

Natural Policy Gradient and Actor Critic Methods for Constrained Multi-Task Reinforcement Learning

Sihan Zeng
JPMorgan AI Research

sihan.zeng@jpmchase.com

Thinh T. Doan
Department of Electrical and Computer Engineering
Virginia Tech

thinhdogan@vt.edu

Justin Romberg
Department of Electrical and Computer Engineering
Georgia Tech

jrom@ece.gatech.edu

Abstract

Multi-task reinforcement learning (RL) aims to find a single policy that effectively solves multiple tasks at the same time. This paper presents a constrained formulation for multi-task RL where the goal is to maximize the average performance of the policy across tasks subject to bounds on the performance in each task. We consider solving this problem both in the centralized setting, where information for all tasks is accessible to a single server, and in the decentralized setting, where a network of agents, each given one task and observing local information, cooperate to find the solution of the globally constrained objective using local communication.

We first propose a primal-dual algorithm that provably converges to the globally optimal solution of this constrained formulation under exact gradient evaluations. When the gradient is unknown, we further develop a sampled-based actor-critic algorithm that finds the optimal policy using online samples of state, action, and reward. Finally, we study the extension of the algorithm to the linear function approximation setting.

1 Introduction

Multi-task reinforcement learning (RL) aims to find a common policy that effectively solves a range of tasks simultaneously, where each task is the policy optimization problem defined over a Markov decision process (MDP). The MDPs can have different state spaces, reward functions, and transition kernels in general, but may be implicitly or explicitly correlated.

The most common mathematical formulation for multi-task RL is to maximize the average cumulative rewards collected by a single policy across all MDPs Zeng et al. (2021); Jiang et al. (2022); Junru et al. (2022). In this paper, we study a generalized formulation in which we maximize the average cumulative rewards subject to constraints on the performance of the policy for each task. This formulation is a special case of the policy optimization problem for a constrained Markov decision process (CMDP) Altman (1999) and is a flexible framework that allows more fine-grained specification of the performance of the optimal policy in each task. In applications where the tasks exhibit major conflicts of interest and/or the magnitude of the rewards varies significantly across tasks Kalashnikov et al. (2021); Guo et al. (2022), the optimal policy under the average-cumulative-reward formulation may excel in some tasks at the cost of compromised performance in others Hayes et al. (2022). The constrained formulation provides a way to mitigate this task imbalance. Illustrative numerical simulations are given in Section 7.

Under the constrained multi-task formulation, we consider centralized and decentralized learning paradigms. “Centralized” in this context means that information of all tasks is available at a single server, while “decen-

tralized” is a scenario where a group of agents, each deployed to one local environment/task, work together to solve the global constrained optimization program without any central coordination. The centralized setting can be regarded as an easier special case of the decentralized problem with a fully connected communication graph. We will initiate our algorithmic development in the centralized setting and extend them to the decentralized scenario.

To solve the constrained multi-task policy optimization problem, we propose a multi-task primal-dual natural policy gradient algorithm (MT-PDNPG) that allows the policy at each local agent to achieve global optimality. The updates in MT-PDNPG require the computation of the exact gradients of the value function, which is impractical in environments with large state spaces and/or unknown transition probability kernel. To extend our algorithm to the case where we do not have perfect knowledge of the environments and can only obtain samples of the state transitions, we propose a sampled-based multi-task primal-dual natural actor-critic algorithm (MT-PDNAC) and study its finite-sample performance. Finally, to tackle problems where the state space is enormous or even infinitely large, we extend MT-PDNAC to the case where the policy and value functions are linearly approximated with pre-determined lower-dimensional feature vectors. Despite the complications introduced to the algorithm and analysis by the linear function approximation, we show that a modified version of the MT-PDNAC algorithm in this setting achieves the same convergence rate as in the tabular case. We note that our proposed sampled-based algorithms are completely online in the sense that they use a single trajectory of continuously generated samples, which makes them convenient to implement in practice.

1.1 Main Contribution

The first contribution of our work is to study the constrained multi-task RL formulation and to propose and analyze a primal-dual natural actor-critic algorithm provably converges to the globally optimal policy in expectation. The algorithm is completely data-driven, single-loop, and relies on a single trajectory of samples. After K iterations, the policy parameter converges in both objective function and constraint violation up to the precision $\mathcal{O}(1/K^{1/6})$. This matches the best-known time and sample complexity of the natural actor-critic algorithm for single-agent (non-constrained) MDPs under comparable assumptions.

Our second contribution is to extend the primal-dual algorithms to the decentralized learning paradigm. The extended algorithm makes each agent compute and update in the direction of a locally observable component of the policy gradient, followed by a parameter averaging step. We show that the decentralized algorithms enjoy finite time and sample complexity matching their centralized counterparts, differing only by a factor that scales inversely with the connectivity of the communication graph.

Finally, we study the setting where both the policy and the critic variables are approximated using linear features. This on-policy linear function approximation setting presents peculiar technical challenges. Specifically, the TD learning target under linear function approximation becomes ill-defined when the policy to evaluate is not completely mixed (i.e. does not have uniformly lower bounded entries). While a well-defined TD learning target can be ensured through careful control of the policy iterates, it is not sufficiently Lipschitz continuous for the direct extension of our algorithm for the tabular setting to converge. We overcome the challenge by dynamically adjusting the number of TD learning iterations devoted to chasing a fixed TD target, which results in a modified algorithm with a sample complexity of $\mathcal{O}(\epsilon^{-6})$, matching that in the tabular case. To our best knowledge, a finite-time and finite-sample analysis for the on-policy natural actor-critic algorithm (even for standard non-constrained MDPs) has been missing from the existing literature, and our work fills in this gap.

1.2 Related Work

This paper presents reliably convergent decentralized algorithm for finding the optimal solution of the new constrained multi-task RL objective. It closely relates to the literature on multi-task RL, decentralized optimization, CMDP, and actor-critic algorithms in RL, which we discuss in this section to give context to our novelty.

Multi-Task Reinforcement Learning. Multi-task RL in general studies efficiently solving the policy optimization tasks for multiple RL environments at the same time by leveraging connections between the tasks. Its most common mathematical formulation is to find a single policy that maximizes the (weighted) average of the cumulative returns collected across all environments, and Zeng et al. (2021); Jiang et al. (2022); Junru et al. (2022); Chen et al. (2022a) study various gradient-based algorithms that provably converge to global or local solutions of this objective. However, as pointed out in Hessel et al. (2019), this average return formulation can be inadequate when modelling practical problems where the tasks have strong conflicting or imbalanced interests. While the authors in Hessel et al. (2019) address this issue by dynamically addressing the weight of each task in the policy updates, we are motivated to propose the constrained multi-task formulation that allows fine-grained control of the performance of the policy in each environment.

It is worth pointing out the large volume of literature that approaches multi-task RL from a more empirical perspective. Some important lines of work include (but are certainly not limited to) policy distillation Rusu et al. (2015); Traoré et al. (2019); Wadhwanian et al. (2019), transfer learning Gupta et al. (2017); D’Eramo et al., and innovated design of the policy representation Yang et al. (2020); Hong et al. (2021). We also note that there exist multi-task RL formulations where rather than learning a single policy, task-specific adaptation is allowed Finn et al. (2017); Raghu et al.

Decentralized Optimization. Closely connected to distributed and decentralized optimization, our problem formulation considers maximizing a global objective function composed of local rewards under local constraints in the case where each agent only has access to its local information. In the unconstrained setting, it has been well-known that decentralized gradient, sub-gradient, and Newton’s methods converge as fast as their centralized counterparts up to a constant that describes the connectivity of the communication graph Nedic & Ozdaglar (2009); Yuan et al. (2016); Nedic (2020); Bullins et al. (2021); Islamov et al. (2021). For constrained convex optimization programs, Chang et al. (2014); Lei et al. (2016) show that primal-dual algorithms provably converge to the globally optimal solution. Our work is inspired by these existing results but focuses on a non-convex optimization program where we establish global convergence using the specific problem structure.

Constrained Markov Decision Process. Our constrained multi-task formulation can be regarded as a special case of policy optimization under a CMDP (if the information for all tasks is centrally available). A common approach to finding the optimal policy for a CMDP is to search for a saddle point of the Lagrangian using primal-dual gradient descent ascent. Variants of this approach are considered in Prashanth & Ghavamzadeh (2016); Chow et al. (2017); Tessler et al. (2018). However, they only establish the asymptotic convergence to a stationary point or locally optimal solution; global optimality is not achieved and the exact finite-time complexity is unknown due to the non-convexity of the underlying optimization program.

For a convex constrained optimization problem under Slater’s condition, it is well known that strong duality holds and the primal-dual gradient descent ascent algorithm efficiently converges to the globally optimal solution. The policy optimization problem for a CMDP is non-convex, but by leveraging the structure of the CMDP, Altman (1999); Paternain et al. (2019) show that the strong duality holds despite the non-convexity. A number of works Ding et al. (2020); Liu et al. (2021b); Ding et al. (2022); Bai et al. (2022) take advantage of this property to design primal-dual natural policy gradient algorithms and establish their finite-time convergence to the globally optimal policy in the tabular case. Some other works focus on deriving regret bounds (rather than finite-time convergence) Zheng & Ratliff (2020); Ding et al. (2021); Agarwal et al. (2022). A limit of these studies is that the policy and dual variable updates rely on an oracle that always returns the exact gradient (or its highly accurate unbiased estimate). In real-life problems where the transition probability kernel is not fully known and/or the state space is large, computing such gradients can become computationally prohibitive. We improve these prior works by presenting a completely sample-based algorithm that solves the CMDP optimization using a single continuously-generated trajectory, both in the tabular setting and under linear function approximation.

There exist non-primal-dual algorithms for finding the optimal policy for a CMDP. For example, Yu et al. (2019) constructs a succession of surrogate convex relaxation to the non-convex CMDP optimization problem and shows that the solutions to these surrogate programs converge to a stationary point of the CMDP

optimization problem. The work Chow et al. (2018) extends value-based methods including value iteration, policy iteration, and Q learning to the context of CMDP. The paper HasanzadeZonuzi et al. (2021) builds an empirical estimate of the probability transition kernel, on which planning is carried out. The authors of Liu et al. (2021a) propose an algorithm driven by the principle of optimistic pessimism that improves the state-of-the-art analysis on the constraint violation. In Ying et al. (2022), a dual-only approach is studied which solves a regularized version of the CMDP.

Actor-Critic Algorithms. The sample-based algorithms presented in our paper fall under the category of actor-critic algorithms, which can be considered as a variant of policy gradient methods where the unknown value function is estimated by an auxiliary variable updated with stochastic approximation. This class of algorithms has been analyzed in various settings of dynamical systems (such as standard MDP, entropy regularized MDP, and linear-quadratic regulator), on/off-policy sample collection, natural/standard policy gradient, and function approximation Yang et al. (2019); Wu et al. (2020); Xu et al. (2020); Zeng et al. (2024); Ju et al. (2022); Khodadadian et al. (2022); Chen et al. (2022b); Barakat et al. (2022). However, as discussed in the previous subsection, analyzing the natural actor-critic algorithm incurs unique technicality that has not been treated in the previous literature under the combined effect of on-policy samples and linear function approximation. Our work thoroughly describes the technical challenges and presents our solution.

Finally, we note that a preliminary version of the work has been presented in Zeng et al. (2022), which analyzes an online actor-critic algorithm for a single-agent CMDP. This current paper significantly generalizes Zeng et al. (2022) along several axes by introducing the multi-task formulation, decentralized learning paradigm, and linear function approximation.

2 Constrained Multi-Task Formulation

Consider a collection of N infinite-horizon discounted-reward MDPs characterized by $\{(\mathcal{S}_i, \mathcal{A}, \mathcal{P}_i, \gamma, r_i)\}_{i=1}^N$, where \mathcal{S}_i is the finite state space, \mathcal{A} is the finite and common action space, $P_i : \mathcal{S}_i \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}_i}$ is the transition probability kernel, $\gamma \in (0, 1)$ is the discount factor, and $r_i : \mathcal{S}_i \times \mathcal{A} \rightarrow [0, 1]$ is the reward function.

Given a policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}_i}$, we define the local value functions for each task $i = 1, \dots, N$

$$V_i^\pi(s) \triangleq \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_i(s_k, a_k) \mid s_0 = s \right], \quad Q_i^\pi(s, a) \triangleq \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_i(s_t, a_t) \mid s_0 = s, a_0 = a \right],$$

$$A_i^\pi(s, a) \triangleq Q_i^\pi(s, a) - V_i^\pi(s), \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

With some abuse of notation, we use $V_i^\pi(\rho)$ to denote the expected cumulative reward collected by policy π under the initial distribution ρ

$$V_i^\pi(\rho) \triangleq \mathbb{E}_{s_0 \sim \rho} [V_i^\pi(s_0)]. \quad (1)$$

For the simplicity of notation, we define $V_0^\pi(\rho)$ to be the averaged cumulative reward over tasks

$$V_0^\pi(\rho) \triangleq \frac{1}{N} \sum_{i=1}^N V_i^\pi(\rho).$$

The common multi-task formulation considered in the literature is to maximize this average value function

$$\max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} V_0^\pi(\rho). \quad (2)$$

Given local performance upper and lower bounds $\{\ell_i \in \mathbb{R}, u_i \in \mathbb{R}\}_{i=1}^N$, our multi-task policy optimization objective is to solve the following constrained program

$$\pi^* \in \operatorname{argmax}_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} V_0^\pi(\rho)$$

$$\text{subject to } \ell_i \leq V_i^\pi(\rho) \leq u_i \quad \forall i = 1, 2, \dots, N. \quad (3)$$

We note that Eq. (3) obviously subsumes the non-constraint multi-task formulation in Zeng et al. (2021); Jiang et al. (2022); Junru et al. (2022) by properly choosing $\{\ell_i, u_i\}$. It can be shown that the multi-task policy optimization problem (even without constraints) does not observe the gradient domination condition in general, which makes it difficult for any gradient-based algorithm to find the globally optimal policy. In this paper, we make the simplifying assumption that the local MDPs have the same state space and transition probability kernel and can only be different in the reward functions (i.e. $\mathcal{S} = \mathcal{S}_1 = \mathcal{S}_2 = \dots$ and $\mathcal{P} = \mathcal{P}_1 = \mathcal{P}_2 = \dots$), under which the gradient domination condition is recovered Zeng et al. (2021).

We consider the following assumption on Eq. (3), which essentially states that the constraint set must have a non-empty interior. This is a standard assumption in the study of constrained MDPs Paternain et al. (2019); Ding et al. (2020); Zeng et al. (2023) and ensures that strong duality holds for the constrained program despite the lack of convexity.

Assumption 1 (Slater’s Condition). *There exists a constant $0 < \xi \leq 1$ and a policy π such that $\ell_i + \xi \leq V_i^\pi(\rho) \leq u_i - \xi$ for all $i = 1, \dots, N$.*

3 Preliminary – Centralized Computation Setting

We start by discussing how Eq. (3) can be solved in the “centralized” setting where information of all tasks are available at a single server. For now, we work with deterministic gradients where we assume having perfect information of the environments to compute the exact value functions. Sampled-based algorithms will be built upon these results in Section 4. Under this simplification, the problem formulation and setting become a special case of those in Ding et al. (2020) where the reward to be optimized is $\frac{1}{N} \sum_{i=1}^N r_i$ and the reward for the definition of the i_{th} constraint is r_i .

Inspired by Ding et al. (2020), we follow a primal-dual approach to solve the constrained program in Eq. (3). The first step is to form the Lagrangian

$$V_L^{\pi, \lambda, \nu}(\rho) \triangleq V_0^\pi(\rho) + \sum_{i=1}^N (\lambda_i (V_i^\pi(\rho) - \ell_i) - \nu_i (V_i^\pi(\rho) - u_i)), \quad (4)$$

where $\lambda = [\lambda_1, \dots, \lambda_N] \in \mathbb{R}_+^N$ and $\nu = [\nu_1, \dots, \nu_N] \in \mathbb{R}_+^N$ are the dual variables associated with the lower and upper bounds. The dual function $V_D^{\lambda, \nu}$ is defined as

$$V_D^{\lambda, \nu}(\rho) = \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} V_L^{\pi, \lambda, \nu}(\rho), \quad (5)$$

and the dual problem is to solve

$$\lambda^*, \nu^* \in \operatorname{argmin}_{\lambda, \nu \in \mathbb{R}_+^N} V_D^{\lambda, \nu}(\rho). \quad (6)$$

Although the program in Eq. (3) is non-convex, it is known that strong duality holds under Slater’s condition in Assumption 1 (see Altman (1999)[Theorem 3.6])

$$V_D^{\lambda^*, \nu^*}(\rho) = V_0^{\pi^*}(\rho), \quad (7)$$

where π^* , λ^* , and ν^* are the (not necessarily unique) optimal solutions to Eq. (3) and Eq. (6), respectively. The optimal dual variables are known to be bounded also as a consequence of Assumption 1, which we present in the following lemma and is a simple extension of Ding et al. (2020)[Lemma 1] to the case of multiple constraints.

Lemma 1. *Under Assumption 1, we have*

$$\|\lambda^*\|_\infty \leq \frac{B_\lambda}{2} \quad \text{and} \quad \|\nu^*\|_\infty \leq \frac{B_\lambda}{2},$$

where $B_\lambda = \frac{1}{\xi(1-\gamma)}$

Motivated by the existence of strong duality, we solve Eq. (3) by finding the minimax saddle point of the Lagrangian. We represent the policy through the softmax parameterization, as that introduces more favorable structure to the optimization landscape Agarwal et al. (2020). Specifically, using a policy parameter $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ that encodes the policy π_θ through the softmax function as follows

$$\pi_\theta(a | s) = \frac{\exp(\theta(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta(s, a'))}, \quad \text{for all } \theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|},$$

we take a gradient-descent-ascent approach to find the saddle point $(\theta^*, \lambda^*, \nu^*)$ such that

$$\theta^*, (\lambda^*, \nu^*) = \underset{\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}}{\operatorname{argmax}} \underset{\lambda, \nu \in \mathbb{R}_+^N}{\operatorname{argmin}} V_L^{\pi_\theta, \lambda, \nu}. \quad (8)$$

Our approach requires computing the gradients of the Lagrangian with respect to both the policy parameter and dual variables, which we now derive.

Gradient of the dual variable. The Lagrangian in Eq. (4) is obviously a linear function of λ and ν , and the gradients have simple closed-form expressions.

$$\begin{aligned} \nabla_{\lambda_i} V_L^{\pi, \lambda, \nu}(\rho) &= V_i^\pi(\rho) - \ell_i = \sum_{s: \rho(s) > 0, a} \rho(s) \pi(a | s) Q_i^\pi(s, a) - \ell_i, \\ \nabla_{\nu_i} V_L^{\pi, \lambda, \nu}(\rho) &= -V_i^\pi(\rho) + u_i = - \sum_{s: \rho(s) > 0, a} \rho(s) \pi(a | s) Q_i^\pi(s, a) + u_i. \end{aligned}$$

This naturally leads to the update in Eq. (21), in which the operator $\Pi_{[0, B_\lambda]} : \mathbb{R}^N \rightarrow \mathbb{R}^N$ denotes the element-wise projection of a vector to the interval $[0, B_\lambda]$. The projection guarantees the stability of the dual variables. We note that the optimal dual variables are in the range $[0, B_\lambda]$ according to Lemma 1.

Gradient of the primal variable. It is known that the natural gradient of the value function under reward r_i with respect to θ , denoted by $\tilde{\nabla}_\theta V_i^{\pi_\theta}(\rho)$, is the advantage function $A_i^{\pi_\theta}$ weighted by $1/(1-\gamma)$ Agarwal et al. (2020). This means that $\tilde{\nabla}_\theta V_L^{\pi_\theta, \lambda, \nu}(\rho)$, the natural policy gradient of the Lagrangian can be expressed as

$$\begin{aligned} &\tilde{\nabla}_{\theta(s, a)} V_L^{\pi_\theta, \lambda, \nu}(\rho) \\ &= \frac{1}{1-\gamma} (A_0^{\pi_\theta}(s, a) + \sum_{i=1}^N (\lambda_i - \nu_i) A_i^{\pi_\theta}(s, a)) = \frac{1}{1-\gamma} \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i - \nu_i \right) A_i^{\pi_\theta}(s, a). \end{aligned} \quad (9)$$

Formally presented in Algorithm 1, the multi-task primal-dual natural policy gradient (MT-PDNPG) algorithm proposed for the centralized setting ascends the primal variable according to Eq. (10) in the direction of the natural gradient, with dual iterates λ_i^k, ν_i^k plugged in as λ_i, ν_i and the Q function as a proxy of the advantage function. On the other hand, dual variables λ_i^k, ν_i^k are iteratively refined according to Eq. (11) with gradient descent.

