, \dots, \xi_{\mathsf{p}}^{(2)} \sim \mu^{(2)}.$$

2. Construct the gradient estimator $h_f^k$ as

$$h_f^k = \nabla_x f(x, y; \xi^{(1)}) -$$
$$\nabla_{xy}^2 g(x, y; \xi_0^{(2)}) \left[ \frac{\mathsf{t}_{\max}(k) \, \mathsf{c}_h}{L_g} \prod_{i=1}^{\mathsf{p}} \left( I - \frac{\mathsf{c}_h}{L_g} \nabla_{yy}^2 g(x, y; \xi_i^{(2)}) \right) \right] \nabla_y f(x, y; \xi^{(1)}),$$

where as a convention, we set $\prod_{i=1}^{0} \left( I - \frac{\mathsf{c}_h}{L_g} \nabla_{yy}^2 g(x, y; \xi_i^{(2)}) \right) = I$.

In the above, the distributions $\mu^{(1)}, \mu^{(2)}$ are defined such that they yield unbiased estimate of the gradients/Jacobians/Hessians as:

$$\nabla_x f(x, y) = \mathbb{E}_{\mu^{(1)}}[\nabla_x f(x, y; \xi^{(1)})], \quad \nabla_y f(x, y) = \mathbb{E}_{\mu^{(1)}}[\nabla_y f(x, y; \xi^{(1)})], \tag{8}$$
$$\nabla_{xy}^2 g(x, y) = \mathbb{E}_{\mu^{(2)}}[\nabla_{xy}^2 g(x, y; \xi^{(2)})], \quad \nabla_{yy}^2 g(x, y) = \mathbb{E}_{\mu^{(2)}}[\nabla_{yy}^2 g(x, y; \xi^{(2)})],$$

and satisfying $\mathbb{E}[\|\nabla_y f(x, y; \xi^{(1)})\|^2] \leq C_y$, $\mathbb{E}[\|\nabla_{xy}^2 g(x, y; \xi^{(2)})\|^2] \leq C_g$,

$$\mathbb{E}[\|\nabla_x f(x, y) - \nabla_x f(x, y; \xi^{(1)})\|^2] \leq \sigma_{fx}^2, \quad \mathbb{E}[\|\nabla_y f(x, y) - \nabla_y f(x, y; \xi^{(1)})\|^2] \leq \sigma_{fy}^2, \tag{9}$$
$$\mathbb{E}[\|\nabla_{xy}^2 g(x, y) - \nabla_{xy}^2 g(x, y; \xi^{(2)})\|_2^2] \leq \sigma_{gxy}^2,$$

note that $\|\cdot\|_2$ is the Schatten-2 norm. For convenience of analysis, we assume $\frac{\mu_g}{\mu_g^2 + \sigma_{gxy}^2} \leq 1$, $L_g \geq 1$. The next lemma shows that $h_f^k$ satisfies Assumption 3.

**Lemma 1.** *Under Assumption 1, 2, (8), (9), and* $\mathsf{c}_h = \mu_g/(\mu_g^2 + \sigma_{gxy}^2)$, *then for any* $x \in X, y \in \mathbb{R}^{d_2}$, $\mathsf{t}_{\max}(k) \geq 1$, *it holds that*

$$\left\| \overline{\nabla}_x f(x^k, y^{k+1}) - \mathbb{E}[h_f^k] \right\| \leq C_{gxy} C_{fy} \cdot \frac{1}{\mu_g} \cdot \left( 1 - \frac{\mu_g^2}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \right)^{\mathsf{t}_{\max}(k)}. \tag{10}$$

*Furthermore, the variance is bounded as*

$$\mathbb{E}[\|h_f^k - \mathbb{E}[h_f^k]\|^2] \leq \sigma_{fx}^2 + \left[ (\sigma_{fy}^2 + C_y^2)\{\sigma_{gxy}^2 + 2C_{gxy}^2\} + \sigma_{fy}^2 C_{gxy}^2 \right] \max\left\{ \frac{3}{\mu_g^2}, \frac{3d_1/L_g}{\mu_g^2 + \sigma_{gxy}^2} \right\}. \tag{11}$$

The proof of the above lemma is relegated to our online appendix, see §E in [24]. Note that the variance bound (11) relies on analyzing the expected norm of product of random matrices using the techniques inspired by [18, 25]. Finally, observe that the upper bounds in (10), (11) correspond to $b_k$, $\sigma_f^2$ in Assumption 3, respectively, and the requirements on $h_f^k$ are satisfied.

To further understand the property of the TTSA algorithm with (6), we borrow the following results from [23] on the Lipschitz continuity of the maps $\nabla\ell(x), y^\star(x)$:

**Lemma 2.** *[23, Lemma 2.2] Under Assumption 1, 2, it holds*

$$\|\overline{\nabla}_x f(x, y) - \nabla\ell(x)\| \leq L\|y^\star(x) - y\|, \quad \|y^\star(x_1) - y^\star(x_2)\| \leq L_y\|x_1 - x_2\|, \tag{12a}$$
$$\|\nabla\ell(x_1) - \nabla\ell(x_2)\| = \|\nabla f(x_1, y^\star(x_1)) - \nabla f(x_2, y^\star(x_2))\| \leq L_f\|x_1 - x_2\|. \tag{12b}$$

*for any $x, x_1, x_2 \in X$ and $y \in \mathbb{R}^{d_2}$, where we have defined*

$$L := L_{f_x} + \frac{L_{f_y} C_{g_{xy}}}{\mu_g} + C_{f_y} \left( \frac{L_{g_{xy}}}{\mu_g} + \frac{L_{g_{yy}} C_{g_{xy}}}{\mu_g^2} \right),$$

$$L_f := L_{f_x} + \frac{(\bar{L}_{f_y} + L) C_{g_{xy}}}{\mu_g} + C_{f_y} \left( \frac{\bar{L}_{g_{xy}}}{\mu_g} + \frac{\bar{L}_{g_{yy}} C_{g_{xy}}}{\mu_g^2} \right), \quad L_y = \frac{C_{g_{xy}}}{\mu_g}. \tag{13}$$

The above properties will be pivotal in establishing the convergence of TTSA. First, we note that (12b) implies that the composite function $\ell(x)$ is weakly convex with a modulus that is at least $(-L_f)$. Furthermore, (7c) in Assumption 3 combined with Lemma 2 leads to the following estimate:

$$\mathbb{E}[\|h_f^k\|^2 | \mathcal{F}_k'] \leq \widetilde{\sigma}_f^2 + 3b_k^2 + 3L^2 \|y^{k+1} - y^\star(x^k)\|^2, \quad \widetilde{\sigma}_f^2 := \sigma_f^2 + 3 \sup_{x \in X} \|\nabla \ell(x)\|^2. \tag{14}$$

Throughout, we assume $\widetilde{\sigma}_f^2$ is bounded, e.g., it can be satisfied if $X$ is bounded, or if $\ell(x)$ has bounded gradient.

## 2.1 Applications

Practical problems such as hyperparameter optimization [22, 41, 50], Stackelberg games [51] can be cast into special cases of bilevel optimization problem (1). To be specific, we discuss three applications of the bilevel optimization problem (1) below.

**Model-Agnostic Meta-Learning** An important paradigm of machine learning is to find model that adapts to multiple training sets in order to achieve the best performance for individual tasks. Among others, a popular formulation is model-agnostic meta learning (MAML) [20] which minimizes an outer objective of empirical risk on all training sets, and the inner objective is the one-step projected gradient. Let $D^{(j)} = \{z_i^{(j)}\}_{i=1}^n$ be the $j$-th ($j \in [J]$) training set with sample size $n$, MAML can be formulated as a bilevel optimization problem [47]:

$$\min_{\theta \in \Theta} \quad \sum_{j=1}^J \sum_{i=1}^n \bar{\ell}(\theta^{\star(j)}(\theta), z_i^{(j)})$$

$$\text{subject to} \quad \theta^{\star(j)}(\theta) \in \arg\min_{\theta^{(j)}} \left\{ \sum_{i=1}^n \langle \theta^{(j)}, \nabla_\theta \bar{\ell}(\theta, z_i^{(j)}) \rangle + \frac{\lambda}{2} \|\theta^{(j)} - \theta\|^2 \right\}. \tag{15}$$

Here $\theta$ is the shared model parameter, $\theta^{(j)}$ is the adaptation of $\theta$ to the $j$th training set, and $\bar{\ell}$ is the loss function. It can be checked that the inner problem is strongly convex. We have Assumption 1, 2, 3 for stochastic gradient updates, assuming $\bar{\ell}$ is sufficiently regular, and the losses are the logistic loss. Moreover, [21] proved that, assuming $\lambda$ is sufficiently large and $\bar{\ell}$ is strongly convex, the outer problem is also strongly convex. In fact, [46] demonstrated that an algorithm with no inner loop achieves a comparable performance to [21].

**Policy Optimization** Another application of the bilevel optimization problem is the policy optimization problem, particularly when combined with an actor-critic scheme. The optimization involved is to find an optimal policy to maximize the expected (discounted) reward. Here, the 'actor' serves as the outer problem and the 'critic' serves as the inner problem which evaluates the

performance of the 'actor' (current policy). To avoid redundancy, we refer our readers to §4 where we present a detailed case study. The latter will also shed lights on the generality of our proof techniques for TTSA algorithms.

**Data hyper-cleaning**    The data hyper-cleaning problem trains a classifier with a dataset of randomly corrupted labels [50]. The problem formulation is given below:

$$\min_{x \in \mathbb{R}^{d_1}} \quad \ell(x) := \sum_{i \in \mathcal{D}_{\text{val}}} L(a_i^\top y^\star(x), b_i) \tag{16}$$
$$\text{s.t.} \quad y^\star(x) = \arg\min_{y \in \mathbb{R}^{d_2}} \left\{ \lambda \|y\|^2 + \sum_{i \in \mathcal{D}_{\text{tr}}} \sigma(x_i) L(a_i^\top y, b_i) \right\}.$$

In this problem, we have $d_1 = |\mathcal{D}_{\text{tr}}|$, and $d_2$ is the dimension of the classifier $y$; $(a_i, b_i)$ is the $i$th data point; $L(\cdot)$ is the loss function; $x_i$ is the parameter that determines the weight for the $i$th data sample; $\sigma : \mathbb{R} \to \mathbb{R}_+$ is the weight function; $\lambda > 0$ is a regularization parameter; $\mathcal{D}_{\text{val}}$ and $\mathcal{D}_{\text{tr}}$ are validation and training sets, respectively. Here, the inner problem finds the classifier $y^*(x)$ with the training set $\mathcal{D}_{\text{tr}}$, while the outer problem finds the best weights $x$ with respect to the validation set $\mathcal{D}_{\text{val}}$.

Before ending this subsection, let us mention that, we do not aware any general sufficient conditions that can be used to verify whether the outer function $\ell(x)$ is (strongly) convex or not. To our knowledge, the convexity of $\ell(x)$ has to be verified in a case-by-case manner; please see [21, Appendix B3] for how this can be done.

# 3    Main Results

This section presents the convergence results for TTSA algorithm for (1). We first summarize a list of important constants in Table 2 to be used in the forthcoming analysis. Next, we discuss a few concepts pivotal to our analysis.

**Tracking Error**    TTSA tackles the inner and outer problems simultaneously using single-loop updates. Due to the coupled nature of the inner and outer problems, in order to obtain an upper bound on the *optimality gap* $\Delta_x^k := \mathbb{E}[\|x^k - x^\star\|^2]$, where $x^\star$ is an optimal solution to (1), we need to estimate the *tracking error* defined as

$$\Delta_y^k := \mathbb{E}[\|y^k - y^\star(x^{k-1})\|^2] \quad \text{where} \quad y^\star(x) = \arg\min_{y \in \mathbb{R}^{d_2}} g(x, y). \tag{17}$$

For any $x \in X$, $y^\star(x)$ is well defined since the inner problem is strongly convex due to Assumption 2. By definition, $\Delta_y^k$ quantifies how close $y^k$ is from the optimal solution to inner problem given $x^{k-1}$.

**Moreau Envelop**    Fix $\rho > 0$, define the Moreau envelop and proximal map as

$$\Phi_{1/\rho}(z) := \min_{x \in X} \left\{ \ell(x) + (\rho/2)\|x - z\|^2 \right\}, \quad \widehat{x}(z) := \arg\min_{x \in X} \left\{ \ell(x) + (\rho/2)\|x - z\|^2 \right\}. \tag{18}$$

For any $\epsilon > 0$, $x^k \in X$ is said to be an $\epsilon$-*nearly stationary solution* [16] if $x^k$ is an approximate fixed point of $\{\widehat{x} - \mathrm{I}\}(\cdot)$, where

$$\widetilde{\Delta}_x^k := \mathbb{E}[\|\widehat{x}(x^k) - x^k\|^2] \leq \rho^{-2} \cdot \epsilon. \tag{19}$$

Table 2: Summary of the Constants for Section 3

| Constant | Description | Reference |
|---|---|---|
| $L_{fx}, L_{fy}$ | Lipschitz constants for $\nabla_x f(x,\cdot)$, $\nabla_y f(x,\cdot)$ w.r.t. $y$, resp. | Assumption 1 |
| $\bar{L}_{fy}$ | Lipschitz constants for $\nabla_y f(\cdot,y)$ w.r.t. $x$ | Assumption 1 |
| $C_{fy}$ | Upper bound on $\|\nabla_y f(x,y)\|$ | Assumption 1 |
| $L_g$ | Lipschitz constant of $\nabla_y g(x,\cdot)$ | Assumption 2 |
| $\mu_g$ | Strong convexity modulus $g(x,\cdot)$ w.r.t. $y$ | Assumption 2 |
| $L_{gxy}, L_{gyy}$ | Lipschitz constants of $\nabla^2_{xy} g(x,\cdot), \nabla^2_{yy} g(x,\cdot)$ w.r.t. $y$, resp. | Assumption 2 |
| $\bar{L}_{gxy}, \bar{L}_{gyy}$ | Lipschitz constants of $\nabla^2_{xy} g(\cdot,y), \nabla^2_{yy} g(\cdot,y)$ w.r.t. $x$, resp. | Assumption 2 |
| $C_{gxy}$ | Upper bound on $\|\nabla^2_{xy} g(x,y)\|$ | Assumption 2 |
| $b_k$ | Bound on the bias of $h_f^k$ at iteration $k$ | Assumption 3 |
| $\sigma_g^2, \sigma_f^2$ | Variance of stochastic estimates $h_g^k$, $h_f^k$, resp. | Assumption 3 |
| $\widetilde{\sigma}_f^2$ | Constant term on the bound for $\mathbb{E}[\|h_f^k\|^2]$ | (14) |
| $L$ | Difference between $\overline{\nabla}_x f(x,y)$, $\nabla\ell(x)$ w.r.t. $\|y^\star(x)-x\|$ | Lemma 2 |
| $L_y$ | Lipschitz constant of $y^\star(x)$ | Lemma 2 |
| $L_f$ | Lipschitz constant of $\nabla\ell(x)$ | Lemma 2 |

We observe that if $\epsilon = 0$, then $x^k \in X$ is a stationary solution to (1) satisfying the second condition in (3). As we will demonstrate next, the *near-stationarity* condition (19) provides an apparatus to quantify the finite-time convergence of TTSA in the case when $\ell(x)$ is non-convex.

## 3.1 Strongly Convex Outer Objective Function

Our first result considers the instance of (1) where $\ell(x)$ is strongly convex. We obtain:

**Theorem 1.** *Under Assumption 1, 2, 3. Assume that $\ell(x)$ is weakly convex with a modulus $\mu_\ell > 0$ (i.e., it is strongly convex), and the step sizes satisfy*

$$\alpha_k \le c_0 \beta_k^{3/2}, \ \beta_k \le c_1 \alpha_k^{2/3}, \ \frac{\beta_{k-1}}{\beta_k} \le 1 + \beta_k \mu_g/8, \ \frac{\alpha_{k-1}}{\alpha_k} \le 1 + 3\alpha_k \mu_\ell/4, \tag{20a}$$

$$\alpha_k \le \frac{1}{\mu_\ell}, \ \beta_k \le \min\left\{\frac{1}{\mu_g}, \frac{\mu_g}{L_g^2(1+\sigma_g^2)}, \frac{\mu_g^2}{48c_0^2 L^2 L_y^2}\right\}, \ 8\mu_\ell \alpha_k \le \mu_g \beta_k, \ \forall \ k \ge 0, \tag{20b}$$

*where the constants $L, L_y$ were defined in Lemma 2 and $c_0, c_1 > 0$ are free parameters. If the bias is bounded as $b_k^2 \le \widetilde{c}_b \alpha_{k+1}$, then for any $k \ge 1$, the TTSA iterates satisfy*

$$\Delta_x^k \lesssim \prod_{i=0}^{k-1}(1-\alpha_i\mu_\ell)\Big[\Delta_x^0 + \frac{L^2}{\mu_\ell^2}\Delta_y^0\Big] + \frac{c_1 L^2}{\mu_\ell^2}\Big[\frac{\sigma_g^2}{\mu_g} + \frac{c_0^2 L_y^2}{\mu_g^2}\widetilde{\sigma}_f^2\Big]\alpha_{k-1}^{2/3},$$

$$\Delta_y^k \lesssim \prod_{i=0}^{k-1}(1-\beta_i\mu_g/4)\Delta_y^0 + \Big[\frac{\sigma_g^2}{\mu_g} + \frac{c_0^2 L_y^2}{\mu_g^2}\widetilde{\sigma}_f^2\Big]\beta_{k-1}, \tag{21}$$

*where the symbol $\lesssim$ denotes that the numerical constants are omitted (see Section 3.3).*

Notice that the bounds in (21) show that the expected optimality gap and tracking error at the $k$th iteration shall compose of a transient and fluctuation terms. For instance, for the bound on $\Delta_x^k$, the first (transient) term decays sub-geometrically as $\prod_{i=0}^{k-1}(1-\alpha_i\mu_\ell)$, while the second (fluctuation)

term scales as $\alpha_{k-1}^{2/3}$. Note that if $\alpha_k, \beta_k \to 0$, then the r.h.s. in (21) converges to zero as $k \to \infty$. While for non-vanishing step sizes, the r.h.s. in (21) may not converge to zero.

The conditions in (20) are satisfied by both *diminishing* and *constant* step sizes. For example, we define the constants:

$$k_\alpha = \max\left\{35\left(\frac{L_g}{\mu_g}\right)^3(1+\sigma_g^2)^{\frac{3}{2}}, \frac{(512)^{\frac{3}{2}}L^2L_y^2}{\mu_\ell^2}\right\}, \; c_\alpha = \frac{8}{3\mu_\ell}, \; k_\beta = \frac{1}{4}k_\alpha, \; c_\beta = \frac{32}{3\mu_g}. \tag{22}$$

Then, for *diminishing step sizes*, we set $\alpha_k = c_\alpha/(k + k_\alpha), \beta_k = c_\beta/(k + k_\beta)^{2/3}$, and for *constant step sizes*, we set $\alpha_k = c_\alpha/k_\alpha, \beta_k = c_\beta/k_\beta^{2/3}$. Both pairs of the step sizes satisfy (20) with $c_0 = \frac{\mu_g^{3/2}}{\mu_\ell}$, $c_1 = 10\frac{\mu_\ell^{2/3}}{\mu_g}$. For *diminishing step sizes*, Theorem 1 shows the last iterate convergence rate for the optimality gap and the tracking error to be $\mathcal{O}(k^{-2/3})$. To compute an $\epsilon$-optimal solution with $\Delta_x^k \le \epsilon$, the TTSA algorithm with diminishing step size requires a total of $\mathcal{O}(\log(1/\epsilon)/\epsilon^{3/2})$ calls of stochastic (gradient/Hessian/Jacobian) oracles of both outer ($f(\cdot, \cdot)$) and inner ($g(\cdot, \cdot)$) functions[1].

