

Federated Natural Policy Gradient Methods for Multi-task Reinforcement Learning

Tong Yang*
CMU

Shicong Cen†
CMU

Yuting Wei‡
UPenn

Yuxin Chen§
UPenn

Yuejie Chi¶
CMU

October 31, 2023

Abstract

Federated reinforcement learning (RL) enables collaborative decision making of multiple distributed agents without sharing local data trajectories. In this work, we consider a multi-task setting, in which each agent has its own private reward function corresponding to different tasks, while sharing the same transition kernel of the environment. Focusing on infinite-horizon tabular Markov decision processes, the goal is to learn a globally optimal policy that maximizes the sum of the discounted total rewards of all the agents in a decentralized manner, where each agent only communicates with its neighbors over some prescribed graph topology.

We develop federated vanilla and entropy-regularized natural policy gradient (NPG) methods under softmax parameterization, where gradient tracking is applied to the global Q-function to mitigate the impact of imperfect information sharing. We establish non-asymptotic global convergence guarantees under exact policy evaluation, which are nearly independent of the size of the state-action space and illuminate the impacts of network size and connectivity. To the best of our knowledge, this is the first time that global convergence is established for federated multi-task RL using policy optimization. Moreover, the convergence behavior of the proposed algorithms is robust against inexactness of policy evaluation.

Keywords: federated reinforcement learning, multi-task reinforcement learning, natural policy gradient methods, entropy regularization, global convergence

Contents

1	Introduction	2
1.1	Our contributions	3
1.2	Related work	4
2	Model and backgrounds	4
2.1	Markov decision processes	4
2.2	Entropy-regularized RL	5
2.3	Natural policy gradient methods	6
3	Federated NPG methods for multi-task RL	7
3.1	Federated multi-task RL	7
3.2	Proposed federated NPG algorithms	8

*Department of Electrical and Computer Engineering, Carnegie Mellon University; email: tongyang@andrew.cmu.edu.

†Department of Electrical and Computer Engineering, Carnegie Mellon University; email: shicongc@andrew.cmu.edu.

‡Department of Statistics and Data Science, Wharton School, University of Pennsylvania; email: ytwei@wharton.upenn.edu.

§Department of Statistics and Data Science, Wharton School, University of Pennsylvania; email: yuxinc@wharton.upenn.edu.

¶Department of Electrical and Computer Engineering, Carnegie Mellon University; email: yuejiechi@cmu.edu.

4	Theoretical guarantees	10
4.1	Global convergence of FedNPG	10
4.2	Global convergence of FedNPG with entropy regularization	11
5	Conclusions	12
A	Convergence analysis	16
A.1	Analysis of entropy-regularized FedNPG with exact policy evaluation	16
A.2	Analysis of entropy-regularized FedNPG with inexact policy evaluation	18
A.3	Analysis of FedNPG with exact policy evaluation	20
A.4	Analysis of FedNPG with inexact policy evaluation	23
B	Proof of key lemmas	25
B.1	Proof of Lemma 1	25
B.2	Proof of Lemma 2	30
B.3	Proof of Lemma 3	32
B.4	Proof of Lemma 4	34
B.5	Proof of Lemma 5	35
C	Proof of auxiliary lemmas	36
C.1	Proof of Lemma 6	36
C.2	Proof of Lemma 8	37
C.3	Proof of Lemma 9	40
C.4	Proof of Lemma 11	41

1 Introduction

Federated reinforcement learning (FRL) is an emerging paradigm that combines the advantages of federated learning (FL) and reinforcement learning (RL) (Qi et al., 2021; Zhuo et al., 2019), allowing multiple agents to learn a shared policy from local experiences, without exposing their private data to a central server nor other agents. FRL is poised to enable collaborative and efficient decision making in scenarios where data is distributed, heterogeneous, and sensitive, which arise frequently in applications such as edge computing, smart cities, and healthcare (Wang et al., 2023, 2020; Zhuo et al., 2019), to name just a few. As has been observed (Lian et al., 2017), decentralized training can lead to performance improvements in FL by avoiding communication congestions at busy nodes such as the server, especially under high-latency scenarios. This motivates us to design algorithms for the fully decentralized setting, a scenario where the agents can only communicate with their local neighbors over a prescribed network topology.

In this work, we study the problem of *federated multi-task reinforcement learning* (Anwar and Raychowdhury, 2021; Qi et al., 2021; Yu et al., 2020), where each agent collects its own reward — possibly unknown to other agents — corresponding to the local task at hand, while having access to the same dynamics (i.e., transition kernel) of the environment. The collective goal is to learn a shared policy that maximizes the total rewards accumulated from all the agents; in other words, one seeks a policy that performs well in terms of overall benefits, rather than biasing towards any individual task, achieving the Pareto frontier in a multi-objective context. There is no shortage of application scenarios where federated multi-task RL becomes highly relevant. For instance, in healthcare (Zerka et al., 2020), different hospitals may be interested in finding an optimal treatment for all patients without disclosing private data, where the effectiveness of the treatment can vary across different hospitals due to demographical differences. As another potential application, to enhance ChatGPT’s performance across different tasks or domains (M Alshater, 2022; Rahman et al., 2023), one might consult domain experts to chat and rate ChatGPT’s outputs for solving different tasks, and train ChatGPT in a federated manner without exposing private data or feedback of each expert.

Nonetheless, despite the promise, provably efficient algorithms for federated multi-task RL remain substantially under-explored, especially in the fully decentralized setting. The heterogeneity of local tasks leads to a higher degree of disagreements between the global value function and local value functions of individual agents. Due to the lack of global information sharing, care needs to be taken to judiciously balance the use

of neighboring information (to facilitate consensus) and local data (to facilitate learning) when updating the policy. To the best of our knowledge, no algorithms are currently available to find the global optimal policy with non-asymptotic convergence guarantees even for tabular infinite-horizon Markov decision processes.

Motivated by the connection with decentralized optimization, it is tempting to take a policy optimization perspective to tackle this challenge. Policy gradient (PG) methods, which seek to learn the policy of interest via first-order optimization methods, play an eminent role in RL due to their simplicity and scalability. In particular, natural policy gradient (NPG) methods (Amari, 1998; Kakade, 2001) are among the most popular variants of PG methods, underpinning default methods used in practice such as trust region policy optimization (TRPO) (Schulman et al., 2015) and proximal policy optimization (PPO) (Schulman et al., 2017). On the theoretical side, it has also been established recently that the NPG algorithm enjoys fast global convergence to the optimal policy in an almost dimension-free manner (Agarwal et al., 2021; Cen et al., 2021), where the iteration complexity is nearly independent of the size of the state-action space. Inspired by the efficacy of NPG methods, it is natural to ask:

*Can we develop **federated** variants of NPG methods that are easy to implement in the fully decentralized setting with **non-asymptotic global convergence** guarantees for multi-task RL?*

1.1 Our contributions

Focusing on infinite-horizon Markov decision processes (MDPs), we provide an affirmative answer to the above question, by developing federated NPG (FedNPG) methods for solving both the vanilla and entropy-regularized multi-task RL problems with finite-time global convergence guarantees. While entropy regularization is often incorporated as an effective strategy to encourage exploration during policy learning, solving the entropy-regularized RL problem is of interest in its own right, as the optimal regularized policy possesses desirable robust properties with respect to reward perturbations (Eysenbach and Levine, 2021; McKelvey and Palfrey, 1995).

Due to the multiplicative update nature of NPG methods under softmax parameterization, it is more convenient to work with the logarithms of local policies in the decentralized setting. In each iteration of the proposed FedNPG method, the logarithms of local policies are updated by a weighted linear combination of two terms (up to normalization): a gossip mixing (Nedic and Ozdaglar, 2009) of the logarithms of neighboring local policies, and a local estimate of the global Q-function tracked via the technique of dynamic average consensus (Zhu and Martínez, 2010), a prevalent idea in decentralized optimization that allows for the use of large constant learning rates (Di Lorenzo and Scutari, 2016; Nedic et al., 2017; Qu and Li, 2017) to accelerate convergence. Our contributions are as follows.

- We propose FedNPG methods for both the vanilla and entropy-regularized multi-task RL problems, where each agent only communicates with its neighbors and performs local computation using its own reward or task information.
- Assuming access to exact policy evaluation, we establish that the average iterate of vanilla FedNPG converges globally at a rate of $\mathcal{O}(1/T^{2/3})$ in terms of the sub-optimality gap for the multi-task RL problem, and that the last iterate of entropy-regularized FedNPG converges globally at a linear rate to the regularized optimal policy. Our convergence theory highlights the impacts of all salient problem parameters (see Table 1 for details), such as the size and connectivity of the communication network. In particular, the iteration complexities of FedNPG are again almost independent of the size of the state-action space, which recover prior results on the centralized NPG methods when the network is fully connected.
- We further demonstrate the stability of the proposed FedNPG methods when policy evaluations are only available in an inexact manner. To be specific, we prove that their convergence rates remain unchanged as long as the approximation errors are sufficiently small in the ℓ_∞ sense.

To the best of our knowledge, the proposed federated NPG methods are the first policy optimization methods for multi-task RL that achieve explicit non-asymptotic global convergence guarantees, allowing for fully decentralized communication without any need to share local reward/task information.

1.2 Related work

Global convergence of NPG methods for tabular MDPs. Agarwal et al. (2021) first establishes a $\mathcal{O}(1/T)$ last-iterate convergence rate of the NPG method under softmax parameterization with constant step size, assuming access to exact policy evaluation. When entropy regularization is in place, Cen et al. (2021) establishes a global linear convergence to the optimal regularized policy for the entire range of admissible constant learning rates using softmax parameterization and exact policy evaluation, which is further shown to be stable in the presence of ℓ_∞ policy evaluation errors. The iteration complexity of NPG methods is nearly independent with the size of the state-action space, which is in sharp contrast to softmax policy gradient methods that may take exponential time to converge (Li et al., 2023c; Mei et al., 2020). Lan (2023) proposed a more general framework through the lens of mirror descent for regularized RL with global linear convergence guarantees, which is further generalized in Zhan et al. (2023); Lan et al. (2023). Earlier analysis of regularized MDPs can be found in Shani et al. (2020). Besides, Xiao (2022) proves that vanilla NPG also achieves linear convergence when geometrically increasing learning rates are used; see also Khodadadian et al. (2021); Bhandari and Russo (2021). Zhou et al. (2022) developed an anchor-changing NPG method for multi-task RL under various optimality criteria in the centralized setting.

Distributed and federated RL. There have been a variety of settings being set forth for distributed and federated RL. Mnih et al. (2016); Espeholt et al. (2018); Assran et al. (2019); Khodadadian et al. (2022); Woo et al. (2023) focused on developing federated versions of RL algorithms to accelerate training, assuming all agents share the same transition kernel and reward function; in particular, Khodadadian et al. (2022); Woo et al. (2023) established the provable benefits of federated learning in terms of linear speedup. More pertinent to our work, Zhao et al. (2023); Anwar and Raychowdhury (2021) considered the federated multi-task framework, allowing different agents having private reward functions. Zhao et al. (2023) proposed an empirically probabilistic algorithm that can seek an optimal policy under the server-client setting, while Anwar and Raychowdhury (2021) developed new attack methods in the presence of adversarial agents. Different from the FRL framework, Chen et al. (2021, 2022b); Omidshafiei et al. (2017); Kar et al. (2012); Chen et al. (2022a); Zeng et al. (2021) considered the distributed multi-agent RL setting where the agents interact with a dynamic environment through a multi-agent Markov decision process, where each agent can have their own state or action spaces. Zeng et al. (2021) developed a decentralized policy gradient method where different agents have different MDPs.

Decentralized first-order optimization algorithms. Early work of consensus-based first-order optimization algorithms for the fully decentralized setting include but are not limited to Lobel and Ozdaglar (2008); Nedic and Ozdaglar (2009); Duchi et al. (2011). Gradient tracking, which leverages the idea of dynamic average consensus (Zhu and Martínez, 2010) to track the gradient of the global objective function, is a popular method to improve the convergence speed (Qu and Li, 2017; Nedic et al., 2017; Di Lorenzo and Scutari, 2016; Pu and Nedić, 2021; Li et al., 2020).

Notation. Boldface small and capital letters denote vectors and matrices, respectively. Sets are denoted with curly capital letters, e.g., \mathcal{S}, \mathcal{A} . We let $(\mathbb{R}^d, \|\cdot\|)$ denote the d -dimensional real coordinate space equipped with norm $\|\cdot\|$. The ℓ^p -norm of \mathbf{v} is denoted by $\|\mathbf{v}\|_p$, where $1 \leq p \leq \infty$, and the spectral norm of a matrix \mathbf{M} is denoted by $\|\mathbf{M}\|_2$. We let $[N]$ denote $\{1, \dots, N\}$, use $\mathbf{1}_N$ to represent the all-one vector of length N , and denote by $\mathbf{0}$ a vector or a matrix consisting of all 0's. We allow the application of functions such as $\log(\cdot)$ and $\exp(\cdot)$ to vectors or matrices, with the understanding that they are applied in an element-wise manner.

2 Model and backgrounds

2.1 Markov decision processes

Markov decision processes. We consider an infinite-horizon discounted Markov decision process (MDP) denoted by $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where \mathcal{S} and \mathcal{A} denote the state space and the action space, respectively, $\gamma \in [0, 1)$ indicates the discount factor, $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition kernel, and $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ stands for the reward function. To be more specific, for each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ and any state

setting	algorithms	iteration complexity	optimality criteria
unregularized	NPG (Agarwal et al., 2021)	$\mathcal{O}\left(\frac{1}{(1-\gamma)^2\varepsilon} + \frac{\log \mathcal{A} }{\eta\varepsilon}\right)$	$V^* - V^{\pi^{(t)}} \leq \varepsilon$
	FedNPG (ours)	$\mathcal{O}\left(\frac{\sqrt{\sigma N} \log \mathcal{A} }{(1-\gamma)^{\frac{9}{2}}(1-\sigma)\varepsilon^{\frac{3}{2}}} + \frac{1}{(1-\gamma)^2\varepsilon}\right)$	$\frac{1}{T} \sum_{t=0}^{T-1} (V^* - V^{\bar{\pi}^{(t)}}) \leq \varepsilon$
regularized	NPG (Cen et al., 2021)	$\mathcal{O}\left(\frac{1}{\tau\eta} \log\left(\frac{1}{\varepsilon}\right)\right)$	$V_\tau^* - V_\tau^{\pi^{(t)}} \leq \varepsilon$
	FedNPG (ours)	$\mathcal{O}\left(\max\left\{\frac{1}{\tau\eta}, \frac{1}{1-\sigma}\right\} \log\left(\frac{1}{\varepsilon}\right)\right)$	$V_\tau^* - V_\tau^{\bar{\pi}^{(t)}} \leq \varepsilon$

Table 1: Iteration complexities of NPG and FedNPG (ours) methods to reach ε -accuracy of the vanilla and entropy-regularized problems, where we assume exact gradient evaluation, and only keep the dominant terms w.r.t. ε . The policy estimates in the t -iteration are $\pi^{(t)}$ and $\bar{\pi}^{(t)}$ for NPG and FedNPG, respectively, where T is the number of iterations. Here, N is the number of agents, $\tau \leq 1$ is the regularization parameter, $\sigma \in [0, 1]$ is the spectral radius of the network, $\gamma \in [0, 1)$ is the discount factor, $|\mathcal{A}|$ is the size of the action space, and $\eta > 0$ is the learning rate. For vanilla FedNPG, the learning rate is set as $\eta = \eta_1 = \mathcal{O}\left(\frac{(1-\gamma)^9(1-\sigma)^2 \log|\mathcal{A}|}{TN\sigma}\right)^{1/3}$; for entropy-regularized FedNPG, the learning rate satisfies $0 < \eta < \eta_0 = \mathcal{O}\left(\frac{(1-\gamma)^7(1-\sigma)^2\tau}{\sigma N}\right)$. The iteration complexities of FedNPG reduce to their centralized counterparts when $\sigma = 0$.