Despite its simplicity, the iterates of Algorithm 1 efficiently converge to the globally optimal solution of Eq. (3). As this algorithm can be regarded as a special case of that in Ding et al. (2020) under a specific choice of reward functions, its analysis directly follows as a corollary of Ding et al. (2020)[Theorem 1], which we state below. In the following sections, we will generalize Algorithm 1 along three dimensions to make it 1) sample-based, 2) compatible with decentralized learning, and 3) effective under linear function approximation.

Algorithm 1 Multi-Task Primal-Dual Natural Policy Gradient Algorithm (Centralized)

- 1: **Initialization:** Initialize $\theta^0 \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} = 0$ and for each task i dual variables $\lambda_i^0, \nu_i^0 \in \mathbb{R}_+ = 0$
- 2: **for** $k = 0, 1, \dots, K - 1$ **do**
- 3: Policy update:

$$\theta^{k+1} = \alpha \sum_{j=1}^N \left(\frac{1}{N} + \lambda_j^k - \nu_j^k \right) Q_j^{\pi_{\theta^k}}, \quad \pi^{k+1}(a | s) = \frac{\exp(\theta^{k+1}(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta^{k+1}(s, a'))} \quad (10)$$

- 4: **for** Each task $i = 1, \dots, N$ **do**
- 5: Dual variable update:

$$\begin{aligned} \lambda_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k - \eta (V_i^{\pi_{\theta^k}}(\rho) - \ell_i) \right) \\ \nu_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k + \eta (V_i^{\pi_{\theta^k}}(\rho) - u_i) \right) \end{aligned} \quad (11)$$

- 6: **end for**
 - 7: **end for**
-

Corollary 1. Consider the iterates $\{\pi^k\}$ obtained from K iterations of Algorithm 1 in the centralized setting. Let the step size sequences be $\alpha = \frac{\alpha_0}{K^{1/2}}, \eta = \frac{\eta_0}{K^{1/2}}$, with $\alpha_0, \eta_0 > 0$. Under Assumption 1, we have

$$\max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} (V_0^{\pi^k}(\rho) - V_0^{\pi^*}(\rho)), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi^k}(\rho)]_+ + [V_i^{\pi^k}(\rho) - u_i]_+ \right) \right\} \leq \mathcal{O}(K^{-1/2}).$$

Corollary 1 guarantees that the policy iterate converges in both objective function value and constraint violation with a rate of $\mathcal{O}(1/\sqrt{K})$. As the centralized setting is a special decentralized case where every agent may communicate with every other agent, this result can also be regarded as a corollary of Theorem 1 to be presented in Section 5.1, which analyzes the MT-PDNPG algorithm under decentralized computation.

4 Sample-Based Setting: Online Primal-Dual Natural Actor-Critic Algorithm

The policy gradient algorithm presented in the previous section employs deterministic gradient updates, which requires evaluating the Q function. In large and/or unknown environments, the Q function cannot be exactly computed instantaneously with samples, which makes Algorithm 1 inapplicable in practice. In Algorithm 2, we introduce a sample-based extension of Algorithm 1, where the key extension is to maintain a variable $\hat{Q}_i^k \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ as an estimate of $Q_i^{\pi^k}$ in each task i , refined over time using samples of the state transition. Instead of relying on the true value function, the updates of the policy and dual variable in Eqs. (23) and (25) use approximate gradients obtained by plugging in the most up-to-date \hat{Q}_i^k estimate. Algorithms of this flavor are usually categorized as actor-critic methods, where the actor refers to the policy iterates and the critic is the value function estimator. We stress that our actor-critic algorithm is single-loop and truly online in the sense that the samples are generated continuously from a single trajectory and updates both actor and critic variables on the run, which makes it very convenient to implement in practice.

While we can generate samples according to the current policy iterate π^k , doing so may cause certain actions to be very infrequently selected, which leads to exploration issues. To guarantee sufficient exploration (visitation of all state-action pairs), we introduce a behavior policy $\hat{\pi}^k$ in Eq. (24), which is the ϵ -exploration version of π^k . We note that ϵ needs to be properly selected with respect to the desired precision and controls an important trade-off: an excessively large ϵ facilitates the exploration of all state-action pairs but lead to a substantial gap between $\hat{\pi}^k$ and π^k , and vice versa. We note that similar behavior policies are also adopted in Borkar (2005); Khodadadian et al. (2022).

In the following subsection, we show that this simple, intuitive, completely sample-based algorithm is guaranteed to converge efficiently.

Algorithm 2 Multi-Task Primal-Dual Natural Actor-Critic Algorithm (Centralized)

- 1: **Initialization:** Initialize $\theta^0 \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} = 0$ and uniform behavior policy $\hat{\pi}^0$. For each task i , initialize dual variables $\lambda_i^0, \nu_i^0 \in \mathbb{R}_+ = 0$ and critic parameters $\hat{Q}_i^0 = 0 \in \mathbb{R}^d$. For each task i , draw the initial sample s_i^0 and $a_i^0 \sim \hat{\pi}_i^0(\cdot | s_i^0)$
- 2: **for** $k = 0, 1, \dots, K - 1$ **do**
- 3: 1) Policy (actor) update: $\forall s \in \mathcal{S}, a \in \mathcal{A}$

$$\theta^{k+1} = \alpha \sum_{j=1}^N \left(\frac{1}{N} + \lambda_j^k - \nu_j^k \right) \hat{Q}_j^k, \quad \pi^{k+1}(a | s) = \frac{\exp(\theta^{k+1}(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta^{k+1}(s, a'))} \quad (12)$$

- 4: 2) Behavior policy update: $\forall s \in \mathcal{S}, a \in \mathcal{A}$

$$\hat{\pi}^{k+1}(a | s) = \frac{\epsilon}{|\mathcal{A}|} + (1 - \epsilon) \pi^{k+1}(a | s). \quad (13)$$

- 5: **for** Each task $i = 1, \dots, N$ **do**
- 6: 1) Observe $s_i^{k+1} \sim P(\cdot | s_i^k, a_i^k)$ and take action $a_i^{k+1} \sim \hat{\pi}^k(\cdot | s_i^{k+1})$
- 7: 2) Value function estimator (critic) update:

$$\hat{Q}_i^{k+1}(s_i^k, a_i^k) = (1 - \beta) \hat{Q}_i^k(s_i^k, a_i^k) + \beta \left(r_i(s_i^k, a_i^k) + \gamma \hat{Q}_i^k(s_i^{k+1}, a_i^{k+1}) \right) \quad (14)$$

- 8: 3) Local dual variable update:

$$\begin{aligned} \lambda_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k - \eta \left(\sum_{s,a} \rho(s) \pi^k(a | s) \hat{Q}_i^k(s, a) - \ell_i \right) \right) \\ \nu_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\nu_i^k + \eta \left(\sum_{s,a} \rho(s) \pi^k(a | s) \hat{Q}_i^k(s, a) - u_i \right) \right) \end{aligned} \quad (15)$$

- 9: **end for**
 - 10: **end for**
-

4.1 Finite-Sample Complexity

We start by introducing some notations and stating our main assumptions.

Given a policy $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$, we use $P^\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{S}}$ to denote the state transition probability under π

$$P^\pi(s' | s) = \sum_{a \in \mathcal{A}} \mathcal{P}(s' | s, a) \pi(a | s),$$

and $\mu_\pi \in \Delta_{\mathcal{S}}$ to denote the stationary distribution of the Markov chain induced by P^π . In addition, we denote by $\tilde{\mu}_\pi \in \Delta_{\mathcal{S} \times \mathcal{A}}$ the stationary distribution of state-action pairs under π , which relates to μ_π as follows

$$\tilde{\mu}_\pi(s, a) = \mu_\pi(s) \pi(a | s).$$

To measure the amount of time a Markov chain takes to approach its stationary distribution, we define the mixing time as follows.

Definition 1. Given policy π , consider the Markov chain $\{s^k\}$ generated according to $s^{k+1} \sim P^\pi(\cdot | s^k)$. For any scalar $c > 0$, the mixing time of $\{s^k\}$ associated with c is

$$\tau_\pi(c) \triangleq \min\{k \geq 0 : \sup_{s \in \mathcal{S}} d_{TV}(P(s^k = \cdot | s^0 = s), \mu_\pi(\cdot)) \leq c\}, \quad (16)$$

where given two probability distributions u_1 and u_2 , d_{TV} denotes their total variation distance

$$d_{TV}(u_1, u_2) = \frac{1}{2} \sup_{\nu: \mathcal{X} \rightarrow [-1, 1]} \left| \int \nu du_1 - \int \nu du_2 \right|. \quad (17)$$

We consider the following assumption on the transition probability kernel \mathcal{P} which states that the Markov chain induced by any fixed policy approaches its stationary distribution geometrically fast.

Assumption 2 (Uniform Ergodicity). *Given any π , the Markov chain $\{s^k\}$ generated by P^π according to $s^{k+1} \sim P^\pi(\cdot | s^k)$ has a unique stationary distribution μ_π and is uniformly geometrically ergodic, i.e., there exist $C_0 \geq 1$ and $\ell \in (0, 1)$ such that*

$$\sup_s d_{TV}(P(s^k = \cdot | s^0 = s), \mu_\pi(\cdot)) \leq C_0 \ell^k, \forall k \geq 0.$$

This assumption is commonly made in RL and optimization papers that study gradient-based algorithms under samples collected from a (time-varying) Markovian chain Wu et al. (2020); Khodadadian et al. (2022); Zeng et al. (2024). Recall the mixing time $\tau_\pi(c)$ defined in Eq. (16). As an obvious result of Assumption 2, there exists a constant $D > 0$ such that

$$\tau_\pi(c) \leq D \log(1/c), \quad \forall c \in (0, 1) \text{ and } \pi. \quad (18)$$

In this work, we denote

$$\tau(c) \triangleq \max_\pi \tau_\pi(c) \leq D \log(1/c). \quad (19)$$

Another consequence of the uniform ergodicity assumption is that the stationary distribution μ_π is uniformly bounded away from 0 for any policy π , and we denote $\underline{\mu} \triangleq \min_{\pi, s} \mu_\pi(s) > 0$.

Algorithm 2 is guaranteed to converge to the optimal multi-task policy π^\star under proper choice of the step sizes. We state the result below, which we take to be a corollary of the analysis of the decentralized MT-PDNAC algorithm to be presented in the next section.

Corollary 2. *Consider the iterates $\{\pi^k\}$ obtained from K iterations of Algorithm 2 in the centralized setting. Let the step size sequences be*

$$\alpha = \frac{\alpha_0}{K^{5/6}}, \quad \beta = \frac{\beta_0}{K^{1/2}}, \quad \eta = \frac{\eta_0}{K^{5/6}}, \quad \epsilon = \frac{\epsilon_0}{K^{1/6}},$$

with $\frac{(1-\gamma)\mu\epsilon_0\beta_0}{|\mathcal{A}|} \leq 1$ and $\alpha_0 = \mathcal{O}(N^{-1/4})$. Then, under Assumptions 1 and 2, we have

$$\max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} (V_0^{\pi^\star}(\rho) - V_0^{\pi^k}(\rho)), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi^k}(\rho)]_+ + [V_i^{\pi^k}(\rho) - u_i]_+ \right) \right\} \\ \leq \mathcal{O}\left(\frac{N^{5/4} \log(K)}{K^{1/6}}\right).$$

Corollary 2 shows a finite-time complexity of $\tilde{\mathcal{O}}(K^{-1/6})$, where $\tilde{\mathcal{O}}$ hides structural constants and logarithm factors. As Algorithm 2 draws exactly N samples in each iteration, it implies that the algorithm will converge in both objective function and constraint violation up to a precision δ with at most $\tilde{\mathcal{O}}(\delta^{-6})$ samples. This result matches the best-known rate of actor-critic algorithms for solving the policy optimization problem under a single-task, unconstrained MDP Khodadadian et al. (2022). Compared with the complexity of Algorithm 1, we have a slightly inferior rate for not exactly knowing the Q function of the current policy iterates.

5 Algorithms Under Decentralized Computation

In many practical applications (for example, each task is associated with an environment physically separated from each other), the information of the multiple tasks may not all be available on a central server. A more realistic learning paradigm in such scenarios is to employ a network of N agents, each placed to explore and learn in one different environment. For generality, we do not assume the existence of a central coordinator that exchanges information with all agents. Rather, the agents are connected according to an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$; agent $i \in \mathcal{V}$ can communicate with agent $j \in \mathcal{N}_i$, where $\mathcal{N}_i = \{j \in \mathcal{V} : (i, j) \in \mathcal{E}\}$ denotes the neighbors of agent i . Since each agent can only access local information, collaboration between the agents is necessary to solve the global objective in Eq. (3).

5.1 Deterministic Gradient Setting

We first present the development of the decentralized MT-PDNPG method, which extends Algorithm 1. While we would still like to perform alternating gradient descent ascent on the dual and primal variables as in Eqs. (11) and (10), the challenge is that the natural gradient in Eq. (10) involves information across all tasks and obviously cannot be computed by any single agent locally.

Algorithm 3 Multi-Task Primal-Dual Natural Policy Gradient Algorithm (Decentralized)

- 1: **Initialization:** Each agent i initializes $\theta_i^0 \in \mathbb{R}^{|S||\mathcal{A}|} = 0$ and dual variables $\lambda_i^0, \nu_i^0 \in \mathbb{R}_+ = 0$
- 2: **for** $k = 0, 1, \dots, K - 1$ **do**
- 3: **for** Each task $i = 1, \dots, N$ **do**
- 4: 1) Policy update:

$$\theta_i^{k+1} = \sum_{j \in \mathcal{N}_i} W_{i,j} \theta_j^k + \alpha \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) Q_i^{\pi_{\theta_i^k}} \quad (20)$$

$$\pi_i^{k+1}(a | s) = \frac{\exp(\theta_i^{k+1}(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta_i^{k+1}(s, a'))}$$

- 5: 2) Dual variable update:

$$\begin{aligned} \lambda_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k - \eta (V_i^{\pi_{\theta_i^k}}(\rho) - \ell_i) \right) \\ \nu_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k + \eta (V_i^{\pi_{\theta_i^k}}(\rho) - u_i) \right) \end{aligned} \quad (21)$$

- 6: **end for**
 - 7: **end for**
-

Our solution to the problem is to make each agent move in the direction of a locally computable portion of this gradient, followed by an averaging step (weighted according to matrix $W \in \mathbb{R}^{N \times N}$) that mixes the agents' policy parameters to achieve consensus. This specific update is shown in Algorithm 3 Eq. (20). In the long run, the local policy parameter obtained by following Eq. (20) behaves almost as if each agent is updated using the global aggregate gradient in Eq. (9).

To analyze the complexity of MT-PDNPG, we make the following assumption on the matrix W , which specifies the averaging weight in Eq. (20).

Assumption 3. *The matrix W is doubly stochastic, i.e. $\sum_{i=1}^N W_{i,j} = \sum_{j=1}^N W_{i,j} = 1$. In addition, $W_{i,j} > 0$ if and only if $(i, j) \in \mathcal{E}$ and $W_{i,j} = 0$ otherwise.*

This assumption is standard in the literature of consensus optimization Yuan et al. (2016); Zhang et al. (2018); Zeng et al. (2021; 2023). For any connected communication graph \mathcal{G} , a weight matrix W that satisfies Assumption 3 can be simply found using the lazy Metropolis method Olshevsky (2015). The largest singular value of W is always 1, and we use $\sigma_2(W) \in [0, 1)$ to denote its second largest singular value. In general, a more densely connected graph \mathcal{G} leads to a smaller $\sigma_2(W)$. We now present the first main theoretical result of the paper, which guarantees the finite-time convergence of Algorithm 3.

Theorem 1. *Consider the iterates $\{\pi_i^k\}$ obtained from K iterations of Algorithm 3. Let the step size sequences be*

$$\alpha = \frac{\alpha_0}{K^{1/2}}, \quad \eta = \frac{\eta_0}{K^{1/2}},$$

with $\alpha_0 = \mathcal{O}(\sqrt{1 - \sigma_2(W)})$. Then, under Assumption 1 and 3, we have for any $j = 1, \dots, N$

$$\max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} (V_0^{\pi^*}(\rho) - V_0^{\pi_j^k}(\rho)), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi_j^k}(\rho)]_+ + [V_i^{\pi_j^k}(\rho) - u_i]_+ \right) \right\}$$

$$\leq \mathcal{O}\left(\frac{N^{5/4}}{\sqrt{1 - \sigma_2(W)}K^{1/2}}\right).$$

Theorem 1 states that the policy iterate at every local agent i converges to the globally optimal multi-task policy π^* in both objective function and constraint violation with rate $\mathcal{O}(1/\sqrt{K})$, which matches the complexity of the algorithm in the centralized setting. The dependency of the bound on N reflects the increasing difficulty as the number of tasks scales up, while the inverse dependency on $\sqrt{1 - \sigma_2(W)}$ captures the impact of the communication graph \mathcal{G} and matrix W .

5.2 Sample-Based Setting

The same principle of distributing computation across the network can be applied to extend Algorithm 2 to the decentralized setting. Policy averaging is performed by each agent with its neighbors, while the critic and dual variables are updated completely locally. We present the updates formally in Algorithm 4 and study its finite-time complexity below.

Algorithm 4 Multi-Task Primal-Dual Natural Actor-Critic Algorithm (Decentralized)

- 1: **Initialization:** Each agent i initializes $\theta_i^0 \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} = 0$, dual variables $\lambda_i^0, \nu_i^0 \in \mathbb{R}_+ = 0$, critic parameters $\hat{Q}_i^0 = 0 \in \mathbb{R}^d$, and uniform behavior policy $\hat{\pi}_i^0$. For each task i , draw the initial sample s_i^0 and $a_i^0 \sim \hat{\pi}_i^0(\cdot | s_i^0)$
- 2: **for** $k = 0, 1, \dots, K - 1$ **do**
- 3: **for** Each task $i = 1, \dots, N$ **do**
- 4: 1) Observe $s_i^{k+1} \sim P(\cdot | s_i^k, a_i^k)$ and take action $a_i^{k+1} \sim \hat{\pi}_i^k(\cdot | s_i^{k+1})$
- 5: 2) Value function estimator (critic) update:

$$\hat{Q}_i^{k+1}(s_i^k, a_i^k) = (1 - \beta)\hat{Q}_i^k(s_i^k, a_i^k) + \beta \left(r_i(s_i^k, a_i^k) + \gamma \hat{Q}_i^k(s_i^{k+1}, a_i^{k+1}) \right) \quad (22)$$

- 6: 3) Policy (actor) update: $\forall s \in \mathcal{S}, a \in \mathcal{A}$

$$\theta_i^{k+1} = \sum_{j \in \mathcal{N}_i} W_{i,j} \theta_j^k + \alpha \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) \hat{Q}_i^k \quad (23)$$

$$\pi_i^{k+1}(a | s) = \frac{\exp(\theta_i^{k+1}(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta_i^{k+1}(s, a'))}$$

- 7: 4) Behavior policy update: $\forall s \in \mathcal{S}, a \in \mathcal{A}$

$$\hat{\pi}_i^{k+1}(a | s) = \frac{\epsilon}{|\mathcal{A}|} + (1 - \epsilon) \pi_i^{k+1}(a | s). \quad (24)$$

- 8: 5) Local dual variable update:

$$\begin{aligned} \lambda_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k - \eta \left(\sum_{s,a} \rho(s) \pi_i^k(a | s) \hat{Q}_i^k(s, a) - \ell_i \right) \right) \\ \nu_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\nu_i^k + \eta \left(\sum_{s,a} \rho(s) \pi_i^k(a | s) \hat{Q}_i^k(s, a) - u_i \right) \right) \end{aligned} \quad (25)$$

- 9: **end for**
 - 10: **end for**
-

Theorem 2. Consider the iterates $\{\pi_i^k\}$ from K iterations of Algorithm 4. Let the step size sequences be

$$\alpha = \frac{\alpha_0}{K^{5/6}}, \quad \beta = \frac{\beta_0}{K^{1/2}}, \quad \eta = \frac{\eta_0}{K^{5/6}}, \quad \epsilon = \frac{\epsilon_0}{K^{1/6}},$$

with $\frac{(1-\gamma)\mu\epsilon_0\beta_0}{|\mathcal{A}|} \leq 1$ and $\alpha_0 = \mathcal{O}(N^{-1/4}\sqrt{1-\sigma_2(W)})$. Then, under Assumptions 1-3, we have for any local agent j

$$\max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} (V_0^{\pi^*}(\rho) - V_0^{\pi_j^k}(\rho)), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi_j^k}(\rho)]_+ + [V_i^{\pi_j^k}(\rho) - u_i]_+ \right) \right\} \\ \leq \mathcal{O}\left(\frac{N^{5/4} \log(K)}{\sqrt{1-\sigma_2(W)} K^{1/6}}\right).$$

Theorem 2 generalizes the result in Section 4 and guarantees that the decentralized MT-PDNAC algorithm finds the globally optimally policy with a convergence rate of $\tilde{\mathcal{O}}(K^{-1/6})$. As Algorithm 4 again draws exactly N sample in each iteration, Theorem 2 implies a sample complexity of $\tilde{\mathcal{O}}(\delta^{-6})$ for converging to an optimal solution up to a precision δ .