While this is arguably the easiest case of (1), we notice that the double-loop algorithm in [23] requires $\mathcal{O}(1/\epsilon)$, $\mathcal{O}(1/\epsilon^2)$ stochastic oracles for the outer (i.e., $f(\cdot, \cdot)$), inner (i.e., $g(\cdot, \cdot)$) functions, respectively. As such, the TTSA algorithm requires less number of stochastic oracles for the inner function.

## 3.2 Smooth (Possibly Non-convex) Outer Objective Function

We focus on the case where $\ell(x)$ is weakly convex. We obtain

**Theorem 2.** *Under Assumption 1, 2, 3, assume that $\ell(\cdot)$ is weakly convex with modulus $\mu_\ell \in \mathbb{R}$. Let $K_{\mathsf{max}} \ge 1$ be the maximum iteration number and set*

$$\alpha = \min\left\{\frac{\mu_g^2}{8L_yLL_g^2(1+\sigma_g^2)}, \frac{1}{4L_yL}K_{\mathsf{max}}^{-3/5}\right\}, \;\; \beta = \min\left\{\frac{\mu_g}{L_g^2(1+\sigma_g^2)}, \frac{2}{\mu_g}K_{\mathsf{max}}^{-2/5}\right\}. \tag{23}$$

*If $b_k^2 \le \alpha$, then for any $K_{\mathsf{max}} \ge 1$, the iterates from the TTSA algorithm satisfy*

$$\mathbb{E}[\widetilde{\Delta}_x^{\mathsf{K}}] = \mathcal{O}(K_{\mathsf{max}}^{-2/5}), \;\; \mathbb{E}[\Delta_y^{\mathsf{K}+1}] = \mathcal{O}(K_{\mathsf{max}}^{-2/5}), \tag{24}$$

*where $\mathsf{K}$ is an independent uniformly distributed random variable on $\{0, ..., K_{\mathsf{max}} - 1\}$; and we recall $\widetilde{\Delta}_x^k := \|\widehat{x}(x^k) - x^k\|^2$. When $K_{\mathsf{max}}$ is large and $\mu_\ell < 0$, setting $\rho = 2|\mu_\ell|$ yields*

$$\mathbb{E}[\widetilde{\Delta}_x^{\mathsf{K}}] \lesssim \left[L^2\left(\Delta^0 + \frac{\sigma_g^2}{\mu_g^2}\right) + \mu_g\widetilde{\sigma}_f^2\right]\frac{K_{\mathsf{max}}^{-\frac{2}{5}}}{|\mu_\ell|^2}, \;\; \mathbb{E}[\Delta_y^{\mathsf{K}+1}] \lesssim \left[\frac{\Delta^0}{\mu_g} + \frac{\sigma_g^2}{\mu_g^2} + \frac{\mu_g\sigma_f^2}{L^2}\right]K_{\mathsf{max}}^{-\frac{2}{5}}, \tag{25}$$

*where we defined $\Delta^0 := \max\{\Phi_{1/\rho}(x^0), \frac{L_y}{L}\mathrm{OPT}^0, \Delta_y^0\}$, and used the conditions $\alpha < 1$, $\mu_g \ll 1$; the symbol $\lesssim$ denotes that the numerical constants are omitted (see Section 3.3).*

The above result uses constant step sizes determined by the maximum number of iterations, $K_{\mathsf{max}}$. Here we set the step sizes as $\alpha_k \asymp K_{\mathsf{max}}^{-3/5}$ and $\beta_k \asymp K_{\mathsf{max}}^{-2/5}$. Similar to the previous case of strongly convex outer objective function, $\alpha_k/\beta_k$ converges to zero as $K_{\mathsf{max}}$ goes to infinity.

---

[1]Notice that as we need $b_k = \mathcal{O}(\sqrt{\alpha_{k+1}})$, from Lemma 1, the polynomial bias decay requires using $\mathsf{t}_{\mathsf{max}}(k) = \mathcal{O}(1 + \log(k))$ samples per iteration, justifying the log factor in the bound.

Nevertheless, Theorem 2 shows that TTSA requires $\mathcal{O}(\epsilon^{-5/2}\log(1/\epsilon))$ calls of stochastic oracle for sampled gradient/Hessian to find an $\epsilon$-nearly stationary solution. In addition, it is worth noting that when $X = \mathbb{R}^{d_1}$, Theorem 2 implies that TTSA achieves $\mathbb{E}[\|\nabla\ell(x^{\mathsf{K}})\|^2] = \mathcal{O}(K_{\mathsf{max}}^{-2/5})$ [16, Sec. 2.2], i.e., $x^{\mathsf{K}}$ is an $\mathcal{O}(K_{\mathsf{max}}^{-2/5})$-approximate (near) stationary point of $\ell(x)$ in expectation.

Let us compare our sampling complexity bounds to the double loop algorithm in [23], which requires $\mathcal{O}(\epsilon^{-3})$ (resp. $\mathcal{O}(\epsilon^{-2})$) stochastic oracle calls for the inner problem (resp. outer problem), to reach an $\epsilon$-stationary solution. The sample complexity of TTSA yields a tradeoff for the inner and outer stochastic oracles. We also observe that a trivial extension to a single-loop algorithm results in a *constant* error bound[2]. Finally, we can extend Theorem 2 to the case where $\ell(\cdot)$ is a convex function.

**Corollary 1.** *Under Assumption 1, 2, 3 and assume that $\ell(x)$ is weakly convex with modulus $\mu_\ell \geq 0$. Consider* (1) *with $X \subseteq \mathbb{R}^{d_1}$, $D_x = \sup_{x,x'\in X}\|x-x'\| < \infty$. Let $K_{\mathsf{max}} \geq 1$ be the maximum iteration number and set*

$$\alpha = \min\Big\{\frac{\mu_g^2}{8L_y L L_g^2(1+\sigma_g^2)}, \frac{1}{4L_y L}K_{\mathsf{max}}^{-3/4}\Big\}, \quad \beta = \min\Big\{\frac{\mu_g}{L_g^2(1+\sigma_g^2)}, \frac{2}{\mu_g}K_{\mathsf{max}}^{-1/2}\Big\}. \quad (26)$$

*If $b_k \leq c_b K_{\mathsf{max}}^{-1/4}$, then for large $K$, the TTSA algorithm satisfies*

$$\mathbb{E}[\ell(x^{\mathsf{K}}) - \ell(x^\star)] = \mathcal{O}(K_{\mathsf{max}}^{-1/4}), \quad \mathbb{E}[\Delta_y^{\mathsf{K}+1}] = \mathcal{O}(K_{\mathsf{max}}^{-1/2}), \quad (27)$$

*where $\mathsf{K}$ is an independent uniform random variable on $\{0,...,K_{\mathsf{max}}-1\}$. By convexity, the above implies $\mathbb{E}[\ell(\frac{1}{K_{\mathsf{max}}}\sum_{k=1}^{K_{\mathsf{max}}} x^k) - \ell(x^\star)] = \mathcal{O}(K_{\mathsf{max}}^{-1/4})$.*

From Corollary 1, the TTSA algorithm requires $\mathcal{O}(\epsilon^{-4}\log(1/\epsilon))$ stochastic oracle calls to find an $\epsilon$-optimal solution (in terms of the optimality gap defined with objective values). This is comparable to the complexity bounds in [23], which requires $\mathcal{O}(\epsilon^{-4}\log(1/\epsilon))$ (resp. $\mathcal{O}(\epsilon^{-2})$) stochastic oracle calls for the inner problem (resp. outer problem). Additionally, we mention that the constant $D_x$, which represents the diameter of the constraint set, appears in the constant of the convergence bounds, therefore it is omitted in the big-O notation in (26). For details, please see the proof in Appendix B.4.

## 3.3 Convergence Analysis

We now present the proofs for Theorem 1, 2. The proof for Corollary 1 is similar to that of Theorem 2. Due to the space limitation, we refer the readers to [24]. We highlight that the proofs of both theorems rely on the similar ideas of tackling coupled inequalities.

**Proof of Theorem 1** Our proof relies on bounding the optimality gap and tracking error *coupled with each other*. First we derive the convergence of the inner problem.

**Lemma 3.** *Under Assumption 1, 2, 3. Suppose that the step size satisfies* (20a), (20b). *For any*

---

[2]To see this, the readers are referred to [23, Theorem 3.1]. If a single inner iteration is performed, $t_k = 1$, so $\bar{A}^k \geq \|y^0 - y^\star(x^k)\|$ which is a constant. Then the r.h.s. of (3.70), (3.73), (3.74) in [23, Theorem 3.1] will all have a constant term.

$k \geq 1$, it holds that

$$\Delta_y^{k+1} \leq \prod_{\ell=0}^{k}(1 - \beta_\ell \mu_g/2)\, \Delta_y^0 + \frac{8}{\mu_g}\Big\{\sigma_g^2 + \frac{4c_0^2 L_y^2}{\mu_g}\big[\widetilde{\sigma}_f^2 + 3b_0^2\big]\Big\}\beta_k. \tag{28}$$

Notice that the bound in (28) relies on the strong convexity of the inner problem and the Lipschitz properties established in Lemma 2 for $y^\star(x)$; see §A.1. We emphasize that the step size condition $\alpha_k \leq c_0 \beta_k^{3/2}$ is crucial in establishing the above bound. As the second step, we bound the convergence of the outer problem.

**Lemma 4.** *Under Assumption 1, 2, 3. Assume that the bias satisfies $b_k^2 \leq \widetilde{c}_b \alpha_{k+1}$. With (20a), (20b), for any $k \geq 1$, it holds that*

$$\Delta_x^{k+1} \leq \prod_{\ell=0}^{k}(1 - \alpha_\ell \mu_\ell)\Delta_x^0 + \Big[\frac{4\widetilde{c}_b}{\mu_\ell^2} + \frac{2\widetilde{\sigma}_f^2 + 6b_0^2}{\mu_\ell}\Big]\alpha_k$$
$$+ \Big[\frac{2L^2}{\mu_\ell} + 3\alpha_0 L^2\Big]\sum_{j=0}^{k}\alpha_j \prod_{\ell=j+1}^{k}(1 - \alpha_\ell \mu_\ell)\Delta_y^{j+1}. \tag{29}$$

see §A.2. We observe that (28), (29) lead to a pair of coupled inequalities. To compute the final bound in the theorem, we substitute (28) into (29). As $\Delta_y^{j+1} = \mathcal{O}(\beta_j) = \mathcal{O}(\alpha_j^{2/3})$, the dominating term in (29) can be estimated as

$$\sum_{j=0}^{k}\alpha_j \prod_{\ell=j+1}^{k}(1 - \alpha_\ell \mu_\ell)\Delta_y^{j+1} = \sum_{j=0}^{k}\mathcal{O}(\alpha_j^{5/3})\prod_{\ell=j+1}^{k}(1 - \alpha_\ell \mu_\ell) = \mathcal{O}(\alpha_k^{2/3}), \tag{30}$$

yielding the desirable rates in the theorem. See §A.3 for details.

**Proof of Theorem 2** Without strong convexity in the outer problem, the analysis becomes more challenging. To this end, we first develop the following lemma on coupled inequalities with numerical sequences, which will be pivotal to our analysis:

**Lemma 5.** *Let $K \geq 1$ be an integer. Consider sequences of non-negative scalars $\{\Omega^k\}_{k=0}^K$, $\{\Upsilon^k\}_{k=0}^K$, $\{\Theta^k\}_{k=0}^K$. Let $c_0, c_1, c_2, d_0, d_1, d_2$ be some positive constants. If the recursion holds*

$$\Omega^{k+1} \leq \Omega^k - c_0 \Theta^{k+1} + c_1 \Upsilon^{k+1} + c_2, \quad \Upsilon^{k+1} \leq (1 - d_0)\Upsilon^k + d_1 \Theta^k + d_2, \tag{31}$$

*for any $k \geq 0$. Then provided that $c_0 - c_1 d_1 (d_0)^{-1} > 0, d_0 - d_1 c_1 (c_0)^{-1} > 0$, it holds*

$$\frac{1}{K}\sum_{k=1}^{K}\Theta^k \leq \frac{\Omega^0 + \frac{c_1}{d_0}\big(\Upsilon^0 + d_1 \Theta^0 + d_2\big)}{\big(c_0 - c_1 d_1 (d_0)^{-1}\big)K} + \frac{c_2 + c_1 d_2 (d_0)^{-1}}{c_0 - c_1 d_1 (d_0)^{-1}}$$
$$\frac{1}{K}\sum_{k=1}^{K}\Upsilon^k \leq \frac{\Upsilon^0 + d_1 \Theta^0 + d_2 + \frac{d_1}{c_0}\Omega^0}{\big(d_0 - d_1 c_1 (c_0)^{-1}\big)K} + \frac{d_2 + d_1 c_2 (c_0)^{-1}}{d_0 - d_1 c_1 (c_0)^{-1}}. \tag{32}$$

The proof of the above lemma is simple and is relegated to §B.1.

We demonstrate that stationarity measures of the TTSA iterates satisfy (31). The conditions $c_0 - c_1 d_1 (d_0)^{-1} > 0, d_0 - d_1 c_1 (c_0)^{-1} > 0$ impose constraints on the step sizes and (32) leads to a

13

finite-time bound on the convergence of TTSA. To begin our derivation of Theorem 2, we observe the following coupled descent lemma:

**Lemma 6.** *Under Assumption 1, 2, 3. If $\mu_g \beta/2 < 1$, $\beta L_g^2(1 + \sigma_g^2) \leq \mu_g$, then the following inequalities hold for any $k \geq 0$:*

$$\mathrm{OPT}^{k+1} \leq \mathrm{OPT}^k - \frac{1 - \alpha L_f}{2\alpha}\mathbb{E}[\|x^{k+1} - x^k\|^2] + \alpha\big[2L^2\Delta_y^{k+1} + 2b_0^2 + \sigma_f^2\big] \tag{33a}$$

$$\Delta_y^{k+1} \leq \big(1 - \mu_g\beta/2\big)\Delta_y^k + \Big(\frac{2}{\mu_g\beta} - 1\Big)L_y^2 \cdot \mathbb{E}[\|x^k - x^{k-1}\|^2] + \beta^2\sigma_g^2. \tag{33b}$$

The proof of (33a) is due to the smoothness of outer function $\ell(\cdot)$ established in Lemma 2, while (33b) follows from the strong convexity of the inner problem. See the details in §B.2. Note that (33a), (33b) together is a special case of (31) with:

$$\Omega^k = \mathrm{OPT}^k, \ \Theta^k = \mathbb{E}[\|x^k - x^{k-1}\|^2], \ c_0 = \frac{1}{2\alpha} - \frac{L_f}{2}, \ c_1 = 2\alpha L^2, \ c_2 = \alpha(2b_0^2 + \sigma_f^2),$$
$$\Upsilon^k = \Delta_y^k, \ d_0 = \mu_g\beta/2, \ d_1 = \Big(\frac{2}{\mu_g\beta} - 1\Big)L_y^2, \ d_2 = \beta^2\sigma_g^2. \tag{34}$$

Notice that $\Theta^0 = 0$. Assuming that $\alpha \leq 1/2L_f$, we notice the following implications:

$$\frac{\alpha}{\beta} \leq \frac{\mu_g}{8L_y L} \implies c_0 - c_1\frac{d_1}{d_0} \geq \frac{1}{8\alpha} > 0, \ d_0 - d_1\frac{c_1}{c_0} \geq \frac{\mu_g\beta}{4} > 0, \tag{35}$$

i.e., if (35) holds, then the conclusion (32) can be applied. It can be shown that the step sizes in (23) satisfy (35). Applying Lemma 5 shows that

$$\frac{1}{K}\sum_{k=1}^K \mathbb{E}[\|x^k - x^{k-1}\|^2] \leq \frac{2\mathrm{OPT}^0 + \frac{L}{L_y}(\Delta_y^0 + 4\sigma_g^2/\mu_g^2)}{L_y L \cdot K^{8/5}} + \frac{(2b_0^2 + \sigma_f^2) + 8\frac{\sigma_g^2 L^2}{\mu_g^2}}{2L_y^2 L^2 \cdot K^{6/5}},$$

$$\frac{1}{K}\sum_{k=1}^K \Delta_y^k \leq \frac{2\Delta_y^0 + \frac{8\sigma_g^2}{\mu_g^2} + \frac{4L_y}{\mu_g L}\mathrm{OPT}^0}{K^{3/5}} + \frac{8\frac{\sigma_g^2}{\mu_g^2} + \frac{\mu_g(2b_0^2 + \sigma_f^2)}{2L^2}}{K^{2/5}}.$$

Again, we emphasize that the two timescales step size design is crucial to establishing the above upper bounds. Now, recalling the properties of the Moreau envelop in (18), we obtain the following descent estimate:

**Lemma 7.** *Under Assumption 1, 2, 3. Set $\rho > -\mu_\ell$, $\rho \geq 0$, then for any $k \geq 0$, it holds that*

$$\mathbb{E}[\Phi_{1/\rho}(x^{k+1}) - \Phi_{1/\rho}(x^k)] \leq \frac{5\rho}{2}\mathbb{E}[\|x^{k+1} - x^k\|^2] + \Big[\frac{2\alpha\rho L^2}{\rho + \mu_\ell} + 3\alpha^2\rho L^2\Big]\Delta_y^{k+1}$$
$$- \frac{(\rho + \mu_\ell)\rho\alpha}{4}\mathbb{E}[\|\widehat{x}^k - x^k\|^2] + \Big[\frac{2\rho}{\rho + \mu_\ell} + \rho(\widetilde{\sigma}_f^2 + 3\alpha)\Big]\alpha^2. \tag{36}$$

See details in §B.3. Summing up the inequality (36) from $k = 0$ to $k = K_{\mathsf{max}} - 1$ gives the following

upper bound:

$$\frac{1}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E}[\|\widehat{x}(x^k) - x^k\|^2] \le \frac{4}{(\rho+\mu_\ell)\rho}\Big[\frac{\Phi_{1/\rho}(x^0)}{\alpha K_{\max}} + \frac{5\rho}{2\alpha K_{\max}}\sum_{k=1}^{K_{\max}}\mathbb{E}[\|x^k - x^{k-1}\|^2]\Big]$$
$$+ \frac{4}{\rho+\mu_\ell}\Big[\frac{\frac{L^2}{\rho+\mu_\ell} + 3\alpha L^2}{K_{\max}}\sum_{k=1}^{K_{\max}}\Delta_y^k + \Big[\frac{2}{\rho+\mu_\ell} + \widetilde{\sigma}_f^2 + 3\alpha\Big]\alpha\Big].$$

Combining the above and $\alpha \asymp K_{\max}^{-3/5}$ yields the desired $\frac{1}{K_{\max}}\sum_{k=0}^{K_{\max}-1}\mathbb{E}[\|\widehat{x}(x^k) - x^k\|^2] = \mathcal{O}(K_{\max}^{-2/5})$. In particular, the asymptotic bound is given by setting $\rho = 2|\mu_\ell|$.