$s' \in \mathcal{S}$, we denote by $P(s'|s, a)$ the transition probability from state s to state s' when action a is taken, and $r(s, a)$ the instantaneous reward received in state s when action a is taken. Furthermore, a policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ specifies an action selection rule, where $\pi(a|s)$ specifies the probability of taking action a in state s for each $(s, a) \in \mathcal{S} \times \mathcal{A}$.

For any given policy π , we denote by $V^\pi : \mathcal{S} \mapsto \mathbb{R}$ the corresponding value function, which is the expected discounted cumulative reward with an initial state $s_0 = s$, given by

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

where the randomness is over the trajectory generated following the policy $a_t \sim \pi(\cdot|s_t)$ and the MDP dynamic $s_{t+1} \sim P(\cdot|s_t, a_t)$. We also overload the notation $V^\pi(\rho)$ to indicate the expected value function of policy π when the initial state follows a distribution ρ over \mathcal{S} , namely, $V^\pi(\rho) := \mathbb{E}_{s \sim \rho} [V^\pi(s)]$. Similarly, the Q-function $Q^\pi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ of policy π is defined by

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad Q^\pi(s, a) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a \right], \quad (2)$$

which measures the expected discounted cumulative reward with an initial state $s_0 = s$ and an initial action $a_0 = a$, with expectation taken over the randomness of the trajectory. The optimal policy π^* refers to the policy that maximizes the value function $V^\pi(s)$ for all states $s \in \mathcal{S}$, which is guaranteed to exist (Puterman, 2014). The corresponding optimal value function and Q-function are denoted as V^* and Q^* , respectively.

2.2 Entropy-regularized RL

Entropy regularization (Williams and Peng, 1991; Ahmed et al., 2019) is a popular technique in practice that encourages stochasticity of the policy to promote exploration, as well as robustness against reward uncertainties. Mathematically, this can be viewed as adjusting the instantaneous reward based the current policy in use as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad r_\tau(s, a) := r(s, a) - \tau \log \pi(a|s), \quad (3)$$

where $\tau \geq 0$ denotes the regularization parameter. Typically, τ should not be too large to outweigh the actual rewards; for ease of presentation, we assume $\tau \leq \min\left\{1, \frac{1}{\log|\mathcal{A}|}\right\}$ (Cen et al., 2022b). Equivalently, this amounts to the entropy-regularized (also known as “soft”) value function, defined as

$$\forall s \in \mathcal{S}: \quad V_\tau^\pi(s) := V^\pi(s) + \tau \mathcal{H}(s, \pi). \quad (4)$$

Here, we define

$$\mathcal{H}(s, \pi) := \mathbb{E} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t) \mid s_0 = s \right] = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^\pi} \left[- \sum_{a \in \mathcal{A}} \pi(a | s') \log \pi(a | s') \right], \quad (5)$$

where $d_{s_0}^\pi$ is the discounted state visitation distribution of policy π given an initial state $s_0 \in \mathcal{S}$, denoted by

$$\forall s \in \mathcal{S}: \quad d_{s_0}^\pi(s) := (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s | s_0), \quad (6)$$

with the trajectory generated by following policy π in the MDP \mathcal{M} starting from state s_0 . Analogously, the regularized (or soft) Q-function Q_τ^π of policy π is related to the soft value function $V_\tau^\pi(s)$ as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad Q_\tau^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \in P(\cdot | s, a)} [V_\tau^\pi(s')], \quad (7a)$$

$$\forall s \in \mathcal{S}: \quad V_\tau^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot | s)} [-\tau \pi(a | s) + Q_\tau^\pi(s, a)]. \quad (7b)$$

The optimal regularized policy, the optimal regularized value function, and the Q-function are denoted by π_τ^* , V_τ^* , and Q_τ^* , respectively.

2.3 Natural policy gradient methods

Natural policy gradient (NPG) methods lie at the heart of policy optimization, serving as the backbone of popular heuristics such as TRPO (Schulman et al., 2015) and PPO (Schulman et al., 2017). Instead of directly optimizing the policy over the probability simplex, one often adopts the softmax parameterization, which parameterizes the policy as

$$\pi_\theta := \text{softmax}(\theta) \quad \text{or} \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \pi_\theta(a | s) := \frac{\exp \theta(s, a)}{\sum_{a' \in \mathcal{A}} \exp \theta(s, a')} \quad (8)$$

for any $\theta: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$.

Vanilla NPG method. In the tabular setting, the update rule of vanilla NPG at the t -th iteration can be concisely represented as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \pi^{(t+1)}(a | s) \propto \pi^{(t)}(a | s) \exp \left(\frac{\eta Q^{(t)}(s, a)}{1-\gamma} \right), \quad (9)$$

where $\eta > 0$ denotes the learning rate, and $Q^{(t)} = Q^{\pi^{(t)}}$ is the Q-function under policy $\pi^{(t)}$. Agarwal et al. (2021) shows that: in order to find an ε -optimal policy, NPG takes at most $\mathcal{O}\left(\frac{1}{(1-\gamma)^2 \varepsilon}\right)$ iterations, assuming exact policy evaluation.

Entropy-regularized NPG method. Turning to the regularized problem, we note that the update rule of entropy-regularized NPG becomes

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \pi^{(t+1)}(a | s) \propto (\pi^{(t)}(a | s))^{1-\frac{\eta\tau}{1-\gamma}} \exp \left(\frac{\eta Q_\tau^{(t)}(s, a)}{1-\gamma} \right), \quad (10)$$

where $\eta \in (0, \frac{1-\gamma}{\tau}]$ is the learning rate, and $Q_\tau^{(t)} = Q_\tau^{\pi^{(t)}}$ is the soft Q-function of policy $\pi^{(t)}$. Cen et al. (2022a) proves that entropy-regularized NPG enjoys fast global linear convergence to the optimal regularized policy: to find an ε -optimal regularized policy, entropy-regularized NPG takes no more than $\mathcal{O}\left(\frac{1}{\eta\tau} \log\left(\frac{1}{\varepsilon}\right)\right)$ iterations.

3 Federated NPG methods for multi-task RL

3.1 Federated multi-task RL

In this paper, we consider the federated multi-task RL setting, where a set of agents learn collaboratively a single policy that maximizes its average performance over all the tasks using only local computation and communication.

Multi-task RL. Each agent $n \in [N]$ has its own private reward function $r_n(s, a)$ — corresponding to different tasks — while sharing the same transition kernel of the environment. The goal is to collectively learn a single policy π that maximizes the global value function given by

$$V^\pi(s) = \frac{1}{N} \sum_{n=1}^N V_n^\pi(s), \quad (11)$$

where V_n^π is the value function of agent $n \in [N]$, defined by

$$\forall s \in \mathcal{S} : \quad V_n^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_n(s_t, a_t) | s_0 = s \right]. \quad (12)$$

Clearly, the global value function (11) corresponds to using the average reward of all agents

$$r(s, a) = \frac{1}{N} \sum_{n=1}^N r_n(s, a). \quad (13)$$

The global Q-function $Q^\pi(s, a)$ and the agent Q-functions $Q_n^\pi(s, a)$ can be defined in a similar manner obeying $Q^\pi(s, a) = \frac{1}{N} \sum_{n=1}^N Q_n^\pi(s, a)$.

In parallel, we are interested in the entropy-regularized setting, where each agent $n \in [N]$ is equipped with a regularized reward function given by

$$r_{\tau,n}(s, a) := r_n(s, a) - \tau \log \pi(a|s), \quad (14)$$

and we define similarly the regularized value function and the global regularized value function as

$$\forall s \in \mathcal{S} : \quad V_{\tau,n}^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_{\tau,n}(s_t, a_t) | s_0 = s \right], \quad \text{and} \quad V_\tau^\pi(s) = \frac{1}{N} \sum_{n=1}^N V_{\tau,n}^\pi(s). \quad (15)$$

The soft Q-function of agent n is given by

$$Q_{\tau,n}^\pi(s, a) = r_n(s, a) + \gamma \mathbb{E}_{s' \in P(\cdot|s,a)} [V_{\tau,n}^\pi(s')], \quad (16)$$

and the global soft Q-function is given by $Q_\tau^\pi(s, a) = \frac{1}{N} \sum_{n=1}^N Q_{\tau,n}^\pi(s, a)$.

Federated policy optimization in the fully decentralized setting. We consider a federated setting with fully decentralized communication, that is, all the agents are synchronized to perform information exchange over some prescribed network topology denoted by an undirected weighted graph $\mathcal{G}([N], E)$. Here, E stands for the edge set of the graph with N nodes — each corresponding to an agent — and two agents can communicate with each other if and only if there is an edge connecting them. The information sharing over the graph is best described by a mixing matrix (Nedic and Ozdaglar, 2009), denoted by $\mathbf{W} = [w_{ij}] \in [0, 1]^{N \times N}$, where w_{ij} is a positive number if $(i, j) \in E$ and 0 otherwise. We also make the following standard assumptions on the mixing matrix.

Assumption 1 (double stochasticity). *The mixing matrix $\mathbf{W} = [w_{ij}] \in [0, 1]^{N \times N}$ is symmetric (i.e., $\mathbf{W}^\top = \mathbf{W}$) and doubly stochastic (i.e., $\mathbf{W}\mathbf{1}_N = \mathbf{1}_N$, $\mathbf{1}_N^\top \mathbf{W} = \mathbf{1}_N^\top$).*

The following standard metric measures how fast information propagates over the graph.

Definition 1 (spectral radius). *The spectral radius of \mathbf{W} is defined as*

$$\sigma := \left\| \mathbf{W} - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \right\|_2 \in [0, 1). \quad (17)$$

The spectral radius σ determines how fast information propagate over the network. For instance, in a fully-connected network, we can achieve $\sigma = 0$ by setting $\mathbf{W} = \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top$. For control of $1/(1 - \sigma)$ regarding different graphs, we refer the readers to paper [Nedić et al. \(2018\)](#). In an Erdős-Rényi random graph, as long as the graph is connected, one has with high probability $\sigma \asymp 1$. Another immediate consequence is that for any $\mathbf{x} \in \mathbb{R}^N$, letting $\bar{x} = \frac{1}{N} \mathbf{1}_N^\top \mathbf{x}$ be its average, we have

$$\| \mathbf{W} \mathbf{x} - \bar{x} \mathbf{1}_N \|_2 \leq \sigma \| \mathbf{x} - \bar{x} \mathbf{1}_N \|_2, \quad (18)$$

where the consensus error contracts by a factor of σ .

3.2 Proposed federated NPG algorithms

Assuming softmax parameterization, the problem can be formulated as decentralized optimization,

$$\text{(unregularized)} \quad \max_{\theta} V^{\pi_{\theta}}(s) = \frac{1}{N} \sum_{n=1}^N V_n^{\pi_{\theta}}(s), \quad (19)$$

$$\text{(regularized)} \quad \max_{\theta} V_{\tau}^{\pi_{\theta}}(s) = \frac{1}{N} \sum_{n=1}^N V_{\tau,n}^{\pi_{\theta}}(s), \quad (20)$$

where $\pi_{\theta} := \text{softmax}(\theta)$ subject to communication constraints. Motivated by the success of NPG methods, we aim to develop federated NPG methods to achieve our goal. For notational convenience, let $\boldsymbol{\pi}^{(t)} := (\pi_1^{(t)}, \dots, \pi_N^{(t)})^\top$ be the collection of policy estimates at all agents in the t -th iteration. Let

$$\bar{\pi}^{(t)} := \text{softmax} \left(\frac{1}{N} \sum_{n=1}^N \log \pi_n^{(t)} \right), \quad (21)$$

which satisfies that $\bar{\pi}^{(t)}(a|s) \propto \left(\prod_{n=1}^N \pi_n^{(t)}(a|s) \right)^{1/N}$ for each $(s, a) \in \mathcal{S} \times \mathcal{A}$. Therefore, $\bar{\pi}^{(t)}$ could be seen as the normalized geometric mean of $\{\pi_n^{(t)}\}_{n \in [N]}$. Define the collection of Q-function estimates as

$$\mathbf{Q}^{(t)} := \left(Q_1^{\pi_1^{(t)}}, \dots, Q_N^{\pi_N^{(t)}} \right)^\top, \quad \mathbf{Q}_{\tau}^{(t)} := \left(Q_{\tau,1}^{\pi_1^{(t)}}, \dots, Q_{\tau,N}^{\pi_N^{(t)}} \right)^\top.$$

We shall often abuse the notation and treat $\boldsymbol{\pi}^{(t)}$, $\mathbf{Q}_{\tau}^{(t)}$ as matrices in $\mathbb{R}^{N \times |\mathcal{S}| |\mathcal{A}|}$, and treat $\pi^{(t)}(a|s)$, $\mathbf{Q}_{\tau}^{(t)}(a|s)$ as vectors in \mathbb{R}^N , for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

Vanilla federated NPG methods. To motivate the algorithm development, observe that the NPG method (cf. (9)) applied to (19) adopts the update rule

$$\pi^{(t+1)}(a|s) \propto \pi^{(t)}(a|s) \exp \left(\frac{\eta Q^{\pi^{(t)}}(s, a)}{1 - \gamma} \right) = \pi^{(t)}(a|s) \exp \left(\frac{\eta \sum_{n=1}^N Q_n^{\pi_n^{(t)}}(s, a)}{N(1 - \gamma)} \right)$$

for all $(s, a) \in \mathcal{S} \times \mathcal{A}$. Two challenges arise when executing this update rule: the policy estimates are maintained locally without consensus, and the global Q-function are unavailable in the decentralized setting. To address these challenges, we apply the idea of dynamic average consensus ([Zhu and Martínez, 2010](#)), where each agent maintains its own estimate $T_n^{(t)}(s, a)$ of the global Q-function, which are collected as vector

$$\mathbf{T}^{(t)} = (T_1^{(t)}, \dots, T_N^{(t)})^\top.$$

Algorithm 1 Federated NPG (FedNPG)

- 1: **Input:** learning rate $\eta > 0$, iteration number $T \in \mathbb{N}_+$, mixing matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$.
- 2: **Initialize:** $\boldsymbol{\pi}^{(0)}, \mathbf{T}^{(0)} = \mathbf{Q}^{(0)}$.
- 3: **for** $t = 0, 1, \dots, T - 1$ **do**
- 4: Update the policy for each $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\log \boldsymbol{\pi}^{(t+1)}(a|s) = \mathbf{W} \log \boldsymbol{\pi}^{(t)}(a|s) + \frac{\eta}{1-\gamma} \mathbf{T}^{(t)}(s, a) - \log \mathbf{z}^{(t)}(s), \quad (22)$$

where $\mathbf{z}^{(t)}(s) = \sum_{a' \in \mathcal{A}} \exp \left\{ \mathbf{W} \log \boldsymbol{\pi}^{(t)}(a'|s) + \frac{\eta}{1-\gamma} \mathbf{T}^{(t)}(s, a') \right\}$.

- 5: Evaluate $\mathbf{Q}^{(t+1)}$.
- 6: Update the global Q-function estimate for each $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\mathbf{T}^{(t+1)}(s, a) = \mathbf{W} \left(\mathbf{T}^{(t)}(s, a) + \underbrace{\mathbf{Q}^{(t+1)}(s, a) - \mathbf{Q}^{(t)}(s, a)}_{\text{Q-tracking}} \right). \quad (23)$$

7: **end for**

At each iteration, each agent updates its policy estimates based on its neighbors' information via gossip mixing, in addition to a correction term that tracks the difference $Q_n^{\pi_n^{(t+1)}}(s, a) - Q_n^{\pi_n^{(t)}}(s, a)$ of the local Q-functions between consecutive policy updates. Note that the mixing is applied linearly to the logarithms of local policies, which translates into a multiplicative mixing of the local policies. Algorithm 1 summarizes the detailed procedure of the proposed algorithm written in a compact matrix form, which we dub as federated NPG (FedNPG). Note that the agents do not need to share their reward functions with others, and agent $n \in [N]$ will only be responsible to evaluate the local policy $\pi_n^{(t)}$ using the local reward r_n .

Entropy-regularized federated NPG methods. Moving onto the entropy regularized case, we adopt similar algorithmic ideas to decentralize (10), and propose the federated NPG (FedNPG) method with entropy regularization, summarized in Algorithm 2. Clearly, the entropy-regularized FedNPG method reduces to the vanilla FedNPG in the absence of the regularization (i.e., when $\tau = 0$).