A main challenge to the analysis of Algorithm 4 lies the coupling between the actor, critic, and dual variables, further complicated by the access to only local information. The actor controls the behavior policy which generates the samples for the update of the critic, which in turn is involved in the dual variable update. In addition, the accuracy of the critic, along with that of the dual variable, affects the update of the actor. To handle this intertwined system of variables, we leverage the fact that the effect of the samples on the actor (and dual variable) is indirect through the critic. This enables us to isolate the analysis of the actor from the evolution of the Markov chain conditioning on the critic. Specifically, as key components of the analysis, we show that the sub-optimality gap converge up to the cumulative error of the critic while the critic converges linearly fast to find the value function of the behavior policy under evolving Markovian samples. Further details of the analysis can be found in Section B.1 of the appendix.

6 Sample-Based Algorithm Under Linear Function Approximation

When we take an optimization approach to solve Eq. (3) in the previous sections, the policy and Q function are both $|\mathcal{S}||\mathcal{A}|$ -dimensional objects. This makes optimizing Eq. (3) (or even evaluating the objective function for a given policy) prohibitively expensive in real-life problems where the state space is huge or even infinitely large.

To reduce the dimensionality of the problem, we consider parameterizing the policy and Q function with linear function approximation. Suppose that the optimal policy parameter and Q functions of all tasks are be (approximately) represented using a given set of d basis vectors $\{\phi_1, \dots, \phi_d \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}\}$ where d can be much smaller than $|\mathcal{S}||\mathcal{A}|$. Each state and action pair (s, a) is associated with the feature vector $\phi(s, a) = [\phi_1(s, a), \dots, \phi_d(s, a)] \in \mathbb{R}^d$. We define the feature matrix $\Phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times d}$ as

$$\Phi = \begin{bmatrix} | & & | \\ \phi_1 & \cdots & \phi_d \\ | & & | \end{bmatrix} = \begin{bmatrix} - & \phi(s_0, a_0)^\top & - \\ \cdots & \cdots & \cdots \\ - & \phi(s_{|\mathcal{S}|}, a_{|\mathcal{A}|})^\top & - \end{bmatrix}.$$

Without loss of generality, we assume that $\|\phi(s, a)\| \leq 1$ for all $s \in \mathcal{S}, a \in \mathcal{A}$ and Φ is full column rank, and use $\sigma_{\min}(\Phi)$ and $\sigma_{\max}(\Phi)$ to denote the smallest and largest singular value of Φ , respectively.

Given the feature vectors, we maintain a parameter $\theta \in \mathbb{R}^d$ which represents the policy through the log-linear parameterization

$$\pi_\theta(a | s) = \frac{\exp(\phi(s, a)^\top \theta)}{\sum_{a' \in \mathcal{A}} \exp(\phi(s, a')^\top \theta)}. \quad (26)$$

We use the same features to approximate the Q function and restrict it to lie in the d -dimensional subspace $\hat{\mathcal{Q}} = \{\Phi\omega : \omega \in \mathbb{R}^d\}$. Given any policy π , its Q function in task i satisfies the Bellman equation

$$Q_i^\pi = T_i^\pi Q_i^\pi,$$

where $T_i^\pi : \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \rightarrow \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ is the Bellman backup operator

$$(T_i^\pi Q_i^\pi)(s, a) = r_i(s, a) + \gamma \sum_{a' \in \mathcal{A}} P(s' | s, a) \pi(a' | s') Q_i^\pi(s, a).$$

Under linear function approximation the feature matrix Φ may not exactly span Q_i^π . In this work, we solve the projected Bellman equation. Letting $M^\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ denote the diagonal matrix such that

$$M_{(s,a),(s,a)}^\pi = \tilde{\mu}_\pi(s, a) = \mu_\pi(s) \pi(a | s), \quad (27)$$

we define the operator $\Pi_\Phi^\pi : \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \rightarrow \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ as the projection to the linear subspace spanned by Φ under the weighted ℓ_2 norm $\|\cdot\|_\pi^2 = \cdot^\top M^\pi \cdot$. If the policy π is in the interior of the probability simplex (has uniformly lower bounded entries), the projected Bellman equation takes the following form

$$\Phi \omega = \Pi_\Phi^\pi T_i^\pi \Phi \omega. \quad (28)$$

A straightforward extension of Tsitsiklis & Vanroy (1997)[Theorem 1] guarantees that a unique solution to Eq. (28) exists, and we use $\omega_i^* : \Delta_{\mathcal{A}}^{\mathcal{S}} \rightarrow d$ to denote the mapping from a policy to its optimal value function parameter, which is the solution of Eq. (28). Note that $\omega_i^*(\pi)$ satisfies

$$\bar{H}^\pi \omega_i^*(\pi) + \bar{b}_i^\pi = 0, \quad (29)$$

where

$$\begin{aligned} \bar{H}^\pi &= \mathbb{E}_{s,a \sim \tilde{\mu}_\pi, s' \sim P(\cdot | s, a), a' \sim \pi(\cdot | s')} [\phi(s, a) (\gamma \phi(s', a') - \phi(s, a))^\top] \\ \bar{b}_i^\pi &= \mathbb{E}_{s,a \sim \tilde{\mu}_\pi} [r_i(s, a) \phi(s, a)]. \end{aligned} \quad (30)$$

6.1 Online Nested-Loop Algorithm

To extend Algorithm 4 to the linear function approximation setting, a most straightforward approach simply replaces the actor and critic updates in Eqs. (23) and (22) with their properly generalized versions (see Eqs.(35) and (34) of Algorithm 5). However, we note that the analysis of the single-loop actor-critic algorithm relies critically on the bounded variation of the TD learning target over iterations. In the tabular setting, the TD learning target is $Q_i^{\pi_i^k}$, the Q function of the current policy iterate π_i^k . The shift $\|Q_i^{\pi_i^{k+1}} - Q_i^{\pi_i^k}\|$ can be easily controlled by $\|\pi_i^{k+1} - \pi_i^k\|$ and eventually by the step size α , due to the following Lipschitz condition (established in Lemma 2)

$$\|Q_i^\pi - Q_i^{\pi'}\| \leq \frac{\gamma |\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2} \|\pi - \pi'\|, \quad \forall \pi, \pi'.$$

In the case of linear function approximation, the TD learning target is ω_i^* , which is nevertheless not a (sufficiently) Lipschitz operator. It can be shown that there exists π, π' and constant $L > 0$ such that

$$\|\omega_i^*(\pi) - \omega_i^*(\pi')\| \geq \frac{L}{\min_{s,a} \{\pi(a | s), \pi'(a | s)\}} \|\pi - \pi'\|.$$

However, to maintain the sample complexity of $\tilde{\mathcal{O}}(\delta^{-6})$ of Algorithm 4, we need at least

$$\|\omega_i^*(\pi_i^k) - \omega_i^*(\pi_i^{k+1})\| \leq \frac{L}{(\min_{s,a} \{\pi_i^k(a | s), \pi_i^{k+1}(a | s)\})^{1/2}} \|\pi_i^k - \pi_i^{k+1}\|, \quad \forall k.$$

This apparent gap in the Lipschitz continuity of ω_i^* means that we cannot apply the analysis established earlier in the paper. To work around the challenge without degrading the sample complexity, we employ nested-loop updates in the linear function approximation setting. Stated formally in Algorithm 5, our method updates the policy, dual variable, and behavior policy in the outer loop and allows more iterations for the critic in the inner loop for it to chase the moving target $\omega_i^*(\pi_i^k)$.

Remark 1. We can still use a single-loop algorithm under linear function approximation despite the degraded Lipschitz condition on the TD learning target. By choosing the step sizes slightly differently from the tabular case, the single-loop algorithm will converge in finite time, but the complexity is slightly worse than $\tilde{\mathcal{O}}(K^{-1/6})$ established in Theorem 2.

Note that Eq. (5) is equivalent to the following direct update on π_i^k Chen et al. (2022b)[Lemma 3.1]

$$\pi_i^{k+1}(a | s) \propto \pi_i^k(a | s) \exp(\alpha(1/N + \lambda_i^k - \nu_i^k)\phi(s, a)^\top)\omega_i^{k,T}. \quad (31)$$

We present the expressions of π_i^{k+1} and $\hat{\pi}_i^{k+1}$ in Eqs. (35) and (36) only for the purpose of clarity. As the policies are large $|\mathcal{S}||\mathcal{A}|$ -dimensional objects, the policies are never explicitly maintained or updated when the algorithm is deployed; we only need to track the policy parameter θ_i^k .

6.2 Finite-Sample Complexity

This section presents the complexity of the nested-loop actor-critic algorithm under linear function approximation. As the function approximation is not necessarily perfect (i.e. the true Q functions may not lie in the column space of the feature matrix Φ), the optimality error does not converge to 0 asymptotically but rather depends on the approximation error associated with Φ , which is the distance between Q_i^π and its approximation in \hat{Q} . In the next assumption, we assume that there is a uniform upper bound on the approximation error.

Assumption 4. There exists a constant $\varepsilon_{\max} > 0$ such that for any π , we have

$$\|\Phi\omega_i^*(\pi) - Q_i^\pi\| \leq \varepsilon_{\max}. \quad (32)$$

A consequence of the assumption is that the target critic parameter always has a bounded norm. More specifically, we have $\|\omega_i^*(\hat{\pi}_i^k)\| \leq B_\omega$ for all $k \geq 0$, where

$$B_\omega = \sigma_{\min}^{-1}(\Phi) \left(\sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{1-\gamma}} + \varepsilon_{\max} \right). \quad (33)$$

This result is stated and proved in Lemma 5 of the appendix. We need to confine the possibly infinite growth of the critic parameter ω_i^k through projection. Knowing the boundedness of $\omega_i^*(\hat{\pi}_i^k)$ allows us to use the operator $\Pi_{B_\omega} : \mathbb{R}^d \rightarrow \mathbb{R}^d$, which projects a vector to the ℓ_2 ball with radius B_ω .

Theorem 3. Let $\delta > 0$ be a desired precision. Under Assumptions 1-4 and properly selected step sizes, the iterates of Algorithm 5 satisfies for any agent j

$$\begin{aligned} \max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} (V_0^{\pi^*}(\rho) - V_0^{\pi_j^k}(\rho)), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi_j^k}(\rho)]_+ + [V_i^{\pi_j^k}(\rho) - u_i]_+ \right) \right\} \\ \leq \frac{N^{5/4}\delta}{\sqrt{1-\sigma_2(W)}} + \mathcal{O}(N\varepsilon_{\max}) \end{aligned}$$

with at most $K = \mathcal{O}(\delta^{-2})$ outer loop iterations and at most $\mathcal{O}(\frac{\log(1/\delta)}{\delta^6})$ total samples.

Theorem 3 again establishes a $\tilde{\mathcal{O}}(\delta^{-6})$ convergence for every agent's local policy parameter, which scales inversely with $1 - \sigma_2(W)$, the spectral gap of the weight matrix. Due to the presence of approximation error, the optimality gap in objective function and constraint violation does not converge to 0, but to a quantity proportional to ε_{\max} . Prior to our work, there are few results on the finite-sample complexity of data-driven algorithms for solving the CMDP optimization problem under linear function approximation, even in the single task setting. The recent work Ding et al. (2022) targets this problem and studies a primal-dual REINFORCE-flavored algorithm. Their overall sample complexity is $\mathcal{O}(\delta^{-8})$, which we significantly improve over. Finally, we note that the sample complexity here matches that of Algorithm 4 for the tabular setting.

Algorithm 5 Multi-Task Primal-Dual Natural Actor-Critic Algorithm under Linear Function Approximation (Decentralized)

- 1: **Initialization:** For each task i , initialize $\theta_i^0 \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} = 0$, dual variables $\lambda_i^0, \nu_i^0 \in \mathbb{R}_+ = 0$, critic parameters $\hat{Q}_i^0 = 0 \in \mathbb{R}^d$, and uniform behavior policy $\hat{\pi}_i^0$. For each task i , draw the initial sample $s_i^{0,0}$ and $a_i^{0,0} \sim \hat{\pi}_i^{0,0}(\cdot | s_i^{0,0})$
- 2: **for** $k = 0, 1, \dots, K-1$ **do**
- 3: **for** Each task $i = 1, \dots, N$ **do**
- 4: **for** $t = 0, 1, \dots, T-1$ **do**
- 5: 1) Observe $s_i^{k,t+1} \sim P(\cdot | s_i^{k,t}, a_i^{k,t})$ and take action $a_i^{k,t+1} \sim \hat{\pi}_i^k(\cdot | s_i^{k,t+1})$
- 6: 2) Critic parameter update:

$$\begin{aligned}\hat{\omega}_i^{k,t+1} &= \omega_i^{k,t} + \beta \phi(s_i^{k,t}, a_i^{k,t}) \left(r_i(s_i^{k,t}, a_i^{k,t}) + (\gamma \phi(s_i^{k,t+1}, a_i^{k,t+1}) - \phi(s_i^{k,t}, a_i^{k,t}))^\top \omega_i^{k,t} \right) \\ \omega_i^{k,t+1} &= \Pi_{B_\omega} \left(\hat{\omega}_i^{k,t+1} \right)\end{aligned}\tag{34}$$

- 7: **end for**
- 8: 3) Policy (actor) update: $\forall s \in \mathcal{S}, a \in \mathcal{A}$

$$\begin{aligned}\theta_i^{k+1} &= \sum_{j \in \mathcal{N}_i} W_{i,j} \theta_j^k + \alpha \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) \omega_i^{k,T} \\ \pi_i^{k+1}(a | s) &= \frac{\exp \left(\theta_i^{k+1}(s, a)^\top \omega_i^{k,T} \right)}{\sum_{a' \in \mathcal{A}} \exp \left(\theta_i^{k+1}(s, a')^\top \omega_i^{k,T} \right)}\end{aligned}\tag{35}$$

- 9: 4) Behavior policy update: $\forall s \in \mathcal{S}, a \in \mathcal{A}$

$$\hat{\pi}_i^{k+1}(a | s) = \frac{\epsilon}{|\mathcal{A}|} + (1 - \epsilon) \pi_i^{k+1}(a | s).\tag{36}$$

- 10: 5) Dual variable update:

$$\begin{aligned}\lambda_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k - \eta \left(\sum_{s,a} \rho(s) \pi_i^k(a | s) \phi(s, a)^\top \omega_i^k - \ell_i \right) \right) \\ \nu_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\nu_i^k + \eta \left(\sum_{s,a} \rho(s) \pi_i^k(a | s) \phi(s, a)^\top \omega_i^k - u_i \right) \right)\end{aligned}\tag{37}$$

- 11: 6) Set $s_i^{k+1,0} = s_i^{k,T}$ and $a_i^{k+1,0} = a_i^{k,T}$
 - 12: **end for**
 - 13: **end for**
-

7 Numerical Simulations

In this section, we use a small-scale GridWorld experiment to show how the constrained formulation allows us to control the multi-task policy in a fine-grained manner. We consider a three-task RL problem. Associated with each task is a 10×10 maze in which an agent aims to navigate to a goal position by crossing bridges that charge different “prices” (i.e. incur different negative rewards). Shown in Figure 1, the starting position for all tasks is the top left corner and the target is the top right corner, both marked in blue. There is a positive reward for reaching the target, which varies in magnitude across tasks.

In the first task, the agent receives the least negative reward for using the first bridge, which is -0.1. The reward of the fourth bridge, -1, is more negative, and the reward of the second and third bridges are the most negative. The reward of any other move in the maze is -0.1. If we are only interested in solving task

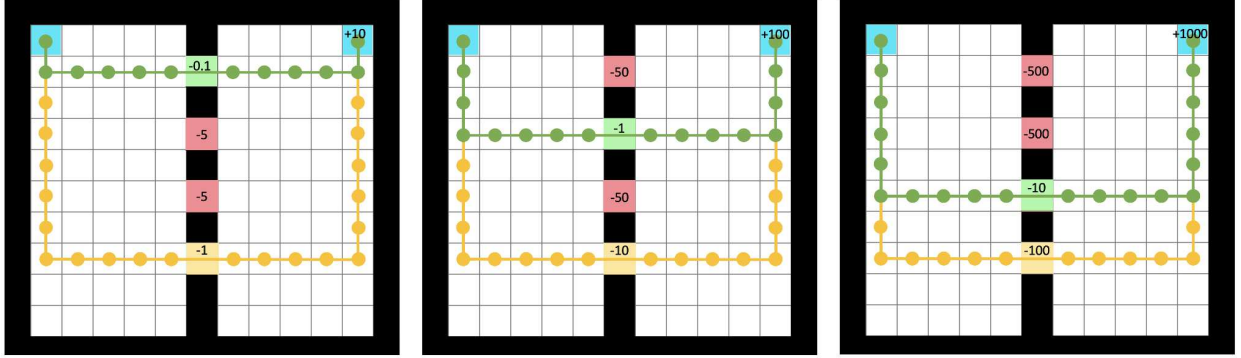


Figure 1: Reward in maze 1 (left): 10 for reaching target, -5 for crossing the second and third bridge, -1 for crossing the fourth bridge, and -0.1 for any other move. Reward in maze 2 (middle): 100 for reaching target, -50 for crossing the first and third bridge, -10 for crossing the fourth bridge, and -1 for any other move. Reward in maze 3 (right): 1000 for reaching target, -500 for crossing the second and the third bridge, -100 for crossing the fourth bridge, and -10 for any other move. Green dotted lines indicate optimal paths for local tasks. The yellow dotted line indicates a sub-optimal but acceptable policy for each local task, which is also the globally optimal policy of the constrained multi-task problem under $\ell_1 = 5$, $\ell_2 = 50$, $\ell_3 = 500$.

1, the optimal policies (one of which is drawn as the green path in the figure) obviously want to use the first bridge to reach the target as soon as possible.

The second task can be regarded as a reward-magnified version of task 1. With the rewards marked in Figure 1, the best bridge becomes the second one, which charges much lower prices than the other bridges. Making any other move not labeled incurs a reward of -1 . One of the optimal policies for this task is drawn as the green dotted path, which uses the second bridge.

The rewards for all moves in the third task are further scaled up. The bonus of reaching the target increases to 1000, while the costs of crossing the bridges becomes 500, 100, and 10. The reward of making any other move also changes to -10 . Since the magnitude of the rewards in this task is dominant over those in the other two tasks, it is easy to verify that the optimal policies of the third task coincide with the globally optimal policy for the unconstrained multi-task RL objective in Eq. (2), one of which is again shown in green in the figure. This optimal path takes the third bridge, which ensures that the greatest reward is collected in the third tasks, at the cost of completely unsatisfactory performance in the other two tasks.

Comparing with solving Eq. (2), our formulation in Eq. (3) allows us to trade off the policy performance in the three tasks by properly specifying the performance lower bounds. In particular, we set $\ell_1 = 5, \ell_2 = 50$, and $\ell_3 = 500$. Calculations show that the optimal policies under this set of constraints will switch to take the fourth bridge, which means a slight but acceptable compromise of the policy performance in all tasks. Numerically, we apply our proposed primal-dual natural gradient algorithm¹, and verify that the local policy at every agent indeed converges to the optimal constrained policy in both objective function and constraint violation (see the first and second plots in Figure 2). In contrast, we also apply the same algorithm with $\ell_1 = \ell_2 = \ell_3 = -\infty$, which essentially means that we run a decentralized natural policy gradient algorithm to solve the unconstrained multi-task problem in Eq. (2). The policy iterates of this unconstrained problem obviously achieves better objective function value, but fails to satisfy the constraints.

8 Conclusion

This paper considers a constrained multi-task RL objective where the goal is to find a policy with the maximum average performance across all tasks and guaranteed minimum performance in each environment. We show that three policy-gradient-based algorithms provably solve the problem, under different policy parameterization and information oracle. In the tabular setting with the access to the exact gradients, we

¹The small dimension and known transition kernel of the environments allow us to derive the exact gradient in this case.