**Proof of Corollary** 1  We observe that Lemma 6 can be applied directly in this setting since convex functions are also weakly convex. With the step size choice (26), similar conclusions hold as:

$$\frac{1}{K}\sum_{k=1}^K \mathbb{E}[\|x^k - x^{k-1}\|^2] \le \frac{2\mathrm{OPT}^0 + \frac{L}{L_y}(\Delta_y^0 + 4\sigma_g^2/\mu_g^2)}{L_y L \cdot K^{7/4}} + \frac{(2b_0^2 + \sigma_f^2) + 8\frac{\sigma_g^2 L^2}{\mu_g^2}}{2L_y^2 L^2 \cdot K^{6/4}},$$
$$\frac{1}{K}\sum_{k=1}^K \Delta_y^k \le \frac{2\Delta_y^0 + \frac{8\sigma_g^2}{\mu_g^2} + \frac{4L_y}{\mu_g L}\mathrm{OPT}^0}{K^{1/2}} + \frac{8\frac{\sigma_g^2}{\mu_g^2} + \frac{\mu_g(2b_0^2 + \sigma_f^2)}{2L^2}}{K^{1/2}}.$$

With the additional property $\mu_\ell \ge 0$, in §B.4 we further derive an alternative descent estimate to (33a) that leads to the desired bound of $K^{-1}\sum_{k=1}^K \mathrm{OPT}^k$.

# 4    Application to Reinforcement Learning

Consider a Markov decision process (MDP) $(S, A, \gamma, P, r)$, where $S$ and $A$ are the state and action spaces, respectively, $\gamma \in [0, 1)$ is the discount factor, $P(s'|s, a)$ is the transition kernel to the next state $s'$ given the current state $s$ and action $a$, and $r(s, a) \in [0, 1]$ is the reward at $(s, a)$. Furthermore, the initial state $s_0$ is drawn from a fixed distribution $\rho_0$. We follow a stationary policy $\pi : S \times A \to \mathbb{R}$. For any $(s, a) \in S \times A$, $\pi(a|s)$ is the probability of the agent choosing action $a \in A$ at state $s \in S$. Note that a policy $\pi \in X$ induces a Markov chain on $S$. Denote the induced Markov transition kernel as $P^\pi$ such that $s_{t+1} \sim P^\pi(\cdot|s_t)$. For any $s, s' \in S$, we have $P^\pi(s'|s) = \sum_{a \in A}\pi(a|s)P(s'|s, a)$. For any $\pi \in X$, $P^\pi$ is assumed to induce a stationary distribution over $S$, denoted by $\mu^\pi$. We assume that $|A| < \infty$ while $|S|$ is possibly infinite (but countable). To simplify our notations, for any distribution $\rho$ on $S$, we let $\langle \cdot, \cdot \rangle_\rho$ be the inner product with respect to $\rho$, and $\|\cdot\|_{\mu^\pi \otimes \pi}$ be the weighted $\ell_2$-norm with respect to the probability measure $\mu^\pi \otimes \pi$ over $S \times A$ (where $f, g$ are measurable functions on $S \times A$)

$$\langle f, g \rangle_\rho = \sum_{s \in S}\langle f(s, \cdot), g(s, \cdot)\rangle\rho(s), \quad \|f\|_{\mu^\pi \otimes \pi} = \sqrt{\sum_{s \in S}\Big\{\sum_{a \in A}\pi(a|s) \cdot [f(s, a)]^2\Big\}\mu^\pi(s)}.$$

In policy optimization, our objective is to maximize the expected total discounted reward re-

ceived by the agent with respect to the policy $\pi$, i.e.,

$$\max_{\pi \in X \subseteq \mathbb{R}^{|S| \times |A|}} -\ell(\pi) = \mathbb{E}_\pi \Big[ \sum_{t \geq 0} \gamma^t \cdot r(s_t, a_t) \,|\, s_0 \sim \rho_0 \Big], \tag{37}$$

where $\mathbb{E}_\pi$ is the expectation with the actions taken according to policy $\pi$. Here we let $\rho_0$ to denote the distribution of the initial state. To see that (37) is approximated as a bilevel problem, set $P^\pi$ as the Markov operator under the policy $\pi$. We let $Q^\pi$ be the unique solution to the following Bellman equation [53]:

$$Q(s, a) = r(s, a) + \gamma (P^\pi Q)(s, a), \ \forall \, s, a \in S \times A. \tag{38}$$

Notice that the following holds:

$$Q^\pi(s, a) = \mathbb{E}_\pi \Big[ \sum_{t \geq 0} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a \Big], \ \mathbb{E}_{a \sim \pi(\cdot|s)}[Q^\pi(s, a)] = \langle Q^\pi(s, \cdot), \pi(\cdot|s) \rangle.$$

Further, we parameterize $Q$ using a linear approximation $Q(s, a) \approx Q_\theta(s, a) := \phi^\top(s, a)\theta$, where $\phi : S \times A \to \mathbb{R}^d$ is a known feature mapping and $\theta \in \mathbb{R}^d$ is a finite-dimensional parameter. Using the fact that $\ell(\pi) = -\mathbb{E}_\pi[Q^\pi(s, a)]$, problem (37) can be approximated as a bilevel optimization problem such that:

$$\begin{aligned} \min_{\pi \in X \subseteq \mathbb{R}^{|S| \times |A|}} \quad & \ell(\pi) = -\langle Q_{\theta^\star(\pi)}, \pi \rangle_{\rho_0} \\ \text{subject to} \quad & \theta^\star(\pi) \in \arg\min_{\theta \in \mathbb{R}^d} \tfrac{1}{2} \|Q_\theta - r - \gamma P^\pi Q_\theta\|_{\mu^\pi \otimes \pi}^2. \end{aligned} \tag{39}$$

**Solving Policy Optimization Problem**  We illustrate how to adopt the TTSA algorithm to solve (39). First, the inner problem is the policy evaluation (a.k.a. 'critic') which minimizes the mean squared Bellman error (MSBE). A standard approach is TD learning [52]. We draw two consecutive state-action pairs $(s, a, s', a')$ satisfying $s \sim \mu^{\pi^k}$, $a \sim \pi^k(\cdot|s)$, $s' \sim P(\cdot|s, a)$, and $a' \sim \pi^k(\cdot|s')$, and update the critic via

$$\theta^{k+1} = \theta^k - \beta h_g^k \quad \text{with} \quad h_g^k = [\phi^\top(s, a)\theta^k - r(s, a) - \gamma \phi^\top(s', a')\theta^k]\phi(s, a), \tag{40}$$

where $\beta$ is the step size. This step resembles (4a) of TTSA except that the mean field $\mathbb{E}[h_g^k | \mathcal{F}_k]$ is a semigradient of the MSBE function.

Secondly, the outer problem searches for the policy (a.k.a. 'actor') that maximizes the expected discounted reward. To develop this step, let us define the visitation measure and the Bregman divergence as:

$$\rho^{\pi^k}(s) := (1 - \gamma)^{-1} \sum_{t \geq 0} \gamma^t \mathbb{P}(s_t = s), \ \bar{D}_{\psi, \rho^{\pi^k}}(\pi, \pi^k) := \sum_{s \in S} D_\psi \big( \pi(\cdot|s), \pi^k(\cdot|s) \big) \rho^{\pi^k}(s),$$

such that $\{s_t\}_{t \geq 0}$ is a trajectory of states obtained by drawing $s_0 \sim \rho_0$ and following the policy $\pi^k$, and $D_\psi$ is the Kullback-Leibler (KL) divergence between probability distributions over $A$. We also define the following gradient surrogate:

$$[\overline{\nabla}_\pi f(\pi^k, \theta^{k+1})](s, a) = -(1 - \gamma)^{-1} Q_{\theta^{k+1}}(s, a) \rho^{\pi^k}(s), \ \forall \, (s, a). \tag{41}$$

Similar to (6) and under the additional assumption that the linear approximation is exact, i.e., $Q_{\theta^\star(\pi^k)} = Q^{\pi^k}$, we can show $\overline{\nabla}_\pi f(\pi^k, \theta^\star(\pi^k)) = \nabla \ell(\pi^k)$ using the policy gradient theorem [53]. In a similar vein as (4b) in TTSA, we consider the mirror descent step for improving the policy

(cf. proximal policy optimization in [49]):

$$\pi^{k+1} = \arg\min_{\pi \in X}\Big\{-(1-\gamma)^{-1}\langle Q_{\theta^{k+1}}, \pi - \pi^k\rangle_{\rho^{\pi^k}} + \frac{1}{\alpha}\bar{D}_{\psi,\rho^{\pi^k}}(\pi, \pi^k)\Big\}, \tag{42}$$

where $\alpha$ is the step size. Note that the above update can be performed as:

$$\pi^{k+1}(\cdot|s) \propto \pi^k(\cdot|s)\exp\big[\alpha_k(1-\gamma)^{-1}Q_{\theta^{k+1}}(s,\cdot)\big] = \pi^0(\cdot|s)\exp\Big[(1-\gamma)^{-1}\phi(s,\cdot)^\top\sum_{i=0}^{k}\alpha\theta^{i+1}\Big].$$

In other words, $\pi^{k+1}$ can be represented using the running sum of critic $\sum_{i=0}^{k}\alpha\theta^{i+1}$. This is similar to the natural policy gradient method [30], and the algorithm requires a low memory footprint. Finally, the recursions (40), (42) give the two-timescale natural actor critic (TT-NAC) algorithm.

## 4.1 Convergence Analysis of TT-NAC

Consider the following assumptions on the MDP model of interest.

**Assumption 4.** *The reward function is uniformly bounded by a constant $\bar{r}$. That is, $|r(s,a)| \leq \bar{r}$ for all $(s,a) \in S \times A$.*

**Assumption 5.** *The feature map $\phi\colon S \times A \to \mathbb{R}^d$ satisfies $\|\phi(s,a)\|_2 \leq 1$ for all $(s,a) \in S \times A$. The action-value function associated with each policy is a linear function of $\phi$. That is, for any policy $\pi \in X$, there exists $\theta^\star(\pi) \in \mathbb{R}^d$ such that $Q^\pi(\cdot,\cdot) = \phi(\cdot,\cdot)^\top\theta^\star(\pi) = Q_{\theta^\star(\pi)}(\cdot,\cdot)$.*

**Assumption 6.** *For each policy $\pi \in X$, the induced Markov chain $P^\pi$ admits a unique stationary distribution $\mu^\pi$ for all $\pi \in X$. Let there exists $\mu_\phi > 0$ such that $\mathbb{E}_{s\sim\mu^\pi, a\sim\pi(\cdot|s)}[\phi(s,a)\phi(s,a)^\top] \succeq \mu_\phi^2\cdot I_d$ for all $\pi \in X$.*

**Assumption 7.** *For any $(s,a) \in S \times A$ and any $\pi \in X$, let $\varrho(s,a,\pi)$ be a probability measure over $S$, defined by*

$$[\varrho(s,a,\pi)](s') = (1-\gamma)^{-1}\sum_{t\geq 0}\gamma^t\cdot\mathbb{P}(s_t = s'), \qquad \forall s' \in S. \tag{43}$$

*That is, $\varrho(s,a,\pi)$ is the visitation measure induced by the Markov chain starting from $(s_0, a_0) = (s,a)$ and follows $\pi$ afterwards. For any $\pi^\star$, there exists $C_\rho > 0$ such that*

$$\mathbb{E}_{s'\sim\rho^\star}\Big[\Big|\frac{\varrho(s,a,\pi)}{\rho^\star}(s')\Big|^2\Big] \leq C_\rho^2, \quad \forall\,(s,a) \in S \times A,\ \pi \in X.$$

*Here we let $\rho^\star$ denote $\rho^{\pi^\star}$ to simplify the notation, which is the visitation measure induced by $\pi^\star$ with $s_0 \sim \rho_0$.*

We remark that Assumption 4 is standard in the reinforcement learning literature [53, 55]. In Assumption 5, we assume that each $Q^\pi$ is linear which implies that the linear function approximation is exact. A sufficient condition for Assumption 5 is that the underlying MDP is a linear MDP [28,62], where both the reward function and Markov transition kernel are linear in $\phi$. Linear MDP contains the tabular MDP as a special case, where the feature mapping $\phi(s,a)$ becomes the canonical vector in $\mathbb{R}^{S\times A}$. Assumption 6 assumes that the stationary distribution $\mu^\pi$ exists for any policy $\pi$, which is a common property for the MDP analyzed in TD learning, e.g., [4,15]. Assumption 6 further requires the smallest eigenvalue of $\Sigma_\pi$ to be bounded uniformly away from zero. Such an assumption is

commonly made in the literature on policy evaluation with linear function approximation, e.g., [4,37]. Finally, Assumption 7 postulates that $\rho^\star$ is regular such that the density ratio between $\varrho(s, a, \pi)$ and $\rho^\star$ has uniformly bounded second-order moments under $\rho^\star$. Such an assumption is closely related to the concentratability coefficient [1,2,43], which characterizes the distribution shift incurred by policy updates and is conjectured essential for the sample complexity analysis of reinforcement learning methods [11]. Assumption 7 is satisfied if the initial distribution $\rho_0$ has lower bounded density over $S \times A$ [1]. For details, please refer to Appendix C.

To state our main convergence results, let us define the quantities of interest:

$$\Delta_Q^{k+1} := \mathbb{E}[\|\theta^{k+1} - \theta^\star(\pi^k)\|_2^2], \quad \mathrm{OPT}^k := \mathbb{E}[\ell(\pi^k) - \ell(\pi^\star)], \tag{44}$$

where the expectations above are taken with respect to the i.i.d. draws of state-action pairs in (40) for TT-NAC. We remark that $\Delta_Q^k$, analogous to $\Delta_y^k$ used in TTSA, is the tracking error that characterizes the performance of TD learning when the target value function, $Q^{\pi^k}$, is time-varying due to policy updates. We obtain:

**Theorem 3.** *Consider the TT-NAC algorithm* (40)-(42) *for the policy optimization problem* (39). *Let $K_{\mathsf{max}} \geq 32^2$ be the maximum number of iterations. Under Assumption 4 – 7, and we set the step sizes as*

$$\alpha = \frac{(1-\gamma)^3 \mu_\phi}{\sqrt{\bar{r} \cdot C_\rho^2}} \min \left\{ \frac{(1-\gamma)^2}{128 \mu_\phi^{-2}}, K_{\mathsf{max}}^{-3/4} \right\}, \ \ \beta = \min \left\{ \frac{(1-\gamma)\mu_\phi^2}{8}, \frac{16}{(1-\gamma)\mu_\phi^2} K_{\mathsf{max}}^{-1/2} \right\}. \tag{45}$$

*Then the following holds*

$$\mathbb{E}[\mathrm{OPT}^{\mathsf{K}}] = \mathcal{O}(K_{\mathsf{max}}^{-1/4}), \quad \mathbb{E}[\Delta_Q^{\mathsf{K}+1}] = \mathcal{O}(K_{\mathsf{max}}^{-1/2}), \tag{46}$$

*where $\mathsf{K}$ is an independent random variable uniformly distributed over $\{0, ..., K_{\mathsf{max}} - 1\}$.*

To shed lights on our analysis, we first observe the following performance difference lemma proven in [29, Lemma 6.1]:

$$\ell(\pi) - \ell(\pi^\star) = (1-\gamma)^{-1} \langle Q^\pi, \pi^\star - \pi \rangle_{\rho^{\pi^\star}}, \qquad \forall \pi \in X, \tag{47}$$

where $\pi^\star$ is an optimal policy solving (39). The above implies a restricted form of convexity, and our analysis uses the insight that (47) plays a similar role as (2) [with $\mu_\ell \geq 0$] and characterizes the loss geometry of the outer problem.

Our result shows that the TT-NAC algorithm finds an optimal policy at the rate of $\mathcal{O}(K^{-1/4})$ in terms of the objective value. This rate is comparable to another variant of the TT-NAC algorithm in [61], which provided a customized analysis for TT-NAC. In contrast, the analysis for our TT-NAC algorithm is rooted in the general TTSA framework developed in §3.3 for tackling bilevel optimization problems. Notice that analysis for the two-timescale actor-critic algorithm can also be found in [59], which provides an $\mathcal{O}(K^{-2/5})$ convergence rate to a stationary solution.

## 5 Numerical Experiments

We consider the data hyper-cleaning task (16), and compare TTSA with several algorithms such as the BSA algorithm [23], the stocBiO [27] for different batch size choices, and the HOAG algorithm
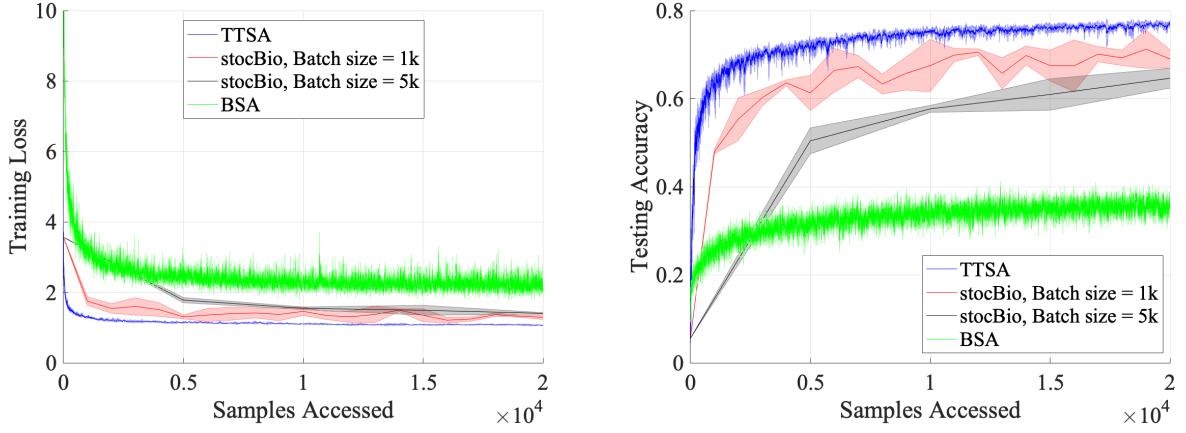
Figure 1: Data hyper-cleaning task on the `FashionMNIST` dataset. We plot the training loss and testing accuracy against the number of gradients evaluated with corruption rate $p = 0.4$.

in [44]. Note that HOAG is a deterministic algorithm and it requires full gradient computation at each iteration. In contrast, stocBiO is a stochastic algorithm but it relies on large batch gradient computations.

We consider problem (16) with $L(\cdot)$ being the cross-entropy loss (i.e., a data cleaning problem for logistic regression); $\sigma(x) := \frac{1}{1+\exp(-x)}$; $c = 0.001$; see [50]. The problem is trained on the `FashionMNIST` dataset [60] with 50k, 10k, and 10k image samples allocated for training, validation and testing purposes, respectively. We consider the setting where each sample in the training dataset is corrupted with probability 0.4. Note that the outer problem $\ell(x)$ is non-convex while the lower level problem is strongly-convex. The simulation results are an average of 3 independent runs. The step sizes for different algorithms are chosen according to their theoretically suggested values. Let the outer iteration be indexed by $t$, for TTSA we choose $\alpha_t = c_\alpha/(1+t)^{3/5}$, $\beta_t = c_\beta/(1+t)^{2/5}$, and tune for $c_\alpha$ and $c_\beta$ in the set $\{10^{-3}, 10^{-2}, 10^{-1}, 10\}$. For BSA [23], we index the outer iteration by $t$ and the inner iteration by $\bar{k} \in \{1, \dots, \bar{k}_t\}$. We set $\bar{k}_t = \lceil \sqrt{t+1} \rceil$ as suggested in [23] and choose the outer and inner step-sizes $\alpha_t$ and $\beta_{\bar{k}}$, respectively, as $\alpha_t = d_\alpha/(1+t)^{1/2}$ and $\beta_{\bar{k}} = d_\beta/(\bar{k}+2)$. We tune for $d_\alpha$ and $d_\beta$ in the set $\{10^{-3}, 10^{-2}, 10^{-1}, 10\}$. Finally, for stocBiO we tune for parameters $\alpha_t$ and $\beta_t$ in the range $[0, 1]$.