Algorithm 2 Federated NPG (FedNPG) with entropy regularization

- 1: **Input:** learning rate $\eta > 0$, iteration number $T \in \mathbb{N}_+$, mixing matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$, regularization coefficient $\tau > 0$.
- 2: **Initialize:** $\boldsymbol{\pi}^{(0)}, \mathbf{T}^{(0)} = \mathbf{Q}_\tau^{(0)}$.
- 3: **for** $t = 0, 1, \dots$ **do**
- 4: Update the policy for each $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\log \boldsymbol{\pi}^{(t+1)}(a|s) = \left(1 - \frac{\eta\tau}{1-\gamma} \right) \mathbf{W} \log \boldsymbol{\pi}^{(t)}(a|s) + \frac{\eta}{1-\gamma} \mathbf{T}^{(t)}(s, a) - \log \mathbf{z}^{(t)}(s), \quad (24)$$

where $\mathbf{z}^{(t)}(s) = \sum_{a' \in \mathcal{A}} \exp \left\{ \left(1 - \frac{\eta\tau}{1-\gamma} \right) \mathbf{W} \log \boldsymbol{\pi}^{(t)}(a'|s) + \frac{\eta}{1-\gamma} \mathbf{T}^{(t)}(s, a') \right\}$.

- 5: Evaluate $\mathbf{Q}_\tau^{(t+1)}$.
- 6: Update the global Q-function estimate for each $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\mathbf{T}^{(t+1)}(s, a) = \mathbf{W} \left(\mathbf{T}^{(t)}(s, a) + \underbrace{\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)}_{\text{Q-tracking}} \right). \quad (25)$$

7: **end for**

4 Theoretical guarantees

4.1 Global convergence of FedNPG

Convergence with exact policy evaluation. We begin with the global convergence of FedNPG (cf. Algorithm 1), stated in the following theorem. The formal statement and proof of this result can be found in Section A.3.

Theorem 1 (Global sublinear convergence of exact FedNPG (informal)). *Suppose $\pi_n^{(0)}, n \in [N]$ are set as the uniform distribution. Then for $0 < \eta \leq \eta_1 := \frac{(1-\sigma)^2(1-\gamma)^3}{16\sqrt{N}\sigma}$, we have*

$$\frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - V^{\bar{\pi}^{(t)}}(\rho) \right) \leq \frac{V^*(d_\rho^{\pi^*})}{(1-\gamma)T} + \frac{\log |\mathcal{A}|}{\eta T} + \frac{32N\sigma\eta^2}{(1-\gamma)^9(1-\sigma)^2}. \quad (26)$$

Theorem 1 characterizes the average-iterate convergence of the average policy $\bar{\pi}^{(t)}$ (cf. (21)) across the agents, which depends logarithmically on the size of the action space, and independently on the size of the state space. When $T \geq \frac{128\sqrt{N}\log |\mathcal{A}|\sigma^2}{(1-\sigma)^4}$, by optimizing the learning rate $\eta = \left(\frac{(1-\gamma)^9(1-\sigma)^2 \log |\mathcal{A}|}{32TN\sigma} \right)^{1/3}$ to balance the latter two terms, we arrive at

$$\frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - V^{\bar{\pi}^{(t)}}(\rho) \right) \lesssim \frac{V^*(d_\rho^{\pi^*})}{(1-\gamma)T} + \frac{N^{1/3}\sigma^{1/3}}{(1-\gamma)^3(1-\sigma)^{2/3}} \left(\frac{\log |\mathcal{A}|}{T} \right)^{2/3}. \quad (27)$$

A few comments are in order.

- *Server-client setting.* When the network is fully connected, i.e., $\sigma = 0$, the convergence rate of FedNPG recovers the $\mathcal{O}(1/T)$ rate, matching that of the centralized NPG established in Agarwal et al. (2021).
- *Well-connected networks.* When the network is relatively well-connected in the sense of $\frac{\sigma}{(1-\sigma)^2} \lesssim \frac{1-\gamma}{N^{1/2}}$, FedNPG first converges at the rate of $\mathcal{O}(1/T)$, and then at the slower $\mathcal{O}(1/T^{2/3})$ rate after $T \gtrsim \frac{(1-\gamma)^3(1-\sigma)^2}{N\sigma}$.
- *Poorly-connected networks.* In addition, when the network is poorly connected in the sense of $\frac{\sigma}{(1-\sigma)^2} \gtrsim \frac{1-\gamma}{N^{1/2}}$, we see that FedNPG converges at the slower $\mathcal{O}(1/T^{2/3})$ rate.

We state the iteration complexity in Corollary 1.

Corollary 1 (Iteration complexity of exact FedNPG). *To reach*

$$\frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - V^{\bar{\pi}^{(t)}}(\rho) \right) \leq \varepsilon,$$

the iteration complexity of FedNPG is at most $\mathcal{O} \left(\left(\frac{\sigma^{1/2}}{(1-\gamma)^9\sigma^{3/2}(1-\sigma)\varepsilon^{3/2}} + \frac{\sigma^2}{(1-\sigma)^4} \right) \sqrt{N} \log |\mathcal{A}| + \frac{1}{\varepsilon(1-\gamma)^2} \right)$.

Convergence with inexact policy evaluation. In practice, the policies need to be evaluated using samples collected by the agents, where the Q-functions are only estimated approximately. We are interested in gauging how the approximation error impacts the performance of FedNPG, as demonstrated in the following theorem.

Theorem 2 (Global sublinear convergence of inexact FedNPG (informal)). *Suppose that an estimate $q_n^{\pi_n^{(t)}}$ are used in replace of $Q_n^{\pi_n^{(t)}}$ in Algorithm 1. Under the assumptions of Theorem 1, we have*

$$\frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - V^{\bar{\pi}^{(t)}}(\rho) \right) \leq \frac{V^*(d_\rho^{\pi^*})}{(1-\gamma)T} + \frac{\log |\mathcal{A}|}{\eta T} + \frac{32N\sigma\eta^2}{(1-\gamma)^9(1-\sigma)^2} + C_3 \max_{n \in [N], t \in [T]} \left\| Q_n^{\pi_n^{(t)}} - q_n^{\pi_n^{(t)}} \right\|_\infty, \quad (28)$$

where $C_3 := \frac{32\sqrt{N}\sigma\eta}{(1-\gamma)^5(1-\sigma)^2} \left(\frac{\eta\sqrt{N}}{(1-\gamma)^3} + 1 \right) + \frac{2}{(1-\gamma)^2}$.

The formal statement and proof of this result is given in Section A.4.

As long as $\max_{n \in [N], t \in [T]} \|Q_n^{\pi_n^{(t)}} - q_n^{\pi_n^{(t)}}\|_\infty \leq \frac{\varepsilon}{C_3}$, inexact FedNPG reaches $\frac{1}{T} \sum_{t=0}^{T-1} (V^*(\rho) - V^{\bar{\pi}^{(t)}}(\rho)) \leq 2\varepsilon$ at the same iteration complexity as predicted in Corollary 1. Equipped with existing sample complexity bounds on policy evaluation, e.g. using a simulator as in Li et al. (2023b) and Li et al. (2023a), this immediate leads to a sample complexity bound for a federated actor-critic type algorithm for multi-task RL. We detail this in the following remark.

Remark 1 (sample complexity bound of inexact FedNPG). *Recall that Li et al. (2023b) shows that for any fixed policy π , model-based policy evaluation achieves $\|q_\tau^\pi - Q_\tau^\pi\|_\infty \leq \varepsilon_{\text{eval}}$ with high probability if the number of samples per state-action pair exceeds the order of $\tilde{\mathcal{O}}\left(\frac{1}{(1-\gamma)^3 \varepsilon_{\text{eval}}^2}\right)$. When $T \gtrsim \frac{\sqrt{N} \log |\mathcal{A}| \sigma^2}{(1-\sigma)^4}$ and $\eta = \left(\frac{(1-\gamma)^9 (1-\sigma)^2 \log |\mathcal{A}|}{32TN\sigma}\right)^{1/3}$, we have $C_3 \asymp 1/(1-\gamma)^2$. By employing fresh samples for the policy evaluation of each agent at every iteration, we can set $\varepsilon_{\text{eval}} := \max_{n \in [N], t \in [T]} \|Q_n^{\pi_n^{(t)}} - q_n^{\pi_n^{(t)}}\|_\infty \asymp \frac{\varepsilon}{C_3} \asymp (1-\gamma)^2 \varepsilon$, and invoke the union bound over all iterations to give a (very loose) upper bound of sample complexity of FedNPG per state-action pair at each agent as follows:*

$$\begin{aligned} & \tilde{\mathcal{O}}\left(\underbrace{\left(\frac{\sigma^{1/2}}{(1-\gamma)^{9/2}(1-\sigma)\varepsilon^{3/2}} + \frac{\sigma^2}{(1-\sigma)^4}\right)\sqrt{N} + \frac{1}{\varepsilon(1-\gamma)^2}}_{\text{iteration complexity}}\right) \cdot \underbrace{\tilde{\mathcal{O}}\left(\frac{1}{(1-\gamma)^7 \varepsilon^2}\right)}_{\text{sample complexity per iteration}} \\ &= \tilde{\mathcal{O}}\left(\frac{1}{(1-\gamma)^7 \varepsilon^2} \cdot \left[\left(\frac{\sigma^{1/2}}{(1-\gamma)^{9/2}(1-\sigma)\varepsilon^{3/2}} + \frac{\sigma^2}{(1-\sigma)^4}\right)\sqrt{N} + \frac{1}{\varepsilon(1-\gamma)^2}\right]\right). \end{aligned}$$

Hence, the total sample complexity scales linearly with respect to the size of the state-action space up to logarithmic factors. When σ is close to 1, which corresponds to the case where the network exhibits a high degree of locality, the above sample complexity becomes

$$\tilde{\mathcal{O}}\left(\frac{\sqrt{N}}{(1-\gamma)^7 \varepsilon^2} \cdot \left[\left(\frac{1}{(1-\gamma)^{9/2}(1-\sigma)\varepsilon^{3/2}} + \frac{1}{(1-\sigma)^4}\right)\right]\right),$$

which further simplifies to $\tilde{\mathcal{O}}\left(\frac{\sqrt{N}}{(1-\gamma)^{11.5}(1-\sigma)\varepsilon^{3.5}}\right)$ for sufficiently small ε .

4.2 Global convergence of FedNPG with entropy regularization

Convergence with exact policy evaluation. Next, we present our global convergence guarantee of entropy-regularized FedNPG with exact policy evaluation (cf. Algorithm 2).

Theorem 3 (Global linear convergence of exact entropy-regularized FedNPG (informal)). *For any $\gamma \in (0, 1)$ and $0 < \tau \leq 1$, there exists $\eta_0 = \min\left\{\frac{1-\gamma}{\tau}, \mathcal{O}\left(\frac{(1-\gamma)^7 (1-\sigma)^2 \tau}{\sigma N}\right)\right\}$, such that if $0 < \eta \leq \eta_0$, then we have*

$$\|\bar{Q}_\tau^{(t)} - Q_\tau^*\|_\infty \leq 2\gamma C_1 \rho(\eta)^t, \quad \|\log \pi_\tau^* - \log \bar{\pi}^{(t)}\|_\infty \leq \frac{2C_1}{\tau} \rho(\eta)^t, \quad (29)$$

where $\bar{Q}_\tau^{(t)} := Q_{\bar{\pi}^{(t)}}^{(t)}$, $\rho(\eta) \leq \max\{1 - \frac{\tau\eta}{2}, \frac{3+\sigma}{4}\} < 1$, and C_1 is some problem-dependent constant.

The exact expressions of C_1 and η_0 are specified in Appendix A.1. Theorem 3 confirms that entropy-regularized FedNPG converges at a linear rate to the optimal regularized policy, which is almost independent of the size of the state-action space, highlighting the positive role of entropy regularization in federated policy optimization. When the network is fully connected, i.e. $\sigma = 0$, the iteration complexity of entropy-regularized FedNPG reduces to $\mathcal{O}\left(\frac{1}{\eta^\tau} \log \frac{1}{\varepsilon}\right)$, matching that of the centralized entropy-regularized NPG established in Cen et al. (2021). When the network is less connected, one needs to be more conservative in the choice of learning rates, leading to a higher iteration complexity, as described in the following corollary.

Corollary 2 (Iteration complexity of exact entropy-regularized FedNPG). *To reach $\|\log \pi_\tau^* - \log \bar{\pi}^{(t)}\|_\infty \leq \varepsilon$, the iteration complexity of entropy-regularized FedNPG is at most*

$$\tilde{O}\left(\max\left\{\frac{2}{\tau\eta}, \frac{4}{1-\sigma}\right\} \log \frac{1}{\varepsilon}\right) \quad (30)$$

up to logarithmic factors. Especially, when $\eta = \eta_0$, the best iteration complexity becomes

$$\tilde{O}\left(\left(\frac{N\sigma}{(1-\gamma)^7(1-\sigma)^2\tau^2} + \frac{1}{1-\gamma}\right) \log \frac{1}{\tau\varepsilon}\right).$$

Convergence with inexact policy evaluation. Last but not the least, we present the informal convergence results of entropy-regularized FedNPG with inexact policy evaluation, whose formal version can be found in Appendix A.2.

Theorem 4 (Global linear convergence of inexact entropy-regularized FedNPG (informal)). *Suppose that an estimate $q_{\tau,n}^{\pi_n^{(t)}}$ are used in replace of $Q_{\tau,n}^{\pi_n^{(t)}}$ in Algorithm 2. Under the assumptions of Theorem 3, we have*

$$\begin{aligned} \|\bar{Q}_\tau^{(t)} - Q_\tau^*\|_\infty &\leq 2\gamma\left(C_1\rho(\eta)^t + C_2 \max_{n\in[N], t\in[T]} \|Q_{\tau,n}^{\pi_n^{(t)}} - q_{\tau,n}^{\pi_n^{(t)}}\|_\infty\right), \\ \|\log \pi_\tau^* - \log \bar{\pi}^{(t)}\|_\infty &\leq \frac{2}{\tau}\left(C_1\rho(\eta)^t + C_2 \max_{n\in[N], t\in[T]} \|Q_{\tau,n}^{\pi_n^{(t)}} - q_{\tau,n}^{\pi_n^{(t)}}\|_\infty\right), \end{aligned} \quad (31)$$

where $\bar{Q}_\tau^{(t)} := Q_\tau^{\bar{\pi}^{(t)}}$, $\rho(\eta) \leq \max\{1 - \frac{\tau\eta}{2}, \frac{3+\sigma}{4}\} < 1$, and C_1, C_2 are problem-dependent constants.

5 Conclusions

This work proposes the first provably efficient federated NPG (FedNPG) methods for solving vanilla and entropy-regularized multi-task RL problems in the fully decentralized setting. The established finite-time global convergence guarantees are almost independent of the size of the state-action space up to some logarithmic factor, and illuminate the impacts of the size and connectivity of the network. Furthermore, the proposed FedNPG methods are robust vis-a-vis inexactness of local policy evaluations, leading to a finite-sample complexity bound of a federated actor-critic method for multi-task RL. When it comes to future directions, it would be of great interest to further explore sample-efficient algorithms and examine if it is possible to go beyond the entrywise approximation error assumption in policy evaluation. Another interesting direction is to extend the analysis of FedNPG to incorporate function approximations.

Acknowledgments

The work of T. Yang, S. Cen and Y. Chi are supported in part by the grants ONR N00014-19-1-2404, NSF CCF-1901199, CCF-2106778, AFRL FA8750-20-2-0504, and a CMU Cylab seed grant. The work of Y. Wei is supported in part by the the NSF grants DMS-2147546/2015447, CAREER award DMS-2143215, CCF-2106778, and the Google Research Scholar Award. The work of Y. Chen is supported in part by the Alfred P. Sloan Research Fellowship, the Google Research Scholar Award, the AFOSR grant FA9550-22-1-0198, the ONR grant N00014-22-1-2354, and the NSF grants CCF-2221009 and CCF-1907661. S. Cen is also gratefully supported by Wei Shen and Xuehong Zhang Presidential Fellowship, Boeing Scholarship, and JP Morgan Chase PhD Fellowship.