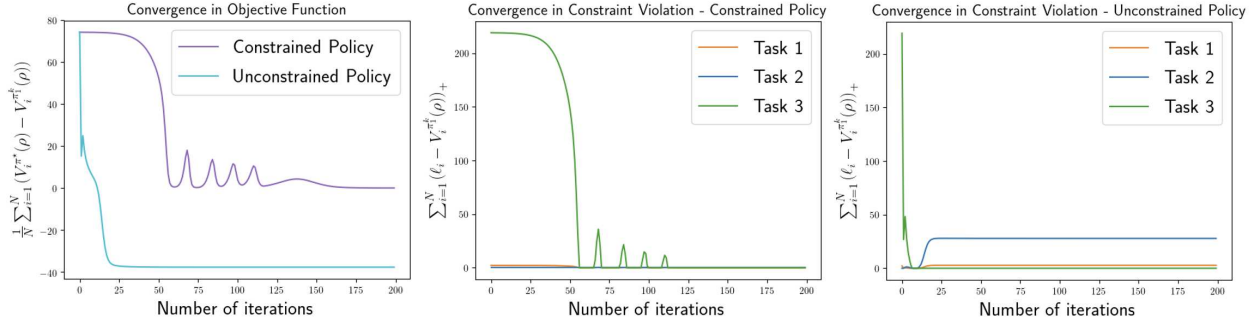


Figure 2: Left – convergence of Algorithm 3 in objective function with $\ell_1 = 5$, $\ell_2 = 50$, $\ell_3 = 500$ (constrained policy), and with $\ell_1 = \ell_2 = \ell_3 = -\infty$ (unconstrained policy). Middle – convergence of Algorithm 3 in constraint violation with $\ell_1 = 5$, $\ell_2 = 50$, $\ell_3 = 500$. Right – convergence of Algorithm 3 in constraint violation with $\ell_1 = \ell_2 = \ell_3 = -\infty$.

study a primal-dual natural-gradient-descent-ascent algorithm that drives each agent’s local policy to the globally optimal solution with finite-time complexity $\mathcal{O}(K^{-1/2})$. When the exact gradient information is not available, we take a completely sampled-based actor-critic approach and show that it converges with rate $\mathcal{O}(K^{-1/6})$. Finally, we extend the results to the linear function approximation setting and establish the same sample complexity as in the tabular case. To our best knowledge, a finite-time and finite-sample analysis for the on-policy natural actor-critic algorithm in the linear function approximation setting (even for standard non-constrained MDPs) has been missing from the existing literature, and our work fills in this gap of knowledge.

Disclaimer

This paper was prepared for informational purposes in part by the Artificial Intelligence Research group of JP Morgan Chase & Co and its affiliates (“JP Morgan”), and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

References

- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. Optimality and approximation with policy gradient methods in markov decision processes. In *Conference on Learning Theory*, pp. 64–66. PMLR, 2020.
- Mridul Agarwal, Qinbo Bai, and Vaneet Aggarwal. Regret guarantees for model-based reinforcement learning with long-term average constraints. In *Uncertainty in Artificial Intelligence*, pp. 22–31. PMLR, 2022.
- Eitan Altman. *Constrained Markov decision processes*, volume 7. Chapman and Hall/CRC Press, 1999.
- Qinbo Bai, Amrit Singh Bedi, Mridul Agarwal, Alec Koppel, and Vaneet Aggarwal. Achieving zero constraint violation for constrained reinforcement learning via primal-dual approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 3682–3689, 2022.
- Anas Barakat, Pascal Bianchi, and Julien Lehmann. Analysis of a target-based actor-critic algorithm with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pp. 991–1040. PMLR, 2022.

-
- Vivek S Borkar. An actor-critic algorithm for constrained markov decision processes. *Systems & control letters*, 54(3):207–213, 2005.
- Brian Bullins, Kshitij Patel, Ohad Shamir, Nathan Srebro, and Blake E Woodworth. A stochastic newton algorithm for distributed convex optimization. *Advances in Neural Information Processing Systems*, 34: 26818–26830, 2021.
- Tsung-Hui Chang, Angelia Nedić, and Anna Scaglione. Distributed constrained optimization by consensus-based primal-dual perturbation method. *IEEE Transactions on Automatic Control*, 59(6):1524–1538, 2014.
- Jinchi Chen, Jie Feng, Weiguo Gao, and Ke Wei. Decentralized natural policy gradient with variance reduction for collaborative multi-agent reinforcement learning. *arXiv preprint arXiv:2209.02179*, 2022a.
- Zaiwei Chen, Sajad Khodadadian, and Siva Theja Maguluri. Finite-sample analysis of off-policy natural actor-critic with linear function approximation. *IEEE Control Systems Letters*, 6:2611–2616, 2022b.
- Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained reinforcement learning with percentile risk criteria. *The Journal of Machine Learning Research*, 18(1):6070–6120, 2017.
- Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. A lyapunov-based approach to safe reinforcement learning. *arXiv preprint arXiv:1805.07708*, 2018.
- Carlo D’Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. Sharing knowledge in multi-task deep reinforcement learning. In *International Conference on Learning Representations*.
- Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo Jovanovic. Natural policy gradient primal-dual method for constrained markov decision processes. *Advances in Neural Information Processing Systems*, 33:8378–8390, 2020.
- Dongsheng Ding, Xiaohan Wei, Zhuoran Yang, Zhaoran Wang, and Mihailo Jovanovic. Provably efficient safe exploration via primal-dual policy optimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 3304–3312. PMLR, 2021.
- Dongsheng Ding, Kaiqing Zhang, Tamer Başar, and Mihailo R Jovanović. Convergence and optimality of policy gradient primal-dual method for constrained markov decision processes. In *2022 American Control Conference (ACC)*, pp. 2851–2856. IEEE, 2022.
- Thinh Doan, Siva Maguluri, and Justin Romberg. Finite-time analysis of distributed td (0) with linear function approximation on multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 1626–1635. PMLR, 2019.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Yijie Guo, Qiucheng Wu, and Honglak Lee. Learning action translator for meta reinforcement learning on sparse-reward tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6792–6800, 2022.
- Abhishek Gupta, Coline Devin, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Learning invariant feature spaces to transfer skills with reinforcement learning. *arXiv preprint arXiv:1703.02949*, 2017.
- Aria HasanzadeZonuzi, Archana Bura, Dileep Kalathil, and Srinivas Shakkottai. Learning with safety constraints: Sample complexity of reinforcement learning for constrained mdps. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 7667–7674, 2021.
- Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane, Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz, et al. A practical guide to multi-objective reinforcement learning and planning. *Autonomous Agents and Multi-Agent Systems*, 36(1):26, 2022.

-
- Matteo Hessel, Hubert Soyer, Lasse Espeholt, Wojciech Czarnecki, Simon Schmitt, and Hado van Hasselt. Multi-task deep reinforcement learning with popart. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3796–3803, 2019.
- Sunghoon Hong, Deunsol Yoon, and Kee-Eung Kim. Structure-aware transformer policy for inhomogeneous multi-task reinforcement learning. In *International Conference on Learning Representations*, 2021.
- Rustem Islamov, Xun Qian, and Peter Richtárik. Distributed second order methods with fast rates and compressed communication. In *International conference on machine learning*, pp. 4617–4628. PMLR, 2021.
- Zhanhong Jiang, Xian Yeow Lee, Sin Yong Tan, Kai Liang Tan, Aditya Balu, Young M Lee, Chinmay Hegde, and Soumik Sarkar. Mdpgt: momentum-based decentralized policy gradient tracking. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 9377–9385, 2022.
- Caleb Ju, Georgios Kotsalis, and Guanghai Lan. A model-free first-order method for linear quadratic regulator with $\tilde{O}(1/\varepsilon)$ sampling complexity. *arXiv preprint arXiv:2212.00084*, 2022.
- Shi Junru, Wang Qiong, Liu Muhua, Ji Zhihang, Zheng Ruijuan, and Wu Qingtao. Decentralized multi-task reinforcement learning policy gradient method with momentum over networks. *Applied Intelligence*, pp. 1–15, 2022.
- Dmitry Kalashnikov, Jacob Varley, Yevgen Chebotar, Benjamin Swanson, Rico Jonschkowski, Chelsea Finn, Sergey Levine, and Karol Hausman. Mt-opt: Continuous multi-task robotic reinforcement learning at scale. *arXiv preprint arXiv:2104.08212*, 2021.
- Sajad Khodadadian, Thinh T Doan, Justin Romberg, and Siva Theja Maguluri. Finite sample analysis of two-time-scale natural actor-critic algorithm. *IEEE Transactions on Automatic Control*, 2022.
- Jinlong Lei, Han-Fu Chen, and Hai-Tao Fang. Primal–dual algorithm for distributed constrained optimization. *Systems & Control Letters*, 96:110–117, 2016.
- Tao Liu, Ruida Zhou, Dileep Kalathil, Panganamala Kumar, and Chao Tian. Learning policies with zero or bounded constraint violation for constrained mdps. *Advances in Neural Information Processing Systems*, 34:17183–17193, 2021a.
- Tao Liu, Ruida Zhou, Dileep Kalathil, PR Kumar, and Chao Tian. Policy optimization for constrained mdps with provable fast global convergence. *arXiv preprint arXiv:2111.00552*, 2021b.
- Angelia Nedic. Distributed gradient methods for convex machine learning problems in networks: Distributed optimization. *IEEE Signal Processing Magazine*, 37(3):92–101, 2020.
- Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- Alex Olshevsky. Linear time average consensus on fixed graphs. *IFAC-PapersOnLine*, 48(22):94–99, 2015.
- Santiago Paternain, Luiz FO Chamon, Miguel Calvo-Fullana, and Alejandro Ribeiro. Constrained reinforcement learning has zero duality gap. *arXiv preprint arXiv:1910.13393*, 2019.
- LA Prashanth and Mohammad Ghavamzadeh. Variance-constrained actor-critic algorithms for discounted and average reward mdps. *Machine Learning*, 105(3):367–417, 2016.
- 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*.
- Andrei A Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. Policy distillation. *arXiv preprint arXiv:1511.06295*, 2015.

-
- Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization. *arXiv preprint arXiv:1805.11074*, 2018.
- René Traoré, Hugo Caselles-Dupré, Timothée Lesort, Te Sun, Guanghang Cai, Natalia Díaz-Rodríguez, and David Filliat. Disorcl: Continual reinforcement learning via policy distillation. *arXiv preprint arXiv:1907.05855*, 2019.
- JN Tsitsiklis and B Vanroy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5):674–690, 1997.
- Samir Wadhwan, Dong-Ki Kim, Shayegan Omidshafiei, and Jonathan P How. Policy distillation and value matching in multiagent reinforcement learning. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 8193–8200. IEEE, 2019.
- Yue Frank Wu, Weitong Zhang, Pan Xu, and Quanquan Gu. A finite-time analysis of two time-scale actor-critic methods. *Advances in Neural Information Processing Systems*, 33:17617–17628, 2020.
- 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.
- Ruihan Yang, Huazhe Xu, Yi Wu, and Xiaolong Wang. Multi-task reinforcement learning with soft modularization. *Advances in Neural Information Processing Systems*, 33:4767–4777, 2020.
- Zhuoran Yang, Yongxin Chen, Mingyi Hong, and Zhaoran Wang. Provably global convergence of actor-critic: A case for linear quadratic regulator with ergodic cost. *Advances in neural information processing systems*, 32, 2019.
- Donghao Ying, Yuhao Ding, and Javad Lavaei. A dual approach to constrained markov decision processes with entropy regularization. In *International Conference on Artificial Intelligence and Statistics*, pp. 1887–1909. PMLR, 2022.
- Ming Yu, Zhuoran Yang, Mladen Kolar, and Zhaoran Wang. Convergent policy optimization for safe reinforcement learning. *Advances in Neural Information Processing Systems*, 32:3127–3139, 2019.
- Kun Yuan, Qing Ling, and Wotao Yin. On the convergence of decentralized gradient descent. *SIAM Journal on Optimization*, 26(3):1835–1854, 2016.
- Sihan Zeng, Malik Aqeel Anwar, Thinh T Doan, Arijit Raychowdhury, and Justin Romberg. A decentralized policy gradient approach to multi-task reinforcement learning. In *Uncertainty in Artificial Intelligence*, pp. 1002–1012. PMLR, 2021.
- Sihan Zeng, Thinh T Doan, and Justin Romberg. Finite-time complexity of online primal-dual natural actor-critic algorithm for constrained markov decision processes. In *2022 IEEE 61st Conference on Decision and Control (CDC)*, pp. 4028–4033. IEEE, 2022.
- Sihan Zeng, Thinh T Doan, and Justin Romberg. Finite-time convergence rates of decentralized stochastic approximation with applications in multi-agent and multi-task learning. *IEEE Transactions on Automatic Control*, 68:2758–2773, 2023.
- Sihan Zeng, Thinh T Doan, and Justin Romberg. A two-time-scale stochastic optimization framework with applications in control and reinforcement learning. *SIAM Journal on Optimization*, 34(1):946–976, 2024.
- Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Basar. Fully decentralized multi-agent reinforcement learning with networked agents. In *International Conference on Machine Learning*, pp. 5872–5881. PMLR, 2018.
- Liyuan Zheng and Lillian Ratliff. Constrained upper confidence reinforcement learning. In *Learning for Dynamics and Control*, pp. 620–629. PMLR, 2020.

A Additional Notations

We introduce some additional shorthand notations frequently used in the appendix. First, we define $\hat{V}_i^k \in \mathbb{R}^{|\mathcal{S}|}$ and $\hat{A}_i^k \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ as the estimated value function and advantage function based on \hat{Q}_i^k such that

$$\begin{aligned}\hat{V}_i^k(s) &\triangleq \sum_{a' \in \mathcal{A}} \pi_i^k(a' | s) \hat{Q}_i^k(s, a'), \quad \forall s \in \mathcal{S}, \\ \hat{A}_i^k(s, a) &\triangleq \hat{Q}_i^k(s, a) - \hat{V}_i^k(s), \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.\end{aligned}$$

We use the subscript g in the value/Q/advantage function to denote the vector formed by stacking the value/Q/advantage functions across agents $i = 1, \dots, N$. Specifically, given a policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$, we write

$$\begin{aligned}V_g^\pi(s) &\triangleq [V_1^\pi(s), \dots, V_N^\pi(s)]^\top \in \mathbb{R}^N \\ Q_g^\pi(s, a) &\triangleq [Q_1^\pi(s, a), \dots, Q_N^\pi(s, a)]^\top \in \mathbb{R}^N \\ A_g^\pi(s, a) &\triangleq [A_1^\pi(s, a), \dots, A_N^\pi(s, a)]^\top \in \mathbb{R}^N\end{aligned}$$

B Proof of Theorems

B.1 Proof of Theorem 1

The proof of Theorem 1 relies on the proposition below, which is an intermediate result that characterizes the convergence of primal and dual variables under a general update rule. This proposition will also be used in the proofs of Theorem 2 and 3.

Proposition 1. *Consider any decentralized algorithm where each agent in iteration k maintains dual variables $\lambda_i^k, \nu_i^k \in \mathbb{R}$, critic variables $Q_i^k \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$, and a primal variable $\theta_i^k \in \mathbb{R}^d$ (for any positive integer d) which parameterizes a policy $\pi_i^k = g_i(\theta_i^k) \in \Delta_{\mathcal{A}}^{\mathcal{S}}$ through some mapping g_i . Suppose the dual variable updates satisfy*

$$\begin{aligned}\lambda_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k - \eta \left(\sum_{s, a} \rho(s) \pi_i^k(a | s) Q_i^k(s, a) - \ell_i \right) \right), \\ \nu_i^{k+1} &= \Pi_{[0, B_\lambda]} \left(\lambda_i^k + \eta \left(\sum_{s, a} \rho(s) \pi_i^k(a | s) Q_i^k(s, a) - u_i \right) \right).\end{aligned}\tag{38}$$

Also suppose that there exists a constant $B > 0$ such that $\|Q_i^k\|_\infty \leq B$ for all i and k , and that there exists a mapping $f : \mathbb{R}^{N \times d} \rightarrow \Delta_{\mathcal{A}}^{\mathcal{S}}$ such that $\bar{\pi}^{k+1} \triangleq f(\{\theta_i^{k+1}\}_i)$ observes the recursive rule

$$\bar{\pi}^{k+1}(a | s) \propto \bar{\pi}^k(a | s) \exp\left(\frac{\alpha}{N} \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k\right) Q_i^k(s, a)\right), \quad \forall k.\tag{39}$$

Then, the parameters $\{\pi_i^k\}_{i,k}$ after K update iterations satisfy

$$\begin{aligned}&\max\left\{\frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho)\right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\bar{\pi}^k}(\rho)]_+ + [V_i^{\bar{\pi}^k}(\rho) - u_i]_+\right)\right\} \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{N}{K\eta} + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|\right).\end{aligned}$$

We also introduce the following lemma which establishes the Lipschitz continuity of the value function and Q function.

Lemma 2 (Lemma 8 of Khodadadian et al. (2022)). *For any policy π_1, π_2 and $i = 1, \dots, N$*

$$\begin{aligned} \|Q_i^{\pi_1} - Q_i^{\pi_2}\| &\leq \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2} \|\pi_1 - \pi_2\|, \quad \|V_i^{\pi_1} - V_i^{\pi_2}\| \leq \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2} \|\pi_1 - \pi_2\|, \\ |Q_i^{\pi_1}(s, a) - Q_i^{\pi_2}(s, a)| &\leq \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{(1-\gamma)^2} \|\pi_1 - \pi_2\|, \quad |V_i^{\pi_1}(s) - V_i^{\pi_2}(s)| \leq \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{(1-\gamma)^2} \|\pi_1 - \pi_2\|. \end{aligned}$$

It is easy to verify that Algorithm 3 observes the update rule in 38 and 39, with $Q_i^k = Q_i^{\pi_i^k}$ and $\bar{\pi}^{k+1} = f(\{\theta_i^{k+1}\}_i)$ defined as

$$\bar{\theta}^{k+1} = \frac{1}{N} \sum_{i=1}^N \theta_i^{k+1}, \quad \bar{\pi}^{k+1} = \frac{\exp(\bar{\theta}^{k+1}(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\bar{\theta}^{k+1}(s, a'))}$$

Due to the bounded update, we can control the distance $\|\bar{\pi}^k - \pi_i^k\|$ by the step size α .

Lemma 3. *The policy iterates $\{\pi_i^k\}$ generated by Algorithm 3 satisfy*

$$\|\bar{\pi}^k - \pi_i^k\| \leq \mathcal{O}\left(\frac{\sqrt{N}\alpha}{1 - \sigma_2(W)}\right), \quad \text{for all } k = 0, \dots, K-1 \text{ and } i = 1, \dots, N.$$

As a result of the proposition and lemma above,

$$\begin{aligned} &\max\left\{\frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho)\right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\bar{\pi}^k}(\rho)]_+ + [V_i^{\bar{\pi}^k}(\rho) - u_i]_+\right)\right\} \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{N}{K\eta} + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|\right) \\ &= \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{N}{K\eta} + \frac{N^{3/2}\alpha}{1 - \sigma_2(W)}\right). \end{aligned}$$

By the bound on consensus error in Lemma 3 and the Lipschitz continuity of the value function in Lemma 2, this implies for any agent $j = 1, \dots, N$

$$\begin{aligned} &\max\left\{\frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\pi_j^k}(\rho)\right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi_j^k}(\rho)]_+ + [V_i^{\pi_j^k}(\rho) - u_i]_+\right)\right\} \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{N}{K\eta} + \frac{N^{3/2}\alpha}{1 - \sigma_2(W)}\right). \end{aligned}$$

Choosing the step sizes as $\alpha = \mathcal{O}\left(\frac{\sqrt{1 - \sigma_2(W)}}{N^{1/4}\sqrt{K}}\right)$ and $\eta = \mathcal{O}(1/\sqrt{K})$ leads to the claimed result.

B.2 Proof of Theorem 2

Since each agent maintains and learns the local value function, the analysis of the critic is the same as in the single agent setting Zeng et al. (2022). Specifically, we define the critic error

$$z_i^k = \hat{Q}_i^k - Q_i^{\bar{\pi}_i^k} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}. \quad (40)$$

and present its convergence rate in the following proposition.