In Figure 1, we compare the performance of different algorithms against the total number of outer samples accessed. As observed, TTSA outperforms BSA, stocBiO and HOAG. We remark that HOAG is a deterministic algorithm and hence requires full batch gradient computations at each iteration. Similarly, stocBio relies on large batch gradients which results in relatively slow convergence.

## 6   Conclusion

This paper develops efficient two-timescale stochastic approximation algorithms for a class of bi-level optimization problems where the inner problem is unconstrained and strongly convex. We

19

show the convergence rates of the proposed TTSA algorithm under the settings where the outer objective function is either strongly convex, convex, or non-convex. Additionally, we show how our theory and analysis can be customized to a two-timescale actor-critic proximal policy optimization algorithm in reinforcement learning, and obtain a comparable convergence rate to existing literature.

# A  Omitted Proofs of Theorem 1

To simplify notations, for any $n, m \in \mathbb{N}$, we define the following quantities for brevity of notations.

$$G^{(1)}_{m:n} = \prod_{i=m}^{n} (1 - \beta_i \mu_g / 4), \quad G^{(2)}_{m:n} = \prod_{i=m}^{n} (1 - \alpha_i \mu_\ell). \tag{48}$$

## A.1  Proof of Lemma 3

Following a direct expansion of the updating rule for $y^{k+1}$ and taking the conditional expectation given filtration $\mathcal{F}_k$ yield that

$$\mathbb{E}[\|y^{k+1} - y^\star(x^k)\|^2 | \mathcal{F}_k] \le (1 - 2\beta_k \mu_g)\|y^k - y^\star(x^k)\|^2 + \beta_k^2 \mathbb{E}[\|h_g^k\|^2 | \mathcal{F}_k],$$

where we used the unbiasedness of $h_g^k$ [cf. Assumption 3] and the strong convexity of $g$. By direct computation and (7b) in Assumption 3, we have

$$\mathbb{E}[\|h_g^k\|^2 | \mathcal{F}_k] = \mathbb{E}\big[\|h_g^k - \nabla g(x^k, y^k)\|^2 \big| \mathcal{F}_k\big] + \|\nabla g(x^k, y^k)\|^2$$
$$\le \sigma_g^2 + (1 + \sigma_g^2)\|\nabla g(x^k, y^k)\|^2 \le \sigma_g^2 + (1 + \sigma_g^2) \cdot L_g^2 \|y^k - y^\star(x^k)\|^2, \tag{49}$$

where the last inequality uses Assumption 2 and the optimality of the inner problem $\nabla_y g(x^k, y^\star(x^k)) = 0$. As $\beta_k L_g^2 (1 + \sigma_g^2) \le \mu_g$, we have

$$\mathbb{E}[\|y^{k+1} - y^\star(x^k)\|^2 | \mathcal{F}_k] \le (1 - \beta_k \mu_g) \cdot \|y^k - y^\star(x^k)\|^2 + \beta_k^2 \sigma_g^2. \tag{50}$$

Using the basic inequality $2ab \le 1/c \cdot a^2 + c \cdot b^2$ for all $c \ge 0$ and $a, b \in \mathbb{R}$, we have

$$\|y^k - y^\star(x^k)\|^2 \le \big(1 + 1/c\big) \cdot \|y^k - y^\star(x^{k-1})\|^2 + \big(1 + c\big)\|y^\star(x^{k-1}) - y^\star(x^k)\|^2. \tag{51}$$

Note we have taken the convention that $x^{-1} = x^0$. Furthermore, we observe that

$$\|y^\star(x^{k-1}) - y^\star(x^k)\|^2 \le L_y^2 \|x^k - x^{k-1}\|^2 \le \alpha_{k-1}^2 L_y^2 \cdot \|h_f^{k-1}\|^2, \tag{52}$$

where the first inequality follows from Lemma 2, and the second inequality follows from the non-expansive property of projection. We have set $h_f^{-1} = 0$ as convention.

Through setting $c = \frac{2(1 - \beta_k \mu_g)}{\beta_k \mu_g}$, we have $\big(1 + 1/c\big)(1 - \beta_k \mu_g) = 1 - \frac{\mu_g}{2}\beta_k$. Substituting the above quantity $c$ into (51) and combining with (50) show that

$$\mathbb{E}[\|y^{k+1} - y^\star(x^k)\|^2 | \mathcal{F}_k]$$
$$\le \big(1 - \frac{\beta_k \mu_g}{2}\big) \cdot \|y^k - y^\star(x^{k-1})\|^2 + \beta_k^2 \cdot \sigma_g^2 + \frac{2 - \mu_g \beta_k}{\mu_g \beta_k} \cdot \alpha_{k-1}^2 L_y^2 \cdot \|h_f^{k-1}\|^2. \tag{53}$$

Taking the total expectation and using (14), we have

$$\Delta_y^{k+1} \le (1 - \beta_k \mu_g/2) \cdot \Delta_y^k + \beta_k^2 \sigma_g^2 + \frac{2 - \mu_g \beta_k}{\mu_g \beta_k} \alpha_{k-1}^2 L_y^2 [\tilde\sigma_f^2 + 3b_{k-1}^2 + 3L^2 \Delta_y^k],$$

with the convention $\alpha_{-1} = 0$. Using $\alpha_{k-1} \le 2\alpha_k$, $\alpha_k \le c_0 \beta_k^{3/2}$, we have

$$\Delta_y^{k+1} \le \left[ 1 - \beta_k \mu_g/2 + \frac{12 c_0^2 L_y^2 L^2}{\mu_g} \beta_k^2 \right] \cdot \Delta_y^k + \beta_k^2 \sigma_g^2 + \frac{4 c_0^2 L_y^2}{\mu_g} \beta_k^2 \cdot [\tilde\sigma_f^2 + 3b_0^2]$$

$$\le \left[ 1 - \beta_k \mu_g/4 \right] \cdot \Delta_y^k + \beta_k^2 \sigma_g^2 + \frac{4 c_0^2 L_y^2}{\mu_g} \beta_k^2 \cdot [\tilde\sigma_f^2 + 3b_0^2],$$

where the last inequality is due to (20a). Solving the recursion leads to

$$\Delta_y^{k+1} \le G_{0:k}^{(1)} \Delta_y^0 + \sum_{j=0}^k \beta_j^2 G_{j+1:k}^{(1)} \left\{ \sigma_g^2 + \frac{4 c_0^2 L_y^2}{\mu_g} [\tilde\sigma_f^2 + 3b_0^2] \right\}. \tag{54}$$

Since $\beta_{k-1}/\beta_k \le 1 + \beta_k \cdot (\mu_g/8)$, applying Lemma 10 to $\{\beta_k\}_{k \ge 0}$ with $a = \mu_g/4$ and $q = 2$, we have $\sum_{j=0}^k \beta_j^2 G_{j+1:k}^{(1)} \le \frac{8\beta_k}{\mu_g}$. Finally, we can simplify (54) as

$$\Delta_y^{k+1} \le G_{0:k}^{(1)} \Delta_y^0 + C_y^{(1)} \beta_k, \quad \text{where} \quad C_y^{(1)} := \frac{8}{\mu_g} \left\{ \sigma_g^2 + \frac{4 c_0^2 L_y^2}{\mu_g} [\tilde\sigma_f^2 + 3b_0^2] \right\}. \tag{55}$$

## A.2  Proof of Lemma 4

Due to the projection property, we get

$$\|x^{k+1} - x^\star\|^2 \le \|x^k - \alpha_k h_f^k - x^\star\|^2 = \|x^k - x^\star\|^2 - 2\alpha_k \langle h_f^k, x^k - x^\star \rangle + \alpha_k^2 \|h_f^k\|^2.$$

Taking the conditional expectation given $\mathcal{F}_k'$ gives

$$\begin{aligned}
\mathbb{E}[\|x^{k+1} - x^\star\|^2 | \mathcal{F}_k'] &\le \|x^k - x^\star\|^2 - 2\alpha_k \langle \nabla\ell(x^k), x^k - x^\star \rangle + \alpha_k^2 \mathbb{E}[\|h_f^k\|^2 | \mathcal{F}_k'] \\
&\quad - 2\alpha_k \langle \overline\nabla_x f(x^k, y^{k+1}) - \nabla\ell(x^k) + B_k, x^k - x^\star \rangle,
\end{aligned} \tag{56}$$

where the inequality follows from (7a). The strong convexity implies $\langle \nabla\ell(x^k), x^k - x^\star \rangle \ge \mu_\ell \|x^k - x^\star\|^2$, we further bound the r.h.s. of (56) via

$$\begin{aligned}
&\mathbb{E}[\|x^{k+1} - x^\star\|^2 | \mathcal{F}_k'] \\
&\le (1 - 2\alpha_k \mu_\ell) \|x^k - x^\star\|^2 - 2\alpha_k \langle \overline\nabla_x f(x^k, y^{k+1}) - \nabla\ell(x^k) + B_k, x^k - x^\star \rangle + \alpha_k^2 \mathbb{E}[\|h_f^k\|^2 | \mathcal{F}_k'] \\
&\le (1 - \alpha_k \mu_\ell) \|x^k - x^\star\|^2 + \frac{\alpha_k}{\mu_\ell} \|\overline\nabla_x f(x^k, y^{k+1}) - \nabla\ell(x^k) + B_k\|^2 + \alpha_k^2 \mathbb{E}[\|h_f^k\|^2 | \mathcal{F}_k'] \\
&\le (1 - \alpha_k \mu_\ell) \cdot \|x^k - x^\star\|^2 + (2\alpha_k/\mu_\ell) \cdot \{ L^2 \|y^{k+1} - y^\star(x^k)\|^2 + b_k^2 \} + \alpha_k^2 \cdot \mathbb{E}[\|h_f^k\|^2 | \mathcal{F}_k'],
\end{aligned}$$

where the last inequality is from Lemma 2. Using (14) and taking total expectation:

$$\begin{aligned}
\Delta_x^{k+1} &\le [1 - \alpha_k \mu_\ell] \Delta_x^k + [2\alpha_k/\mu_\ell] L^2 \Delta_y^{k+1} + 2\alpha_k b_k^2/\mu_\ell + \alpha_k^2 [\tilde\sigma_f^2 + 3b_0^2 + 3L^2 \Delta_y^{k+1}] \\
&\le [1 - \alpha_k \mu_\ell] \cdot \Delta_x^k + [2\alpha_k/\mu_\ell + 3\alpha_k^2] \cdot L^2 \cdot \Delta_y^{k+1} + \alpha_k^2 [2\tilde c_b/\mu_\ell + \tilde\sigma_f^2 + 3b_0^2],
\end{aligned}$$

where we have used $b_k^2 \leq \widetilde{c}_b \alpha_k$. Solving the recursion above leads to

$$\Delta_x^{k+1} \leq G_{0:k}^{(2)} \Delta_x^0 + \sum_{j=0}^{k} \left\{ \left[ \frac{2\widetilde{c}_b}{\mu_\ell} + \widetilde{\sigma}_f^2 + 3b_0^2 \right] \alpha_j^2 G_{j+1:k}^{(2)} + \left[ \frac{2L^2}{\mu_\ell} + 3\alpha_0 L^2 \right] \alpha_j G_{j+1:k}^{(2)} \Delta_y^{j+1} \right\}$$

$$\leq G_{0:k}^{(2)} \Delta_x^0 + \frac{2}{\mu_\ell} \left[ \frac{2\widetilde{c}_b}{\mu_\ell} + \widetilde{\sigma}_f^2 + 3b_0^2 \right] \cdot \alpha_k + \left( \frac{2L^2}{\mu_\ell} + 3\alpha_0 L^2 \right) \sum_{j=0}^{k} \alpha_j G_{j+1:k}^{(2)} \Delta_y^{j+1}.$$

The last inequality follows from applying Lemma 10 with $q = 2$ and $a = \mu_\ell$.

## A.3 Bounding $\Delta_x^k$ by coupling with $\Delta_y^k$

Using (55), we observe that

$$\sum_{j=0}^{k} \alpha_j G_{j+1:k}^{(2)} \Delta_y^{j+1} \leq \sum_{j=0}^{k} \alpha_j G_{j+1:k}^{(2)} \left\{ G_{0:j}^{(1)} \Delta_y^0 + C_y^{(1)} \beta_j \right\}. \tag{57}$$

We bound each term on the right-hand side of (57). For the first term, as $\mu_\ell \alpha_i \leq \mu_g \beta_i / 8$ [cf. (20b)], applying Lemma 11 with $a = \mu_\ell$, $b = \mu_g/4$, $\gamma_i = \alpha_i$, $\rho_i = \beta_i$ gives

$$\sum_{j=0}^{k} \alpha_j G_{j+1:k}^{(2)} G_{0:j}^{(1)} \Delta_y^0 \leq \frac{1}{\mu_\ell} G_{0:k}^{(2)} \Delta_y^0. \tag{58}$$

Recall that $\beta_j \leq c_1 \cdot \alpha_j^{2/3}$. Applying Lemma 10 with $q = 5/3$, $a = \mu_\ell$ yields

$$\sum_{j=0}^{k} \alpha_j \beta_j G_{j+1:k}^{(2)} \leq c_1 \sum_{j=0}^{k} \alpha_j^{5/3} G_{j+1:k}^{(2)} \leq c_1 \frac{2}{\mu_\ell} \alpha_k^{2/3}, \tag{59}$$

We obtain a bound on the optimality gap as

$$\Delta_x^{k+1} \leq G_{0:k}^{(2)} \left\{ \Delta_x^0 + \left[ \frac{2L^2}{\mu_\ell^2} + \frac{3\alpha_0 L^2}{\mu_\ell} \right] \Delta_y^0 \right\} + \frac{2}{\mu_\ell} \left[ \frac{2\widetilde{c}_b}{\mu_\ell} + \widetilde{\sigma}_f^2 + 3b_0^2 \right] \alpha_k$$

$$+ \frac{2c_1}{\mu_\ell} \left[ \frac{2L^2}{\mu_\ell} + 3\alpha_0 L^2 \right] C_y^{(1)} \alpha_k^{\frac{2}{3}}.$$

To simplify the notation, we define the constants

$$C_x^{(0)} = \Delta_x^0 + \left[ \frac{2L^2}{\mu_\ell^2} + \frac{3\alpha_0 L^2}{\mu_\ell} \right] \Delta_y^0, \quad C_x^{(1)} = \frac{2}{\mu_\ell} \left[ \frac{2\widetilde{c}_b}{\mu_\ell} + \widetilde{\sigma}_f^2 + 3b_0^2 \right] + \frac{2c_1}{\mu_\ell} \left[ \frac{2L^2}{\mu_\ell} + 3\alpha_0 L^2 \right] C_y^{(1)}$$

Then, as long as $\alpha_k < 1/\mu_\ell$ and we use the step size parameters in (22), we have

$$\Delta_x^{k+1} \leq G_{0:k}^{(2)} C_x^{(0)} + C_x^{(1)} \alpha_k^{2/3} = \mathcal{O}\left( \left[ \frac{L^2}{\mu_\ell^2 \mu_g^2} + \frac{L^2 L_y^2}{\mu_\ell^4} \right] \frac{\widetilde{\sigma}_f^2}{k^{2/3}} + \frac{L^2}{\mu_\ell^2 \mu_g} \frac{\sigma_g^2}{k^{2/3}} \right),$$

and we recall that $\widetilde{\sigma}_f^2 = \sigma_f^2 + 3 \sup_{x \in X} \| \nabla \ell(x) \|^2$.

# B Omitted Proofs of Theorem 2 and Corollary 1

## B.1 Proof of Lemma 5

We observe that summing the first and the second inequalities in (31) from $k = 0$ to $k = K - 1$ gives:

$$c_0 \sum_{k=1}^{K} \Theta^k \leq \Omega^0 + c_1 \sum_{k=1}^{K} \Upsilon^k + c_2 \cdot K. \tag{60}$$

$$d_0 \sum_{k=1}^{K} \Upsilon^k \leq \Upsilon^1 + d_1 \sum_{k=1}^{K} \Theta^k + d_2 \cdot K. \tag{61}$$

Substituting (60) into (61) gives

$$d_0 \sum_{k=1}^{K} \Upsilon^k \leq \Upsilon^1 + d_2 \cdot K + \frac{d_1}{c_0} \Big[ \Omega^0 + c_1 \sum_{k=1}^{K} \Upsilon^k + c_2 \cdot K \Big]. \tag{62}$$

Therefore, if $d_0 - d_1 \frac{c_1}{c_0} > 0$, a simple computation yields the second inequality in (32). Similarly, we substitute (61) into (60) to yield

$$c_0 \sum_{k=1}^{K} \Theta^k \leq \Omega^0 + c_2 \cdot K + \frac{c_1}{d_0} \Big[ \Upsilon^1 + d_1 \sum_{k=1}^{K} \Theta^k + d_2 \cdot K \Big]. \tag{63}$$

Under $c_0 - c_1 \frac{d_1}{d_0} > 0$, simple computation yields the first inequality in (32).