References

Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. (2021). On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *The Journal of Machine Learning Research*, 22(1):4431–4506.

- Ahmed, Z., Le Roux, N., Norouzi, M., and Schuurmans, D. (2019). Understanding the impact of entropy on policy optimization. In *International Conference on Machine Learning*, pages 151–160.
- Amari, S.-I. (1998). Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276.
- Anwar, A. and Raychowdhury, A. (2021). Multi-task federated reinforcement learning with adversaries. *arXiv preprint arXiv:2103.06473*.
- Assran, M., Romoff, J., Ballas, N., Pineau, J., and Rabbat, M. (2019). Gossip-based actor-learner architectures for deep reinforcement learning. *Advances in Neural Information Processing Systems*, 32.
- Bhandari, J. and Russo, D. (2021). On the linear convergence of policy gradient methods for finite MDPs. In *International Conference on Artificial Intelligence and Statistics*, pages 2386–2394. PMLR.
- Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. (2022a). Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 70(4):2563–2578.
- Cen, S., Chi, Y., Du, S. S., and Xiao, L. (2022b). Faster last-iterate convergence of policy optimization in zero-sum Markov games. In *The Eleventh International Conference on Learning Representations*.
- Cen, S., Wei, Y., and Chi, Y. (2021). Fast policy extragradient methods for competitive games with entropy regularization. *Advances in Neural Information Processing Systems*, 34:27952–27964.
- Chen, J., Feng, J., Gao, W., and Wei, K. (2022a). Decentralized natural policy gradient with variance reduction for collaborative multi-agent reinforcement learning. *arXiv preprint arXiv:2209.02179*.
- Chen, T., Zhang, K., Giannakis, G. B., and Başar, T. (2021). Communication-efficient policy gradient methods for distributed reinforcement learning. *IEEE Transactions on Control of Network Systems*, 9(2):917–929.
- Chen, Z., Zhou, Y., and Chen, R.-R. (2022b). Multi-agent off-policy tdc with near-optimal sample and communication complexities.
- Di Lorenzo, P. and Scutari, G. (2016). Next: In-network nonconvex optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2):120–136.
- Duchi, J. C., Agarwal, A., and Wainwright, M. J. (2011). Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606.
- Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., et al. (2018). Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International conference on machine learning*, pages 1407–1416. PMLR.
- Eysenbach, B. and Levine, S. (2021). Maximum entropy RL (provably) solves some robust RL problems. In *International Conference on Learning Representations*.
- Horn, R. A. and Johnson, C. R. (2012). *Matrix analysis*. Cambridge university press.
- Kakade, S. M. (2001). A natural policy gradient. *Advances in neural information processing systems*, 14.
- Kar, S., Moura, J. M., and Poor, H. V. (2012). Qd-learning: A collaborative distributed strategy for multi-agent reinforcement learning through consensus. *arXiv preprint arXiv:1205.0047*.
- Khodadadian, S., Jhunjhunwala, P. R., Varma, S. M., and Maguluri, S. T. (2021). On the linear convergence of natural policy gradient algorithm. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 3794–3799. IEEE.
- Khodadadian, S., Sharma, P., Joshi, G., and Maguluri, S. T. (2022). Federated reinforcement learning: Linear speedup under Markovian sampling. In *International Conference on Machine Learning*, pages 10997–11057. PMLR.

- Lan, G. (2023). Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *Mathematical programming*, 198(1):1059–1106.
- Lan, G., Li, Y., and Zhao, T. (2023). Block policy mirror descent. *SIAM Journal on Optimization*, 33(3):2341–2378.
- Li, B., Cen, S., Chen, Y., and Chi, Y. (2020). Communication-efficient distributed optimization in networks with gradient tracking and variance reduction. *The Journal of Machine Learning Research*, 21(1):7331–7381.
- Li, G., Cai, C., Chen, Y., Wei, Y., and Chi, Y. (2023a). Is q-learning minimax optimal? a tight sample complexity analysis. *Operations Research*.
- Li, G., Wei, Y., Chi, Y., and Chen, Y. (2023b). Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Operations Research*.
- Li, G., Wei, Y., Chi, Y., and Chen, Y. (2023c). Softmax policy gradient methods can take exponential time to converge. *Mathematical Programming*, pages 1–96.
- Lian, X., Zhang, C., Zhang, H., Hsieh, C.-J., Zhang, W., and Liu, J. (2017). Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. *Advances in neural information processing systems*, 30.
- Lobel, I. and Ozdaglar, A. (2008). Convergence analysis of distributed subgradient methods over random networks. In *2008 46th Annual Allerton Conference on Communication, Control, and Computing*, pages 353–360. IEEE.
- M Alshater, M. (2022). Exploring the role of artificial intelligence in enhancing academic performance: A case study of chatgpt. *Available at SSRN*.
- McKelvey, R. D. and Palfrey, T. R. (1995). Quantal response equilibria for normal form games. *Games and economic behavior*, 10(1):6–38.
- Mei, J., Xiao, C., Szepesvari, C., and Schuurmans, D. (2020). On the global convergence rates of softmax policy gradient methods. In *International Conference on Machine Learning*, pages 6820–6829. PMLR.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937.
- Nachum, O., Norouzi, M., Xu, K., and Schuurmans, D. (2017). Bridging the gap between value and policy based reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 2775–2785.
- Nedić, A., Olshevsky, A., and Rabbat, M. G. (2018). Network topology and communication-computation tradeoffs in decentralized optimization. *Proceedings of the IEEE*, 106(5):953–976.
- Nedic, A., Olshevsky, A., and Shi, W. (2017). Achieving geometric convergence for distributed optimization over time-varying graphs. *SIAM Journal on Optimization*, 27(4):2597–2633.
- Nedic, A. and Ozdaglar, A. (2009). Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61.
- Omidshafiei, S., Pazis, J., Amato, C., How, J. P., and Vian, J. (2017). Deep decentralized multi-task multi-agent reinforcement learning under partial observability. In *International Conference on Machine Learning*, pages 2681–2690. PMLR.
- Petersen, K. B. and Pedersen, M. S. (2008). The matrix cookbook. *Technical University of Denmark*, 7(15):510.
- Pu, S. and Nedić, A. (2021). Distributed stochastic gradient tracking methods. *Mathematical Programming*, 187:409–457.

- Puterman, M. L. (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Qi, J., Zhou, Q., Lei, L., and Zheng, K. (2021). Federated reinforcement learning: Techniques, applications, and open challenges. *arXiv preprint arXiv:2108.11887*.
- Qu, G. and Li, N. (2017). Harnessing smoothness to accelerate distributed optimization. *IEEE Transactions on Control of Network Systems*, 5(3):1245–1260.
- Rahman, M. M., Terano, H. J., Rahman, M. N., Salamzadeh, A., and Rahaman, M. S. (2023). Chatgpt and academic research: a review and recommendations based on practical examples. *Rahman, M., Terano, HJR, Rahman, N., Salamzadeh, A., Rahaman, S.(2023). ChatGPT and Academic Research: A Review and Recommendations Based on Practical Examples. Journal of Education, Management and Development Studies*, 3(1):1–12.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. (2015). Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Shani, L., Efroni, Y., and Mannor, S. (2020). Adaptive trust region policy optimization: Global convergence and faster rates for regularized MDPs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5668–5675.
- Wang, H., Kaplan, Z., Niu, D., and Li, B. (2020). Optimizing federated learning on non-iid data with reinforcement learning. In *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*, pages 1698–1707. IEEE.
- Wang, J., Hu, J., Mills, J., Min, G., Xia, M., and Georgalas, N. (2023). Federated ensemble model-based reinforcement learning in edge computing. *IEEE Transactions on Parallel and Distributed Systems*.
- Williams, R. J. and Peng, J. (1991). Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268.
- Woo, J., Joshi, G., and Chi, Y. (2023). The blessing of heterogeneity in federated q-learning: Linear speedup and beyond. *arXiv preprint arXiv:2305.10697*.
- Xiao, L. (2022). On the convergence rates of policy gradient methods. *The Journal of Machine Learning Research*, 23(1):12887–12922.
- Yu, T., Li, T., Sun, Y., Nanda, S., Smith, V., Sekar, V., and Seshan, S. (2020). Learning context-aware policies from multiple smart homes via federated multi-task learning. In *2020 IEEE/ACM Fifth International Conference on Internet-of-Things Design and Implementation (IoTDI)*, pages 104–115. IEEE.
- Zeng, S., Anwar, M. A., Doan, T. T., Raychowdhury, A., and Romberg, J. (2021). A decentralized policy gradient approach to multi-task reinforcement learning. In *Uncertainty in Artificial Intelligence*, pages 1002–1012. PMLR.
- Zerka, F., Barakat, S., Walsh, S., Bogowicz, M., Leijenaar, R. T., Jochems, A., Miraglio, B., Townend, D., and Lambin, P. (2020). Systematic review of privacy-preserving distributed machine learning from federated databases in health care. *JCO clinical cancer informatics*, 4:184–200.
- Zhan, W., Cen, S., Huang, B., Chen, Y., Lee, J. D., and Chi, Y. (2023). Policy mirror descent for regularized reinforcement learning: A generalized framework with linear convergence. *SIAM Journal on Optimization*, 33(2):1061–1091.
- Zhao, F., Ren, X., Yang, S., Zhao, P., Zhang, R., and Xu, X. (2023). Federated multi-objective reinforcement learning. *Information Sciences*, 624:811–832.

Zhou, R., Liu, T., Kalathil, D., Kumar, P., and Tian, C. (2022). Anchor-changing regularized natural policy gradient for multi-objective reinforcement learning. *Advances in Neural Information Processing Systems*, 35:13584–13596.

Zhu, M. and Martínez, S. (2010). Discrete-time dynamic average consensus. *Automatica*, 46(2):322–329.

Zhuo, H. H., Feng, W., Lin, Y., Xu, Q., and Yang, Q. (2019). Federated deep reinforcement learning. *arXiv preprint arXiv:1901.08277*.

A Convergence analysis

For technical convenience, we present first the analysis for entropy-regularized FedNPG and then for vanilla FedNPG.

A.1 Analysis of entropy-regularized FedNPG with exact policy evaluation

To facilitate analysis, we introduce several notation below. For all $t \geq 0$, we recall $\bar{\pi}^{(t)}$ as the normalized geometric mean of $\{\pi_n^{(t)}\}_{n \in [N]}$:

$$\bar{\pi}^{(t)} := \text{softmax} \left(\frac{1}{N} \sum_{n=1}^N \log \pi_n^{(t)} \right), \quad (32)$$

from which we can easily see that for each $(s, a) \in \mathcal{S} \times \mathcal{A}$, $\bar{\pi}^{(t)}(a|s) \propto \left(\prod_{n=1}^N \pi_n^{(t)}(a|s) \right)^{\frac{1}{N}}$. We denote the soft Q -functions of $\bar{\pi}^{(t)}$ by $\bar{Q}_\tau^{(t)}$:

$$\bar{Q}_\tau^{(t)} := \begin{pmatrix} Q_{\tau,1}^{\bar{\pi}^{(t)}} \\ \vdots \\ Q_{\tau,N}^{\bar{\pi}^{(t)}} \end{pmatrix}. \quad (33)$$

In addition, we define $\hat{Q}_\tau^{(t)}, \bar{Q}_\tau^{(t)} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ and $\bar{V}_\tau^{(t)} \in \mathbb{R}^{|\mathcal{S}|}$ as follows

$$\hat{Q}_\tau^{(t)} := \frac{1}{N} \sum_{n=1}^N Q_{\tau,n}^{\pi_n^{(t)}}, \quad (34a)$$

$$\bar{Q}_\tau^{(t)} := Q_\tau^{\bar{\pi}^{(t)}} = \frac{1}{N} \sum_{n=1}^N Q_{\tau,n}^{\bar{\pi}^{(t)}}. \quad (34b)$$

$$\bar{V}_\tau^{(t)} := V_\tau^{\bar{\pi}^{(t)}} = \frac{1}{N} \sum_{n=1}^N V_{\tau,n}^{\bar{\pi}^{(t)}}. \quad (34c)$$

For notational convenience, we also denote

$$\alpha := 1 - \frac{\eta\tau}{1-\gamma}. \quad (35)$$

Following [Cen et al. \(2022a\)](#), we introduce the following auxiliary sequence $\{\boldsymbol{\xi}^{(t)} = (\xi_1^{(t)}, \dots, \xi_N^{(t)})^\top \in \mathbb{R}^{N \times |\mathcal{S}| \times |\mathcal{A}|}\}_{t=0,1,\dots}$, each recursively defined as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad \boldsymbol{\xi}^{(0)}(s, a) := \frac{\|\exp(Q_\tau^*(s, \cdot)/\tau)\|_1}{\left\| \exp\left(\frac{1}{N} \sum_{n=1}^N \log \pi_n^{(0)}(\cdot|s)\right) \right\|_1} \cdot \boldsymbol{\pi}^{(0)}(a|s), \quad (36a)$$

$$\log \boldsymbol{\xi}^{(t+1)}(s, a) = \alpha \mathbf{W} \log \boldsymbol{\xi}^{(t)}(s, a) + (1 - \alpha) \mathbf{T}^{(t)}(s, a)/\tau, \quad (36b)$$

where $\mathbf{T}^{(t)}(s, a)$ is updated via (23). Similarly, we introduce an averaged auxiliary sequence $\{\bar{\xi}^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}\}$ given by

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \quad \bar{\xi}^{(0)}(s, a) := \|\exp(Q_\tau^*(s, \cdot)/\tau)\|_1 \cdot \bar{\pi}^{(0)}(a|s), \quad (37a)$$

$$\log \bar{\xi}^{(t+1)}(s, a) = \alpha \log \bar{\xi}^{(t)}(s, a) + (1 - \alpha) \widehat{Q}_\tau^{(t)}(s, a)/\tau. \quad (37b)$$

We introduces four error metrics defined as

$$\Omega_1^{(t)} := \|u^{(t)}\|_\infty, \quad (38a)$$

$$\Omega_2^{(t)} := \|v^{(t)}\|_\infty, \quad (38b)$$

$$\Omega_3^{(t)} := \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty, \quad (38c)$$

$$\Omega_4^{(t)} := \max \left\{ 0, -\min_{s,a} \left(\widehat{Q}_\tau^{(t)}(s, a) - \tau \log \bar{\xi}^{(t)}(s, a) \right) \right\}, \quad (38d)$$

where $u^{(t)}, v^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ are defined as

$$u^{(t)}(s, a) := \|\log \xi^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N\|_2, \quad (39)$$

$$v^{(t)}(s, a) := \|\mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a) \mathbf{1}_N\|_2. \quad (40)$$

We collect the error metrics above in a vector $\boldsymbol{\Omega}^{(t)} \in \mathbb{R}^4$:

$$\boldsymbol{\Omega}^{(t)} := \left(\Omega_1^{(t)}, \Omega_2^{(t)}, \Omega_3^{(t)}, \Omega_4^{(t)} \right)^\top. \quad (41)$$

With the above preparation, we are ready to state the convergence guarantee of Algorithm 2 in Theorem 5 below, which is the formal version of Theorem 3.

Theorem 5. *For any $N \in \mathbb{N}_+, \tau > 0, \gamma \in (0, 1)$, there exists $\eta_0 > 0$ which depends only on $N, \gamma, \tau, \sigma, |\mathcal{A}|$, such that if $0 < \eta \leq \eta_0$ and $1 - \sigma > 0$, then the updates of Algorithm 2 satisfy*

$$\|\widehat{Q}_\tau^{(t)} - Q_\tau^*\|_\infty \leq 2\gamma \rho(\eta)^t \|\boldsymbol{\Omega}^{(0)}\|_2, \quad (42)$$

$$\|\log \pi_\tau^* - \log \bar{\pi}^{(t)}\|_\infty \leq \frac{2}{\tau} \rho(\eta)^t \|\boldsymbol{\Omega}^{(0)}\|_2, \quad (43)$$

where

$$\rho(\eta) \leq \max \left\{ 1 - \frac{\tau\eta}{2}, \frac{3 + \sigma}{4} \right\} < 1.$$

The dependency of η_0 on $N, \gamma, \tau, \sigma, |\mathcal{A}|$ is made clear in Lemma 2 that will be presented momentarily in this section. The rest of this section is dedicated to the proof of Theorem 5. We first state a key lemma that tracks the error recursion of Algorithm 2.