Proposition 2. *Let Assumption 2 hold and $\frac{(1-\gamma)\mu\epsilon\beta}{|\mathcal{A}|} \leq 1$, then the iterates $\{\hat{Q}_i^k\}$ and $\{\bar{\pi}_i^k\}$ satisfy*

$$\frac{1}{K} \sum_{k=0}^{K-1} \|z_i^k\| \leq \mathcal{O}\left(\frac{\tau^{1/2}}{\epsilon^{1/2}\beta^{1/2}K^{1/2}} + \frac{N\beta^{1/2}\tau}{\epsilon^{1/2}} + \frac{N\alpha}{\epsilon\beta}\right).$$

Lemma 4. *The policy iterates $\{\pi_i^k\}$ generated by Algorithm 4 satisfy*

$$\|\bar{\pi}^k - \pi_i^k\| \leq \mathcal{O}\left(\frac{\sqrt{N}\alpha}{1 - \sigma_2(W)}\right), \quad \text{for all } k = 0, \dots, K-1 \text{ and } i = 1, \dots, N.$$

Similar to the proof of Theorem 1, we can verify that Algorithm 4 observes the update rule in 38 and 39, with $Q_i^k = \hat{Q}_i^k$ and $\bar{\pi}^{k+1} = f(\{\theta_i^{k+1}\}_i)$ being defined as

$$\bar{\theta}^{k+1} = \frac{1}{N} \sum_{i=1}^N \theta_i^{k+1}, \quad \bar{\pi}^{k+1} = \frac{\exp(\bar{\theta}^{k+1}(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\bar{\theta}^{k+1}(s, a'))}$$

The bounds in Proposition 1 depend on the error $\frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=0}^N \|\hat{Q}_i^k - Q_i^{\pi_i^k}\|$, which can be decomposed. From the definition of z_i^k in Eq. (40), we have for any $i = 0, 1, \dots, N$

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \|\hat{Q}_i^k - Q_i^{\pi_i^k}\| &\leq \frac{1}{K} \sum_{k=0}^{K-1} \|Q_i^{\hat{\pi}_i^k} - Q_i^{\pi_i^k}\| + \frac{1}{K} \sum_{k=0}^{K-1} \|z_i^k\| \\ &\leq \frac{1}{K} \sum_{k=0}^{K-1} \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2} \|\hat{\pi}_i^k - \pi_i^k\| + \frac{1}{K} \sum_{k=0}^{K-1} \|z_i^k\| \\ &\leq \frac{2|\mathcal{S}|^{3/2}|\mathcal{A}|\epsilon}{(1-\gamma)^2} + \frac{1}{K} \sum_{k=0}^{K-1} \|z_i^k\| \\ &= \mathcal{O}\left(\epsilon + \frac{\tau^{1/2}}{\epsilon^{1/2}\beta^{1/2}K^{1/2}} + \frac{N\beta^{1/2}\tau}{\epsilon^{1/2}} + \frac{N\alpha}{\epsilon\beta}\right), \end{aligned} \quad (41)$$

where the second inequality uses Lemma 2 and the third inequality follows from

$$\begin{aligned} \|\hat{\pi}_i^k - \pi_i^k\| &= \left\| \frac{\epsilon}{|\mathcal{A}|} \mathbf{1}_{|\mathcal{S}||\mathcal{A}|} + (1-\epsilon)\pi_i^k - \pi_i^k \right\| \\ &\leq \epsilon \left\| \frac{1}{|\mathcal{A}|} \mathbf{1}_{|\mathcal{S}||\mathcal{A}|} \right\| + \epsilon \|\pi_i^k\| \leq \epsilon \frac{|\mathcal{S}|^{1/2}}{|\mathcal{A}|^{1/2}} + \epsilon |\mathcal{S}|^{1/2} \leq 2\epsilon |\mathcal{S}|^{1/2}. \end{aligned}$$

Using Eq. (41) and Lemma 4 in Proposition 1, we obtain the following bound on the optimality gap

$$\begin{aligned} &\max\left\{\frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho)\right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\bar{\pi}^k}(\rho)]_+ + [V_i^{\bar{\pi}^k}(\rho) - u_i]_+\right)\right\} \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - \hat{Q}_i^k\|\right) \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{N}{K\eta} + \frac{N^{3/2}\alpha}{1 - \sigma_2(W)} + \epsilon + \frac{\tau^{1/2}}{\epsilon^{1/2}\beta^{1/2}K^{1/2}} + \frac{N\beta^{1/2}\tau}{\epsilon^{1/2}} + \frac{N\alpha}{\epsilon\beta}\right). \end{aligned}$$

By the bound on consensus error in Lemma 4 and the Lipschitz continuity of the value function in Lemma 2, this implies for any agent $j = 1, \dots, N$

$$\begin{aligned} &\max\left\{\frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\pi_j^k}(\rho)\right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi_j^k}(\rho)]_+ + [V_i^{\pi_j^k}(\rho) - u_i]_+\right)\right\} \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{N}{K\eta} + \frac{N^{3/2}\alpha}{1 - \sigma_2(W)} + \epsilon + \frac{\tau^{1/2}}{\epsilon^{1/2}\beta^{1/2}K^{1/2}} + \frac{N\beta^{1/2}\tau}{\epsilon^{1/2}} + \frac{N\alpha}{\epsilon\beta}\right). \end{aligned}$$

Plugging in the step sizes to the two inequalities above and recognizing from Eq. (18) that

$$\tau^{1/2} \leq \tau \leq D \log(1/\alpha) = D \log\left(\frac{K^{5/6}}{\alpha_0}\right) = \mathcal{O}(\log(K))$$

completes the proof. \square

B.3 Proof of Theorem 3

In the context of linear function approximation, we denote

$$\hat{Q}_i^k = \Phi \omega_i^{k,T}, \quad \hat{V}_i^k(s), \quad (42)$$

and adopt the rest of the notations from Section A.

We note that Algorithm 5 observes the update rule in 38 and 39, with $Q_i^k = \hat{Q}_i^k$ and $\bar{\pi}^{k+1} = f(\{\theta_i^{k+1}\}_i)$ being defined as

$$\bar{\theta}^{k+1} = \frac{1}{N} \sum_{i=1}^N \theta_i^{k+1}, \quad \bar{\pi}^{k+1} = \frac{\exp(\phi(s, a)^\top \bar{\theta}^{k+1}(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\phi(s, a')^\top \bar{\theta}^{k+1}(s, a'))}.$$

As a result, we can apply Proposition 1, which implies

$$\begin{aligned} & \max\left\{\frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho)\right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\bar{\pi}^k}(\rho)]_+ + [V_i^{\bar{\pi}^k}(\rho) - u_i]_+\right)\right\} \\ & \leq \mathcal{O}\left(\frac{N}{K\alpha} + N\eta + \frac{N}{K\eta} + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - \hat{Q}_i^k\|\right). \end{aligned} \quad (43)$$

A straightforward consequence of the assumption is the boundedness of \hat{Q}_i^k , \hat{V}_i^k , $Q_i^{\hat{\pi}_i^k}$, and $\omega_i^*(\hat{\pi}_i^k)$, which we state in the lemma below.

Lemma 5. *Recall the definition of B_ω in Eq. (33) For all $i = 1, \dots, N$ and $k \geq 0$, we have*

$$\|\omega_i^*(\hat{\pi}_i^k)\| \leq B_\omega, \quad \max\{\|\hat{Q}_i^k\|, \|\hat{V}_i^k\|, \|Q_i^{\hat{\pi}_i^k}\|\} \leq Q_{\max},$$

where $Q_{\max} = \sigma_{\max}(\Phi)B_\omega$.

We treat the policy space consensus error in the following lemma.

Lemma 6. *The policy iterates $\{\pi_i^k\}$ generated by Algorithm 5 satisfy*

$$\|\bar{\pi}^k - \pi_i^k\| \leq \mathcal{O}\left(\frac{\sqrt{N}\alpha}{1 - \sigma_2(W)}\right), \quad \text{for all } k = 0, \dots, K-1 \text{ and } i = 1, \dots, N.$$

We will later decompose the error $\frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - \hat{Q}_i^k\|$ and bound a component with the following proposition.

Proposition 3. *Define $z_i^{k,t} = \omega_i^{k,t} - \omega_i^*(\hat{\pi}_i^k)$. Under Assumptions 1-4, we have*

$$\mathbb{E}[\|z_i^{k,T}\|^2] \leq 4B_\omega^2 \left(1 - \frac{2(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{T-\tau} + \frac{4(1+18B_\omega)^2|\mathcal{A}|\beta\tau}{(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon}.$$

By Jensen's inequality, Proposition 3 implies for all i and k

$$\mathbb{E}[\|z_i^{k,T}\|] \leq \sqrt{\mathbb{E}[\|z_i^{k,T}\|^2]}$$

$$\leq 2B_\omega \left(1 - \frac{2(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{\frac{T-\tau}{2}} + \frac{2(1+18B_\omega)\sqrt{|\mathcal{A}|\beta\tau}}{\sqrt{(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon}}. \quad (44)$$

By the definition of \hat{Q}_i^k in Eq. (42), we have

$$\begin{aligned} \|\hat{Q}_i^k - Q_i^{\pi_i^k}\| &\leq \|\Phi(\omega_i^{k,T} - \omega_i^*(\hat{\pi}_i^k))\| + \|\Phi\omega_i^*(\hat{\pi}_i^k) - Q_i^{\hat{\pi}_i^k}\| + \|Q_i^{\hat{\pi}_i^k} - Q_i^{\pi_i^k}\| \\ &\leq \sigma_{\max}(\Phi) \|\omega_i^{k,T} - \omega_i^*(\hat{\pi}_i^k)\| + \varepsilon_{\max} + \frac{\gamma|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2} \|\hat{\pi}_i^k - \pi_i^k\| \\ &\leq \sigma_{\max}(\Phi) \|z_i^{k,T}\| + \varepsilon_{\max} + \frac{2\gamma|\mathcal{S}|^{3/2}|\mathcal{A}|\epsilon}{(1-\gamma)^2}, \end{aligned} \quad (45)$$

where the second inequality employs Assumption 4 and the last inequality follows from

$$\begin{aligned} \|\hat{\pi}_i^k - \pi_i^k\| &= \left\| \frac{\epsilon}{|\mathcal{A}|} \mathbf{1}_{|\mathcal{S}||\mathcal{A}|} + (1-\epsilon)\pi_i^k - \pi_i^k \right\| \leq \epsilon \left\| \frac{1}{|\mathcal{A}|} \mathbf{1}_{|\mathcal{S}||\mathcal{A}|} \right\| + \epsilon \|\pi_i^k\| \\ &\leq \epsilon \frac{|\mathcal{S}|^{1/2}}{|\mathcal{A}|^{1/2}} + \epsilon |\mathcal{S}|^{1/2} \leq 2\epsilon |\mathcal{S}|^{1/2}. \end{aligned}$$

As a result of Eqs. (44) and (45),

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=0}^N \mathbb{E}[\|\hat{Q}_i^k - Q_i^{\pi_i^k}\|] &\leq \frac{\sigma_{\max}(\Phi)}{K} \sum_{k=0}^{K-1} \sum_{i=0}^N \mathbb{E}[\|z_i^{k,T}\|] + N\varepsilon_{\max} + \frac{2\gamma|\mathcal{S}|^{3/2}|\mathcal{A}|N\epsilon}{(1-\gamma)^2} \\ &\leq \mathcal{O}\left(N\left(1 - \frac{2(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{\frac{T-\tau}{2}} + \frac{\sqrt{\beta\tau}}{\sqrt{\epsilon}} + N\varepsilon_{\max} + N\epsilon\right). \end{aligned} \quad (46)$$

Plugging Eq. (46) and the result of Lemma 6 into Eq. (43),

$$\begin{aligned} &\max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\pi^k}(\rho) \right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi^k}(\rho)]_+ + [V_i^{\pi^k}(\rho) - u_i]_+ \right) \right\} \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N(\eta + \epsilon + \varepsilon_{\max}) + \frac{N}{K\eta} + \frac{N^{3/2}\alpha}{1 - \sigma_2(W)} + N\left(1 - \frac{2(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{\frac{T-\tau}{2}} + \frac{\sqrt{\beta\tau}}{\sqrt{\epsilon}} \right). \end{aligned}$$

By the bound on consensus error in Lemma 6 and the Lipschitz continuity of the value function in Lemma 2, this implies for any agent $j = 1, \dots, N$

$$\begin{aligned} &\max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\pi_j^k}(\rho) \right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi_j^k}(\rho)]_+ + [V_i^{\pi_j^k}(\rho) - u_i]_+ \right) \right\} \\ &\leq \mathcal{O}\left(\frac{N}{K\alpha} + N(\eta + \epsilon + \varepsilon_{\max}) + \frac{N}{K\eta} + \frac{N^{3/2}\alpha}{1 - \sigma_2(W)} + N\left(1 - \frac{2(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{\frac{T-\tau}{2}} + \frac{\sqrt{\beta\tau}}{\sqrt{\epsilon}} \right). \end{aligned}$$

In order to have

$$\begin{aligned} &\max \left\{ \frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\pi_j^k}(\rho) \right), \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \left([\ell_i - V_i^{\pi_j^k}(\rho)]_+ + [V_i^{\pi_j^k}(\rho) - u_i]_+ \right) \right\} \\ &\leq \frac{N^{5/4}\delta}{\sqrt{1 - \sigma_2(W)}} + \mathcal{O}(N\varepsilon_{\max}), \end{aligned}$$

we can choose

$$\alpha \sim \mathcal{O}\left(\frac{(\sqrt{1 - \sigma_2(W)})\delta}{N^{1/4}}\right), \quad \beta \sim \mathcal{O}\left(\frac{\delta^3}{\log(1/\delta)}\right), \quad \epsilon \sim \mathcal{O}(\delta), \quad \eta \sim \mathcal{O}(\delta),$$

which implies

$$K \sim \mathcal{O}(\delta^{-2}), \quad T \sim \mathcal{O}\left(\frac{\log(N/\delta)}{\beta\epsilon}\right) = \mathcal{O}\left(\frac{\log(N/\delta)}{\delta^4}\right), \quad TK \sim \mathcal{O}\left(\frac{\log(1/\delta)}{\delta^6}\right).$$

□

C Proof of Propositions

C.1 Proof of Proposition 1

Without loss of generality, we assume that $B \geq \frac{1}{1-\gamma}$, since otherwise we can safely set B to be $\frac{1}{1-\gamma}$ in Proposition 1. This helps us simplify notations later in the analysis.

We define the following notation on the aggregate policy over the network.

$$\boldsymbol{\pi}^k = [\pi_1^k, \dots, \pi_N^k] \quad (47)$$

We have from Eq. (39)

$$\bar{\pi}^{k+1}(a | s) = \bar{\pi}^k(a | s) \frac{\exp(\frac{\alpha}{N} \sum_{i=1}^N (\frac{1}{N} + \lambda_i^k - \nu_i^k) Q_i^k(s, a))}{Z^k(s)}, \quad (48)$$

where $Z^k(s) = \sum_{a' \in \mathcal{A}} \bar{\pi}^k(a' | s) \exp(\frac{\alpha}{N} \sum_{i=1}^N (\frac{1}{N} + \lambda_i^k - \nu_i^k) Q_i^k(s, a'))$.

Let $V_i^k(s) = \sum_a \pi_i^k(a | s) Q_i^k(s, a)$.

Define $Q_g^{\boldsymbol{\pi}^k} = [Q_1^{\boldsymbol{\pi}^k}, \dots, Q_N^{\boldsymbol{\pi}^k}]$ and $V_g^{\boldsymbol{\pi}^k} = [V_1^{\boldsymbol{\pi}^k}, \dots, V_N^{\boldsymbol{\pi}^k}]$. Define $Q_g^k = [Q_1^k, \dots, Q_N^k]$ and $V_g^k = [V_1^k, \dots, V_N^k]$. Define $Q_{L,k}^k = \sum_{i=1}^N (\frac{1}{N} + \lambda_i^k - \nu_i^k) Q_i^k$ and $V_{L,k}^k = \sum_{i=1}^N (\frac{1}{N} + \lambda_i^k - \nu_i^k) V_i^k$.

Objective function convergence. From the dual update in Eq. (21), we have

$$\begin{aligned} 0 \leq \|\lambda^K\|^2 &= \sum_{k=0}^{K-1} (\|\lambda^{k+1}\|^2 - \|\lambda^k\|^2) \\ &= \sum_{k=0}^{K-1} \left(\left\| \Pi_{[0, B\lambda]} \left(\lambda^k - \eta \left(\sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) Q_g^k(s, a) - \ell \right) \right) \right\|^2 - \|\lambda^k\|^2 \right) \\ &\leq \sum_{k=0}^{K-1} \left(\left\| \lambda^k - \eta \left(\sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) Q_g^k(s, a) - \ell \right) \right\|^2 - \|\lambda^k\|^2 \right) \\ &= -2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top \left(\sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) Q_g^{\boldsymbol{\pi}^k}(s, a) - \ell \right) \\ &\quad + 2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top \left(\sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) \left(Q_g^{\boldsymbol{\pi}^k}(s, a) - Q_g^k(s, a) \right) \right) \\ &\quad + \eta^2 \sum_{k=0}^{K-1} \left\| \sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) Q_g^k(s, a) - \ell \right\|^2 \\ &= -2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top \left(V_g^{\boldsymbol{\pi}^k}(\rho) - \ell \right) + \eta^2 \sum_{k=0}^{K-1} \left\| \sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) Q_g^k(s, a) - \ell \right\|^2 \\ &\quad + 2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top \left(\sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) \left(Q_g^{\boldsymbol{\pi}^k}(s, a) - Q_g^k(s, a) \right) \right). \end{aligned} \quad (49)$$

Since the value function and constant ℓ_i are within $[0, \frac{1}{1-\gamma}]$, the second term of Eq. (49) obeys

$$\begin{aligned}
\sum_{k=0}^{K-1} \left\| \sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) Q_g^k(s, a) - \ell \right\|^2 &= \sum_{k=0}^{K-1} \sum_{i=1}^N \left(\sum_{s,a} \rho(s) \pi_i^k(a | s) Q_i^k(s, a) - \ell_i \right)^2 \\
&\leq 2 \sum_{k=0}^{K-1} \sum_{i=1}^N \left(\left(\sum_{s,a} \rho(s) \pi_i^k(a | s) Q_i^k(s, a) \right)^2 + (\ell_i)^2 \right) \\
&\leq 2 \sum_{k=0}^{K-1} \sum_{i=1}^N \left(B^2 + \frac{1}{(1-\gamma)^2} \right) \leq \frac{4KN}{(1-\gamma)^2}.
\end{aligned} \tag{50}$$

The third term of Eq. (49) can be treated as

$$\begin{aligned}
&2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top \left(\sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) (Q_g^{\pi^k}(s, a) - Q_g^k(s, a)) \right) \\
&= 2\eta \sum_{k=0}^{K-1} \sum_{i=1}^N \lambda_i^k \left(\sum_{s,a} \rho(s) \pi_i^k(a | s) (Q_i^{\pi^k}(s, a) - Q_i^k(s, a)) \right) \\
&\leq 2B_\lambda \eta \sum_{k=0}^{K-1} \sum_{i=1}^N \left(\sum_{s,a} \rho(s)^2 \pi_i^k(a | s)^2 \right)^{1/2} \|Q_i^{\pi^k} - Q_i^k\| \\
&\leq 2B_\lambda \eta \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi^k} - Q_i^k\|,
\end{aligned} \tag{51}$$

where the first inequality follows from the Cauchy-Schwarz inequality, and the second inequality follows from the fact that the ℓ_2 norm of a vector is upper bounded by its ℓ_1 norm.

Using Eqs. (50) and (51) in Eq. (49), we get

$$\begin{aligned}
0 &\leq -2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top (V_g^{\pi^k}(\rho) - \ell) + \eta^2 \sum_{k=0}^{K-1} \left\| \sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) Q_g^k(s, a) - b \right\|^2 \\
&\quad + 2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top \left(\sum_{s,a} \rho(s) \text{diag}(\boldsymbol{\pi}^k(a | s)) (Q_g^{\pi^k}(s, a) - Q_g^k(s, a)) \right) \\
&\leq 2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top (V_g^{\pi^*}(\rho) - V_g^{\bar{\pi}^k}(\rho)) + 2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top (V_g^{\bar{\pi}^k}(\rho) - V_g^{\pi^k}(\rho)) \\
&\quad + \frac{4KN\eta^2}{(1-\gamma)^2} + 2B_\lambda \eta \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi^k} - Q_i^k\| \\
&\leq 2\eta \sum_{k=0}^{K-1} (\lambda^k)^\top (V_g^{\pi^*}(\rho) - V_g^{\bar{\pi}^k}(\rho)) + \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda\eta}{(1-\gamma)^2} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
&\quad + \frac{4KN\eta^2}{(1-\gamma)^2} + 2B_\lambda \eta \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi^k} - Q_i^k\|,
\end{aligned}$$

where the second inequality follows from the fact that the optimal policy satisfies the constraints, i.e. $V_i^{\pi^*}(\rho) \geq \ell_i$ for all $i = 1, \dots, N$, and the third inequality is applies Lemma 2.