## B.2 Proof of Lemma 6

Recall that we defined $\mathrm{OPT}^k := \mathbb{E}[\ell(x^k) - \ell(x^\star)]$ for each $k \geq 0$. To begin with, we have the following descent estimate

$$\ell(x^{k+1}) \leq \ell(x^k) + \langle \nabla \ell(x^k), x^{k+1} - x^k \rangle + (L_f/2) \|x^{k+1} - x^k\|^2. \tag{64}$$

The optimality condition of step (4b) leads to the following bound

$$\langle \nabla \ell(x^k), x^{k+1} - x^k \rangle \leq \langle \nabla \ell(x^k) - \overline{\nabla}_x f(x^k, y^{k+1}) - B_k, x^{k+1} - x^k \rangle$$
$$+ \langle B_k + \overline{\nabla}_x f(x^k, y^{k+1}) - h_f^k, x^{k+1} - x^k \rangle - \frac{1}{\alpha} \|x^{k+1} - x^k\|^2,$$

where we obtained the inequality by adding and subtracting $B_k + \overline{\nabla}_x f(x^k, y^{k+1}) - h_f^k$. Then, taking the conditional expectation on $\mathcal{F}_k'$, for any $c, d > 0$, we obtain

$$\mathbb{E}[\langle \nabla \ell(x^k), x^{k+1} - x^k \rangle | \mathcal{F}_k']$$
$$\leq \mathbb{E}\big[ \|\nabla \ell(x^k) - \overline{\nabla}_x f(x^k, y^{k+1}) - B_k\| \cdot \|x^{k+1} - x^k\| \big| \mathcal{F}_k' \big]$$
$$+ \mathbb{E}\big[ \|B_k + \overline{\nabla}_x f(x^k, y^{k+1}) - h_f^k\| \|x^{k+1} - x^k\| \big| \mathcal{F}_k' \big] - \frac{1}{\alpha} \mathbb{E}[\|x^{k+1} - x^k\|^2 | \mathcal{F}_k']$$
$$\leq \frac{1}{2c} \mathbb{E}[\|\nabla \ell(x^k) - \overline{\nabla}_x f(x^k, y^{k+1}) - B_k\|^2 | \mathcal{F}_k'] + \frac{c}{2} \mathbb{E}[\|x^{k+1} - x^k\|^2 | \mathcal{F}_k']$$
$$+ \frac{\sigma_f^2}{2d} + \frac{d}{2} \mathbb{E}[\|x^{k+1} - x^k\|^2 | \mathcal{F}_k'] - \frac{1}{\alpha} \mathbb{E}[\|x^{k+1} - x^k\|^2 | \mathcal{F}_k'],$$

where the second inequality follows from the Young's inequality and Assumption 3. Simplifying the terms above leads to

$$\mathbb{E}[\langle \nabla \ell(x^k), x^{k+1} - x^k \rangle | \mathcal{F}_k']$$
$$\leq \frac{1}{2c} \mathbb{E}[\|\nabla \ell(x^k) - \overline{\nabla}_x f(x^k, y^{k+1}) - B_k\|^2 | \mathcal{F}_k'] + \frac{\sigma_f^2}{2d} + \left( \frac{c+d}{2} - \frac{1}{\alpha} \right) \cdot \mathbb{E}[\|x^{k+1} - x^k\|^2 | \mathcal{F}_k'].$$

Setting $d = c = \frac{1}{2\alpha}$, plugging the above to (64), and taking the full expectation:

$$\text{OPT}^{k+1} \leq \text{OPT}^k - \left( \frac{1}{2\alpha} - \frac{L_f}{2} \right) \cdot \mathbb{E}[\|x^{k+1} - x^k\|^2] + \alpha \Delta^{k+1} + \alpha \sigma_f^2, \tag{65}$$

where we have denoted $\Delta^{k+1}$ as follows

$$\Delta^{k+1} := \mathbb{E}[\|\overline{\nabla}_x f(x^k; y^{k+1}) - \nabla \ell(x^k) - B_k\|^2] \overset{(12a)}{\leq} 2L^2 \mathbb{E}[\|y^{k+1} - y^\star(x^k)\|^2] + 2b_k^2,$$

where the last inequality follows from Lemma 2 and (7a) in Assumption 3. Next, following from the standard SGD analysis [cf. (50)] and using $\beta \leq \mu_g/(L_g^2(1 + \sigma_g^2))$, we have

$$\mathbb{E}[\|y^{k+1} - y^\star(x^k)\|^2 | \mathcal{F}_k] \leq (1 - \mu_g \beta) \mathbb{E}[\|y^k - y^\star(x^k)\|^2 | \mathcal{F}_k] + \beta^2 \sigma_g^2$$
$$\leq (1+c)(1 - \mu_g \beta) \mathbb{E}[\|y^k - y^\star(x^{k-1})\|^2 | \mathcal{F}_k] + (1 + 1/c) \mathbb{E}[\|y^\star(x^k) - y^\star(x^{k-1})\|^2 | \mathcal{F}_k] + \beta^2 \sigma_g^2$$
$$\leq \left( 1 - \mu_g \beta/2 \right) \mathbb{E}[\|y^k - y^\star(x^{k-1})\|^2 | \mathcal{F}_k] + \left( \frac{2}{\mu_g \beta} - 1 \right) \cdot \mathbb{E}[\|y^\star(x^k) - y^\star(x^{k-1})\|^2 | \mathcal{F}_k] + \beta^2 \sigma_g^2,$$
$$\leq \left( 1 - \mu_g \beta/2 \right) \mathbb{E}[\|y^k - y^\star(x^{k-1})\|^2 | \mathcal{F}_k] + \left( \frac{2}{\mu_g \beta} - 1 \right) L_y^2 \cdot \mathbb{E}[\|x^k - x^{k-1}\|^2 | \mathcal{F}_k] + \beta^2 \sigma_g^2,$$

where the last inequality is due to the Lipschitz continuity property (12a) and $\mu_g \beta < 1$. Furthermore, we have picked $c = \mu_g \beta \cdot [2(1 - \mu_g \beta)]^{-1}$, so that

$$(1+c)(1 - \mu_g \beta) = 1 - \mu_g \beta/2, \qquad 1/c + 1 = 2/(\mu_g \beta) - 1.$$

Taking a full expectation on both sides leads to the desired result.

## B.3  Proof of Lemma 7

For simplicity, we let $\widehat{x}^{k+1}$ and $\widehat{x}$ denote $\widehat{x}(x^{k+1})$ and $\widehat{x}(x)$, respectively. For any $x \in X$, letting $x_1 = \widehat{x}$ and $x_2 = x$ in (2), we get

$$\ell(\widehat{x}) \geq \ell(x) + \langle \nabla \ell(x), \widehat{x} - x \rangle + \frac{\mu_\ell}{2} \|\widehat{x} - x\|^2. \tag{66}$$

Moreover, by the definition of $\widehat{x}$, for any $x \in X$, we have

$$\ell(x) + \frac{\rho}{2} \|x - x\|^2 - \left[ \ell(\widehat{x}) + \frac{\rho}{2} \|\widehat{x} - x\|^2 \right] = \ell(x) - \left[ \ell(\widehat{x}) + \frac{\rho}{2} \|\widehat{x} - x\|^2 \right] \geq 0. \tag{67}$$

Adding the two inequalities above, we obtain

$$-\frac{\mu_\ell + \rho}{2} \cdot \|\widehat{x} - x\|^2 \geq \langle \nabla \ell(x), \widehat{x} - x \rangle. \tag{68}$$

Note that we choose $\rho$ such that $\rho + \mu_\ell > 0$. To proceed, combining the definitions of the Moreau envelop and $\widehat{x}$ in (18), for $x^{k+1}$, we have

$$\Phi_{1/\rho}(x^{k+1}) \overset{(18)}{=} \ell(\widehat{x}^{k+1}) + \frac{\rho}{2} \cdot \|x^{k+1} - \widehat{x}^{k+1}\|^2 \leq \ell(\widehat{x}^k) + \frac{\rho}{2} \cdot \|x^{k+1} - \widehat{x}^k\|^2$$

$$\leq \ell(\widehat{x}^k) + \frac{\rho}{2} \cdot \|x^k - \widehat{x}^k\|^2 + \frac{\rho}{2} \cdot \|x^{k+1} - x^k\|^2 + \rho\alpha \cdot \langle \widehat{x}^k - x^k, h_f^k \rangle$$

$$+ \alpha\rho \cdot \langle h_f^k, x^k - x^{k+1}\rangle + \rho\|x^{k+1} - x^k\|^2$$

$$\overset{(18)}{=} \Phi_{1/\rho}(x^k) + \frac{5\rho}{2} \cdot \|x^{k+1} - x^k\|^2 + \rho\alpha\langle \widehat{x}^k - x^k, h_f^k\rangle + \alpha^2\rho\|h_f^k\|^2, \tag{69}$$

where the first equality and the first inequality follow from the optimality of $\widehat{x}^{k+1} = \widehat{x}(x^{k+1})$, and the second term is from the optimality condition in (4b). For any $x^\star$ that is a global optimal solution for the original problem $\min_{x \in X} \ell(x)$, we must have

$$\Phi_{1/\rho}(x^\star) = \min_{x \in X}\left\{\ell(x) + \frac{\rho}{2}\|x - x^\star\|^2\right\} = \ell(x^\star),$$

where the last equality holds because

$$\Phi_{1/\rho}(x^\star) = \min_{x \in X}\left\{\ell(x) + \frac{\rho}{2}\|x - x^\star\|^2\right\} \leq \ell(x^\star) + \frac{\rho}{2}\|x^\star - x^\star\|^2 = \ell(x^\star), \tag{70}$$

$$\Phi_{1/\rho}(z) = \min_{x \in X}\left\{\ell(x) + \frac{\rho}{2}\|x - z\|^2\right\} \geq \min_{x \in X} \ell(x) = \ell(x^\star), \ \forall \ z \in X. \tag{71}$$

Taking expectation of $\langle \widehat{x}^k - x^k, h_f^k\rangle$ while conditioning on $\mathcal{F}'_k$, we have:

$$\mathbb{E}[\langle \widehat{x}^k - x^k, h_f^k\rangle | \mathcal{F}'_k]$$

$$= \mathbb{E}[\langle \widehat{x}^k - x^k, h_f^k - \overline{\nabla}_x f(x^k, y^{k+1}) + \overline{\nabla}_x f(x^k, y^{k+1}) - \nabla\ell(x^k) + \nabla\ell(x^k)\rangle | \mathcal{F}'_k]$$

$$= \langle \widehat{x}^k - x^k, B_k\rangle + \mathbb{E}[\langle \widehat{x}^k - x^k, \overline{\nabla}_x f(x^k, y^{k+1}) - \nabla\ell(x^k)\rangle + \langle \widehat{x}^k - x^k, \nabla\ell(x^k)\rangle | \mathcal{F}'_k], \tag{72}$$

where the second equality follows from (7a) in Assumption 3. By Young's inequality, for any $c > 0$, we have

$$\langle \widehat{x}^k - x^k, B_k\rangle \leq \frac{c}{4}\|\widehat{x}^k - x^k\|^2 + \frac{1}{c}b_k^2, \tag{73}$$

$$\mathbb{E}[\langle \widehat{x}^k - x^k, \overline{\nabla}_x f(x^k, y^{k+1}) - \nabla\ell(x^k)\rangle | \mathcal{F}'_k] \leq \frac{1}{c}\|\overline{\nabla}_x f(x^k, y^{k+1}) - \nabla\ell(x^k)\|^2 + \frac{c}{4}\|\widehat{x}^k - x^k\|^2,$$

where we also use (7a) in deriving (73). Combining (68), (72), (73), and setting $c = (\rho + \mu_\ell)/2$, we obtain that

$$\mathbb{E}[\langle \widehat{x}^k - x^k, h_f^k\rangle | \mathcal{F}'_k] \tag{74}$$

$$\leq \frac{c}{2}\|\widehat{x}^k - x^k\|^2 + \frac{1}{c}b_k^2 + \frac{1}{c}\mathbb{E}[\|\overline{\nabla}_x f(x^k, y^{k+1}) - \nabla\ell(x^k)\|^2] - \frac{\rho + \mu_\ell}{2}\|\widehat{x}^k - x^k\|^2$$

$$= \frac{2}{\rho + \mu_\ell} \cdot \mathbb{E}[\|\overline{\nabla}_x f(x^k, y^{k+1}) - \nabla\ell(x^k)\|^2] - \frac{(\rho + \mu_\ell)}{4} \cdot \|\widehat{x}^k - x^k\|^2 + \frac{2}{\rho + \mu_\ell} \cdot b_k^2$$

$$\leq \frac{2L^2}{\rho + \mu_\ell} \cdot \mathbb{E}[\|y^{k+1} - y^\star(x^k)\|^2] - \frac{(\rho + \mu_\ell)}{4} \cdot \|\widehat{x}^k - x^k\|^2 + \frac{2}{\rho + \mu_\ell} \cdot b_k^2,$$

where the last step follows from the first inequality of Lemma 2. Plugging the above into (69), and taking a full expectation, we obtain

$$
\begin{aligned}
&\mathbb{E}[\Phi_{1/\rho}(x^{k+1})] - \mathbb{E}[\Phi_{1/\rho}(x^k)] \\
&\leq \frac{5\rho}{2}\mathbb{E}[\|x^{k+1} - x^k\|^2] + \frac{2\rho\alpha L^2}{\rho + \mu_\ell}\Delta_y^{k+1} - \frac{(\rho + \mu_\ell)\rho\alpha}{4}\mathbb{E}[\|\widehat{x}^k - x^k\|^2] + \frac{2\rho\alpha b_k^2}{\rho + \mu_\ell} + \alpha^2\rho\mathbb{E}[\|h_f^k\|^2] \\
&\leq \frac{5\rho}{2}\mathbb{E}[\|x^{k+1} - x^k\|^2] + \frac{2\rho\alpha L^2}{\rho + \mu_\ell}\Delta_y^{k+1} - \frac{(\rho + \mu_\ell)\rho\alpha}{4}\mathbb{E}[\|\widehat{x}^k - x^k\|^2] + \frac{2\rho\alpha b_k^2}{\rho + \mu_\ell} \\
&\quad + \alpha^2\rho(\widetilde{\sigma}_f^2 + 3b_k^2 + 3L^2\Delta_y^{k+1}) \\
&\leq \frac{5\rho}{2}\mathbb{E}[\|x^{k+1} - x^k\|^2] + \Big[\frac{2\alpha\rho L^2}{\rho + \mu_\ell} + 3\alpha^2\rho L^2\Big]\Delta_y^{k+1} - \frac{(\rho + \mu_\ell)\rho\alpha}{4}\mathbb{E}[\|\widehat{x}^k - x^k\|^2] \\
&\quad + \Big[\frac{2\rho}{\rho + \mu_\ell} + \rho(\widetilde{\sigma}_f^2 + 3b_0^2)\Big]\alpha^2,
\end{aligned}
$$

where the last inequality is due to the assumption $b_k^2 \leq \alpha$.

## B.4  Proof of Corollary 1

Our proof departs from that of Theorem 2 through manipulating the descent estimate (64) in an alternative way. The key is to observe the following three-point inequality [3]:

$$
\langle h_f^k, x^{k+1} - x^\star\rangle \leq \frac{1}{2\alpha}\Big\{\|x^\star - x^k\|^2 - \|x^\star - x^{k+1}\|^2 - \|x^k - x^{k+1}\|^2\Big\}, \tag{75}
$$

where $x^\star$ is an optimal solution to (1). Observe that

$$
\langle\nabla\ell(x^k), x^{k+1} - x^k\rangle = \langle\nabla\ell(x^k) - h_f^k, x^{k+1} - x^\star\rangle + \langle h_f^k, x^{k+1} - x^\star\rangle + \langle\nabla\ell(x^k), x^\star - x^k\rangle. \tag{76}
$$

Notice that due to the convexity of $\ell(x)$, we have $\langle\nabla\ell(x^k), x^\star - x^k\rangle \leq -\text{OPT}^k$. Furthermore,

$$
\begin{aligned}
&\langle\nabla\ell(x^k) - h_f^k, x^{k+1} - x^\star\rangle \\
&= \langle\nabla\ell(x^k) - h_f^k + B_k + \overline{\nabla}_x f(x^k, y^{k+1}) - B_k - \overline{\nabla}_x f(x^k, y^{k+1}), x^{k+1} - x^\star\rangle \\
&\leq D_x\big\{b_k + L\|y^{k+1} - y^\star(x^k)\|\big\} + \langle B_k + \overline{\nabla}_x f(x^k, y^{k+1}) - h_f^k, x^{k+1} - x^k + x^k - x^\star\rangle.
\end{aligned}
$$

We notice that $\mathbb{E}[\langle B_k + \overline{\nabla}_x f(x^k, y^{k+1}) - h_f^k, x^k - x^\star\rangle | \mathcal{F}_k'] = 0$. Thus, taking the total expectation on both sides and applying Young's inequality on the last inner product lead to

$$
\mathbb{E}[\langle\nabla\ell(x^k) - h_f^k, x^{k+1} - x^\star\rangle \leq D_x\big\{b_k + L\mathbb{E}[\|y^{k+1} - y^\star(x^k)\|]\big\} + \frac{\alpha}{2}\sigma_f^2 + \frac{1}{2\alpha}\mathbb{E}[\|x^{k+1} - x^k\|^2]
$$

Substituting the above observations into (64) and using the three-point inequality (75) give

$$
\begin{aligned}
\mathbb{E}[\ell(x^{k+1}) - \ell(x^k)] &\leq D_x\big\{b_k + L\mathbb{E}[\|y^{k+1} - y^\star(x^k)\|]\big\} + \frac{1}{2\alpha}\Big\{\|x^\star - x^k\|^2 - \|x^\star - x^{k+1}\|^2\Big\} \\
&\quad + \frac{\alpha}{2}\sigma_f^2 - \text{OPT}^k + \frac{L_f}{2}\mathbb{E}[\|x^{k+1} - x^k\|^2].
\end{aligned}
$$

Summing up both sides from $k = 0$ to $k = K_{\max} - 1$ and dividing by $K_{\max}$ gives

$$\frac{1}{K_{\max}} \sum_{k=1}^{K_{\max}} \text{OPT}^k \leq D_x b_0 + \frac{\alpha \sigma_f^2}{2} + \frac{\|x^\star - x^0\|^2}{2\alpha K_{\max}} + \frac{D_x L}{K} \sum_{k=1}^{K_{\max}} \mathbb{E}[\|y^k - y^\star(x^{k-1})\|]$$

$$+ \frac{L_f}{2K_{\max}} \sum_{k=1}^{K_{\max}} \mathbb{E}[\|x^k - x^{k-1}\|^2].$$

Applying Cauchy-Schwartz inequality and Lemma 5, 6 with $\alpha = \mathcal{O}(K_{\max}^{-3/4})$, $\beta = \mathcal{O}(K_{\max}^{-1/2})$ as in (26) show that $\frac{1}{K_{\max}} \sum_{k=1}^{K_{\max}} \mathbb{E}[\|y^k - y^\star(x^{k-1})\|] \leq \sqrt{\frac{1}{K_{\max}} \sum_{k=1}^{K_{\max}} \Delta_y^k} = \mathcal{O}(K_{\max}^{-1/4})$; cf. (34). The proof is concluded.

## C  Justifications to Assumption 4–Assumption 7

In the following, we list these assumptions and provide explanations for when the assumptions are satisfied.

- (Assumption 4) The reward function is uniformly bounded by a constant $\bar{r}$. That is, $|r(s,a)| \leq \bar{r}$ for all $(s,a) \in S \times A$.

  This assumption merely states that the reward functions are uniformly bounded. This is a standard assumption used in MDP and reinforcement learning community. See, e.g., Chapter 2.2 of [55]. In practice, the reward functions are usually hand-crafted by the problem solver. They often encode the scores earned by the agent in each step, or whether some desired goal is reached.

- (Assumption 5) The feature map $\phi \colon S \times A \to \mathbb{R}^d$ satisfies $\|\phi(s,a)\|_2 \leq 1$ for all $(s,a) \in S \times A$. The action-value function associated with each policy is a linear function of $\phi$. That is, for any policy $\pi \in X$, there exists $\theta^\star(\pi) \in \mathbb{R}^d$ such that $Q^\pi(\cdot, \cdot) = \phi(\cdot, \cdot)^\top \theta^\star(\pi) = Q_{\theta^\star(\pi)}(\cdot, \cdot)$.

  This assumption assumes that the action-value function $Q^\pi$ is a linear function in a known feature mapping $\phi$ and $\phi$ is bounded. Such an assumption is standard in the literature on reinforcement learning with linear function approximation. See, e.g., Chapter 3.2 of [55]. In this line of research, it is oftentime postulated that $V^\pi(\cdot)$ or $Q^\pi(\cdot, \cdot)$ are linear functions of a known feature mapping.

  As for the feature mapping $\phi$, it is usually constructed based on domain knowledge. Some of the common examples include polynomial functions, on $[0, 1]$ radial basis function, and random features, which are all bounded. Here we assume that $\sup_{(s,a) \in S \times A} \|\phi(s,a)\|_2$ is bounded by one for simplicity, which can be replaced by any fixed parameter.