Lemma 1. *The following linear system holds for all $t \geq 0$:*

$$\boldsymbol{\Omega}^{(t+1)} \leq \underbrace{\begin{pmatrix} \sigma\alpha & \frac{\eta}{1-\gamma} & 0 & 0 \\ S\sigma & \left(1 + \frac{\eta M \sqrt{N}}{1-\gamma}\right) \sigma & \frac{(2+\gamma)\eta MN}{1-\gamma} \sigma & \frac{\gamma \eta MN}{1-\gamma} \sigma \\ (1-\alpha)M & 0 & (1-\alpha)\gamma + \alpha & (1-\alpha)\gamma \\ \frac{2\gamma + \eta\tau}{1-\gamma} M & 0 & 0 & \alpha \end{pmatrix}}_{=: \mathbf{A}(\eta)} \boldsymbol{\Omega}^{(t)}, \quad (44)$$

where we let

$$S := M\sqrt{N} \left(2\alpha + (1-\alpha) \cdot \sqrt{2N} + \frac{1-\alpha}{\tau} \cdot \sqrt{NM} \right), \quad (45)$$

and

$$M := \frac{1 + \gamma + 2\tau(1 - \gamma) \log |\mathcal{A}|}{(1 - \gamma)^2} \cdot \gamma.$$

In addition, it holds for all $t \geq 0$ that

$$\left\| \bar{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty \leq \gamma \Omega_3^{(t)} + \gamma \Omega_4^{(t)}, \quad (46)$$

$$\left\| \log \bar{\pi}^{(t)} - \log \pi_\tau^* \right\|_\infty \leq \frac{2}{\tau} \Omega_3^{(t)}. \quad (47)$$

Proof. See Appendix B.1. □

Let $\rho(\eta)$ denote the spectral norm of $\mathbf{A}(\eta)$. As $\Omega^{(t)} \geq 0$, it is immediate from (44) that

$$\left\| \Omega^{(t)} \right\|_2 \leq \rho(\eta)^t \left\| \Omega^{(0)} \right\|_2,$$

and therefore we have

$$\left\| \bar{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty \leq 2\gamma \left\| \Omega^{(t)} \right\|_\infty \leq 2\gamma \rho(\eta)^t \left\| \Omega^{(0)} \right\|_2,$$

and

$$\left\| \log \bar{\pi}^{(t)} - \log \pi_\tau^* \right\|_\infty \leq \frac{2}{\tau} \left\| \Omega^{(t)} \right\|_\infty \leq \frac{2}{\tau} \rho(\eta)^t \left\| \Omega^{(0)} \right\|_2.$$

It remains to bound the spectral radius $\rho(\eta)$, which is achieved by the following lemma.

Lemma 2 (Bounding the spectral norm of $\mathbf{A}(\eta)$). *Let*

$$\zeta := \frac{(1 - \gamma)(1 - \sigma)^2 \tau}{8(\tau S_0 \sigma + 10M c \sigma / (1 - \gamma) + (1 - \sigma)^2 \tau^2 / 16)}, \quad (48)$$

where $S_0 := M\sqrt{N} \left(2 + \sqrt{2N} + \frac{M\sqrt{N}}{\tau} \right)$. For any $N \in \mathbb{N}_+$, $\tau > 0$, $\gamma \in (0, 1)$, if

$$0 < \eta \leq \eta_0 := \min \left\{ \frac{1 - \gamma}{\tau}, \zeta \right\}, \quad (49)$$

then we have

$$\rho(\eta) \leq \max \left\{ \frac{3 + \sigma}{4}, \frac{1 + (1 - \alpha)\gamma + \alpha}{2} \right\} < 1. \quad (50)$$

Proof. See Appendix B.2. □

A.2 Analysis of entropy-regularized FedNPG with inexact policy evaluation

We define the collection of *inexact* Q-function estimates as

$$\mathbf{q}_\tau^{(t)} := \left(q_{\tau,1}^{\pi_1^{(t)}}, \dots, q_{\tau,N}^{\pi_N^{(t)}} \right)^\top,$$

and then the update rule (25) should be understood as

$$\mathbf{T}^{(t+1)}(s, a) = \mathbf{W} \left(\mathbf{T}^{(t)}(s, a) + \mathbf{q}_\tau^{(t+1)}(s, a) - \mathbf{q}_\tau^{(t)}(s, a) \right) \quad (51)$$

in the inexact setting. For notational simplicity, we define $e_n \in \mathbb{R}$ as

$$e_n := \max_{t \in [T]} \left\| Q_{\tau,n}^{\pi_n^{(t)}} - q_{\tau,n}^{\pi_n^{(t)}} \right\|_\infty, \quad n \in [N], \quad (52)$$

and let $\mathbf{e} = (e_1, \dots, e_n)^\top$. Define $\hat{q}_\tau^{(t)}$, the approximation of $\hat{Q}_\tau^{(t)}$ as

$$\hat{q}_\tau^{(t)} := \frac{1}{N} \sum_{n=1}^N q_{\tau,n}^{\pi_n^{(t)}}. \quad (53)$$

With slight abuse of notation, we adapt the auxiliary sequence $\{\bar{\xi}^{(t)}\}_{t=0,\dots}$ to the inexact updates as

$$\bar{\xi}^{(0)}(s, a) := \|\exp(Q_\tau^*(s, \cdot)/\tau)\|_1 \cdot \bar{\pi}^{(0)}(a|s), \quad (54a)$$

$$\bar{\xi}^{(t+1)}(s, a) := \left[\bar{\xi}^{(t)}(s, a)\right]^\alpha \exp\left((1-\alpha)\frac{\hat{q}_\tau^{(t)}(s, a)}{\tau}\right), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}, t \geq 0. \quad (54b)$$

In addition, we define

$$\Omega_1^{(t)} := \|u^{(t)}\|_\infty, \quad (55a)$$

$$\Omega_2^{(t)} := \|v^{(t)}\|_\infty, \quad (55b)$$

$$\Omega_3^{(t)} := \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty, \quad (55c)$$

$$\Omega_4^{(t)} := \max\left\{0, -\min_{s,a} \left(\bar{q}_\tau^{(t)}(s, a) - \tau \log \bar{\xi}^{(t)}(s, a)\right)\right\}, \quad (55d)$$

where

$$u^{(t)}(s, a) := \left\| \log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N \right\|_2, \quad (56)$$

$$v^{(t)}(s, a) := \left\| \mathbf{T}^{(t)}(s, a) - \hat{q}_\tau^{(t)}(s, a) \mathbf{1}_N \right\|_2. \quad (57)$$

We let $\boldsymbol{\Omega}^{(t)}$ be

$$\boldsymbol{\Omega}^{(t)} := \left(\Omega_1^{(t)}, \Omega_2^{(t)}, \Omega_3^{(t)}, \Omega_4^{(t)}\right)^\top. \quad (58)$$

With the above preparation, we are ready to state the inexact convergence guarantee of Algorithm 2 in Theorem 6 below, which is the formal version of Theorem 4.

Theorem 6. *Suppose that $q_{\tau,n}^{\pi_n^{(t)}}$ are used in replace of $Q_{\tau,n}^{\pi_n^{(t)}}$ in Algorithm 2. For any $N \in \mathbb{N}_+, \tau > 0, \gamma \in (0, 1)$, there exists $\eta_0 > 0$ which depends only on $N, \gamma, \tau, \sigma, |\mathcal{A}|$, such that if $0 < \eta \leq \eta_0$ and $1 - \sigma > 0$, we have*

$$\left\| \bar{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty \leq 2\gamma \left(\rho(\eta)^t \|\boldsymbol{\Omega}^{(0)}\|_2 + C_2 \max_{n \in [N], t \in [T]} \left\| Q_{\tau,n}^{\pi_n^{(t)}} - q_{\tau,n}^{\pi_n^{(t)}} \right\|_\infty \right), \quad (59)$$

$$\left\| \log \pi_\tau^* - \log \bar{\pi}^{(t)} \right\|_\infty \leq \frac{2}{\tau} \left(\rho(\eta)^t \|\boldsymbol{\Omega}^{(0)}\|_2 + C_2 \max_{n \in [N], t \in [T]} \left\| Q_{\tau,n}^{\pi_n^{(t)}} - q_{\tau,n}^{\pi_n^{(t)}} \right\|_\infty \right), \quad (60)$$

where $\rho(\eta) \leq \max\{1 - \frac{\tau\eta}{2}, \frac{3+\sigma}{4}\} < 1$ is the same as in Theorem 5, and $C_2 := \frac{\sigma\sqrt{N}(2(1-\gamma)+M\sqrt{N}\eta)+2\gamma^2+\eta\tau}{(1-\gamma)(1-\rho(\eta))}$.

From Theorem 6, we can conclude that if

$$\max_{n \in [N], t \in [T]} \left\| Q_{\tau,n}^{\pi_n^{(t)}} - q_{\tau,n}^{\pi_n^{(t)}} \right\|_\infty \leq \frac{(1-\gamma)(1-\rho(\eta))\varepsilon}{2\gamma \left(\sigma\sqrt{N}(2(1-\gamma)+M\sqrt{N}\eta) + 2\gamma^2 + \eta\tau \right)}, \quad (61)$$

then inexact entropy-regularized FedNPG could still achieve 2ε -accuracy (i.e. $\left\| \bar{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty \leq 2\varepsilon$) within $\max\left\{\frac{2}{\tau\eta}, \frac{4}{1-\sigma}\right\} \log \frac{2\gamma\|\boldsymbol{\Omega}^{(0)}\|_2}{\varepsilon}$ iterations.

Remark 2. *When $\eta = \eta_0$ (cf. (49) and (48)) and $\tau \leq 1$, the RHS of (61) is of the order*

$$\mathcal{O}\left(\frac{(1-\gamma)\tau\eta_0\varepsilon}{\gamma(\gamma^2 + \sigma\sqrt{N}(1-\gamma))}\right) = \mathcal{O}\left(\frac{(1-\gamma)^8\tau^2(1-\sigma)^2\varepsilon}{\gamma(\gamma^2 + \sigma\sqrt{N}(1-\gamma))(\gamma^2 N\sigma + (1-\sigma)^2\tau^2(1-\gamma)^6)}\right),$$

which can be translated into a crude sample complexity bound when using fresh samples to estimate the soft Q -functions in each iteration.

The rest of this section outlines the proof of Theorem 6. We first state a key lemma that tracks the error recursion of Algorithm 2 with inexact policy evaluation, which is a modified version of Lemma 1.

Lemma 3. *The following linear system holds for all $t \geq 0$:*

$$\mathbf{\Omega}^{(t+1)} \leq \mathbf{A}(\eta)\mathbf{\Omega}^{(t)} + \underbrace{\begin{pmatrix} 0 \\ \sigma\sqrt{N}\left(2 + \frac{M\sqrt{N}\eta}{1-\gamma}\right) \\ \frac{\eta\tau}{1-\gamma} \\ \frac{2\gamma^2}{1-\gamma} \end{pmatrix}}_{=: \mathbf{b}(\eta)} \|e\|_\infty, \quad (62)$$

where $\mathbf{A}(\eta)$ is provided in Lemma 1. In addition, it holds for all $t \geq 0$ that

$$\left\| \overline{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty \leq \gamma\Omega_3^{(t)} + \gamma\Omega_4^{(t)}, \quad (63)$$

$$\left\| \log \overline{\pi}^{(t)} - \log \pi_\tau^* \right\|_\infty \leq \frac{2}{\tau}\Omega_3^{(t)}. \quad (64)$$

Proof. See Appendix B.3. □

By (62), we have

$$\forall t \in N_+ : \quad \mathbf{\Omega}^{(t)} \leq \mathbf{A}(\eta)^t \mathbf{\Omega}^{(0)} + \sum_{s=1}^t \mathbf{A}(\eta)^{t-s} \mathbf{b}(\eta),$$

which gives

$$\begin{aligned} \left\| \mathbf{\Omega}^{(t)} \right\|_2 &\leq \rho(\eta)^t \left\| \mathbf{\Omega}^{(0)} \right\|_2 + \sum_{s=1}^t \rho(\eta)^{t-s} \|\mathbf{b}(\eta)\|_2 \|e\|_\infty \\ &\leq \rho(\eta)^t \left\| \mathbf{\Omega}^{(0)} \right\|_2 + \frac{\sigma\sqrt{N}(2(1-\gamma) + M\sqrt{N}\eta) + 2\gamma^2 + \eta\tau}{(1-\gamma)(1-\rho(\eta))} \|e\|_\infty. \end{aligned} \quad (65)$$

Here, (65) follows from $\|\mathbf{b}(\eta)\|_2 \leq \|\mathbf{b}(\eta)\|_1 = \frac{\sigma\sqrt{N}(2(1-\gamma) + M\sqrt{N}\eta) + 2\gamma^2 + \eta\tau}{1-\gamma} \|e\|_\infty$ and $\sum_{s=1}^t \rho(\eta)^{t-s} \leq 1/(1-\rho(\eta))$. Recall that the bound on $\rho(\eta)$ has already been established in Lemma 2. Therefore we complete the proof of Theorem 6 by combining the above inequality with (63) and (64) in a similar fashion as before. We omit further details for conciseness.

A.3 Analysis of FedNPG with exact policy evaluation

We state the formal version of Theorem 1 below.

Theorem 7. *Suppose all $\pi_n^{(0)}$ in Algorithm 1 are initialized as uniform distribution. When*

$$0 < \eta \leq \eta_1 := \frac{(1-\sigma)^2(1-\gamma)^3}{8(1+\gamma)\gamma\sqrt{N}\sigma},$$

we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - V^{\overline{\pi}^{(t)}}(\rho) \right) \leq \frac{V^*(d_\rho^{\pi^*})}{(1-\gamma)T} + \frac{\log |\mathcal{A}|}{\eta T} + \frac{8(1+\gamma)^2\gamma^2 N\sigma}{(1-\gamma)^9(1-\sigma)^2} \eta^2 \quad (66)$$

for any fixed state distribution ρ .

The rest of this section is dedicated to prove Theorem 7. Similar to (33), we denote the Q -functions of $\overline{\pi}^{(t)}$ by $\overline{Q}^{(t)}$:

$$\overline{Q}^{(t)} := \begin{pmatrix} Q_1^{\overline{\pi}^{(t)}} \\ \vdots \\ Q_N^{\overline{\pi}^{(t)}} \end{pmatrix}. \quad (67)$$

In addition, similar to (34), we define $\widehat{Q}^{(t)}, \overline{Q}^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\overline{V}^{(t)} \in \mathbb{R}^{|\mathcal{S}|}$ as follows

$$\widehat{Q}^{(t)} := \frac{1}{N} \sum_{n=1}^N Q_n^{\pi_n^{(t)}}, \quad (68a)$$

$$\overline{Q}^{(t)} := Q^{\overline{\pi}^{(t)}} = \frac{1}{N} \sum_{n=1}^N Q_n^{\overline{\pi}^{(t)}}. \quad (68b)$$

$$\overline{V}^{(t)} := V^{\overline{\pi}^{(t)}} = \frac{1}{N} \sum_{n=1}^N V_n^{\overline{\pi}^{(t)}}. \quad (68c)$$

Following the same strategy in the analysis of entropy-regularized FedNPG, we introduce the auxiliary sequence $\{\boldsymbol{\xi}^{(t)} = (\xi_1^{(t)}, \dots, \xi_N^{(t)})^\top \in \mathbb{R}^{N \times |\mathcal{S}||\mathcal{A}|}\}$ recursively:

$$\boldsymbol{\xi}^{(0)}(s, a) := \frac{1}{\left\| \exp \left(\frac{1}{N} \sum_{n=1}^N \log \pi_n^{(0)}(\cdot|s) \right) \right\|_1} \cdot \boldsymbol{\pi}^{(0)}(a|s), \quad (69a)$$

$$\log \boldsymbol{\xi}^{(t+1)}(s, a) = \mathbf{W} \log \boldsymbol{\xi}^{(t)}(s, a) + \frac{\eta}{1-\gamma} \mathbf{T}^{(t)}(s, a), \quad (69b)$$

as well as the averaged auxiliary sequence $\{\bar{\xi}^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}\}$:

$$\bar{\xi}^{(0)}(s, a) := \overline{\pi}^{(0)}(a|s), \quad (70a)$$

$$\log \bar{\xi}^{(t+1)}(s, a) := \log \bar{\xi}^{(t)}(s, a) + \frac{\eta}{1-\gamma} \widehat{Q}^{(t)}(s, a), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}, t \geq 0. \quad (70b)$$

As usual, we collect the consensus errors in a vector $\boldsymbol{\Omega}^{(t)} = (\|u^{(t)}\|_\infty, \|v^{(t)}\|_\infty)^\top$, where $u^{(t)}, v^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ are defined as:

$$u^{(t)}(s, a) := \left\| \log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N \right\|_2, \quad (71)$$

$$v^{(t)}(s, a) := \left\| \mathbf{T}^{(t)}(s, a) - \widehat{Q}^{(t)}(s, a) \mathbf{1}_N \right\|_2. \quad (72)$$

Step 1: establishing the error recursion. The next key lemma establishes the error recursion of Algorithm 1.