Re-arranging this inequality and dividing by $2K\eta$ lead to

$$\frac{1}{K} \sum_{k=0}^{K-1} (\lambda^k)^\top (V_g^{\pi^*}(\rho) - V_g^{\bar{\pi}^k}(\rho))$$

$$\geq -\frac{\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| - \frac{2N\eta}{(1-\gamma)^2} - \frac{B_\lambda}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|. \quad (52)$$

A similar analysis on ν^k implies

$$\begin{aligned} & -\frac{1}{K} \sum_{k=0}^{K-1} (\nu^k)^\top \left(V_g^{\pi^*}(\rho) - V_g^{\bar{\pi}^k}(\rho) \right) \\ & \geq -\frac{\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| - \frac{2N\eta}{(1-\gamma)^2} - \frac{B_\lambda}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|. \end{aligned} \quad (53)$$

Lemma 7. *The iterates of Eq. (38) satisfy for all $k = 0, \dots, K-1$*

$$\begin{aligned} V_{L,k}^{\bar{\pi}^{k+1}}(\zeta) - V_{L,k}^{\bar{\pi}^k}(\zeta) & \geq \frac{N}{\alpha} \mathbb{E}_{s \sim \zeta} \left[\log Z_k(s) - \frac{\alpha}{N} V_{L,k}^{\bar{\pi}^k}(s) \right] \\ & \quad - \frac{2(B_\lambda + 1/N)}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| - \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|. \end{aligned}$$

Lemma 8. *The iterates of Eq. (38) satisfy for all $k = 0, \dots, K-1$*

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \left(V_{L,k}^{\pi^*}(\rho) - V_{L,k}^{\bar{\pi}^k}(\rho) \right) & \leq \frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{3(B_\lambda + 1/N)}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\ & \quad + \frac{3\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{4N\eta}{(1-\gamma)^3K}. \end{aligned}$$

Combining Eq. (52), Eq. (53), and Lemma 8,

$$\begin{aligned} & \frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) \right) \\ & \leq \frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{3(B_\lambda + 1/N)}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\ & \quad + \frac{3\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{4N\eta}{(1-\gamma)^3K} \\ & \quad + \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{4N\eta}{(1-\gamma)^2} + \frac{2B_\lambda}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\ & \leq \frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{8N\eta}{(1-\gamma)^3} + \frac{5\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\ & \quad + \frac{8(B_\lambda + 1/N)}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|. \end{aligned}$$

Constraint violation convergence. For any $\lambda \in [0, B_\lambda]^N$, since the projection operator $\Pi_{[0, B_\lambda]}$ is non-expansive, we have

$$\begin{aligned} \|\lambda^{k+1} - \lambda\|^2 & = \|\Pi_{[0, B_\lambda]}(\lambda^k - \eta(V_g^k(\rho) - \ell)) - \lambda\|^2 \\ & \leq \|\lambda^k - \eta(V_g^k(\rho) - \ell) - \lambda\|^2 \\ & = \|\lambda^k - \lambda\|^2 - 2\eta(\lambda^k - \lambda)^\top (V_g^k(\rho) - \ell) + \eta^2 \sum_{i=1}^N \left\| \sum_{s,a} \rho(s) \pi_i^k(a | s) Q_i^k(\rho) - \ell_i \right\|^2 \end{aligned}$$

$$\leq \|\lambda^k - \lambda\|^2 - 2\eta(\lambda^k - \lambda)^\top (V_g^k(\rho) - \ell) + \frac{4N\eta^2}{(1-\gamma)^2},$$

where the last inequality bounds the quadratic term using an approach similar to Eq. (50).

Re-arranging the terms and summing up from $k = 0$ to $k = K - 1$, we get

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} (\lambda^k - \lambda)^\top (V_g^k(\rho) - b) &\leq \frac{1}{K} (\|\lambda^0 - \lambda\|^2 - \|\lambda^K - \lambda\|^2) + \frac{2N\eta}{(1-\gamma)^2} \\ &\leq \frac{1}{2K\eta} \|\lambda^0 - \lambda\|^2 + \frac{2N\eta}{(1-\gamma)^2}, \end{aligned}$$

which implies

$$\begin{aligned} &\frac{1}{K} \sum_{k=0}^{K-1} (\lambda^k - \lambda)^\top (V_g^{\bar{\pi}^k}(\rho) - \ell) \\ &= \frac{1}{K} \sum_{k=0}^{K-1} (\lambda^k - \lambda)^\top (V_g^k(\rho) - \ell) + \frac{1}{K} \sum_{k=0}^{K-1} (\lambda^k - \lambda)^\top (V_g^{\bar{\pi}^k}(\rho) - V_g^{\pi^k}(\rho)) \\ &\quad + \frac{1}{K} \sum_{k=0}^{K-1} (\lambda^k - \lambda)^\top (V_g^{\pi^k}(\rho) - V_g^k(\rho)) \\ &\leq \frac{1}{2K\eta} \|\lambda^0 - \lambda\|^2 + \frac{2N\eta}{(1-\gamma)^2} + \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{2B_\lambda}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^k - Q_i^{\pi_i^k}\|. \end{aligned}$$

Similarly, we can show for any $\nu \in [0, B_\lambda]^N$

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} (\lambda^k - \lambda)^\top (u - V_g^{\bar{\pi}^k}(\rho)) &\leq \frac{1}{2K\eta} \|\nu^0 - \nu\|^2 + \frac{2N\eta}{(1-\gamma)^2} + \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\ &\quad + \frac{2B_\lambda}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^k - Q_i^{\pi_i^k}\|. \end{aligned}$$

Since λ^k, ν^k are non-negative, we have from Lemma 8 and the two inequalities above

$$\begin{aligned} &\frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) + \lambda^\top (\ell - V_g^{\bar{\pi}^k}(\rho)) + \nu^\top (V_g^{\bar{\pi}^k}(\rho) - u) \right) \\ &\leq \frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) + (\lambda^k)^\top (V_g^{\pi^*}(\rho) - \ell) + \lambda^\top (\ell - V_g^{\bar{\pi}^k}(\rho)) + (\nu^k)^\top (u - V_g^{\pi^*}(\rho)) + \nu^\top (V_g^{\bar{\pi}^k}(\rho) - u) \right) \\ &= \frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) + (\lambda^k - \nu^k)^\top (V_g^{\pi^*}(\rho) - V_g^{\bar{\pi}^k}(\rho)) \right. \\ &\quad \left. + (\lambda^k - \lambda)^\top (V_g^{\bar{\pi}^k}(\rho) - \ell) + (\nu^k - \nu)^\top (u - V_g^{\bar{\pi}^k}(\rho)) \right) \\ &\leq \frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{3(B_\lambda + 1/N)}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\ &\quad + \frac{3\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{4N\eta}{(1-\gamma)^3K} \\ &\quad + \frac{1}{2K\eta} \|\lambda^0 - \lambda\|^2 + \frac{2N\eta}{(1-\gamma)^2} + \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{2B_\lambda}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^k - Q_i^{\pi_i^k}\| \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2K\eta} \|\nu^0 - \nu\|^2 + \frac{2N\eta}{(1-\gamma)^2} + \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}B_\lambda}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{2B_\lambda}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^k - Q_i^{\pi_i^k}\| \\
& \leq \frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{8N\eta}{(1-\gamma)^3} + \frac{\|\lambda^0 - \lambda\|^2 + \|\nu^0 - \nu\|^2}{2K\eta} \\
& + \frac{7\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{7(B_\lambda + 1/N)}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|.
\end{aligned} \tag{54}$$

Now, choosing λ and ν such that

$$\lambda_i = \begin{cases} B_\lambda, & \text{if } \ell_i - V_i^{\pi^k}(\rho) \geq 0 \\ 0, & \text{else} \end{cases} \quad \nu_i = \begin{cases} B_\lambda, & \text{if } V_i^{\pi^k}(\rho) - u_i \geq 0 \\ 0, & \text{else} \end{cases}$$

Then, Eq. (54) leads to

$$\begin{aligned}
& \frac{1}{K} \sum_{k=0}^{K-1} \left(V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) \right) + \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^N B_\lambda \left([\ell_i - V_i^{\bar{\pi}^k}(\rho)]_+ + [V_i^{\bar{\pi}^k}(\rho) - u_i]_+ \right) \\
& \leq \frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{8N\eta}{(1-\gamma)^3} + \frac{NB_\lambda^2}{K\eta} \\
& + \frac{7\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{7(B_\lambda + 1/N)}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|.
\end{aligned} \tag{55}$$

Note that there always exists a policy $\tilde{\pi}^K$ such that $d_{\rho}^{\tilde{\pi}^K} = \frac{1}{K} \sum_{k=0}^{K-1} d_{\rho}^{\bar{\pi}^k}$, which implies

$$V_i^{\tilde{\pi}^K} = \frac{1}{K} \sum_{k=0}^{K-1} V_i^{\bar{\pi}^k} \quad \forall i = 0, 1, \dots, N.$$

As a result, Eq. (55) becomes

$$\begin{aligned}
& \left(V_0^{\pi^*}(\rho) - V_0^{\tilde{\pi}^K}(\rho) \right) + B_\lambda \sum_{i=1}^N \left([\ell_i - V_i^{\tilde{\pi}^K}(\rho)]_+ + [V_i^{\tilde{\pi}^K}(\rho) - u_i]_+ \right) \\
& \leq \frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{8N\eta}{(1-\gamma)^3} + \frac{NB_\lambda^2}{K\eta} \\
& + \frac{7\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{7(B_\lambda + 1/N)}{(1-\gamma)^2K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|.
\end{aligned} \tag{56}$$

Lemma 9 (Theorem 6 of Ding et al. (2020)). *Suppose that Assumption 1 holds. Let the constant C obey $C \geq 2\|\lambda^*\|_\infty$ and $C \geq 2\|\nu^*\|_\infty$. Then, given a policy π , if there exists a constant $\delta > 0$ such that*

$$V_0^{\pi^*}(\rho) - V_0^{\pi}(\rho) + C \sum_{i=1}^N ([\ell_i - V_i^{\pi}(\rho)]_+ + [V_i^{\pi}(\rho) - u_i]_+) \leq \delta,$$

then we have

$$\sum_{i=1}^N ([\ell_i - V_i^{\pi}(\rho)]_+ + [V_i^{\pi}(\rho) - u_i]_+) \leq \frac{2\delta}{C}.$$

Recall that Lemma 1 states that $2\|\lambda^*\|_\infty \leq B_\lambda$ and $2\|\nu^*\|_\infty \leq B_\lambda$. Applying Lemma 9 with $C = B_\lambda$ and δ being the terms on the left hand side of Eq. (56), we have

$$\sum_{i=1}^N \left([\ell_i - V_i^{\tilde{\pi}^K}(\rho)]_+ + [V_i^{\tilde{\pi}^K}(\rho) - u_i]_+ \right) \leq \frac{2}{B_\lambda} \left(\frac{N \log |\mathcal{A}|}{(1-\gamma)K\alpha} + \frac{2NB_\lambda}{(1-\gamma)^2K} + \frac{8N\eta}{(1-\gamma)^3} + \frac{NB_\lambda^2}{K\eta} \right)$$

$$+ \frac{7\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4 K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| + \frac{7(B_\lambda + 1/N)}{(1-\gamma)^2 K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|.$$

□

C.2 Proof of Proposition 2

We omit the proof and note that this result is adapted from Zeng et al. (2022)[Proposition 2].

C.3 Proof of Proposition 3

This section presents the proof of Proposition 2. We use $O_i^{k,t}$ to denote the data observation used for variable updates in iteration (t, k) , i.e.

$$O_i^{k,t} = (s_i^{k,t}, a_i^{k,t}, s_i^{k,t+1}, a_i^{k,t+1}).$$

We use $R_i : |\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}| \times |\mathcal{A}| \rightarrow \mathbb{R}^d$ to denote the feature-reward composite operator such that

$$R_i(s, a, s', a') = r_i(s, a)\phi(s, a). \quad (57)$$

We also define the operator $H : |\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}| \times |\mathcal{A}| \rightarrow \mathbb{R}^{d \times d}$ such that

$$H(s, a, s', a') = \phi(s, a)(\gamma\phi(s', a') - \phi(s, a))^\top, \quad (58)$$

which means that \bar{H}^π defined in Eq. (30) satisfies

$$\bar{H}^\pi = \mathbb{E}_{s \sim \mu_\pi, a \sim \pi(\cdot|s), s' \sim P(\cdot|s, a), a' \sim \pi(\cdot|s')} [H(s, a, s', a')].$$

Finally, we denote

$$\begin{aligned} \Gamma_i(\pi, z, O) &\triangleq z^\top (R_i(O) + H(O)\omega_i^*(\pi)) + z^\top (H(O) - \bar{H}^\pi)z, \\ p_i^{k,t} &\triangleq \omega_i^{k,t+1} - \hat{\omega}_i^{k,t+1} \end{aligned} \quad (59)$$

We introduce a few technical lemmas to support the analysis.

Lemma 10. *For all $k \geq 0$, we have*

$$\begin{aligned} \|p_i^{k,t}\| &\leq 2(1 + 2B_\omega)\beta, \\ \|z_i^{k,t+1} - z_i^{k,t}\| &= \|\omega_i^{k,t+1} - \omega_i^{k,t}\| \leq (1 + 2B_\omega)\beta. \end{aligned}$$

Lemma 11. *The matrix $\bar{H}^{\hat{\pi}_i^k}$ is negative definite*

$$\omega^\top \bar{H}^{\hat{\pi}_i^k} \omega \leq -\frac{(1-\gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon}{|\mathcal{A}|} \|\omega\|^2, \quad \forall \omega \in \mathbb{R}^d, k \geq 0,$$

where $\underline{\mu} = \min_{\pi, s} \mu_\pi(s)$ is a positive constant due to the uniform ergodicity of the Markov chain under any policy.

Lemma 12. *Under Assumption 2, we have for all $t \geq \tau$ and $i = 0, \dots, N$*

$$\mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t}, O_i^{k,t})] \leq 2(1 + 18B_\omega)^2 \beta \tau.$$

With the definition of R_i and H in Eqs. (57) and (58), the critic update in Eq. (34) can be re-expressed as

$$\omega_i^{k,t+1} = \omega_i^{k,t} + \beta(R_i(O_i^{k,t}) + H(O_i^{k,t})\omega_i^{k,t}) + p_i^{k,t},$$

which implies

$$\begin{aligned}
z_i^{k,t+1} - z_i^{k,t} &= \omega_i^{k,t+1} - \omega_i^{k,t} \\
&= \beta(R_i(O_i^{k,t}) + H(O_i^{k,t})\omega_i^{k,t}) + p_i^{k,t} \\
&= \beta(R_i(O_i^{k,t}) + H(O_i^{k,t})\omega_i^*(\hat{\pi}_i^k) + H(O_i^{k,t})z_i^{k,t}) + p_i^{k,t}.
\end{aligned} \tag{60}$$

Then, straightforward manipulations yield

$$\begin{aligned}
\|z_i^{k,t+1}\|^2 - \|z_i^{k,t}\|^2 &= 2(z_i^{k,t})^\top(z_i^{k,t+1} - z_i^{k,t}) + \|z_i^{k,t+1} - z_i^{k,t}\|^2 \\
&= 2(z_i^{k,t})^\top(z_i^{k,t+1} - z_i^{k,t} - \beta\bar{H}\hat{\pi}_i^k z_i^{k,t}) + \|z_i^{k,t+1} - z_i^{k,t}\|^2 + 2\beta(z_i^{k,t})^\top\bar{H}\hat{\pi}_i^k z_i^{k,t} \\
&= 2\beta(z_i^{k,t})^\top(R_i(O_i^{k,t}) + H(O_i^{k,t})(\omega_i^*(\hat{\pi}_i^k) + z_i^{k,t}) - \bar{H}\hat{\pi}_i^k z_i^{k,t}) \\
&\quad + 2(z_i^{k,t})^\top p_i^{k,t} + \|z_i^{k,t+1} - z_i^{k,t}\|^2 + 2\beta(z_i^{k,t})^\top\bar{H}\hat{\pi}_i^k z_i^{k,t} \\
&\leq 2\beta\Gamma_i(\hat{\pi}_i^k, z_i^{k,t}, O_i^{k,t}) + 2(z_i^{k,t})^\top p_i^{k,t} \\
&\quad + \|z_i^{k,t+1} - z_i^{k,t}\|^2 + 2\beta(z_i^{k,t})^\top\bar{H}\hat{\pi}_i^k z_i^{k,t},
\end{aligned} \tag{61}$$

where the third equality applies Eq. (60).

To bound $(z_i^{k,t})^\top p_i^{k,t}$,

$$\begin{aligned}
(z_i^{k,t})^\top p_i^{k,t} &= \langle \omega_i^{k,t} - \omega_i^*(\hat{\pi}_i^k), p_i^{k,t} \rangle \\
&\leq \langle \omega_i^{k,t+1} - \omega_i^*(\hat{\pi}_i^k), \omega_i^{k,t+1} - \hat{\omega}_i^{k,t+1} \rangle \\
&\quad + \|z_i^{k,t+1} - z_i^{k,t}\| \|p_i^{k,t}\| \\
&= \langle \hat{\omega}_i^{k,t+1} - \omega_i^*(\hat{\pi}_i^k), \omega_i^{k,t+1} - \hat{\omega}_i^{k,t+1} \rangle \\
&\quad + \|\omega_i^{k,t+1} - \hat{\omega}_i^{k,t+1}\|^2 + \|z_i^{k,t+1} - z_i^{k,t}\| \|p_i^{k,t}\| \\
&\leq 2(1 + 2B_\omega)^2 \beta^2,
\end{aligned} \tag{62}$$

where the second inequality is due to Doan et al. (2019)[Lemma 3(a)].

Taking the expectation in Eq. (61) and applying Eq. (62) and Lemmas 11 and 12, we have

$$\begin{aligned}
\mathbb{E}[\|z_i^{k,t+1}\|^2 - \|z_i^{k,t}\|^2] &\leq 2\beta\mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t}, O_i^{k,t})] + 2\mathbb{E}[(z_i^{k,t})^\top p_i^{k,t}] \\
&\quad + \mathbb{E}[\|z_i^{k,t+1} - z_i^{k,t}\|^2] + 2\beta\mathbb{E}[(z_i^{k,t})^\top\bar{H}\hat{\pi}_i^k z_i^{k,t}] \\
&\leq 4(1 + 18B_\omega)^2 \beta^2 \tau + 2(1 + 2B_\omega)^2 \beta^2 + 2(1 + 2B_\omega)^2 \beta^2 - \frac{2(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|} \mathbb{E}[\|z_i^{k,t}\|^2].
\end{aligned}$$

Re-arranging the terms,

$$\mathbb{E}[\|z_i^{k,t+1}\|^2] \leq \left(1 - \frac{2(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right) \mathbb{E}[\|z_i^{k,t}\|^2] + 8(1 + 18B_\omega)^2 \beta^2 \tau. \tag{63}$$

Recursively applying Eq. (63), we get

$$\begin{aligned}
&\mathbb{E}[\|z_i^{k,T}\|^2] \\
&\leq \left(1 - \frac{2(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{T-\tau} \mathbb{E}[\|z_i^{k,\tau}\|^2] + \sum_{t=0}^{T-\tau-1} 8(1 + 18B_\omega)^2 \beta^2 \tau \left(1 - \frac{2(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^t \\
&\leq \left(1 - \frac{2(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{T-\tau} \mathbb{E}[\|z_i^{k,\tau}\|^2] + 8(1 + 18B_\omega)^2 \beta^2 \tau \sum_{t=0}^{\infty} \left(1 - \frac{2(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^t \\
&= 4B_\omega^2 \left(1 - \frac{2(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon\beta}{|\mathcal{A}|}\right)^{T-\tau} + \frac{4(1 + 18B_\omega)^2 |\mathcal{A}| \beta \tau}{(1 - \gamma)\underline{\mu}\sigma_{\min}(\Phi)\epsilon}.
\end{aligned}$$

□

D Proof of Technical Lemmas

D.1 Proof of Lemma 1

We skip the proof as it is a simple extension of Ding et al. (2020)[Lemma 1].

D.2 Proof of Lemma 3

We denote $g_i^k = (\frac{1}{N} + \lambda_i^k - \nu_i^k)Q_i^{\pi_i^k} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $g^k = [(g_1^k)^\top, \dots, (g_N^k)^\top]^\top \in \mathbb{R}^{N|\mathcal{S}||\mathcal{A}|}$. It is easy to see

$$\|g_i^k\| \leq \left| \frac{1}{N} + \lambda_i^k - \nu_i^k \right| \|Q_i^{\pi_i^k}\| \leq \frac{(B_\lambda + \frac{1}{N})\sqrt{|\mathcal{S}||\mathcal{A}|}}{1-\gamma},$$

which implies $\|g^k\| \leq \frac{(B_\lambda + \frac{1}{N})\sqrt{N|\mathcal{S}||\mathcal{A}|}}{1-\gamma}$ for all k . Then, using an argument similar to the one in Yuan et al. (2016)[Lemma 1], we can get

$$\|\bar{\theta}^k - \theta_i^k\| \leq \frac{(B_\lambda + \frac{1}{N})\sqrt{N|\mathcal{S}||\mathcal{A}|}\alpha}{(1-\gamma)(1-\sigma_2(W))}. \quad (64)$$

The softmax function is Lipschitz with constant 1, i.e.

$$\|\pi_\theta - \pi_{\theta'}\| \leq \|\theta - \theta'\|, \quad \forall \theta, \theta',$$

Recall the definition of $\bar{\pi}^k$ in Eq. (47). The Lipschitz continuity and Eq. (64) imply the claimed result. \square

D.3 Proof of Lemma 4

We skip the proof and note that it is almost identical to the proof of Lemma 3.