  Moreover, a concrete mathematical model that satisfies such a model is known as the linear MDP (see [28]), which assumes that both the reward function and the Markov transition kernel are linear in the given feature mapping $\phi$. Specifically, it is assumed that there exist $\lambda \in \mathbb{R}^d$ and $\mu \colon S \to \mathbb{R}^d$ such that

$$r(s,a) = \phi(s,a)^\top \lambda, \qquad P(s' \,|\, s,a) = \phi(s,a)^\top \mu(s') \qquad \forall (s,a,s') \in S \times A \times S. \tag{77}$$

Such a model includes the finite tabular MDP as a special case with $\phi(s,a)$ being the canonical vector $\boldsymbol{e}_{s,a}$ in $\mathbb{R}^{S \times A}$. Under the linear MDP assumption, for any policy $\pi$, the value functions $Q^\pi$ and $V^\pi$ exist and satisfy

$$V^\pi(s) = \sum_{a \in A} \pi(a \,|\, s) Q^\pi(s,a)$$

$$Q^\pi(s,a) = r(s,a) + \gamma \cdot \sum_{s' \in S} P(s' \,|\, s,a) \cdot V^\pi(s') = \phi(s,a)^\top \underbrace{\left( \lambda + \sum_{\in S} \mu(s') \cdot V^\pi(s') \right)}_{\theta^\star(\pi)}.$$

Thus (77) serves as a sufficient condition for the assumption.

- (Assumption 6) For each policy $\pi \in X$, the induced Markov chain $P^\pi$ admits a unique stationary distribution $\mu^\pi$ for all $\pi \in X$. Moreover, there exists $\mu_\phi > 0$ such that

$$\mathbb{E}_{s \sim \mu^\pi, a \sim \pi(\cdot|s)}[\phi(s,a)\phi^\top(s,a)] \succeq \mu_\phi^2 \cdot I_d, \ \forall \ \pi \in X.$$

The assumption that the Markov chain $P^\pi$ induced by any policy $\pi$ has a unique stationary distribution $\mu^\pi$ is a common assumption made in the literature on policy gradient. A sufficient condition ensures this is that all deterministic (stationary) policies visit all states eventually with probability one, i.e., the MDP is *unichain* (see Section 4.2.4 of [55]).

Furthermore, for asymptotic convergence analysis, classical RL literature often assumes that $\mathbb{E}_{s \sim \mu^\pi, a \sim \pi(\cdot|s)}[\phi(s,a)\phi(s,a)^\top]$ is invertible (see Section 4.4.2 of [55]; page 70). Here we additionally assumes that such a matrix is well-conditioned in the sense the smallest eigenvalue is lower bounded for nonasymptotic analysis. Such an assumption is also required for establishing statistical rates in linear regression.

In the tabular setting, a sufficient condition that justifying such an assumption is that the transition model is sufficient stochastic such that every policy induces $\pi$ induces a stationary distribution $\mu^\pi$ over $S$ such that the mass of $\mu^\pi$ on each state $s$ is lower bounded by $\mu_\phi > 0$.

- (Assumption 7) For any $(s,a) \in S \times A$ and any $\pi \in X$, let $\varrho(s,a,\pi)$ be a probability measure over $S$, defined by

$$[\varrho(s,a,\pi)](s') = (1-\gamma)^{-1} \sum_{t \geq 0} \gamma^t \cdot \mathbb{P}(s_t = s'), \qquad \forall s' \in S. \tag{78}$$

That is, $\varrho(s,a,\pi)$ is the visitation measure induced by the Markov chain starting from $(s_0, a_0) = (s,a)$ and follows $\pi$ afterwards. For any $\pi^\star$, there exists $C_\rho > 0$ such that

$$\mathbb{E}_{s' \sim \rho^{\pi^\star}}\left[ \left| \frac{\varrho(s,a,\pi)}{\rho^\star}(s') \right|^2 \right] \leq C_\rho^2, \quad \forall \ (s,a) \in S \times A, \ \pi \in X.$$

This assumption postulates that the distribution shift between the visitation measure induced by any policy $\pi$ and that induced by the optimal policy $\pi^*$ is bounded. Here the distribution shift is defined by the second-order moment of the density ratio.

Such an assumption is commonly made in reinforcement learning literature with various forms, which are referred to *concentrability coefficients* in general. It is conjectured in [11] that such an assumption is necessary for theoretical analysis. Moreover, our version is slightly weaker than that in [11], which essentially assumes the $\ell_\infty$-norm of the density ratio between $\varrho(s,a,\pi)$ and $\rho^\star$ is upper bounded. Moreover, a sufficient condition of Assumption 7 is that the initial

distribution $\rho_0$ is lower bounded everywhere over $S \times A$. Such a condition also appears in existing work, e.g., [1]. Note that $\rho^* \geq (1-\gamma) \cdot \rho_0$. Thus when the probability mass function of $\rho_0$ is lower bounded by $c_0$, Assumption 7 is satisfied with $C_\rho = (1-\gamma)^{-1} \cdot c_0^{-1}$.

# D Proof of Theorem 3

Hereafter, we let $\langle \cdot, \cdot \rangle$ and $\| \cdot \|$ denote the inner product and $\ell_1$-norm on $\mathbb{R}^{|A|}$, respectively. For any two policies $\pi_1$ and $\pi_2$, for any $s \in S$, $\|\pi_1(\cdot|s) - \pi_2(\cdot|s)\|_1$ is the total variation distance between $\pi_1(\cdot|s)$ and $\pi_2(\cdot|s)$. For any $f, f' \colon S \times A \to \mathbb{R}$, define the following norms:

$$\|f\|_{\rho,1} = \big[ \sum_{s \in S} \|f(s, \cdot)\|_1^2 \rho(s) \big]^{1/2}, \quad \|f\|_{\rho,\infty} = \big[ \sum_{s \in S} \|f(s, \cdot)\|_\infty^2 \rho(s) \big]^{1/2}$$

The following result can be derived from the Hölder's inequality:

$$\big| \langle f, f' \rangle_\rho \big| \leq \sum_{s \in S} \big| \langle f(s, \cdot), f'(s, \cdot) \rangle \big| \rho(s) \leq \|f\|_{\rho,1} \|f'\|_{\rho,\infty}. \tag{79}$$

Lastly, it can be shown that $\|\pi\|_{\rho,1} = 1, \|\pi\|_{\rho,\infty} \leq 1$.

Under Assumption 5, $\theta^\star(\pi)$ is the solution to the inner problem with $Q^\pi(\cdot, \cdot) = \phi(\cdot, \cdot)^\top \theta^\star(\pi)$. Below we first show that $\theta^\star(\pi)$ and $Q^\pi$ are Lipschitz continuous maps with respect to $\| \cdot \|_{\rho^\star, 1}$, where $\rho^\star$ is the visitation measure of an optimal policy $\pi^\star$.

**Lemma 8.** *Under Assumption 4–7, for any two policies $\pi_1, \pi_2 \in X$,*

$$\|Q^{\pi_1} - Q^{\pi_2}\|_{\rho^\star, \infty} \leq (1-\gamma)^{-2} \cdot \overline{r} \cdot C_\rho \cdot \|\pi_1 - \pi_2\|_{\rho^\star, 1},$$
$$\big\| \theta^\star(\pi_1) - \theta^\star(\pi_2) \big\|_2 \leq (1-\gamma)^{-2} \cdot \overline{r} \cdot C_\rho / \mu_\phi \cdot \|\pi_1 - \pi_2\|_{\rho^\star, 1}, \tag{80}$$

*where $\overline{r}$ is an upper bound on the reward function, $\mu_\phi$ is specified in Assumption 6, and $C_\rho$ is defined in Assumption 7.*

The proof of the above lemma is relegated to §D.1.

In the sequel, we first derive coupled inequalities on the non-negative sequences $\text{OPT}^k := \mathbb{E}[\ell(\pi^k) - \ell(\pi^\star)]$, $\Delta_Q^k := \mathbb{E}[\|\theta^k - \theta^\star(\pi^{k-1})\|^2]$, $\mathbb{E}[\|\pi^k - \pi^{k+1}\|_{\rho^\star, 1}^2]$, then we apply Lemma 5 to derive the convergence rates of TT-NAC. Using the performance difference lemma [cf. (47)], we obtain the following

$$\ell(\pi^{k+1}) - \ell(\pi^k) = -(1-\gamma)^{-1} \langle Q^{\pi^{k+1}}, \pi^{k+1} - \pi^\star \rangle_{\rho^\star} + (1-\gamma)^{-1} \langle Q^{\pi^k}, \pi^k - \pi^\star \rangle_{\rho^\star}$$
$$= (1-\gamma)^{-1} \langle -Q^{\pi^k}, \pi^{k+1} - \pi^k \rangle_{\rho^\star} + (1-\gamma)^{-1} \langle Q^{\pi^k} - Q^{\pi^{k+1}}, \pi^{k+1} - \pi^\star \rangle_{\rho^\star}. \tag{81}$$

Applying the inequality (79), we further have

$$\langle Q^{\pi^k} - Q^{\pi^{k+1}}, \pi^{k+1} - \pi^\star \rangle_{\rho^\star} \leq \|Q^{\pi^k} - Q^{\pi^{k+1}}\|_{\rho^\star, \infty} \|\pi^{k+1} - \pi^\star\|_{\rho^\star, 1} \tag{82}$$

$$\leq 2L_Q \|\pi^k - \pi^{k+1}\|_{\rho^\star, 1} \leq \frac{1-\gamma}{4\alpha} \|\pi^{k+1} - \pi^k\|_{\rho^\star, 1}^2 + \frac{4L_Q^2 \alpha}{1-\gamma}, \tag{83}$$

where $L_Q := (1-\gamma)^{-2} \overline{r} \cdot C_\rho$. The above inequality follows from $\|\pi^k - \pi^{k+1}\|_{\rho^\star, 1} \leq 2$ for any $\pi_1, \pi_2 \in X$ and applying Lemma 8. Then, combining (81), (83) leads to

$$\ell(\pi^{k+1}) - \ell(\pi^k) \leq \frac{-1}{1-\gamma} \langle Q^{\pi^k}, \pi^{k+1} - \pi^k \rangle_{\rho^\star} + \frac{1}{4\alpha} \|\pi^{k+1} - \pi^k\|_{\rho^\star, 1}^2 + 4 \frac{\alpha L_Q^2}{(1-\gamma)^2}. \tag{84}$$

Let us bound the first term in the right-hand side of (84). To proceed, note that the policy update (42) can be implemented for each state individually as below:

$$\pi^{k+1}(\cdot|s) = \underset{\nu:\ \sum_a \nu(a)=1, \nu(a)\geq 0}{\arg\min} \left\{ -(1-\gamma)^{-1} \cdot \langle Q_{\theta^{k+1}}(s,\cdot), \nu \rangle + 1/\alpha_k \cdot D_\psi\big(\nu, \pi^k(\cdot|s)\big) \right\}, \qquad (85)$$

for all $s \in S$. Observe that we can modify $\rho^{\pi^k}$ in (42) to $\rho^\star$ without changing the optimal solution for this subproblem. Specifically, (42) can be written as

$$\pi^{k+1} = \underset{\pi \in X}{\arg\min}\left\{ -(1-\gamma)^{-1}\langle Q_{\theta^{k+1}}, \pi - \pi^k \rangle_{\rho^\star} + \frac{1}{\alpha}\bar{D}_{\psi,\rho^\star}(\pi, \pi^k) \right\}. \qquad (86)$$

We have

$$-(1-\gamma)^{-1}\langle Q^{\pi^k}, \pi^{k+1} - \pi^k \rangle_{\rho^*} = (1-\gamma)^{-1}\big[\langle Q_{\theta^{k+1}} - Q^{\pi^k}, \pi^{k+1} - \pi^k \rangle_{\rho^\star} - \langle Q_{\theta^{k+1}}, \pi^{k+1} - \pi^k \rangle_{\rho^\star}\big].$$

Furthermore, from (86), we obtain

$$\frac{\langle Q_{\theta^{k+1}}, \pi^{k+1} - \pi^k \rangle_{\rho^\star}}{1-\gamma} \geq \frac{1}{\alpha}\sum_{s\in S}\langle \nabla D_\psi(\pi^{k+1}(\cdot|s), \pi^k(\cdot|s)), \pi^{k+1}(\cdot|s) - \pi^k(\cdot|s)\rangle\rho^\star(s), \qquad (87)$$

where the inequality follows from the optimality condition of the mirror descent step. Meanwhile, the 1-strong convexity of $D_\psi(\cdot, \cdot)$ implies that

$$\big\langle \nabla D_\psi\big(\pi^{k+1}(\cdot|s), \pi^k(\cdot|s)\big), \pi^{k+1}(\cdot|s) - \pi^k(\cdot|s)\big\rangle \geq \|\pi^{k+1}(\cdot|s) - \pi^k(\cdot|s)\|^2. \qquad (88)$$

Thus, combining (87) and (88), and applying Young's inequality, we further have

$$-(1-\gamma)^{-1}\langle Q^{\pi^k}, \pi^{k+1} - \pi^k\rangle_{\rho^*}$$
$$\leq \frac{1}{4\alpha}\|\pi^{k+1} - \pi^k\|^2_{\rho^\star,1} + \alpha(1-\gamma)^{-2}\cdot\|Q^{\pi^k} - Q_{\theta^{k+1}}\|^2_{\rho^\star,\infty} - \frac{1}{\alpha}\|\pi^{k+1} - \pi^k\|^2_{\rho^\star,1}$$
$$= \alpha(1-\gamma)^{-2}\cdot\|Q^{\pi^k} - Q_{\theta^{k+1}}\|^2_{\rho^\star,\infty} - \frac{3}{4\alpha}\|\pi^{k+1} - \pi^k\|^2_{\rho^\star,1}. \qquad (89)$$

By direct computation and using $\|\phi(s,a)\| \leq 1$ [cf. Assumption 5], we have

$$\|Q^{\pi^k} - Q_{\theta^{k+1}}\|^2_{\rho^\star,\infty} = \sum_{s\in S}\left\{\max_{a\in A}|\phi(s,a)^\top[\theta^\star(\pi^k) - \theta^{k+1}]|\right\}^2\rho^\star(s)$$
$$\leq \sum_{s\in S}\max_{a\in A}\{\|\phi(s,a)\|^2\}\|\theta^\star(\pi^k) - \theta^{k+1}\|^2\rho^\star(s) \leq \|\theta^\star(\pi^k) - \theta^{k+1}\|^2. \qquad (90)$$

Combining (84), (89), and (90), we obtain

$$\ell(\pi^{k+1}) - \ell(\pi^k) \leq \alpha(1-\gamma)^{-2}\|\theta^\star(\pi^k) - \theta^{k+1}\|^2 - \frac{1}{2\alpha}\|\pi^{k+1} - \pi^k\|^2_{\rho^\star,1} + 4(1-\gamma)^{-2}L_Q^2\alpha.$$

Taking full expectation leads to

$$\mathrm{OPT}^{k+1} - \mathrm{OPT}^k \leq \alpha(1-\gamma)^{-2}\Delta_Q^{k+1} - \frac{1}{2\alpha}\mathbb{E}[\|\pi^{k+1} - \pi^k\|^2_{\rho^\star,1}] + 4(1-\gamma)^{-2}L_Q^2\alpha. \qquad (91)$$

Next, we consider the convergence of $\Delta_Q^k$. Let $\mathcal{F}_k = \sigma\{\theta^0, \pi^0, \ldots, \theta^k, \pi^k\}$ be the $\sigma$-algebra generated by the first $k+1$ actor and critic updates. Under Assumption 5, we can write the

conditional expectation of $h_g^k$ as

$$\mathbb{E}[h_g^k|\mathcal{F}_k] = \mathbb{E}_{\mu^{\pi^k}}\left[\{Q_{\theta^k}(s,a) - r(s,a) - \gamma Q_{\theta^k}(s',a')\}\phi(s,a)|\mathcal{F}_k\right]$$
$$= \mathbb{E}_{\mu^{\pi^k}}\left[\phi(s,a)\{\phi(s,a) - \gamma\phi(s',a')\}^\top\right][\theta^k - \theta^\star(\pi^k)], \tag{92}$$

where $\mathbb{E}_{\mu^{\pi^k}}[\cdot]$ denotes the expectation taken with $s \sim \mu^{\pi^k}$, $a \sim \pi^k(\cdot|s)$, $s' \sim P(\cdot|s,a)$, $a' \sim \pi^k(\cdot|s')$. Under Assumption 5 and 6, Lemma 3 of [4] shows that $\mathbb{E}[h_g^k|\mathcal{F}_k]$ is a semigradient of the MSBE function $\|Q_{\theta^k} - Q_{\theta^\star(\pi^k)}\|_{\mu^{\pi^k}\otimes\pi^k}^2$. Particularly, we obtain

$$\mathbb{E}[h_g^k|\mathcal{F}_k]^\top[\theta^k - \theta^\star(\pi^k)] \geq (1-\gamma)\|Q_{\theta^k} - Q_{\theta^\star(\pi^k)}\|_{\mu^{\pi^k}\otimes\pi^k}^2 \geq \mu_{\mathsf{td}}\|\theta^k - \theta^\star(\pi^k)\|_2^2, \tag{93}$$

where we have let $\mu_{\mathsf{td}} = (1-\gamma)\mu_\phi^2$. Moreover, Lemma 5 of [4] demonstrates that the second order moment $\mathbb{E}[\|h_g^k\|_2^2|\mathcal{F}_k]$ is bounded as

$$\mathbb{E}[\|h_g^k\|_2^2|\mathcal{F}_k] \leq 8\|Q_{\theta^k} - Q_{\theta^\star(\pi^k)}\|_{\mu^{\pi^k}\otimes\pi^k}^2 + \sigma_{\mathsf{td}}^2 \leq 8\|\theta^k - \theta^\star(\pi^k)\|_2^2 + \sigma_{\mathsf{td}}^2, \tag{94}$$

where $\sigma_{\mathsf{td}}^2 = 4\bar{r}^2(1-\gamma)^{-2}$. Combining (93), (94) and recalling $\beta \leq \mu_{\mathsf{td}}/8$, it holds

$$\mathbb{E}[\|\theta^{k+1} - \theta^\star(\pi^k)\|_2^2|\mathcal{F}_k] = \|\theta^k - \theta^\star(\pi^k)\|_2^2 - 2\beta\mathbb{E}[h_g^k|\mathcal{F}_k]^\top[\theta^k - \theta^\star(\pi^k)] + \beta^2\mathbb{E}[\|h_g^k\|_2^2|\mathcal{F}_k]$$
$$\leq \left(1 - 2\mu_{\mathsf{td}}\beta + 8\beta^2\right) \cdot \|\theta^k - \theta^\star(\pi^k)\|_2^2 + \beta^2 \cdot \sigma_{\mathsf{td}}^2$$
$$\leq \left(1 - \mu_{\mathsf{td}}\beta\right) \cdot \|\theta^k - \theta^\star(\pi^k)\|_2^2 + \beta^2 \cdot \sigma_{\mathsf{td}}^2. \tag{95}$$

By Young's inequality and Lemma 8, we further have

$$\mathbb{E}[\|\theta^{k+1} - \theta^\star(\pi^k)\|_2^2|\mathcal{F}_k]$$
$$\leq (1+c)(1-\mu_{\mathsf{td}}\beta)\|\theta^k - \theta^\star(\pi^{k-1})\|_2^2 + (1+1/c)\|\theta^\star(\pi^k) - \theta^\star(\pi^{k-1})\|_2^2 + \beta^2\sigma_{\mathsf{td}}^2$$
$$\leq (1-\mu_{\mathsf{td}}\beta/2)\|\theta^k - \theta^\star(\pi^{k-1})\|_2^2 + \left(\frac{2}{\mu_{\mathsf{td}}\beta} - 1\right)\overline{L}_Q\|\pi^k - \pi^{k+1}\|_{\rho^\star,1}^2 + \beta^2\sigma_{\mathsf{td}}^2, \tag{96}$$

where we have chosen $c > 0$ such that $(1+c)(1-\mu_{\mathsf{td}}\beta) = 1 - \mu_{\mathsf{td}}\beta/2$, which implies that $1/c + 1 = 2/(\mu_{\mathsf{td}}\beta) - 1 > 0$ [cf. (45)]; The last inequality comes from Lemma 8 with the constant $\overline{L}_Q = (1-\gamma)^{-4} \cdot \bar{r} \cdot C_\rho^2 \cdot \mu_\phi^{-2}$.