Lemma 4. *The updates of FedNPG satisfy*

$$\boldsymbol{\Omega}^{(t+1)} \leq \underbrace{\begin{pmatrix} \sigma \\ J\sigma \end{pmatrix} \sigma \left(1 + \frac{\frac{\eta}{1-\gamma}}{\frac{(1+\gamma)\gamma\sqrt{N}\eta}{(1-\gamma)^3}} \right)}_{=: \mathbf{B}(\eta)} \boldsymbol{\Omega}^{(t)} + \underbrace{\begin{pmatrix} 0 \\ \frac{(1+\gamma)\gamma N\sigma}{(1-\gamma)^4} \eta \end{pmatrix}}_{=: \mathbf{d}(\eta)} \quad (73)$$

for all $t \geq 0$, where

$$J := \frac{2(1+\gamma)\gamma}{(1-\gamma)^2} \sqrt{N}. \quad (74)$$

In addition, we have

$$\phi^{(t+1)}(\eta) \leq \phi^{(t)}(\eta) + \frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta \|u^{(t)}\|_\infty - \eta \left(V^*(\rho) - \overline{V}^{(t)}(\rho) \right), \quad (75)$$

where

$$\phi^{(t)}(\eta) := \mathbb{E}_{s \sim d_{\rho^*}} \left[\text{KL}(\pi^*(\cdot|s) \| \overline{\pi}^{(t)}(\cdot|s)) \right] - \frac{\eta}{1-\gamma} \overline{V}^{(t)}(d_{\rho^*}), \quad \forall t \geq 0. \quad (76)$$

Proof. See Appendix B.4. □

Step 2: bounding the value functions. Let $\mathbf{p} \in \mathbb{R}^2$ be defined as:

$$\mathbf{p}(\eta) = \begin{pmatrix} p_1(\eta) \\ p_2(\eta) \end{pmatrix} := \frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \begin{pmatrix} \frac{(1-\gamma)(1-\sigma-(1+\gamma)\gamma\sqrt{N\sigma\eta}/(1-\gamma)^3)\eta}{(1-\gamma)(1-\sigma-(1+\gamma)\gamma\sqrt{N\sigma\eta}/(1-\gamma)^3)(1-\sigma)-J\sigma\eta} \\ \frac{\eta^2}{(1-\gamma)(1-\sigma-(1+\gamma)\gamma\sqrt{N\sigma\eta}/(1-\gamma)^3)(1-\sigma)-J\sigma\eta} \end{pmatrix}; \quad (77)$$

the rationale for this choice will be made clear momentarily. We define the following Lyapunov function

$$\Phi^{(t)}(\eta) = \phi^{(t)}(\eta) + \mathbf{p}(\eta)^\top \boldsymbol{\Omega}^{(t)}, \quad \forall t \geq 0, \quad (78)$$

which satisfies

$$\begin{aligned} \Phi^{(t+1)}(\eta) &= \phi^{(t+1)}(\eta) + \mathbf{p}(\eta)^\top \boldsymbol{\Omega}^{(t+1)} \\ &\leq \phi^{(t)}(\eta) + \frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta \|u^{(t)}\|_\infty - \eta \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right) + \mathbf{p}(\eta)^\top \left(\mathbf{B}(\eta) \boldsymbol{\Omega}^{(t)} + \mathbf{d}(\eta) \right) \\ &= \Phi^{(t)}(\eta) + \left[\mathbf{p}(\eta)^\top (\mathbf{B}(\eta) - \mathbf{I}) + \left(\frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta, 0 \right) \right] \boldsymbol{\Omega}^{(t)} - \eta \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right) \\ &\quad + p_2(\eta) \frac{(1+\gamma)\gamma N \sigma}{(1-\gamma)^4} \eta. \end{aligned} \quad (79)$$

Here, the second inequality follows from (75). One can verify that the second term vanishes due to the choice of $\mathbf{p}(\eta)$:

$$\mathbf{p}(\eta)^\top (\mathbf{B}(\eta) - \mathbf{I}) + \left(\frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta, 0 \right) = (0, 0). \quad (80)$$

Therefore, we conclude that

$$V^*(\rho) - \bar{V}^{(t)}(\rho) \leq \frac{\Phi^{(t)}(\eta) - \Phi^{(t+1)}(\eta)}{\eta} + p_2(\eta) \frac{(1+\gamma)\gamma N \sigma}{(1-\gamma)^4}.$$

Averaging over $t = 0, \dots, T-1$,

$$\begin{aligned} &\frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right) \\ &\leq \frac{\Phi^{(0)}(\eta) - \Phi^{(T)}(\eta)}{\eta T} + \frac{2(1+\gamma)^2 \gamma^2}{(1-\gamma)^8} \cdot \frac{N \sigma \eta^2}{(1-\gamma)(1-\sigma-(1+\gamma)\gamma\sqrt{N\sigma\eta}/(1-\gamma)^3)(1-\sigma)-\sigma J \eta}. \end{aligned} \quad (81)$$

Step 3: simplifying the expression. We first upper bound the first term in the RHS of (81). Assuming uniform initialization for all $\pi_n^{(0)}$ in Algorithm 1, we have $\|u^{(0)}\|_\infty = \|v^{(0)}\|_\infty = 0$, and

$$\mathbb{E}_{s \sim d_{\rho^*}} \left[\text{KL}(\pi^*(\cdot|s) \parallel \bar{\pi}^{(0)}(\cdot|s)) \right] \leq \log |\mathcal{A}|.$$

Therefore, putting together relations (78) and (158) we have

$$\frac{\Phi^{(0)}(\eta) - \Phi^{(T)}(\eta)}{\eta T} \leq \frac{\log |\mathcal{A}|}{T \eta} + \frac{1}{T} \left(\mathbf{p}(\eta)^\top \boldsymbol{\Omega}^{(0)} / \eta + \frac{V^*(d_{\rho^*})}{1-\gamma} \right) = \frac{\log |\mathcal{A}|}{T \eta} + \frac{V^*(d_{\rho^*})}{T(1-\gamma)}, \quad (82)$$

To continue, we upper bound the second term in the RHS of (81). Note that

$$\eta \leq \eta_1 \leq \frac{(1-\sigma)(1-\gamma)^3}{2(1+\gamma)\gamma\sqrt{N\sigma}},$$

which gives

$$\frac{(1+\gamma)\gamma\sqrt{N\sigma}}{(1-\gamma)^3} \eta \leq \frac{1-\sigma}{2}. \quad (83)$$

Thus we have

$$\begin{aligned}
& (1-\gamma)(1-\sigma - (1+\gamma)\gamma\sqrt{N}\sigma\eta/(1-\gamma)^3)(1-\sigma) - J\sigma\eta \\
& \geq (1-\gamma)(1-\sigma)^2/2 - J\sigma\eta_1 \\
& \geq (1-\gamma)(1-\sigma)^2/4,
\end{aligned} \tag{84}$$

where the first inequality follows from (83) and the second inequality follows from the definition of η_1 and J . By (84), we deduce

$$\frac{2(1+\gamma)^2\gamma^2}{(1-\gamma)^8} \cdot \frac{N\sigma\eta^2}{(1-\gamma)(1-\sigma - (1+\gamma)\gamma\sqrt{N}\sigma\eta/(1-\gamma)^3)(1-\sigma) - J\sigma\eta} \leq \frac{8(1+\gamma)^2\gamma^2N\sigma}{(1-\gamma)^9(1-\sigma)^2}\eta^2, \tag{85}$$

and our advertised bound (66) thus follows from plugging (82) and (85) into (81).

A.4 Analysis of FedNPG with inexact policy evaluation

We state the formal version of Theorem 2 below.

Theorem 8. *Suppose that $q_n^{\pi_n^{(t)}}$ are used in replace of $Q_n^{\pi_n^{(t)}}$ in Algorithm 1. Suppose all $\pi_n^{(0)}$ in Algorithm 1 set to uniform distribution. Let*

$$0 < \eta \leq \eta_1 := \frac{(1-\sigma)^2(1-\gamma)^3}{8(1+\gamma)\gamma\sqrt{N}\sigma},$$

we have

$$\begin{aligned}
& \frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - V^{\bar{\pi}^{(t)}}(\rho) \right) \\
& \leq \frac{V^*(d_\rho^*)}{(1-\gamma)T} + \frac{\log |\mathcal{A}|}{\eta T} + \frac{8(1+\gamma)^2\gamma^2N\sigma}{(1-\gamma)^9(1-\sigma)^2}\eta^2 \\
& \quad + \left[\frac{8(1+\gamma)\gamma}{(1-\gamma)^5(1-\sigma)^2}\sqrt{N}\sigma\eta \left(\frac{(1+\gamma)\gamma\eta\sqrt{N}}{(1-\gamma)^3} + 2 \right) + \frac{2}{(1-\gamma)^2} \right] \max_{n \in [N], t \in [T]} \left\| Q_n^{\pi_n^{(t)}} - q_n^{\pi_n^{(t)}} \right\|_\infty
\end{aligned}$$

for any fixed state distribution ρ .

We next outline the proof of Theorem 8. With slight abuse of notation, we again define $e_n \in \mathbb{R}$ as

$$e_n := \max_{t \in [T]} \left\| Q_n^{\pi_n^{(t)}} - q_n^{\pi_n^{(t)}} \right\|_\infty, \quad n \in [N], \tag{86}$$

and let $\mathbf{e} = (e_1, \dots, e_N)^\top$. We define the collection of *inexact* Q-function estimates as

$$\mathbf{q}^{(t)} := \left(q_1^{\pi_1^{(t)}}, \dots, q_N^{\pi_N^{(t)}} \right)^\top,$$

and then the update rule (23) should be understood as

$$\mathbf{T}^{(t+1)}(s, a) = \mathbf{W} \left(\mathbf{T}^{(t)}(s, a) + \mathbf{q}^{(t+1)}(s, a) - \mathbf{q}^{(t)}(s, a) \right) \tag{87}$$

in the inexact setting. Define $\hat{q}^{(t)}$, the approximation of $\hat{Q}^{(t)}$ as

$$\hat{q}^{(t)} := \frac{1}{N} \sum_{n=1}^N q_n^{\pi_n^{(t)}}, \tag{88}$$

we adapt the averaged auxiliary sequence $\{\bar{\xi}^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}\}$ to the inexact updates as follows:

$$\bar{\xi}^{(0)}(s, a) := \bar{\pi}^{(0)}(a|s), \tag{89a}$$

$$\bar{\xi}^{(t+1)}(s, a) := \bar{\xi}^{(t)}(s, a) \exp\left(\frac{\eta}{1-\gamma} \hat{q}^{(t)}(s, a)\right), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}, t \geq 0. \quad (89b)$$

As usual, we define the consensus error vector as $\boldsymbol{\Omega}^{(t)} = (\|u^{(t)}\|_\infty, \|v^{(t)}\|_\infty)^\top$, where $u^{(t)}, v^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ are given by

$$u^{(t)}(s, a) := \left\| \log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N \right\|_2, \quad (90)$$

$$v^{(t)}(s, a) := \left\| \mathbf{T}^{(t)}(s, a) - \hat{q}^{(t)}(s, a) \mathbf{1}_N \right\|_2. \quad (91)$$

The following lemma characterizes the dynamics of the error vector $\boldsymbol{\Omega}^{(t)}$, perturbed by additional approximation error.

Lemma 5. *The updates of inexact FedNPG satisfy*

$$\boldsymbol{\Omega}^{(t+1)} \leq \mathbf{B}(\eta) \boldsymbol{\Omega}^{(t)} + \mathbf{d}(\eta) + \underbrace{\left(\sqrt{N} \sigma \left(\frac{0}{(1-\gamma)^3} + 2 \right) \right)}_{=: \mathbf{c}(\eta)} \|\mathbf{e}\|_\infty. \quad (92)$$

In addition, we have

$$\phi^{(t+1)}(\eta) \leq \phi^{(t)}(\eta) + \frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta \|u^{(t)}\|_\infty + \frac{2\eta}{(1-\gamma)^2} \|\mathbf{e}\|_\infty - \eta \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right), \quad (93)$$

where $\phi^{(t)}(\eta)$ is defined in (76).

Proof. See Appendix B.5. □

Similar to (79), we can recursively bound $\Phi^{(t)}(\eta)$ (defined in (78)) as

$$\begin{aligned} \Phi^{(t+1)}(\eta) &= \phi^{(t+1)}(\eta) + \mathbf{p}(\eta)^\top \boldsymbol{\Omega}^{(t+1)} \\ &\stackrel{(93)}{\leq} \phi^{(t)}(\eta) + \frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta \|u^{(t)}\|_\infty + \frac{2\eta}{(1-\gamma)^2} \|\mathbf{e}\|_\infty - \eta \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right) \\ &\quad + \mathbf{p}(\eta)^\top \left(\mathbf{B}(\eta) \boldsymbol{\Omega}^{(t)} + \mathbf{d}(\eta) + \mathbf{c}(\eta) \right) \\ &= \Phi^{(t)}(\eta) + \underbrace{\left[\mathbf{p}(\eta)^\top (\mathbf{B}(\eta) - \mathbf{I}) + \left(\frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta, 0 \right) \right]}_{=(0,0) \text{ via (80)}} \boldsymbol{\Omega}^{(t)} - \eta \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right) \\ &\quad + p_2(\eta) \frac{(1+\gamma)\gamma N \sigma}{(1-\gamma)^4} \eta + \left[p_2(\eta) \sqrt{N} \sigma \left(\frac{(1+\gamma)\gamma \eta \sqrt{N}}{(1-\gamma)^3} + 2 \right) + \frac{2\eta}{(1-\gamma)^2} \right] \|\mathbf{e}\|_\infty. \quad (94) \end{aligned}$$

From the above expression we know that

$$V^*(\rho) - \bar{V}^{(t)}(\rho) \leq \frac{\Phi^{(t)}(\eta) - \Phi^{(t+1)}(\eta)}{\eta} + p_2(\eta) \frac{(1+\gamma)\gamma N \sigma}{(1-\gamma)^4} + \left[p_2(\eta) \sqrt{N} \sigma \left(\frac{(1+\gamma)\gamma \eta \sqrt{N}}{(1-\gamma)^3} + \frac{2}{\eta} \right) + \frac{2}{(1-\gamma)^2} \right] \|\mathbf{e}\|_\infty,$$

which gives

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right) &\leq \frac{\Phi^{(0)}(\eta) - \Phi^{(T)}(\eta)}{\eta T} + p_2(\eta) \frac{(1+\gamma)\gamma N \sigma}{(1-\gamma)^4} \\ &\quad + \left[p_2(\eta) \sqrt{N} \sigma \left(\frac{(1+\gamma)\gamma \eta \sqrt{N}}{(1-\gamma)^3} + \frac{2}{\eta} \right) + \frac{2}{(1-\gamma)^2} \right] \|\mathbf{e}\|_\infty \quad (95) \end{aligned}$$

via telescoping. Combining the above expression with (82), (84) and (85), we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \left(V^*(\rho) - \bar{V}^{(t)}(\rho) \right) &\leq \frac{\log |\mathcal{A}|}{T\eta} + \frac{V^*(d_\rho^{\pi^*})}{T(1-\gamma)} + \frac{8(1+\gamma)^2\gamma^2 N\sigma}{(1-\gamma)^9(1-\sigma)^2} \eta^2 \\ &\quad + \left[\frac{8(1+\gamma)\gamma}{(1-\gamma)^5(1-\sigma)^2} \sqrt{N}\sigma\eta \left(\frac{(1+\gamma)\gamma\eta\sqrt{N}}{(1-\gamma)^3} + 2 \right) + \frac{2}{(1-\gamma)^2} \right] \|e\|_\infty, \end{aligned} \quad (96)$$

which establishes (86).