D.4 Proof of Lemma 5

Due to the boundedness of the reward function, it is easy to see that $|Q_i^\pi(s, a)| \leq \frac{1}{1-\gamma}$, which implies

$$\|Q_i^\pi\| \leq \sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{1-\gamma}}.$$

Since $\|\Phi\omega_i^*(\hat{\pi}_i^k) - Q_i^{\hat{\pi}_i^k}\| \leq \varepsilon_{\max}$ due to Assumption 4, we have

$$\|\Phi\omega_i^*(\hat{\pi}_i^k)\| \leq \|Q_i^{\hat{\pi}_i^k}\| + \varepsilon_{\max} \leq \sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{1-\gamma}} + \varepsilon_{\max}, \quad (65)$$

which implies

$$\|\omega_i^*(\hat{\pi}_i^k)\| \leq \sigma_{\min}^{-1}(\Phi) \left(\sqrt{\frac{|\mathcal{S}||\mathcal{A}|}{1-\gamma}} + \varepsilon_{\max} \right) = B_\omega. \quad (66)$$

To show the bound on \hat{Q}_i^k , note that $\|\omega_i^{k,T}\| \leq B_\omega$ due to the projection in Eq. (34). As a result,

$$\|\hat{Q}_i^k\| = \|\Phi\omega_i^{k,T}\| \leq \sigma_{\max}(\Phi)B_\omega.$$

The bound on $\|\hat{V}_{i,t}\|$ easily follows from Jensen's inequality

$$\|\hat{V}_i^k\|^2 = \sum_s \left(\sum_a \pi_i^k(a | s) \hat{Q}_i^k(s, a) \right)^2 \leq \sum_s \sum_a \pi_i^k(a | s) \left(\hat{Q}_i^k(s, a) \right)^2 \leq \sum_{s,a} \left(\hat{Q}_i^k(s, a) \right)^2 = \|\hat{Q}_i^k\|^2.$$

\square

D.5 Proof of Lemma 6

We denote

$$g_i^k = \left(\frac{1}{N} + \lambda_i^k - \nu_i^k\right) \omega_i^{k,T} \in \mathbb{R}^d, \quad \text{and} \quad g^k = [(g_1^k)^\top, \dots, (g_N^k)^\top]^\top \in \mathbb{R}^{Nd}.$$

Then,

$$\|g_i^k\| \leq \left| \frac{1}{N} + \lambda_i^k - \nu_i^k \right| \|\omega_i^{k,T}\| \leq (B_\lambda + \frac{1}{N}) B_\omega,$$

which implies $\|g^k\| \leq (B_\lambda + \frac{1}{N}) B_\omega \sqrt{N}$ for all k . Then, using an argument similar to the one in Yuan et al. (2016)[Lemma 1], we can get

$$\|\bar{\theta}^k - \theta_i^k\| \leq \frac{(B_\lambda + \frac{1}{N}) B_\omega \sqrt{N} \alpha}{1 - \sigma_2(W)}. \quad (67)$$

The softmax function is Lipschitz continuous with constant 1, which implies

$$\|\bar{\pi}^k - \pi_i^k\| \leq \|\Phi(\bar{\theta}^k - \theta_i^k)\| \leq \sigma_{\max}(\Phi) \|\bar{\theta}^k - \theta_i^k\| \leq \mathcal{O}\left(\frac{\sqrt{N} \alpha}{1 - \sigma_2(W)}\right).$$

□

D.6 Proof of Lemma 7

The performance difference lemma states that for any policies π_1, π_2 , initial distribution ζ , and $i = 0, \dots, N$

$$V_i^{\pi_1}(\zeta) - V_i^{\pi_2}(\zeta) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\pi^*}, a \sim \pi^*(\cdot|s)} [A_0^{\pi_k}(s, a)]. \quad (68)$$

By this lemma,

$$\begin{aligned} & V_0^{\bar{\pi}^{k+1}}(\zeta) - V_0^{\bar{\pi}^k}(\zeta) \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} [A_0^{\bar{\pi}^k}(s, a)] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} [Q_0^{\bar{\pi}^k}(s, a)] - \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}} [V_0^{\bar{\pi}^k}(s)] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} [Q_{L,k}^k(s, a)] + \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} [Q_{L,k}^{\pi^k}(s, a) - Q_{L,k}^k(s, a)] \\ &\quad + \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} [Q_{L,k}^{\bar{\pi}^k}(s, a) - Q_{L,k}^{\pi^k}(s, a)] \\ &\quad - \frac{(\lambda^k - \nu^k)^\top}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} [Q_g^{\bar{\pi}^k}(s, a)] - \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}} [V_0^{\bar{\pi}^k}(s)]. \end{aligned}$$

Note that the actor update rule in Eq. (48) implies

$$Q_{L,k}^k(s, a) = \frac{N}{\alpha} \log \left(\frac{\bar{\pi}^{k+1}(a | s)}{\bar{\pi}^k(a | s)} Z_k(s) \right).$$

Combining the two equalities above, we have

$$\begin{aligned} & V_0^{\bar{\pi}^{k+1}}(\zeta) - V_0^{\bar{\pi}^k}(\zeta) \\ &= \frac{N}{\alpha(1 - \gamma)} \mathbb{E}_{s \sim d_\zeta^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} \left[\log \left(\frac{\bar{\pi}^{k+1}(a | s)}{\bar{\pi}^k(a | s)} Z_k(s) \right) \right] \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} \left[Q_{L,k}^{\bar{\pi}^k}(s, a) - Q_{L,k}^k(s, a) \right] \\
& + \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^k}, a \sim \bar{\pi}^k(\cdot|s)} \left[Q_{L,k}^{\bar{\pi}^k}(s, a) - Q_{L,k}^{\pi^k}(s, a) \right] \\
& - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} \left[Q_g^{\bar{\pi}^k}(s, a) \right] - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} \left[V_0^{\bar{\pi}^k}(s) \right] \\
& \geq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} \left[D_{KL}(\bar{\pi}^{k+1}(\cdot|s) \parallel \bar{\pi}^k(\cdot|s)) \right] + \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} [\log Z_k(s)] \\
& \quad - \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi^k} - Q_i^k\| - \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
& \quad - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}, a \sim \bar{\pi}^{k+1}(\cdot|s)} \left[A_g^{\bar{\pi}^k}(s, a) \right] - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} \left[V_{L,k}^{\bar{\pi}^k}(s) \right] \\
& \geq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} [\log Z_k(s)] - \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi^k} - Q_i^k\| - \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
& \quad - (\lambda^k - \nu^k)^\top \left(V_g^{\bar{\pi}^{k+1}}(\zeta) - V_g^{\bar{\pi}^k}(\zeta) \right) - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} \left[V_{L,k}^{\bar{\pi}^k}(s) \right],
\end{aligned}$$

where the last inequality applies the performance difference lemma. Rearranging this inequality leads to

$$\begin{aligned}
V_{L,k}^{\bar{\pi}^{k+1}}(\zeta) - V_{L,k}^{\bar{\pi}^k}(\zeta) & \geq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} [\log Z_k(s)] - \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
& \quad - \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi^k} - Q_i^k\| - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} \left[V_{L,k}^{\bar{\pi}^k}(s) \right]. \tag{69}
\end{aligned}$$

From the definition of Z^k and Jensen's inequality,

$$\begin{aligned}
\log Z^k(s) & = \log \left(\sum_{a' \in \mathcal{A}} \bar{\pi}^k(a' | s) \exp \left(\frac{\alpha}{N} \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) Q_i^k(s, a') \right) \right) \\
& \geq \sum_{a' \in \mathcal{A}} \bar{\pi}^k(a' | s) \log \left(\exp \left(\frac{\alpha}{N} \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) Q_i^k(s, a') \right) \right) \\
& = \frac{\alpha}{N} \sum_{a' \in \mathcal{A}} \bar{\pi}^k(a' | s) \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) Q_i^k(s, a') \\
& = \frac{\alpha}{N} \sum_{a' \in \mathcal{A}} \bar{\pi}^k(a' | s) \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) Q_i^{\bar{\pi}^k}(s, a') \\
& \quad + \frac{\alpha}{N} \sum_{a' \in \mathcal{A}} \bar{\pi}^k(a' | s) \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) (Q_i^k(s, a') - Q_i^{\bar{\pi}^k}(s, a')) \\
& \quad + \frac{\alpha}{N} \sum_{a' \in \mathcal{A}} \bar{\pi}^k(a' | s) \sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) (Q_i^{\pi^k}(s, a') - Q_i^{\bar{\pi}^k}(s, a')) \\
& \geq \frac{\alpha}{N} V_{L,k}^{\bar{\pi}^k}(s) - \frac{(B_\lambda + 1/N)\alpha}{N} \sum_{i=1}^N \|Q_i^{\pi^k} - Q_i^k\| - \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)\alpha}{N(1-\gamma)^2} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|.
\end{aligned}$$

This bound on $\log Z_k(s)$ implies

$$\frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} [\log Z_k(s)] - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\zeta}^{\bar{\pi}^{k+1}}} \left[V_{L,k}^{\bar{\pi}^k}(s) \right]$$

$$\begin{aligned}
&= \frac{N}{\alpha(1-\gamma)} \sum_s d_\zeta^{\bar{\pi}^{k+1}}(s) \left(\log Z_k(s) - \frac{\alpha}{N} V_{L,k}^{\bar{\pi}^k}(s) \right) \\
&= \frac{N}{\alpha(1-\gamma)} \sum_s d_\zeta^{\bar{\pi}^{k+1}}(s) \left(\log Z_k(s) - \frac{\alpha}{N} V_{L,k}^{\bar{\pi}^k}(s) + \frac{(B_\lambda + 1/N)\alpha}{N} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \right. \\
&\quad \left. + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)\alpha}{N(1-\gamma)^2} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \right) \\
&\quad - \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| - \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
&\geq \frac{N}{\alpha} \sum_s \zeta(s) \left(\log Z_k(s) - \frac{\alpha}{N} V_{L,k}^{\bar{\pi}^k}(s) \right) - \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\
&\quad - \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|,
\end{aligned}$$

where the inequality follows from the fact that $d_\zeta^\pi \geq (1-\gamma)\zeta$ elementwise for any policy π . Plugging this bound into Eq. (69), we have

$$\begin{aligned}
V_{L,k}^{\bar{\pi}^{k+1}}(\zeta) - V_{L,k}^{\bar{\pi}^k}(\zeta) &\geq \frac{N}{\alpha} \mathbb{E}_{s \sim \zeta} \left[\log Z_k(s) - \frac{\alpha}{N} V_{L,k}^{\bar{\pi}^k}(s) \right] - \frac{2(B_\lambda + 1/N)}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\
&\quad - \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|.
\end{aligned}$$

□

D.7 Proof of Lemma 8

By the performance difference lemma in Eq. (68),

$$\begin{aligned}
V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) &= \frac{1}{N} \sum_{i=1}^N (V_i^{\pi^*}(\rho) - V_i^{\bar{\pi}^k}(\rho)) \\
&= \frac{1}{N} \sum_{i=1}^N (V_i^{\pi^*}(\rho) - V_i^{\pi_i^k}(\rho)) + \frac{1}{N} \sum_{i=1}^N (V_i^{\pi_i^k}(\rho) - V_i^{\bar{\pi}^k}(\rho)) \\
&\leq \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho^*}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[A_i^{\pi_i^k}(s, a) \right] + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
&= \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho^*}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[Q_i^{\pi_i^k}(s, a) \right] - \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho^*}^{\pi^*}} \left[V_i^{\pi_i^k}(s) \right] \\
&\quad + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|.
\end{aligned}$$

Plugging in the update rule of the policy,

$$\begin{aligned}
V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) &\leq \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho^*}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[Q_i^{\pi_i^k}(s, a) \right] - \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho^*}^{\pi^*}} \left[V_i^{\pi_i^k}(s) \right] \\
&\quad + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[\sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) Q_i^{\pi_i^k}(s, a) \right] - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[Q_g^{\pi^k}(s, a) \right] \\
&\quad - \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[V_i^{\pi_i^k}(s) \right] + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
&= \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[\log \left(\frac{\bar{\pi}^{k+1}(a|s)}{\bar{\pi}^k(a|s)} Z_k(s) \right) \right] \\
&\quad + \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[\sum_{i=1}^N \left(\frac{1}{N} + \lambda_i^k - \nu_i^k \right) (Q_i^{\pi_i^k}(s, a) - Q_i^k(s, a)) \right] \\
&\quad - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[A_g^{\pi^k}(s, a) \right] - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[V_g^{\pi^k}(s) \right] \\
&\quad - \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[V_i^{\pi_i^k}(s) \right] + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
&\leq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[D_{\text{KL}}(\pi^*(\cdot|s) \|\bar{\pi}^k(\cdot|s)) \right] - \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[D_{\text{KL}}(\pi^*(\cdot|s) \|\bar{\pi}^{k+1}(\cdot|s)) \right] \\
&\quad + \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[\log Z^k(s) \right] - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[A_g^{\pi^k}(s, a) \right] \\
&\quad - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[V_g^{\pi^k}(s) \right] - \frac{1}{N(1-\gamma)} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[V_i^{\pi_i^k}(s) \right] \\
&\quad + \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^k - Q_i^{\pi_i^k}\| + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|.
\end{aligned}$$

Re-grouping the terms,

$$\begin{aligned}
&V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) \\
&\leq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[D_{\text{KL}}(\pi^*(\cdot|s) \|\bar{\pi}^k(\cdot|s)) \right] - \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[D_{\text{KL}}(\pi^*(\cdot|s) \|\bar{\pi}^{k+1}(\cdot|s)) \right] \\
&\quad + \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[\log Z^k(s) \right] - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[A_g^{\pi^k}(s, a) \right] \\
&\quad + \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}, a \sim \pi^*(\cdot|s)} \left[A_g^{\bar{\pi}^k}(s, a) - A_g^{\pi^k}(s, a) \right] - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[V_{L,k}^{\pi^k}(s) \right] \\
&\quad + \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^k - Q_i^{\pi_i^k}\| + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| \\
&\leq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[D_{\text{KL}}(\pi^*(\cdot|s) \|\bar{\pi}^k(\cdot|s)) \right] - \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[D_{\text{KL}}(\pi^*(\cdot|s) \|\bar{\pi}^{k+1}(\cdot|s)) \right] \quad (70) \\
&\quad + \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[\log Z^k(s) \right] - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \left(V_g^{\pi^*}(s) - V_g^{\bar{\pi}^k}(s) \right) - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[V_{L,k}^{\pi^k}(s) \right] \\
&\quad + \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^k - Q_i^{\pi_i^k}\| + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|,
\end{aligned}$$

where the second inequality follows from the performance difference lemma and the Lipschitz continuity of the advantage.

Applying Lemma 7 with $\zeta = d_{\rho}^{\pi^*}$,

$$\frac{N}{\alpha} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} \left[\log Z_k(s) - \frac{\alpha}{N} V_{L,k}^{\bar{\pi}^k}(s) \right] \leq V_{L,k}^{\bar{\pi}^{k+1}}(d_{\rho}^{\pi^*}) - V_{L,k}^{\bar{\pi}^k}(d_{\rho}^{\pi^*}) + \frac{2(B_\lambda + 1/N)}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\|$$

$$+ \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|. \quad (71)$$

Combining Eqs. (70) and (71),

$$\begin{aligned} & V_0^{\pi^*}(\rho) - V_0^{\bar{\pi}^k}(\rho) \\ & \leq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_\rho^{\pi^*}} [D_{\text{KL}}(\pi^*(\cdot | s) || \bar{\pi}^k(\cdot | s))] - \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_\rho^{\pi^*}} [D_{\text{KL}}(\pi^*(\cdot | s) || \bar{\pi}^{k+1}(\cdot | s))] \\ & \quad + \frac{1}{1-\gamma} \left(V_{L,k}^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) - V_{L,k}^{\bar{\pi}^k}(d_\rho^{\pi^*}) \right) + \frac{2(B_\lambda + 1/N)}{(1-\gamma)^2} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\ & \quad + \frac{2\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\| - \frac{(\lambda^k - \nu^k)^\top}{1-\gamma} \left(V_g^{\pi^*}(s) - V_g^{\bar{\pi}^k}(s) \right) \\ & \quad + \frac{B_\lambda + 1/N}{1-\gamma} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| + \frac{\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{N(1-\gamma)^3} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|, \end{aligned}$$

which implies

$$\begin{aligned} & V_{L,k}^{\pi^*}(\rho) - V_{L,k}^{\bar{\pi}^k}(\rho) \leq \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_\rho^{\pi^*}} [D_{\text{KL}}(\pi^*(\cdot | s) || \bar{\pi}^k(\cdot | s))] \\ & \quad - \frac{N}{\alpha(1-\gamma)} \mathbb{E}_{s \sim d_\rho^{\pi^*}} [D_{\text{KL}}(\pi^*(\cdot | s) || \bar{\pi}^{k+1}(\cdot | s))] + \frac{1}{1-\gamma} \left(V_{L,k}^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) - V_{L,k}^{\bar{\pi}^k}(d_\rho^{\pi^*}) \right) \\ & \quad + \frac{3(B_\lambda + 1/N)}{(1-\gamma)^2} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| + \frac{3\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|. \end{aligned}$$

Taking the average from $k = 0$ to $k = K - 1$, we have

$$\begin{aligned} & \frac{1}{K} \sum_{k=0}^{K-1} \left(V_{L,k}^{\pi^*}(\rho) - V_{L,k}^{\bar{\pi}^k}(\rho) \right) \\ & \leq \frac{N}{(1-\gamma)K\alpha} \mathbb{E}_{s \sim d_\rho^{\pi^*}} [D_{\text{KL}}(\pi^*(\cdot | s) || \pi_0(\cdot | s))] + \frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} \left(V_{L,k}^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) - V_{L,k}^{\bar{\pi}^k}(d_\rho^{\pi^*}) \right) \\ & \quad + \frac{3(B_\lambda + 1/N)}{(1-\gamma)^2 K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| + \frac{3\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4 K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\bar{\pi}^k - \pi_i^k\|. \quad (72) \end{aligned}$$

The second term on the right hand side can be decomposed as follows

$$\begin{aligned} & \frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} \left(V_{L,k}^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) - V_{L,k}^{\bar{\pi}^k}(d_\rho^{\pi^*}) \right) \\ & \leq \frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} \left(V_0^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) - V_0^{\bar{\pi}^k}(d_\rho^{\pi^*}) \right) \\ & \quad + \frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} (\lambda^k - \nu^k)^\top \left(V_g^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) - V_g^{\bar{\pi}^k}(d_\rho^{\pi^*}) \right) \\ & = \frac{V_0^{\bar{\pi}^K}(d_\rho^{\pi^*})}{(1-\gamma)K} + \frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} \left((\lambda^{k+1} - \nu^{k+1})^\top V_g^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) - (\lambda^k - \nu^k)^\top V_g^{\bar{\pi}^k}(d_\rho^{\pi^*}) \right) \\ & \quad + \frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} (\lambda^k - \nu^k - \lambda^{k+1} + \nu^{k+1})^\top V_g^{\bar{\pi}^{k+1}}(d_\rho^{\pi^*}) \end{aligned}$$

$$\begin{aligned}
&= \frac{V_0^{\pi^K}(d_\rho^{\pi^\star})}{(1-\gamma)K} + \frac{1}{(1-\gamma)K} \sum_{i=1}^N (\lambda_i^K - \nu_i^K) V_i^{\pi^K}(d_\rho^{\pi^\star}) \\
&\quad + \frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} \sum_{i=1}^N (\lambda_i^k - \nu_i^k - \lambda_i^{k+1} + \nu_i^{k+1}) V_i^{\pi^{k+1}}(d_\rho^{\pi^\star}).
\end{aligned} \tag{73}$$

We know that the value functions are bounded between $[0, \frac{1}{1-\gamma}]$. The projection in the update of the dual variable in Eq. (38) guarantees $\lambda_i^k \in [0, B_\lambda]$. It is also straightforward to see that

$$|\lambda_{i,k} - \lambda_{i,k+1}| \leq \frac{\eta}{1-\gamma} + B\eta, |\nu_{i,k} - \nu_{i,k+1}| \leq \frac{\eta}{1-\gamma} + B\eta, \forall i = 1, 2, \dots, N, k = 0, 1, \dots, K-1.$$