From (91), (96), we identify that condition (31) of Lemma 5 holds with:

$$\Omega^k = \mathrm{OPT}^k, \ \Theta^k = \mathbb{E}[\|\pi^k - \pi^{k-1}\|_{\rho^*,1}^2], \ \mathrm{c}_0 = \frac{1}{2\alpha}, \ \mathrm{c}_1 = \frac{\alpha}{(1-\gamma)^2}, \ \mathrm{c}_2 = \frac{4L_Q^2\alpha}{(1-\gamma)^2},$$

$$\Upsilon^k = \mathbb{E}[\|\theta^k - \theta^\star(\pi^{k-1})\|^2], \ \mathrm{d}_0 = \mu_{\mathsf{td}}\beta/2, \ \mathrm{d}_1 = \left(\frac{2}{\mu_{\mathsf{td}}\beta} - 1\right)\overline{L}_Q > 0, \ \mathrm{d}_2 = \beta^2 \cdot \sigma_{\mathsf{td}}^2.$$

Selecting the step sizes as in (45), one can verify that $\frac{\alpha}{\beta} < \frac{\mu_{\mathsf{td}}(1-\gamma)}{16\sqrt{\overline{L}_Q}}$. This ensures

$$\mathrm{c}_0 - \mathrm{c}_1\mathrm{d}_1(\mathrm{d}_0)^{-1} > 1/(4\alpha), \ \ \mathrm{d}_0 - \mathrm{c}_1\mathrm{d}_1(\mathrm{c}_0)^{-1} > \mu_{\mathsf{td}}\beta/4. \tag{97}$$

Applying Lemma 5, we obtain for any $K \geq 1$ that

$$\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\pi^k - \pi^{k+1}\|_{\rho^\star,1}^2] \leq \frac{\text{OPT}^0 \cdot 4\alpha + \frac{8\alpha^2(1-\gamma)^{-2}}{\mu_{\text{td}}\beta}(\Delta_Q^0 + \beta^2\sigma_{\text{td}}^2)}{K} + \frac{8\alpha^2(4L_Q^2 + \beta\sigma_{\text{td}}^2/\mu_{\text{td}})}{(1-\gamma)^2}$$

$$\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\Delta_Q^{k+1}] \leq \frac{\mathbb{E}[\Delta_Q^0] + \beta^2\sigma_{\text{td}}^2 + \frac{4\alpha}{\mu_{\text{td}}\beta}\bar{L}_Q\text{OPT}^0}{\mu_{\text{td}}\beta K/4} + \frac{\beta^2\sigma_{\text{td}}^2 + \frac{16\alpha^2}{\mu_{\text{td}}\beta}(1-\gamma)^{-2}L_Q^2\alpha^2}{\mu_{\text{td}}\beta/4}.$$

Particularly, plugging in $\alpha \asymp K_{\text{max}}^{-3/4}$, $\beta \asymp K_{\text{max}}^{-1/2}$ shows that the convergence rates are $K_{\text{max}}^{-1}\sum_{k=1}^{K_{\text{max}}}\mathbb{E}[\|\pi^k - \pi^{k+1}\|_{\rho^\star,1}^2] = \mathcal{O}(K_{\text{max}}^{-3/2})$, $K_{\text{max}}^{-1}\sum_{k=1}^{K_{\text{max}}}\mathbb{E}[\Delta_Q^{k+1}] = \mathcal{O}(K_{\text{max}}^{-1/2})$.

Our last step is to analyze the convergence rate of the objective value $\text{OPT}^k$. To this end, we observe the following three-point inequality [3]

$$\frac{-\langle Q_{\theta^{k+1}}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star}}{1-\gamma} \leq \frac{1}{\alpha}\left[\bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^k) - \bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^{k+1}) - \bar{D}_{\psi,\rho^\star}(\pi^{k+1}, \pi^k)\right]. \qquad (98)$$

Meanwhile, by the inequalities (79), (81), (82), we have

$$\ell(\pi^{k+1}) - \ell(\pi^k) = -\frac{1}{1-\gamma}\langle Q^{\pi^k}, \pi^{k+1} - \pi^k\rangle_{\rho^\star} + \frac{1}{1-\gamma}\langle Q^{\pi^k} - Q^{\pi^{k+1}}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star}$$

$$\leq \frac{1}{1-\gamma}\left[\langle -Q^{\pi^k}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star} - \langle Q^{\pi^k}, \pi^\star - \pi^k\rangle_{\rho^\star} + \|\pi^{k+1} - \pi^\star\|_{\rho^\star,1}\|Q^{\pi^k} - Q^{\pi^{k+1}}\|_{\rho^\star,\infty}\right]$$

$$\leq \frac{1}{1-\gamma}\left[\langle -Q^{\pi^k}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star} - \langle Q^{\pi^k}, \pi^\star - \pi^k\rangle_{\rho^\star} + 2L_Q\|\pi^{k+1} - \pi^k\|_{\rho^\star,1}\right],$$

where the last inequality follows from Lemma 8. Now, with the performance difference lemma $\ell(\pi^*) - \ell(\pi^k) = (1-\gamma)^{-1}\langle -Q^{\pi^k}, \pi^\star - \pi^k\rangle_{\rho^\star}$, the above simplifies to

$$\ell(\pi^{k+1}) - \ell(\pi^\star) \leq (1-\gamma)^{-1}\left[-\langle Q^{\pi^k}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star} + 2L_Q\|\pi^{k+1} - \pi^k\|_{\rho^\star,1}\right]$$

With $\langle Q^{\pi^k}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star} = \langle Q^{\pi^k} - Q_{\theta^{k+1}}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star} + \langle Q_{\theta^{k+1}}, \pi^{k+1} - \pi^\star\rangle_{\rho^\star}$ and applying the three-point inequality (98), we have

$$\ell(\pi^{k+1}) - \ell(\pi^\star) \leq \frac{2}{1-\gamma}\left[L_Q\|\pi^{k+1} - \pi^k\|_{\rho^\star,1} + \|Q^{\pi^k} - Q_{\theta^{k+1}}\|_{\rho^\star,\infty}\right]$$

$$+ \frac{1}{\alpha}\left[\bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^k) - \bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^{k+1}) - \bar{D}_{\psi,\rho^\star}(\pi^{k+1}, \pi^k)\right].$$

$$\leq \frac{2}{1-\gamma}\left[L_Q\|\pi^{k+1} - \pi^k\|_{\rho^\star,1} + \|\theta^\star(\pi^k) - \theta^{k+1}\|_2\right] + \frac{1}{\alpha}\left[\bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^k) - \bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^{k+1})\right]$$

where the last inequality uses (90) and the fact that $\bar{D}_{\psi,\rho^\star}$ is non-negative. Finally, taking the full expectation on both sides of the inequality, we obtain

$$\text{OPT}^{k+1} \leq 2(1-\gamma)^{-1}\mathbb{E}[\|\theta^\star(\pi^k) - \theta^{k+1}\|_2] + 2(1-\gamma)^{-1}L_Q\mathbb{E}[\|\pi^{k+1} - \pi^k\|_{\rho^\star,1}]$$

$$+ \frac{1}{\alpha}\mathbb{E}\left[\bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^k) - \bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^{k+1})\right]. \qquad (99)$$

Summing up both sides from $k = 0$ to $k = K_{\mathsf{max}} - 1$ and dividing by $K_{\mathsf{max}}$ yields

$$\frac{1}{K_{\mathsf{max}}} \sum_{k=1}^{K_{\mathsf{max}}} \mathrm{OPT}^k \leq \frac{1}{\alpha K_{\mathsf{max}}} \big\{ \bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^0) - \bar{D}_{\psi,\rho^\star}(\pi^\star, \pi^{K_{\mathsf{max}}}) \big\}$$

$$+ \frac{2}{(1-\gamma)} \frac{1}{K_{\mathsf{max}}} \sum_{k=1}^{K_{\mathsf{max}}} \big\{ \mathbb{E}[\|\theta^\star(\pi^{k-1}) - \theta^k\|_2] + L_Q \cdot \mathbb{E}[\|\pi^k - \pi^{k-1}\|_{\rho^\star,1}] \big\}. \tag{100}$$

Using Cauchy-Schwarz's inequality, it can be easily seen that the right-hand side is $\mathcal{O}(K_{\mathsf{max}}^{-1/4})$. This concludes the proof of the theorem.

## D.1   Proof of Lemma 8

We first bound $|Q^{\pi_1}(s,a) - Q^{\pi_2}(s,a)|$. By the Bellman equation (38) and the performance difference lemma (47), we have

$$Q^{\pi_1}(s,a) - Q^{\pi_2}(s,a) = \sum_{s' \in S} P(s'|s,a) \cdot \big[ V^{\pi_1}(s') - V^{\pi_2}(s') \big]$$

$$= (1-\gamma)^{-1} \sum_{s' \in S} P(s'|s,a) \cdot \mathbb{E}_{\widetilde{s} \sim \widetilde{\varrho}(s',\pi_1)} \big[ \langle Q^{\pi_2}(\widetilde{s}, \cdot), \pi_1(\cdot|\widetilde{s}) - \pi_2(\cdot|\widetilde{s}) \rangle \big], \tag{101}$$

where $\widetilde{\varrho}(s', \pi_1)$ is the visitation measure obtained by the Markov chain induced by $\pi_1$ with the initial state fixed to $s'$. Recall the definition of the visitation measure $\varrho(s,a,\pi)$ in (78). We rewrite (101) as

$$Q^{\pi_1}(s,a) - Q^{\pi_2}(s,a) = (1-\gamma)^{-1} \cdot \mathbb{E}_{\widetilde{s} \sim \varrho(s,a,\pi_1)} \big[ \langle Q^{\pi_2}(\widetilde{s}, \cdot), \pi_1(\cdot|\widetilde{s}) - \pi_2(\cdot|\widetilde{s}) \rangle \big]. \tag{102}$$

Moreover, notice that $\sup_{(s,a) \in S \times A} |Q^\pi(s,a)| \leq (1-\gamma)^{-1} \cdot \bar{r}$ under Assumption 4. Then, applying Hölder's inequality to (102), we obtain

$$\big| Q^{\pi_1}(s,a) - Q^{\pi_2}(s,a) \big| \leq (1-\gamma)^{-1} \cdot \mathbb{E}_{\widetilde{s} \sim \widetilde{\varrho}(s',\pi_1)} \big[ \|Q^{\pi_2}(\widetilde{s}, \cdot)\|_\infty \cdot \|\pi_1(\cdot|\widetilde{s}) - \pi_2(\cdot|\widetilde{s})\|_1 \big]$$

$$\leq (1-\gamma)^{-2} \cdot \bar{r} \cdot \mathbb{E}_{\widetilde{s} \sim \rho^\star} \left[ \frac{\varrho(s,a,\pi)}{\rho^\star}(\widetilde{s}) \cdot \|\pi_1(\cdot|\widetilde{s}) - \pi_2(\cdot|\widetilde{s})\|_1 \right]$$

$$\leq (1-\gamma)^{-2} \cdot \bar{r} \cdot \left\{ \mathbb{E}_{\widetilde{s} \sim \rho^\star} \left[ \Big| \frac{\varrho(s,a,\pi)}{\rho^\star}(\widetilde{s}) \Big|^2 \right] \mathbb{E}_{\widetilde{s} \sim \rho^\star} \big[ \|\pi_1(\cdot|\widetilde{s}) - \pi_2(\cdot|\widetilde{s})\|_1^2 \big] \right\}^{1/2}$$

$$\leq (1-\gamma)^{-2} \cdot \bar{r} \cdot C_\rho \cdot \|\pi_1 - \pi_2\|_{\rho^\star,1}, \tag{103}$$

where the second inequality is from the boundedness of $Q^\pi$, the third one is the Cauchy-Schwarz inequality, and the last one is from Assumption 7. Finally, we have

$$\|Q^{\pi_1} - Q^{\pi_2}\|_{\rho^\star,\infty} \leq (1-\gamma)^{-2} \cdot \bar{r} \cdot C_\rho \cdot \|\pi_1 - \pi_2\|_{\rho^\star,1}.$$

It remains to bound $\|\theta^\star(\pi_1) - \theta^\star(\pi_2)\|^2$. Under Assumption 5, we have

$$\|Q^{\pi_1} - Q^{\pi_2}\|_{\mu^{\pi^\star} \otimes \pi^\star}^2 = \mathbb{E}_{s \sim \mu^{\pi^\star}, a \sim \pi^\star(\cdot|s)} \big\{ \big[ Q^{\pi_1}(s,a) - Q^{\pi_2}(s,a) \big]^2 \big\}$$

$$= \mathbb{E}_{s \sim \mu^{\pi^\star}, a \sim \pi^\star(\cdot|s)} \Big( \big\{ \phi(s,a)^\top [\theta^\star(\pi_1) - \theta^\star(\pi_2)] \big\}^2 \Big)$$

$$= [\theta^\star(\pi_1) - \theta^\star(\pi_2)]^\top \Sigma_{\pi^\star} [\theta^\star(\pi_1) - \theta^\star(\pi_2)]. \tag{104}$$

Then, combining Assumption 6 and (103), we have

$$\mu_\phi^2 \|\theta^\star(\pi_1) - \theta^\star(\pi_2)\|^2 \leq \|Q^{\pi_1} - Q^{\pi_2}\|_{\mu^{\pi^\star} \otimes \pi^\star}^2 \leq (1-\gamma)^{-4} \cdot \bar{r}^2 \cdot C_\rho^2 \cdot \|\pi_1 - \pi_2\|_{\rho^\star, 1}^2, \qquad (105)$$

which yields the second inequality in Lemma 8. We conclude the proof.

# E   Auxiliary Lemmas

The proofs for the lemmas below can be found in the online appendix [24].

**Lemma 9.** *[31, Lemma 12] Let $a > 0$, $\{\gamma_j\}_{j \geq 0}$ be a non-increasing, non-negative sequence such that $\gamma_0 < 1/a$, it holds for any $k \geq 0$ that*

$$\sum_{j=0}^{k} \gamma_j \prod_{\ell=j+1}^{k} (1 - \gamma_\ell a) \leq \frac{1}{a}. \qquad (106)$$

**Lemma 10.** *Fix a real number $1 < q \leq 2$. Let $a > 0$, $\{\gamma_j\}_{j \geq 0}$ be a non-increasing, non-negative sequence such that $\gamma_0 < 1/(2a)$. Suppose that $\frac{\gamma_{\ell-1}}{\gamma_\ell} \leq 1 + \frac{a}{2(q-1)}\gamma_\ell$. Then, it holds for any $k \geq 0$ that*

$$\sum_{j=0}^{k} \gamma_j^q \prod_{\ell=j+1}^{k} (1 - \gamma_\ell a) \leq \frac{2}{a}\gamma_k^{q-1}. \qquad (107)$$

**Lemma 11.** *Fix the real numbers $a, b > 0$. Let $\{\gamma_j\}_{j \geq 0}, \{\rho_j\}_{j \geq 0}$ be nonincreasing, non-negative sequences such that $2a\gamma_j \leq b\rho_j$ for all $j$, and $\rho_0 < 1/b$. Then, it holds that*

$$\sum_{j=0}^{k} \gamma_j \prod_{\ell=j+1}^{k} (1 - \gamma_\ell a) \prod_{i=0}^{j} (1 - \rho_i b) \leq \frac{1}{a} \prod_{\ell=0}^{k} (1 - \gamma_\ell a), \ \forall \ k \geq 0. \qquad (108)$$

# F   Technical Results Omitted from the Main Paper

## F.1   Proof of Lemma 10

To derive this result, we observe that

$$\sum_{j=0}^{k} \gamma_j^q \prod_{\ell=j+1}^{k} (1 - \gamma_\ell a) \leq \gamma_k^{q-1} \sum_{j=0}^{k} \gamma_j \frac{\gamma_j^{q-1}}{\gamma_k^{q-1}} \prod_{\ell=j+1}^{k} (1 - \gamma_\ell a)$$

$$= \gamma_k^{q-1} \sum_{j=0}^{k} \gamma_j \prod_{\ell=j+1}^{k} \left(\frac{\gamma_{\ell-1}}{\gamma_\ell}\right)^{q-1} (1 - \gamma_\ell a).$$

Furthermore, from the conditions on $\gamma_\ell$,

$$\left(\frac{\gamma_{\ell-1}}{\gamma_\ell}\right)^{q-1} (1 - \gamma_\ell a) \leq \left(1 + \frac{a}{2(q-1)}\gamma_\ell\right)^{q-1} (1 - \gamma_\ell a) \leq 1 - \frac{a}{2}\gamma_\ell.$$

Therefore,

$$\sum_{j=0}^{k} \gamma_j^q \prod_{\ell=j+1}^{k} (1 - \gamma_\ell a) \leq \gamma_k^{q-1} \sum_{j=0}^{k} \gamma_j \prod_{\ell=j+1}^{k} (1 - \frac{a}{2} \gamma_\ell) \leq \frac{2}{a} \gamma_k^{q-1}.$$

This concludes the proof.

## F.2 Proof of Lemma 11

First of all, the condition $2a\gamma_j \leq b\rho_j$ implies

$$\frac{1 - \rho_i b}{1 - \gamma_i a} \leq 1 - \rho_i b/2, \ \forall \ i \geq 0.$$

As such, we observe $\prod_{i=0}^{j} \frac{1 - \rho_i b}{1 - \gamma_i a} \leq \prod_{i=0}^{j} (1 - \rho_i b/2)$ and subsequently,

$$\sum_{j=0}^{k} \gamma_j \prod_{\ell=j+1}^{k} (1 - \gamma_\ell a) \prod_{i=0}^{j} (1 - \rho_i b) \leq \Big[ \prod_{\ell=0}^{k} (1 - \gamma_\ell a) \Big] \sum_{j=0}^{k} \gamma_j \prod_{i=0}^{j} (1 - \rho_i b/2).$$

Furthermore, for any $j = 0, ..., k$, it holds

$$\rho_j \prod_{i=0}^{j-1} (1 - \rho_i b/2) = \frac{2}{b} \Big[ \prod_{i=0}^{j-1} (1 - \rho_i b/2) - \prod_{i=0}^{j} (1 - \rho_i b/2) \Big], \tag{109}$$

where we have taken the convention $\prod_{i=0}^{-1} (1 - \rho_i b/2) = 1$. We obtain that

$$\sum_{j=0}^{k} \gamma_j \prod_{i=0}^{j} (1 - \rho_i b/2) \leq \frac{b}{2a} \sum_{j=0}^{k} \rho_j \prod_{i=0}^{j-1} (1 - \rho_i b/2)$$

$$= \frac{1}{a} \sum_{j=0}^{k} \Big[ \prod_{i=0}^{j-1} (1 - \rho_i b/2) - \prod_{i=0}^{j} (1 - \rho_i b/2) \Big] \leq \frac{1}{a},$$

where the last inequality follows from the bound $1 - \prod_{i=0}^{k} (1 - \rho_i \mu_g b/2) \leq 1$. Combining with the above inequality yields the desired results.