B Proof of key lemmas

B.1 Proof of Lemma 1

Before proceeding, we summarize several useful properties of the auxiliary sequences (cf. (36) and (37)), whose proof is postponed to Appendix C.1.

Lemma 6 (Properties of auxiliary sequences $\{\bar{\xi}^{(t)}\}$ and $\{\xi^{(t)}\}$). *$\{\bar{\xi}^{(t)}\}$ and $\{\xi^{(t)}\}$ have the following properties:*

1. $\xi^{(t)}$ can be viewed as an unnormalized version of $\pi^{(t)}$, i.e.,

$$\pi_n^{(t)}(\cdot|s) = \frac{\xi_n^{(t)}(s, \cdot)}{\|\xi_n^{(t)}(s, \cdot)\|_1}, \quad \forall n \in [N], s \in \mathcal{S}. \quad (97)$$

2. For any $t \geq 0$, $\log \bar{\xi}^{(t)}$ keeps track of the average of $\log \xi^{(t)}$, i.e.,

$$\frac{1}{N} \mathbf{1}_N^\top \log \xi^{(t)} = \log \bar{\xi}^{(t)}. \quad (98)$$

It follows that

$$\forall s \in \mathcal{S}, t \geq 0: \quad \bar{\pi}^{(t)}(\cdot|s) = \frac{\bar{\xi}^{(t)}(s, \cdot)}{\|\bar{\xi}^{(t)}(s, \cdot)\|_1}. \quad (99)$$

Lemma 7 ((Cen et al., 2022a, Appendix. A.2)). *For any vector $\theta = [\theta_a]_{a \in \mathcal{A}} \in \mathbb{R}^{|\mathcal{A}|}$, we denote by $\pi_\theta \in \mathbb{R}^{|\mathcal{A}|}$ the softmax transform of θ such that*

$$\pi_\theta(a) = \frac{\exp(\theta_a)}{\sum_{a' \in \mathcal{A}} \exp(\theta_{a'})}, \quad a \in \mathcal{A}. \quad (100)$$

For any $\theta_1, \theta_2 \in \mathbb{R}^{|\mathcal{A}|}$, we have

$$\left| \log(\|\exp(\theta_1)\|_1) - \log(\|\exp(\theta_2)\|_1) \right| \leq \|\theta_1 - \theta_2\|_\infty, \quad (101)$$

$$\|\log \pi_{\theta_1} - \log \pi_{\theta_2}\|_\infty \leq 2 \|\theta_1 - \theta_2\|_\infty. \quad (102)$$

Step 1: bound $u^{(t+1)}(s, a) = \|\log \xi^{(t+1)}(s, a) - \log \bar{\xi}^{(t+1)}(s, a) \mathbf{1}_N\|_2$. By (36b) and (37b) we have

$$\begin{aligned} u^{(t+1)}(s, a) &= \|\log \xi^{(t+1)}(s, a) - \log \bar{\xi}^{(t+1)}(s, a) \mathbf{1}_N\|_2 \\ &= \left\| \alpha \left(\mathbf{W} \log \xi^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N \right) + (1-\alpha) \left(\mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a) \mathbf{1}_N \right) / \tau \right\|_2 \\ &\leq \sigma\alpha \|\log \xi^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N\|_2 + \frac{1-\alpha}{\tau} \|\mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a) \mathbf{1}_N\|_2 \\ &\leq \sigma\alpha \|u^{(t)}\|_\infty + \frac{1-\alpha}{\tau} \|v^{(t)}\|_\infty, \end{aligned} \quad (103)$$

where the penultimate step results from the averaging property of \mathbf{W} (property (18)). Taking maximum over $(s, a) \in \mathcal{S} \times \mathcal{A}$ establishes the bound on $\Omega_1^{(t+1)}$ in (44).

Step 2: bound $v^{(t+1)}(s, a) = \|\mathbf{T}^{(t+1)}(s, a) - \widehat{Q}_\tau^{(t+1)}(s, a)\mathbf{1}_N\|_2$. By (25) we have

$$\begin{aligned}
& \|\mathbf{T}^{(t+1)}(s, a) - \widehat{Q}_\tau^{(t+1)}(s, a)\mathbf{1}_N\|_2 \\
&= \left\| \mathbf{W} \left(\mathbf{T}^{(t)}(s, a) + \mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a) \right) - \widehat{Q}_\tau^{(t+1)}(s, a)\mathbf{1}_N \right\|_2 \\
&= \left\| \left(\mathbf{W}\mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a)\mathbf{1}_N \right) + \mathbf{W} \left(\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a) \right) + \left(\widehat{Q}_\tau^{(t)}(s, a) - \widehat{Q}_\tau^{(t+1)}(s, a) \right) \mathbf{1}_N \right\|_2 \\
&\leq \sigma \|\mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a)\mathbf{1}_N\|_2 + \sigma \left\| \left(\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a) \right) + \left(\widehat{Q}_\tau^{(t)}(s, a) - \widehat{Q}_\tau^{(t+1)}(s, a) \right) \mathbf{1}_N \right\|_2 \\
&\leq \sigma \|\mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a)\mathbf{1}_N\|_2 + \sigma \|\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)\|_2, \tag{104}
\end{aligned}$$

where the penultimate step uses property (18), and the last step is due to

$$\begin{aligned}
& \left\| \left(\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a) \right) + \left(\widehat{Q}_\tau^{(t)}(s, a) - \widehat{Q}_\tau^{(t+1)}(s, a) \right) \mathbf{1}_N \right\|_2^2 \\
&= \|\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)\|_2^2 + N(\widehat{Q}_\tau^{(t)}(s, a) - \widehat{Q}_\tau^{(t+1)}(s, a))^2 \\
&\quad - 2 \sum_{n=1}^N \left(Q_{\tau, n}^{\pi_n^{(t+1)}}(s, a) - Q_{\tau, n}^{\pi_n^{(t)}}(s, a) \right) \left(\widehat{Q}_\tau^{(t+1)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a) \right) \\
&= \|\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)\|_2^2 - N(\widehat{Q}_\tau^{(t)}(s, a) - \widehat{Q}_\tau^{(t+1)}(s, a))^2 \\
&\leq \|\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)\|_2^2.
\end{aligned}$$

Step 3: bound $\|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty$. We decompose the term of interest as

$$\begin{aligned}
Q_\tau^* - \tau \log \bar{\xi}^{(t+1)} &= Q_\tau^* - \tau \alpha \log \bar{\xi}^{(t)} - (1 - \alpha) \widehat{Q}_\tau^{(t)} \\
&= \alpha(Q_\tau^* - \tau \log \bar{\xi}^{(t)}) + (1 - \alpha)(Q_\tau^* - \bar{Q}_\tau^{(t)}) + (1 - \alpha)(\bar{Q}_\tau^{(t)} - \widehat{Q}_\tau^{(t)}),
\end{aligned}$$

which gives

$$\|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty \leq \alpha \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty + (1 - \alpha) \|Q_\tau^* - \bar{Q}_\tau^{(t)}\|_\infty + (1 - \alpha) \|\bar{Q}_\tau^{(t)} - \widehat{Q}_\tau^{(t)}\|_\infty. \tag{105}$$

Note that we can upper bound $\|\bar{Q}_\tau^{(t)} - \widehat{Q}_\tau^{(t)}\|_\infty$ by

$$\begin{aligned}
\|\bar{Q}_\tau^{(t)} - \widehat{Q}_\tau^{(t)}\|_\infty &= \left\| \frac{1}{N} \sum_{n=1}^N Q_{\tau, n}^{\pi_n^{(t)}} - \frac{1}{N} \sum_{n=1}^N Q_{\tau, n}^{\bar{\pi}_n^{(t)}} \right\|_\infty \\
&\leq \frac{1}{N} \sum_{n=1}^N \|Q_{\tau, n}^{\pi_n^{(t)}} - Q_{\tau, n}^{\bar{\pi}_n^{(t)}}\|_\infty \\
&\leq \frac{M}{N} \sum_{n=1}^N \|\log \xi_n^{(t)} - \log \bar{\xi}^{(t)}\|_\infty \leq M \|u^{(t)}\|_\infty. \tag{106}
\end{aligned}$$

The last step is due to $|\log \xi_n^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a)| \leq u^{(t)}(s, a)$, while the penultimate step results from writing

$$\begin{aligned}
\bar{\pi}^{(t)}(\cdot | s) &= \text{softmax} \left(\log \bar{\xi}^{(t)}(s, \cdot) \right), \\
\pi_n^{(t)}(\cdot | s) &= \text{softmax} \left(\log \xi_n^{(t)}(s, \cdot) \right),
\end{aligned}$$

and applying the following lemma.

Lemma 8 (Lipschitz constant of soft Q-function). *Assume that $r(s, a) \in [0, 1], \forall (s, a) \in \mathcal{S} \times \mathcal{A}$ and $\tau \geq 0$. For any $\theta, \theta' \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$, we have*

$$\|Q_\tau^{\pi_{\theta'}} - Q_\tau^{\pi_\theta}\|_\infty \leq \underbrace{\frac{1 + \gamma + 2\tau(1 - \gamma) \log |\mathcal{A}|}{(1 - \gamma)^2}}_{=: M} \cdot \gamma \|\theta' - \theta\|_\infty. \quad (107)$$

Plugging (106) into (105) gives

$$\|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty \leq \alpha \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty + (1 - \alpha) \|Q_\tau^* - \bar{Q}_\tau^{(t)}\|_\infty + (1 - \alpha) M \|u^{(t)}\|_\infty. \quad (108)$$

Step 4: bound $\|Q_\tau^{(t+1)}(s, a) - Q_\tau^{(t)}(s, a)\|_2$. Let $w^{(t)} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be defined as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad w^{(t)}(s, a) := \left\| \log \xi^{(t+1)}(s, a) - \log \xi^{(t)}(s, a) - (1 - \alpha) V_\tau^*(s) \mathbf{1}_N / \tau \right\|_2. \quad (109)$$

Again, we treat $w^{(t)}$ as vectors in $\mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ whenever it is clear from context. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $n \in [N]$, by Lemma 8 it follows that

$$\begin{aligned} \left| Q_{\tau, n}^{\pi_n^{(t+1)}}(s, a) - Q_{\tau, n}^{\pi_n^{(t)}}(s, a) \right| &\leq M \max_{s \in \mathcal{S}} \left\| \log \xi_n^{(t+1)}(s, \cdot) - \log \xi_n^{(t)}(s, \cdot) - (1 - \alpha) V_\tau^*(s) \mathbf{1}_{|\mathcal{A}|} / \tau \right\|_\infty \\ &\leq M \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} w^{(t)}(s, a) \leq M \|w^{(t)}\|_\infty, \end{aligned} \quad (110)$$

and consequently

$$\|Q_\tau^{(t+1)}(s, a) - Q_\tau^{(t)}(s, a)\|_2 \leq M \sqrt{N} \|w^{(t)}\|_\infty. \quad (111)$$

It boils down to control $\|w^{(t)}\|_\infty$. To do so, we first note that for each $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have

$$\begin{aligned} w^{(t)}(s, a) &= \left\| \alpha \mathbf{W} \log \xi^{(t)}(s, a) + (1 - \alpha) \mathbf{T}^{(t)}(s, a) / \tau - \log \xi^{(t)}(s, a) - (1 - \alpha) V_\tau^*(s) \mathbf{1}_N / \tau \right\|_2 \\ &\stackrel{(a)}{=} \left\| \alpha (\mathbf{W} - \mathbf{I}_N) \left(\log \xi^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N \right) + (1 - \alpha) \left(\mathbf{T}^{(t)}(s, a) / \tau - \log \xi^{(t)}(s, a) - V_\tau^*(s) \mathbf{1}_N / \tau \right) \right\|_2 \\ &\stackrel{(b)}{\leq} 2\alpha \left\| \log \xi^{(t)}(s, a) - \log \bar{\xi}^{(t)}(s, a) \mathbf{1}_N \right\|_2 + \frac{1 - \alpha}{\tau} \left\| \mathbf{T}^{(t)}(s, a) - \tau \log \xi^{(t)}(s, a) - V_\tau^*(s) \mathbf{1}_N \right\|_2 \end{aligned} \quad (112)$$

where (a) is due to the doubly stochasticity property of \mathbf{W} and (b) is from the fact $\|\mathbf{W} - \mathbf{I}_N\|_2 \leq 2$. We further bound the second term as follows:

$$\begin{aligned} &\left\| \mathbf{T}^{(t)}(s, a) - \tau \log \xi^{(t)}(s, a) - V_\tau^*(s) \mathbf{1}_N \right\|_2 \\ &= \left\| \mathbf{T}^{(t)}(s, a) - \tau \log \xi^{(t)}(s, a) - (Q_\tau^*(s, a) - \tau \log \pi_\tau^*(a|s)) \mathbf{1}_N \right\|_2 \\ &\leq \left\| \mathbf{T}^{(t)}(s, a) - Q_\tau^*(s, a) \mathbf{1}_N \right\|_2 + \tau \left\| \log \xi^{(t)}(s, a) - \log \pi_\tau^*(a|s) \mathbf{1}_N \right\|_2 \\ &\leq \left\| \mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau(s, a) \mathbf{1}_N \right\|_2 + \left\| \widehat{Q}_\tau(s, a) \mathbf{1}_N - Q_\tau^*(s, a) \mathbf{1}_N \right\|_2 \\ &\quad + \tau \left\| \log \xi^{(t)}(s, a) - \log \bar{\pi}^{(t)}(a|s) \mathbf{1}_N \right\|_2 + \tau \left\| \log \bar{\pi}^{(t)}(a|s) \mathbf{1}_N - \log \pi_\tau^*(a|s) \mathbf{1}_N \right\|_2 \\ &= \left\| \mathbf{T}^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a) \mathbf{1}_N \right\|_2 + \sqrt{N} \left| \widehat{Q}_\tau^{(t)}(s, a) - Q_\tau^*(s, a) \right| \\ &\quad + \tau \left\| \log \xi^{(t)}(s, a) - \log \bar{\pi}^{(t)}(a|s) \mathbf{1}_N \right\|_2 + \tau \sqrt{N} \left| \log \bar{\pi}^{(t)}(a|s) - \log \pi_\tau^*(a|s) \right|. \end{aligned} \quad (113)$$

Here, the first step results from the following relation established in Nachum et al. (2017):