Using these bounds in Eq. (73), we get

$$\begin{aligned}
\frac{1}{(1-\gamma)K} \sum_{k=0}^{K-1} \left(V_{L,k}^{\pi^{k+1}}(d_\rho^{\pi^\star}) - V_{L,k}^{\pi^k}(d_\rho^{\pi^\star}) \right) &\leq \frac{1}{(1-\gamma)^2 K} + \frac{NB_\lambda}{(1-\gamma)^2 K} + \frac{2N\eta}{(1-\gamma)^3 K} + \frac{2NB\eta}{(1-\gamma)^2 K} \\
&\leq \frac{2NB_\lambda}{(1-\gamma)^2 K} + \frac{4N\eta}{(1-\gamma)^3 K}.
\end{aligned} \tag{74}$$

Finally, combining Eqs. (72) and (74) yields

$$\begin{aligned}
&\frac{1}{K} \sum_{k=0}^{K-1} \left(V_{L,k}^{\pi^\star}(\rho) - V_{L,k}^{\pi^k}(\rho) \right) \\
&\leq \frac{N}{(1-\gamma)K\alpha} \mathbb{E}_{s \sim d_\rho^{\pi^\star}} [D_{\text{KL}}(\pi^\star(\cdot | s) || \pi_0(\cdot | s))] + \frac{3(B_\lambda + 1/N)}{(1-\gamma)^2 K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|Q_i^{\pi_i^k} - Q_i^k\| \\
&\quad + \frac{3\sqrt{|\mathcal{S}||\mathcal{A}|}(B_\lambda + 1/N)}{(1-\gamma)^4 K} \sum_{k=0}^{K-1} \sum_{i=1}^N \|\pi^k - \pi_i^k\| + \frac{2NB_\lambda}{(1-\gamma)^2 K} + \frac{4N\eta}{(1-\gamma)^3 K},
\end{aligned}$$

which leads to the claimed result by recognizing the fact that for $D_{\text{KL}}(p_1 || p_2) \leq \log |\mathcal{A}|$ for $p_1, p_2 \in \Delta_{\mathcal{A}}$ if p_2 is a uniform distribution. \square

D.8 Proof of Lemma 10

From the definition of $p_i^{k,t}$ in Eq. (59),

$$\begin{aligned}
\|p_i^{k,t}\| &= \|\omega_i^{k,t+1} - \hat{\omega}_i^{k,t+1}\| \\
&\leq \|\omega_i^{k,t+1} - \omega_i^{k,t}\| + \|\omega_i^{k,t} - \hat{\omega}_i^{k,t+1}\| \\
&\leq 2\|\omega_i^{k,t} - \hat{\omega}_i^{k,t+1}\| \\
&\leq 2\beta \|r_i(s_i^{k,t}, a_i^{k,t}) + (\gamma\phi(s_i^{k,t+1}, a_i^{k,t+1}) - \phi(s_i^{k,t}, a_i^{k,t}))^\top \omega_i^{k,t}\| \\
&\leq 2(1 + 2B_\omega)\beta,
\end{aligned}$$

where the second inequality is due to the fact that Π_{B_ω} is the projection to a convex set and $\Pi_{B_\omega} \omega_i^{k,t} = \omega_i^{k,t}$.

Similarly,

$$\|z_i^{k,t+1} - z_i^{k,t}\| = \|\omega_i^{k,t+1} - \omega_i^{k,t}\| \leq \|\omega_i^{k,t} - \hat{\omega}_i^{k,t+1}\| \leq (1 + 2B_\omega)\beta.$$

\square

D.9 Proof of Lemma 11

Recall the definition of matrix M^π in Eq. (27). Given any π , we define the matrix $\tilde{P}^\pi \in \mathbb{R}^{|S||\mathcal{A}| \times |S||\mathcal{A}|}$ such that $\tilde{P}^\pi(s', a' | s, a) = P(s' | s, a)\pi(a' | s')$. Then, we have for any two vectors $\omega_1, \omega_2 \in \mathbb{R}^d$

$$\begin{aligned}\mathbb{E}_{\hat{\pi}_i^k}[\omega_1^\top \phi(s, a)\phi(s', a')^\top \omega_2] &= \sum_{s, s', a, a'} \tilde{\mu}_{\hat{\pi}_i^k}(s, a) \tilde{P}^{\hat{\pi}_i^k}(s', a' | s, a) \omega_1^\top \phi(s, a)\phi(s', a')^\top \omega_2 \\ &= \omega_1^\top \Phi^\top M^{\hat{\pi}_i^k} \tilde{P}^{\hat{\pi}_i^k} \Phi \omega_2.\end{aligned}$$

Since this is true for any ω_1, ω_2 , we know

$$\mathbb{E}_{\hat{\pi}_i^k}[\phi(s, a)\phi(s', a')^\top] = \Phi^\top M^{\hat{\pi}_i^k} \tilde{P}^{\hat{\pi}_i^k} \Phi.$$

Similarly, we can show

$$\begin{aligned}\mathbb{E}_{\hat{\pi}_i^k}[\phi(s, a)\phi(s, a)^\top] &= \Phi^\top M^{\hat{\pi}_i^k} \Phi, \\ \mathbb{E}_{\hat{\pi}_i^k}[r(s, a)\phi(s, a)] &= \Phi^\top M^{\hat{\pi}_i^k} r_i,\end{aligned}$$

where $r_i = [r_i(s_0, a_0), \dots, r_i(s_{|S|}, a_{|\mathcal{A}|})]^\top \in \mathbb{R}^{|S||\mathcal{A}|}$.

We define $e_i^{\hat{\pi}_i^k}(\omega) = \bar{H}^{\hat{\pi}_i^k} \omega + \bar{b}_i^{\hat{\pi}_i^k}$ and note that $e_i^{\hat{\pi}_i^k}(\omega_i^*(\hat{\pi}_i^k)) = 0$. From the equations above and the definition of \bar{H}^π and \bar{b}_i^π ,

$$\begin{aligned}e_i^{\hat{\pi}_i^k}(\omega) &= \bar{H}^{\hat{\pi}_i^k} \omega + \bar{b}_i^{\hat{\pi}_i^k} \\ &= \mathbb{E}_{\hat{\pi}_i^k}[\phi(s, a)((\gamma\phi(s', s')^\top - \phi(s, a)^\top)\omega - r_i(s, a))] \\ &= (\gamma\Phi^\top M^{\hat{\pi}_i^k} \tilde{P}^{\hat{\pi}_i^k} \Phi - \Phi^\top M^{\hat{\pi}_i^k} \Phi)\omega + \Phi^\top M^{\hat{\pi}_i^k} r_i \\ &= \Phi^\top M^{\hat{\pi}_i^k} (\gamma\tilde{P}^{\hat{\pi}_i^k} \Phi\omega - \Phi\omega + r_i) \\ &= \Phi^\top M^{\hat{\pi}_i^k} (T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega),\end{aligned}$$

where $T_i^\pi : \mathbb{R}^{|S||\mathcal{A}|} \rightarrow \mathbb{R}^{|S||\mathcal{A}|}$ for any policy π is the Bellman operator defined such that

$$(T_i^\pi Q)(s, a) = \mathbb{E}_\pi[r_i(s, a) + \gamma Q(s', a')], \quad \forall Q \in \mathbb{R}^{|S||\mathcal{A}|}. \quad (75)$$

The projection of a vector to the span of Φ under the weighted $M^{\hat{\pi}_i^k}$ norm is carried out through the projection matrix $\Pi_\Phi^{\hat{\pi}_i^k}$

$$\Pi_\Phi^{\hat{\pi}_i^k} = \Phi(\Phi^\top M^{\hat{\pi}_i^k} \Phi)^{-1} \Phi^\top M^{\hat{\pi}_i^k}.$$

It is obvious that $\Phi^\top M^{\hat{\pi}_i^k} \Pi_\Phi^{\hat{\pi}_i^k} = \Phi^\top M^{\hat{\pi}_i^k}$.

Eq. (75) implies that for any $Q_1, Q_2 \in \mathbb{R}^{|S||\mathcal{A}|}$

$$\begin{aligned}\|T_i^{\hat{\pi}_i^k} Q_1 - T_i^{\hat{\pi}_i^k} Q_2\|_{M^{\hat{\pi}_i^k}}^2 &= \gamma(Q_1 - Q_2)^\top (\tilde{P}^{\hat{\pi}_i^k})^\top M^{\hat{\pi}_i^k} \tilde{P}^{\hat{\pi}_i^k} (Q_1 - Q_2) \\ &= \sum_{s, a} \tilde{\mu}_{\hat{\pi}_i^k}(s, a) \left(\sum_{s', a'} P(s' | s, a) \hat{\pi}_i^k(a' | s') (Q_1 - Q_2)(s', a') \right)^2 \\ &\leq \sum_{s, a} \tilde{\mu}_{\hat{\pi}_i^k}(s, a) \sum_{s', a'} P(s' | s, a) \hat{\pi}_i^k(a' | s') (Q_1(s', a') - Q_2(s', a'))^2 \\ &\leq \sum_{s', a'} \tilde{\mu}_{\hat{\pi}_i^k}(s', a') (Q_1(s', a') - Q_2(s', a'))^2\end{aligned}$$

$$= \|Q_1 - Q_2\|_{M^{\hat{\pi}_i^k}}^2, \quad (76)$$

where the inequality follows from Jensen's inequality. Eq. (76) implies another property of $\Pi_{\Phi}^{\hat{\pi}_i^k}$, which is the contraction of $\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}$ under the weighted $M^{\hat{\pi}_i^k}$ norm. Specifically, we have for any $\omega \in \mathbb{R}^d$

$$\begin{aligned} \|\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega_i^*(\hat{\pi}_i^k)\|_{M^{\hat{\pi}_i^k}} &= \|\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega_i^*(\hat{\pi}_i^k))\|_{M^{\hat{\pi}_i^k}} \\ &\leq \|T_i^{\hat{\pi}_i^k}(\Phi\omega) - T_i^{\hat{\pi}_i^k}(\Phi\omega_i^*(\hat{\pi}_i^k))\|_{M^{\hat{\pi}_i^k}} \\ &\leq \gamma \|\Phi\omega - \Phi\omega_i^*(\hat{\pi}_i^k)\|_{M^{\hat{\pi}_i^k}}. \end{aligned}$$

Then, we have for any $\omega \in \mathbb{R}^d$,

$$\begin{aligned} &(\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} \bar{H}^{\hat{\pi}_i^k} (\omega - \omega_i^*(\hat{\pi}_i^k)) \\ &= (\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} (e_i^{\hat{\pi}_i^k}(\omega) - e_i^{\hat{\pi}_i^k}(\omega_i^*(\hat{\pi}_i^k))) \\ &= (\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} e_i^{\hat{\pi}_i^k}(\omega) \\ &= (\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} \Phi^{\top} M^{\hat{\pi}_i^k} (T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega) \\ &= (\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} \Phi^{\top} M^{\hat{\pi}_i^k} ((I - \Pi_{\Phi}^{\hat{\pi}_i^k}) T_i^{\hat{\pi}_i^k}(\Phi\omega) + \Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega) \\ &= (\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} \Phi^{\top} M^{\hat{\pi}_i^k} (\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega), \end{aligned}$$

where the last equality follows from $\Phi^{\top} M^{\hat{\pi}_i^k} \Pi_{\Phi}^{\hat{\pi}_i^k} = \Phi^{\top} M^{\hat{\pi}_i^k}$.

Using the contraction of $\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}$ and the Cauchy-Schwarz inequality, we have for any $\omega \in \mathbb{R}^d$

$$\begin{aligned} &(\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} \bar{H}^{\hat{\pi}_i^k} (\omega - \omega_i^*(\hat{\pi}_i^k)) \\ &= (\omega - \omega_i^*(\hat{\pi}_i^k))^{\top} \Phi^{\top} M^{\hat{\pi}_i^k} (\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega) \\ &= (\Phi\omega - \Phi\omega_i^*(\hat{\pi}_i^k))^{\top} M^{\hat{\pi}_i^k} (\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega_i^*(\hat{\pi}_i^k)) + (\Phi\omega - \Phi\omega_i^*(\hat{\pi}_i^k))^{\top} M^{\hat{\pi}_i^k} (\Phi\omega_i^*(\hat{\pi}_i^k) - \Phi\omega) \\ &\leq \|\Phi\omega - \Phi\omega_i^*(\hat{\pi}_i^k)\|_{M^{\hat{\pi}_i^k}} \|\Pi_{\Phi}^{\hat{\pi}_i^k} T_i^{\hat{\pi}_i^k}(\Phi\omega) - \Phi\omega_i^*(\hat{\pi}_i^k)\|_{M^{\hat{\pi}_i^k}} - \|\Phi\omega - \Phi\omega_i^*(\hat{\pi}_i^k)\|_{M^{\hat{\pi}_i^k}}^2 \\ &\leq (\gamma - 1) \|\Phi\omega - \Phi\omega_i^*(\hat{\pi}_i^k)\|_{M^{\hat{\pi}_i^k}}^2 \\ &\leq (\gamma - 1) \sigma_{\min}(\Phi) \sigma_{\min}(M^{\hat{\pi}_i^k}) \|\omega - \omega_i^*(\hat{\pi}_i^k)\|^2 \\ &\leq (\gamma - 1) \sigma_{\min}(\Phi) \min_{s,a} \mu_{\hat{\pi}_i^k}(s) \hat{\pi}_t(a | s) \|\omega - \omega_i^*(\hat{\pi}_i^k)\|^2, \end{aligned}$$

where σ_{\min} denotes the singular value with the smallest magnitude, and the last inequality follows from the fact that the singular values of a diagonal matrix are the diagonal entries. \square

D.10 Proof of Lemma 12

Recall that the Markovian samples generated by our algorithm are

$$s_i^{k,t-\tau} \xrightarrow{\hat{\pi}_i^k} a_i^{k,t-\tau} \xrightarrow{P} s_i^{k,t-\tau+1} \xrightarrow{\hat{\pi}_i^k} a_i^{k,t-\tau+1} \xrightarrow{P} \dots \xrightarrow{P} s_i^{k,t} \xrightarrow{\hat{\pi}_i^k} a_i^{k,t} \xrightarrow{P} s_i^{k,t+1} \xrightarrow{\hat{\pi}_i^k} a_i^{k,t+1}.$$

Let $\bar{O}_i^k = (\bar{s}_i^k, \bar{a}_i^k, \bar{s}_i^{k'}, \bar{a}_i^{k'})$ where $\bar{s}_i^k \sim \mu_{\hat{\pi}_i^k}$, $\bar{a}_i^k \sim \hat{\pi}_i^k(\cdot | \bar{s}_i^k)$, $\bar{s}_i^{k'} \sim P(\cdot | \bar{s}_i^k, \bar{a}_i^k)$, and $\bar{a}_i^{k'} \sim \hat{\pi}_i^k(\cdot | \bar{s}_i^{k'})$. We can decompose the term of interest as

$$\mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t}, O_i^{k,t})] = \mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t}, O_i^{k,t}) - \Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, O_i^{k,t})]$$

$$+ \mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, O_i^{k,t}) - \Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, \bar{O}_i^k)] + \mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, \bar{O}_i^k)]. \quad (77)$$

To bound the first term of Eq. (77),

$$\begin{aligned} & \Gamma_i(\hat{\pi}_i^k, z_i^{k,t}, O_i^{k,t}) - \Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, O_i^{k,t}) \\ & \leq (z_i^{k,t} - z_i^{k,t-\tau})^\top \left(R_i(O_i^{k,t}) + H(O_i^{k,t})\omega_i^*(\hat{\pi}_i^k) \right) \\ & \quad + (z_i^{k,t})^\top (H(O_i^{k,t}) - \bar{H}^{\hat{\pi}_i^k})z_i^{k,t} - (z_i^{k,t-\tau})^\top (H(O_i^{k,t}) - \bar{H}^{\hat{\pi}_i^k})z_i^{k,t-\tau} \\ & \leq \|z_i^{k,t} - z_i^{k,t-\tau}\| (1 + 2B_\omega) + (z_i^{k,t} - z_i^{k,t-\tau})^\top (H(O_i^{k,t}) - \bar{H}^{\hat{\pi}_i^k})z_i^{k,t-\tau} \\ & \quad + (z_i^{k,t})^\top (H(O_i^{k,t}) - \bar{H}^{\hat{\pi}_i^k})(z_i^{k,t} - z_i^{k,t-\tau}) \\ & \leq (1 + 2B_\omega)\|z_i^{k,t} - z_i^{k,t-\tau}\| + 16B_\omega \sum_{t'=t-\tau}^{t-1} \|z_i^{k,t'+1} - z_i^{k,t'}\| \\ & \leq (1 + 18B_\omega) \cdot \tau(1 + 2B_\omega)\beta \\ & \leq (1 + 18B_\omega)^2\beta\tau. \end{aligned} \quad (78)$$

Let $\mathcal{F}_i^{k,t}$ denote the past randomness in outer loop iteration k up to inner loop iteration t at agent i , i.e. $\mathcal{F}_i^{k,t} = \{O_i^{k,0}, O_i^{k,1}, \dots, O_i^{k,t}\}$. To treat the second term of Eq. (77), we have for all $t \geq \tau$

$$\begin{aligned} & \mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, O_i^{k,t}) \mid \mathcal{F}_i^{k,t-\tau} - \Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, \bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}] \\ & = (z_i^{k,t-\tau})^\top \mathbb{E}[R_i(O_i^{k,t}) - R_i(\bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}] + (z_i^{k,t-\tau})^\top \mathbb{E}[H(O_i^{k,t}) - H(\bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}]\omega_i^*(\hat{\pi}_i^k) \\ & \quad + (z_i^{k,t-\tau})^\top \mathbb{E}[H(O_i^{k,t}) - H(\bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}]z_i^{k,t-\tau} \\ & \leq \|z_i^{k,t-\tau}\| \|\mathbb{E}[R_i(O_i^{k,t}) - R_i(\bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}]\| + \|z_i^{k,t-\tau}\| \mathbb{E}[H(O_i^{k,t}) - H(\bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}] \|\omega_i^*(\hat{\pi}_i^k)\| \\ & \quad + \|z_i^{k,t-\tau}\|^2 \|\mathbb{E}[H(O_i^{k,t}) - H(\bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}]\| \\ & \leq B_\omega \cdot d_{TV}(O_i^{k,t}, \bar{O}_i^k \mid \mathcal{F}_i^{k,t-\tau}) + 2B_\omega \cdot B_\omega \cdot 2d_{TV}(O_i^{k,t}, \bar{O}_i^k \mid \mathcal{F}_i^{k,t-\tau}) \\ & \quad + 4B_\omega^2 \cdot 2d_{TV}(O_i^{k,t}, \bar{O}_i^k \mid \mathcal{F}_i^{k,t-\tau}) \\ & \leq (B_\omega + 12B_\omega^2)d_{TV}(O_i^{k,t}, \bar{O}_i^k \mid \mathcal{F}_i^{k,t-\tau}) \\ & \leq (B_\omega + 12B_\omega^2)C_0\ell^k \\ & \leq (B_\omega + 12B_\omega^2)\beta, \end{aligned} \quad (79)$$

where the fourth inequality follows from an argument similar to the one in Wu et al. (2020)[Lemma D.11], the last inequality uses Assumption 2, and the last inequality is a result of Eq. (19).

By the definition of $H(\bar{O}_i^k)$, $\bar{H}^{\hat{\pi}_i^k}$, and $\omega_i^*(\hat{\pi}_i^k)$, we have for the last term of Eq. (77)

$$\begin{aligned} \mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t-\tau}, \bar{O}_i^k) \mid \mathcal{F}_i^{k,t-\tau}] & = \mathbb{E}[(z_i^{k,t-\tau})^\top (R_i(\bar{O}_i^k) + H(\bar{O}_i^k)\omega_i^*(\hat{\pi}_i^k)) \mid \mathcal{F}_i^{k,t-\tau}] \\ & \quad + \mathbb{E}[(z_i^{k,t-\tau})^\top (H(\bar{O}_i^k) - \bar{H}^{\hat{\pi}_i^k})z_i^{k,t-\tau} \mid \mathcal{F}_i^{k,t-\tau}] \\ & = 0 + 0 = 0. \end{aligned} \quad (80)$$

Plugging Eqs. (78) and (79) into Eq. (77), we have

$$\begin{aligned} \mathbb{E}[\Gamma_i(\hat{\pi}_i^k, z_i^{k,t}, O_i^{k,t})] & \leq (1 + 18B_\omega)^2\beta\tau + (B_\omega + 12B_\omega^2)\beta \\ & \leq 2(1 + 18B_\omega)^2\beta\tau. \end{aligned}$$

□