## F.3 Proof of Lemma 1

*Proof.* Since the samples are drawn independently, the expected value of $h_f^k$ is

$$\mathbb{E}[h_f^k] = \nabla_x f(x, y) - \nabla_{xy}^2 g(x, y) \mathbb{E} \Big[ \frac{\mathsf{t_{max}} \mu_g}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \prod_{i=1}^{\mathsf{p}} \big( I - \frac{\mu_g}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \nabla_{yy}^2 g(x, y; \xi_i^{(2)}) \big) \Big] \nabla_y f(x, y). \tag{110}$$

We have
$$\|\overline{\nabla}_x f(x, y) - \mathbb{E}[h_f^k]\|$$
$$= \Big\| \nabla_{xy}^2 g(x, y) \big\{ \mathbb{E} \big[ \frac{\mathsf{t_{max}} \mu_g}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \prod_{i=1}^{\mathsf{p}} \big( I - \frac{\mu_g}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \nabla_{yy}^2 g(x, y; \xi_i^{(2)}) \big) \big] - [\nabla_{yy} g(x, y)]^{-1} \big\} \nabla_y f(x, y) \Big\|$$
$$\leq C_{gxy} C_{fy} \Big\| \mathbb{E} \big[ \frac{\mathsf{t_{max}} \mu_g}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \prod_{i=1}^{\mathsf{p}} \big( I - \frac{\mu_g}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \nabla_{yy}^2 g(x, y; \xi_i^{(2)}) \big) \big] - [\nabla_{yy} g(x, y)]^{-1} \Big\|,$$

where the last inequality follows from Assumption 1-3 and Assumption 2-5.

Applying [23, Lemma 3.2], the latter norm can be bounded by $\frac{1}{\mu_g}\left(1 - \frac{\mu_g^2}{L_g(\mu_g^2+\sigma_{gxy}^2)}\right)^{\mathsf{t}_{\max}}$. This concludes the proof for the first part of the lemma.

It remains to bound the variance of $h_f^k$. We first let

$$H_{yy} = \frac{\mathsf{t}_{\max}\,\mu_g}{L_g(\mu_g^2+\sigma_{gxy}^2)}\prod_{i=1}^{\mathsf{p}}\left(I - \frac{\mu_g}{L_g(\mu_g^2+\sigma_{gxy}^2)}\nabla_{yy}^2 g(x,y;\xi_i^{(2)})\right).$$

To estimate the variance of $h_f^k$, using (110), we observe that

$$\mathbb{E}[\|h_f^k - \mathbb{E}[h_f^k]\|^2] = \mathbb{E}[\|\nabla_x f(x,y;\xi^{(1)}) - \nabla_x f(x,y)\|^2]$$
$$+ \mathbb{E}\left[\left\|\nabla_{xy}^2 g(x,y;\xi_0^{(2)})H_{yy}\nabla_y f(x,y;\xi^{(1)}) - \nabla_{xy}^2 g(x,y)\mathbb{E}[H_{yy}]\nabla_y f(x,y)\right\|^2\right].$$

The first term on the right hand side can be bounded by $\sigma_{fx}^2$. Furthermore

$$\nabla_{xy}^2 g(x,y;\xi_0^{(2)})H_{yy}\nabla_y f(x,y;\xi^{(1)}) - \nabla_{xy}^2 g(x,y)\mathbb{E}[H_{yy}]\nabla_y f(x,y)$$
$$= \left\{\nabla_{xy}^2 g(x,y;\xi_0^{(2)}) - \nabla_{xy}^2 g(x,y)\right\}H_{yy}\nabla_y f(x,y;\xi^{(1)})$$
$$+ \nabla_{xy}^2 g(x,y)\{H_{yy} - \mathbb{E}[H_{yy}]\}\nabla_y f(x,y;\xi^{(1)})$$
$$+ \nabla_{xy}^2 g(x,y)\mathbb{E}[H_{yy}]\{\nabla_y f(x,y;\xi^{(1)}) - \nabla_y f(x,y)\}.$$

We also observe

$$\mathbb{E}[\|\nabla_y f(x,y;\xi^{(1)})\|^2] = \mathbb{E}[\|\nabla_y f(x,y;\xi^{(1)}) - \nabla_y f(x,y)\|^2] + \mathbb{E}[\|\nabla_y f(x,y)\|^2] \leq \sigma_{fy}^2 + C_y^2.$$

Using $(a+b+c)^2 \leq 3(a^2+b^2+c^2)$ and the Cauchy-Schwarz's inequality, we have

$$\mathbb{E}\left[\left\|\nabla_{xy}^2 g(x,y;\xi_0^{(2)})H_{yy}\nabla_y f(x,y;\xi^{(1)}) - \nabla_{xy}^2 g(x,y)\mathbb{E}[H_{yy}]\nabla_y f(x,y)\right\|^2\right]$$
$$\leq 3\Big\{(\sigma_{fy}^2 + C_y^2)\{\sigma_{gxy}^2\mathbb{E}[\|H_{yy}\|^2] + C_{gxy}^2\mathbb{E}[\|H_{yy} - \mathbb{E}[H_{yy}]\|^2]\} + \sigma_{fy}^2 C_{gxy}^2 \cdot \|\mathbb{E}[H_{yy}]\|^2\Big\}.$$

Next, we observe that

$$\mathbb{E}[\|H_{yy}\|^2] = \frac{\mu_g^2}{L_g^2(\mu_g^2+\sigma_{gxy}^2)^2}\sum_{p=0}^{\mathsf{t}_{\max}-1}\mathbb{E}\left[\left\|\prod_{i=1}^{p}\left(I - \frac{\mu_g}{L_g(\mu_g^2+\sigma_{gxy}^2)}\nabla_{yy}^2 g(x,y;\xi_i^{(2)})\right)\right\|^2\right] \quad (111)$$

Observe that the product of random matrices satisfies the conditions in Lemma 12 with $\mu \equiv \frac{\mu_g^2}{L_g(\mu_g^2+\sigma_{gxy}^2)}$, $\sigma^2 \equiv \frac{\mu_g^2\sigma_{gxy}^2}{L_g^2(\mu_g^2+\sigma_{gxy}^2)^2}$. Under the condition $L_g \geq 1$, it can be seen that

$$(1-\mu)^2 + \sigma^2 \leq 1 - \mu_g^2/(L_g(\mu_g^2+\sigma_{gxy}^2)) < 1.$$

Applying Lemma 12 shows that

$$\mathbb{E}\left[\left\|\prod_{i=1}^{p}\left(I - \frac{\mu_g}{L_g(\mu_g^2+\sigma_{gxy}^2)}\nabla_{yy}^2 g(x,y;\xi_i^{(2)})\right)\right\|^2\right] \leq d_1\left(1 - \frac{\mu_g^2}{L_g(\mu_g^2+\sigma_{gxy}^2)}\right)^p.$$

Subsequently,

$$\mathbb{E}[\|H_{yy}\|^2] \leq \frac{d_1}{L_g(\mu_g^2+\sigma_{gxy}^2)}.$$

36

[3]Furthermore, it is easy to derive that $\|\mathbb{E}[H_y y]\| \le 1/\mu_g$. Together, the above gives the following estimate on the variance:

$$\mathbb{E}[\|h_f^k - \mathbb{E}[h_f^k]\|^2] \le \sigma_{fx}^2 + \left\{ (\sigma_{fy}^2 + C_y^2)\{\sigma_{gxy}^2 + 2C_{gxy}^2\} + \sigma_{fy}^2 C_{gxy}^2 \right\} \max \left\{ \frac{3}{\mu_g^2}, \frac{3d_1}{L_g(\mu_g^2 + \sigma_{gxy}^2)} \right\}. \quad (112)$$

This concludes the proof for the second part of the lemma. $\square$

We observe the following lemma on the product of (possibly non-PSD) matrices, which is inspired by [18, 25]:

**Lemma 12.** *Let $Z_i, i = 0, 1, ...$ be a sequence of random matrices defined recursively as $Z_i = Y_i Z_{i-1}$, $i \ge 1$, with $Z_0 = I \in \mathbb{R}^{d \times d}$, and $Y_i, i = 0, 1, ...$ are independent, symmetric, random matrices satisfying $\|\mathbb{E}[Y_i]\| \le 1 - \mu$ and $\mathbb{E}[\|Y_i - \mathbb{E}[Y_i]\|_2^2] \le \sigma^2$. If $(1 - \mu)^2 + \sigma^2 < 1$, then for any $t \ge 0$, it holds*

$$\mathbb{E}[\|Z_t\|^2] \le \mathbb{E}[\|Z_t\|_2^2] \le d\left( (1 - \mu)^2 + \sigma^2 \right)^t, \quad (113)$$

*where $\|X\|_2$ denotes the Schatten-2 norm of the matrix $X$.*

*Proof.* We note from the norm equivalence between spectral norm and Schatten-2 norm which yields $\|Z_t\| \le \|Z_t\|_2$ and thus $\mathbb{E}[\|Z_t\|^2] \le \mathbb{E}[\|Z_t\|_2^2]$. For any $i \ge 1$, we observe that

$$Z_i = \underbrace{(Y_i - \mathbb{E}[Y_i])Z_{i-1}}_{=A_i} + \underbrace{\mathbb{E}[Y_i]Z_{i-1}}_{=B_i}.$$

Notice that as $\mathbb{E}[A_i | B_i] = 0$, applying [25, Proposition 4.3] yields

$$\mathbb{E}[\|Z_t\|_2^2] \le \mathbb{E}[\|A_t\|_2^2] + \mathbb{E}[\|B_t\|_2^2]. \quad (114)$$

Furthermore, using the fact that $Y_i$s are independent random matrices and Hölder's inequality for matrices, we observe that

$$\mathbb{E}[\|A_t\|_2^2] \le \mathbb{E}[\|Y_t - \mathbb{E}[Y_t]\|_2^2 \|Z_{t-1}\|_2^2] = \mathbb{E}[\|Y_t - \mathbb{E}[Y_t]\|_2^2] \mathbb{E}[\|Z_{t-1}\|_2^2] \le \sigma^2 \mathbb{E}[\|Z_{t-1}\|_2^2]$$

and using [25, (4.1)],

$$\mathbb{E}[\|B_t\|_2^2] \le \|\mathbb{E}[Y_t]\|^2 \mathbb{E}[\|Z_{t-1}\|_2^2] \le (1 - \mu)^2 \mathbb{E}[\|Z_{t-1}\|_2^2]$$

Substituting the above into (114) yields $\mathbb{E}[\|Z_t\|_2^2] \le ((1 - \mu)^2 + \sigma^2)\mathbb{E}[\|Z_{t-1}\|_2^2]$. Repeating the same arguments for $t$ times and using $\|I\|_2^2 = d$ yields the upper bound. $\square$

# References

[1] Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 2021.

[2] András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.

[3] Amir Beck. *First-order methods in optimization*, volume 25. SIAM, 2017.

---

[3]Notice that a slight modification of the proof of [23, Lemma 3.2] yields $(\mathbb{E}[\|H_{yy}\|])^2 \le \mu_g^{-2}$. However, the latter lemma requires $I - (1/L_g)\nabla_{yy}^2 g(x, y, \xi_i^{(2)})$ to be a PSD matrix almost surely, which is not required in our analysis.

[4] Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite time analysis of temporal difference learning with linear function approximation. In *Conference On Learning Theory*, pages 1691–1692, 2018.

[5] Shalabh Bhatnagar, Mohammad Ghavamzadeh, Mark Lee, and Richard S Sutton. Incremental natural actor-critic algorithms. In *Advances in Neural Information Processing Systems*, pages 105–112, 2008.

[6] Shalabh Bhatnagar, Doina Precup, David Silver, Richard S Sutton, Hamid R Maei, and Csaba Szepesvári. Convergent temporal-difference learning with arbitrary smooth function approximation. In *Advances in Neural Information Processing Systems*, pages 1204–1212, 2009.

[7] Vivek S Borkar. Stochastic approximation with two time scales. *Systems & Control Letters*, 29(5):291–294, 1997.

[8] Vivek S Borkar and Sarath Pattathil. Concentration bounds for two time scale stochastic approximation. In *Allerton Conference on Communication, Control, and Computing*, pages 504–511, 2018.

[9] Jerome Bracken, James E. Falk, and James T. McGill. Technical note—the equivalence of two mathematical programs with optimization problems in the constraints. *Operations Research*, 22(5):1102–1104, 1974.

[10] Jerome Bracken and James T. McGill. Mathematical programs with optimization problems in the constraints. *Operations Research*, 21(1):37–44, 1973.

[11] Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In *ICML*, 2019.

[12] Benoît Colson, Patrice Marcotte, and Gilles Savard. An overview of bilevel optimization. *Annals of Operations Research*, 153(1):235–256, 2007.

[13] Nicolas Couellan and Wenjuan Wang. On the convergence of stochastic bi-level gradient methods. 2016.

[14] Gal Dalal, Balazs Szorenyi, and Gugan Thoppe. A tale of two-timescale reinforcement learning with the tightest finite-time bound. *arXiv preprint arXiv:1911.09157*, 2019.

[15] Christoph Dann, Gerhard Neumann, Jan Peters, et al. Policy evaluation with temporal differences: A survey and comparison. *Journal of Machine Learning Research*, 15:809–883, 2014.

[16] Damek Davis and Dmitriy Drusvyatskiy. Stochastic subgradient method converges at the rate $o(k^{-1/4})$ on weakly convex functions. *arXiv preprint arXiv:1802.02988*, 2018.

[17] Thinh T Doan. Nonlinear two-time-scale stochastic approximation: Convergence and finite-time performance. *arXiv preprint arXiv:2011.01868*, 2020.

[18] Alain Durmus, Eric Moulines, Alexey Naumov, Sergey Samsonov, Kevin Scaman, and Hoi-To Wai. Tight high probability bounds for linear stochastic approximation with fixed stepsize. In *NeurIPS*, 2021.

[19] James E. Falk and Jiming Liu. On bilevel programming, part I: General nonlinear cases. *Mathematical Programming*, 70:47–72, 1995.

[20] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pages 1126–1135, 2017.

[21] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *International Conference on Machine Learning*, 2019.

[22] L Franceschi, P Frasconi, S Salzo, R Grazzi, and M Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning*, pages 1563–1572, 2018.

[23] Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.

[24] M. Hong, H. Wai, Z. Wang, and Z. Yang. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint arXiv:2007.05170*, 2020.

[25] De Huang, Jonathan Niles-Weed, Joel A Tropp, and Rachel Ward. Matrix concentration for products. *Foundations of Computational Mathematics*, pages 1–33, 2021.

[26] Y. Ishizuka and E. Aiyoshi. Double penalty method for bilevel optimization problems. *Ann Oper Res*, 34:73–88, 1992.

[27] Kaiyi Ji, Junjie Yang, and Yingbin Liang. Provably faster algorithms for bilevel optimization and applications to meta-learning. *arXiv preprint arXiv:2010.07962*, 2020.

[28] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR, 2020.

[29] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, pages 267–274, 2002.

[30] Sham M Kakade. A natural policy gradient. In *Advances in Neural Information Processing Systems*, pages 1531–1538, 2002.

[31] Maxim Kaledin, Eric Moulines, Alexey Naumov, Vladislav Tadic, and Hoi-To Wai. Finite time analysis of linear two-timescale stochastic approximation with Markovian noise. In *COLT*, 2020.

[32] Prasenjit Karmakar and Shalabh Bhatnagar. Two time-scale stochastic approximation with controlled Markov noise and off-policy temporal-difference learning. *Mathematics of Operations Research*, 43(1):130–151, 2018.

[33] Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in Neural Information Processing Systems*, pages 1008–1014, 2000.

[34] Vijay R Konda, John N Tsitsiklis, et al. Convergence rate of linear two-time-scale stochastic approximation. *Annals of Applied Probability*, 14(2):796–819, 2004.

[35] Junyi Li, Bin Gu, and Heng Huang. Improved bilevel model: Fast and optimal algorithm with theoretical guarantee. *arXiv preprint arXiv:2009.00690*, 2020.

[36] Valerii Likhosherstov, Xingyou Song, Krzysztof Choromanski, Jared Davis, and Adrian Weller. Ufo-blo: Unbiased first-order bilevel optimization. *arXiv preprint arXiv:2006.03631*, 2020.

[37] Bo Liu, Ian Gemp, Mohammad Ghavamzadeh, Ji Liu, Sridhar Mahadevan, and Marek Petrik. Proximal gradient temporal difference learning: Stable reinforcement learning with polynomial sample complexity. *Journal of Artificial Intelligence Research*, 63:461–494, 2018.

[38] Risheng Liu, Pan Mu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A generic first-order algorithmic framework for bi-level programming beyond lower-level singleton. *arXiv preprint arXiv:2006.04045*, 2020.

[39] Zhi-Quan Luo, Jong-Shi Pang, and Daniel Ralph. *Mathematical Programs with Equilibrium Constraints*. Cambridge University Press, 1996.

[40] Hamid Reza Maei, Csaba Szepesvári, Shalabh Bhatnagar, and Richard S Sutton. Toward off-policy learning control with function approximation. In *International Conference on Machine Learning*, 2010.

[41] Akshay Mehra and Jihun Hamm. Penalty method for inversion-free deep bilevel optimization. In *Asian Conference on Machine Learning*, 2019.

[42] Abdelkader Mokkadem, Mariane Pelletier, et al. Convergence rate and averaging of nonlinear two-time-scale stochastic approximation algorithms. *Annals of Applied Probability*, 16(3):1671–1702, 2006.

[43] Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9(May):815–857, 2008.

[44] Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In *International conference on machine learning*, pages 737–746. PMLR, 2016.

[45] Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.

[46] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of MAML. In *International Conference on Learning Representations*, 2019.

[47] Aravind Rajeswaran, Chelsea Finn, Sham Kakade, and Sergey Levine. Meta-learning with implicit gradients. 2019.

[48] Shoham Sabach and Shimrit Shtern. A first order method for solving convex bilevel optimization problems. *SIAM Journal on Optimization*, 27(2):640–660, 2017.

[49] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[50] Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation for bilevel optimization. In *Artificial Intelligence and Statistics*, pages 1723–1732, 2019.

[51] H. Van Stackelberg. *The theory of market economy*. Oxford University Press, 1952.

[52] Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–44, 1988.

[53] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[54] Richard S Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *International Conference on Machine Learning*, pages 993–1000, 2009.

[55] Csaba Szepesvári. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.

[56] L. Vicente, G. Savard, and J. Judice. Descent approaches for quadratic bilevel programming. *Journal of Optimization Theory and Applications*, 81:379–399, 1994.

[57] Mengdi Wang, Ethan X Fang, and Han Liu. Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions. *Mathematical Programming*, 161(1-2):419–449, 2017.

[58] D. J. White and G. Anandalingam. A penalty function approach for solving bi-level linear programs. *Journal of Global Optimization*, 3:397–419, 1993.

[59] Yue Wu, Weitong Zhang, Pan Xu, and Quanquan Gu. A finite time analysis of two time-scale actor-critic methods. In *NeurIPS*, 2020.

[60] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

[61] Tengyu Xu, Zhe Wang, and Yingbin Liang. Non-asymptotic convergence analysis of two time-scale (natural) actor-critic algorithms. *arXiv preprint arXiv:2005.03557*, 2020.

[62] Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004. PMLR, 2019.