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad V_\tau^*(s) = -\tau \log \pi_\tau^*(a|s) + Q_\tau^*(s, a), \quad (114)$$

which also leads to

$$\left\| \log \bar{\pi}^{(t)} - \log \pi_\tau^* \right\|_\infty \leq \frac{2}{\tau} \left\| Q_\tau^* - \tau \log \bar{\xi}^{(t)} \right\|_\infty \quad (115)$$

by Lemma 7. For the remaining terms in (113), we have

$$\left| \widehat{Q}_\tau^{(t)}(s, a) - Q_\tau^*(s, a) \right| \leq \left\| \widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)} \right\|_\infty + \left\| \overline{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty, \quad (116)$$

and

$$\begin{aligned} \left\| \log \boldsymbol{\xi}^{(t)}(s, a) - \log \overline{\pi}^{(t)}(a|s) \mathbf{1}_N \right\|_2 &= \sqrt{\sum_{n=1}^N \left(\log \xi_n^{(t)}(s, a) - \log \overline{\xi}^{(t)}(a|s) \right)^2} \\ &\leq \sqrt{\sum_{n=1}^N 2 \left\| \log \xi_n^{(t)} - \log \overline{\xi}^{(t)} \right\|_\infty^2} \\ &\leq \sqrt{\sum_{n=1}^N 2 \left\| u^{(t)} \right\|_\infty^2} = \sqrt{2N} \left\| u^{(t)} \right\|_\infty, \end{aligned} \quad (117)$$

where the first inequality again results from Lemma 7. Plugging (115), (116), (117) into (113) and using the definition of $u^{(t)}, v^{(t)}$, we arrive at

$$\begin{aligned} w^{(t)}(s, a) &\leq \left(2\alpha + (1 - \alpha) \cdot \sqrt{2N} \right) \left\| u^{(t)} \right\|_\infty + \frac{1 - \alpha}{\tau} \left\| v^{(t)} \right\|_\infty + \frac{1 - \alpha}{\tau} \cdot \sqrt{N} \left(\left\| \widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)} \right\|_\infty + \left\| \overline{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty \right) \\ &\quad + \frac{1 - \alpha}{\tau} \cdot 2\sqrt{N} \left\| Q_\tau^* - \tau \log \overline{\xi}^{(t)} \right\|_\infty. \end{aligned}$$

Using previous display, we can write (111) as

$$\begin{aligned} &\left\| \mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a) \right\|_2 \\ &\leq M\sqrt{N} \left\{ \left(2\alpha + (1 - \alpha) \cdot \sqrt{2N} \right) \left\| u^{(t)} \right\|_\infty + \frac{1 - \alpha}{\tau} \left\| v^{(t)} \right\|_\infty \right. \\ &\quad \left. + \frac{1 - \alpha}{\tau} \cdot \sqrt{N} \left(M \left\| u^{(t)} \right\|_\infty + \left\| \overline{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty \right) + \frac{1 - \alpha}{\tau} \cdot 2\sqrt{N} \left\| Q_\tau^* - \tau \log \overline{\xi}^{(t)} \right\|_\infty \right\}. \end{aligned} \quad (118)$$

Combining (104) with the above expression (118), we get

$$\begin{aligned} \left\| v^{(t+1)} \right\|_\infty &\leq \sigma \left(1 + \frac{\eta M \sqrt{N}}{1 - \gamma} \right) \left\| v^{(t)} \right\|_\infty + \sigma M \sqrt{N} \left\{ \left(2\alpha + (1 - \alpha) \cdot \sqrt{2N} + \frac{1 - \alpha}{\tau} \cdot \sqrt{NM} \right) \left\| u^{(t)} \right\|_\infty \right. \\ &\quad \left. + \frac{1 - \alpha}{\tau} \cdot \sqrt{N} \left\| \overline{Q}_\tau^{(t)} - Q_\tau^* \right\|_\infty + \frac{1 - \alpha}{\tau} \cdot 2\sqrt{N} \left\| Q_\tau^* - \tau \log \overline{\xi}^{(t)} \right\|_\infty \right\}. \end{aligned} \quad (119)$$

Step 5: bound $\left\| \overline{Q}_\tau^{(t+1)} - Q_\tau^* \right\|_\infty$. For any state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, we observe that

$$\begin{aligned} &Q_\tau^*(s, a) - \overline{Q}_\tau^{(t+1)}(s, a) \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V_\tau^*(s')] - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V_\tau^{\overline{\pi}^{(t+1)}}(s')] \right) \\ &= \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\tau \log \left(\left\| \exp \left(\frac{Q_\tau^*(s', \cdot)}{\tau} \right) \right\|_1 \right) \right] - \gamma \mathbb{E}_{\substack{s' \sim P(\cdot|s, a), \\ a' \sim \overline{\pi}^{(t+1)}(\cdot|s')}} [\overline{Q}_\tau^{(t+1)}(s', a') - \tau \log \overline{\pi}^{(t+1)}(a'|s')], \end{aligned} \quad (120)$$

where the first step invokes the definition of Q_τ^* (cf. (7a)), and the second step is due to the following expression of V_τ^* established in Nachum et al. (2017):

$$V_\tau^*(s) = \tau \log \left(\left\| \exp \left(\frac{Q_\tau^*(s, \cdot)}{\tau} \right) \right\|_1 \right). \quad (121)$$

To continue, note that by (99) and (37b) we have

$$\begin{aligned}\log \bar{\pi}^{(t+1)}(a|s) &= \log \bar{\xi}^{(t+1)}(s, a) - \log \left(\|\bar{\xi}^{(t+1)}(s, \cdot)\|_1 \right) \\ &= \alpha \log \bar{\xi}^{(t)}(s, a) + (1 - \alpha) \frac{\widehat{Q}_\tau^{(t)}(s, a)}{\tau} - \log \left(\|\bar{\xi}^{(t+1)}(s, \cdot)\|_1 \right).\end{aligned}\quad (122)$$

Plugging (122) into (120) and (118) establishes the bounds on

$$\begin{aligned}Q_\tau^*(s, a) - \bar{Q}_\tau^{(t+1)}(s, a) &= \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\tau \log \left(\left\| \exp \left(\frac{Q_\tau^*(s', \cdot)}{\tau} \right) \right\|_1 \right) - \tau \log \left(\|\bar{\xi}^{(t+1)}(s', \cdot)\|_1 \right) \right] \\ &\quad - \gamma \mathbb{E}_{\substack{s' \sim P(\cdot|s, a), \\ a' \sim \bar{\pi}^{(t+1)}(\cdot|s')}} \left[\underbrace{\bar{Q}_\tau^{(t+1)}(s', a') - \tau \left(\alpha \log \bar{\xi}^{(t)}(s', a') + (1 - \alpha) \frac{\widehat{Q}_\tau^{(t)}(s', a')}{\tau} \right)}_{= \log \bar{\xi}^{(t+1)}(s', a')} \right]\end{aligned}\quad (123)$$

for any $(s, a) \in \mathcal{S} \times \mathcal{A}$. In view of property (101), the first term on the right-hand side of (123) can be bounded by

$$\tau \log \left(\left\| \exp \left(\frac{Q_\tau^*(s', \cdot)}{\tau} \right) \right\|_1 \right) - \tau \log \left(\|\bar{\xi}^{(t+1)}(s', \cdot)\|_1 \right) \leq \|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty.$$

Plugging the above expression into (123), we have

$$0 \leq Q_\tau^*(s, a) - \bar{Q}_\tau^{(t+1)}(s, a) \leq \gamma \|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty - \gamma \min_{s, a} \left(\bar{Q}_\tau^{(t+1)}(s, a) - \tau \log \bar{\xi}^{(t+1)}(s, a) \right),$$

which gives

$$\|Q_\tau^* - \bar{Q}_\tau^{(t+1)}\|_\infty \leq \gamma \|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty + \gamma \max \left\{ 0, -\min_{s, a} \left(\bar{Q}_\tau^{(t+1)}(s, a) - \tau \log \bar{\xi}^{(t+1)}(s, a) \right) \right\}.\quad (124)$$

Plugging the above inequality into (108) and (119) establishes the bounds on $\Omega_3^{(t+1)}$ and $\Omega_2^{(t+1)}$ in (44), respectively.

Step 6: bound $-\min_{s, a} \left(\bar{Q}_\tau^{(t+1)}(s, a) - \tau \log \bar{\xi}^{(t+1)}(s, a) \right)$. We need the following lemma which is adapted from Lemma 1 in Cen et al. (2022a):

Lemma 9 (Performance improvement of FedNPG with entropy regularization). *Suppose $0 < \eta \leq (1 - \gamma)/\tau$. For any state-action pair $(s_0, a_0) \in \mathcal{S} \times \mathcal{A}$, one has*

$$\begin{aligned}\bar{V}_\tau^{(t+1)}(s_0) - \bar{V}_\tau^{(t)}(s_0) &\geq \frac{1}{\eta} \mathbb{E}_{s \sim d_{s_0}^{(t+1)}} \left[\alpha \text{KL}(\bar{\pi}^{(t+1)}(\cdot|s_0) \| \bar{\pi}^{(t)}(\cdot|s_0)) + \text{KL}(\bar{\pi}^{(t)}(\cdot|s_0) \| \bar{\pi}^{(t+1)}(\cdot|s_0)) \right] \\ &\quad - \frac{2}{1 - \gamma} \|\widehat{Q}_\tau^{(t)} - \bar{Q}_\tau^{(t)}\|_\infty,\end{aligned}\quad (125)$$

$$\bar{Q}_\tau^{(t+1)}(s_0, a_0) - \bar{Q}_\tau^{(t)}(s_0, a_0) \geq -\frac{2\gamma}{1 - \gamma} \|\widehat{Q}_\tau^{(t)} - \bar{Q}_\tau^{(t)}\|_\infty.\quad (126)$$

Proof. See Appendix C.3. □

Using (126), we have

$$\bar{Q}_\tau^{(t+1)}(s, a) - \tau \left(\alpha \log \bar{\xi}^{(t)}(s, a) + (1 - \alpha) \frac{\widehat{Q}_\tau^{(t)}(s, a)}{\tau} \right)$$

$$\begin{aligned}
&\geq \bar{Q}_\tau^{(t)}(s, a) - \tau \left(\alpha \log \bar{\xi}^{(t)}(s, a) + (1 - \alpha) \frac{\widehat{Q}_\tau^{(t)}(s, a)}{\tau} \right) - \frac{2\gamma}{1 - \gamma} \|\widehat{Q}_\tau^{(t)} - \bar{Q}_\tau^{(t)}\|_\infty \\
&\geq \alpha \left(\bar{Q}_\tau^{(t)}(s, a) - \tau \log \bar{\xi}^{(t)}(s, a) \right) - \frac{2\gamma + \eta\tau}{1 - \gamma} \|\widehat{Q}_\tau^{(t)} - \bar{Q}_\tau^{(t)}\|_\infty,
\end{aligned} \tag{127}$$

which gives

$$\begin{aligned}
& - \min_{s, a} \left(\bar{Q}_\tau^{(t+1)}(s, a) - \tau \log \bar{\xi}^{(t+1)}(s, a) \right) \\
& \leq -\alpha \min_{s, a} \left(\bar{Q}_\tau^{(t)}(s, a) - \tau \log \bar{\xi}^{(t)}(s, a) \right) + \frac{2\gamma + \eta\tau}{1 - \gamma} M \|u^{(t)}\|_\infty \\
& \leq \alpha \max \left\{ 0, \min_{s, a} \left(\bar{Q}_\tau^{(t)}(s, a) - \tau \log \bar{\xi}^{(t)}(s, a) \right) \right\} + \frac{2\gamma + \eta\tau}{1 - \gamma} M \|u^{(t)}\|_\infty.
\end{aligned} \tag{128}$$

This establishes the bounds on $\Omega_4^{(t+1)}$ in (44).

B.2 Proof of Lemma 2

Let $f(\lambda)$ denote the characteristic function. In view of some direct calculations, we obtain

$$\begin{aligned}
f(\lambda) &= (\lambda - \alpha) \left\{ \underbrace{(\lambda - \sigma\alpha)(\lambda - \sigma(1 + b\eta))(\lambda - (1 - \alpha)\gamma - \alpha)}_{=: f_0(\lambda)} \right. \\
& \quad \left. - \frac{\eta\sigma}{1 - \gamma} \underbrace{[S(\lambda - (1 - \alpha)\gamma - \alpha) + \gamma cdM\eta + (1 - \alpha)(2 + \gamma)Mcd\eta]}_{=: f_1(\lambda)} \right\} \\
& \quad - \frac{\tau\eta^3\gamma}{(1 - \gamma)^2} \cdot 2cdM\sigma,
\end{aligned} \tag{129}$$

where, for the notation simplicity, we let

$$b := \frac{M\sqrt{N}}{1 - \gamma}, \tag{130a}$$

$$c := \frac{MN}{1 - \gamma} = \sqrt{N}b, \tag{130b}$$

$$d := \frac{2\gamma + \eta\tau}{1 - \gamma}. \tag{130c}$$

Note that among all these new notation we introduce, S, d are dependent of η . To decouple the dependence, we give their upper bounds as follows

$$d_0 := \frac{1 + \gamma}{1 - \gamma} M \geq d, \tag{131}$$

$$S_0 := M\sqrt{N} \left(2 + \sqrt{2N} + \frac{M\sqrt{N}}{\tau} \right) \geq S, \tag{132}$$

where (131) follows from $\eta \leq (1 - \gamma)/\tau$, and (132) uses the fact that $\alpha \leq 1$ and $1 - \alpha \leq 1$.

Let

$$\lambda^* := \max \left\{ \frac{3 + \sigma}{4}, \frac{1 + (1 - \alpha)\gamma + \alpha}{2} \right\}. \tag{133}$$

Since $\mathbf{A}(\rho)$ is a nonnegative matrix, by Perron-Frobenius Theorem (see [Horn and Johnson \(2012\)](#), Theorem 8.3.1), $\rho(\eta)$ is an eigenvalue of $\mathbf{A}(\rho)$. So to verify (50), it suffices to show that $f(\lambda) > 0$ for any $\lambda \in [\lambda^*, \infty)$. To do so, in the following we first show that $f(\lambda^*) > 0$, and then we prove that f is non-decreasing on $[\lambda^*, \infty)$.

- *Showing $f(\lambda^*) > 0$.* We first lower bound $f_0(\lambda^*)$. Since $\lambda^* \geq \frac{3+\sigma}{4}$, we have

$$\lambda^* - \sigma(1 + b\eta) \geq \frac{1 - \sigma}{4}, \quad (134)$$

and from $\lambda^* \geq \frac{1+(1-\alpha)\gamma+\alpha}{2}$ we deduce

$$\lambda^* - (1 - \alpha)\gamma - \alpha \geq \frac{(1 - \gamma)(1 - \alpha)}{2} \quad (135)$$

and

$$\lambda^* > \frac{1 + \alpha}{2}, \quad (136)$$

which gives

$$\lambda^* - \sigma\alpha \geq \frac{1 + \alpha}{2} - \sigma\alpha. \quad (137)$$

Combining (137), (134), (135), we have that

$$f_0(\lambda^*) \geq \frac{1 - \sigma}{8} \left(\frac{1 + \alpha}{2} - \sigma\alpha \right) \eta\tau. \quad (138)$$

To continue, we upper bound $f_1(\lambda^*)$ as follows.

$$\begin{aligned} f_1(\lambda^*) &\leq S\tau\eta + \gamma cdM\eta + \frac{2 + \gamma}{1 - \gamma} cM\tau\eta^2 \\ &= \eta \left(\tau \left(S + \frac{2 + \gamma}{1 - \gamma} Mc\eta \right) + \gamma cdM \right). \end{aligned} \quad (139)$$

Plugging (138),(139) into (129) and using (136), we have

$$\begin{aligned} f(\lambda^*) &> \frac{1 - \alpha}{2} \left(f_0(\lambda^*) - \frac{\eta\sigma}{1 - \gamma} f_1(\lambda^*) \right) - \frac{\tau\eta^3\gamma}{(1 - \gamma)^2} \cdot 2cdM\sigma \\ &\geq \frac{\tau\eta^2}{2(1 - \gamma)} \left[\frac{1 - \sigma}{8} \tau \left(1 - \sigma + (1 - \alpha)(\sigma - \frac{1}{2}) \right) - \frac{\eta\sigma}{1 - \gamma} \left(\tau \left(S + \frac{2 + \gamma}{1 - \gamma} Mc\eta \right) + 5\gamma cdM \right) \right] \\ &= \frac{\tau\eta^2}{2(1 - \gamma)} \left[\frac{(1 - \sigma)^2}{8} \tau - \frac{\eta}{1 - \gamma} \left(S\tau\sigma + \frac{2 + \gamma}{1 - \gamma} Mc\sigma\tau\eta + \tau^2 \left(\frac{1}{2} - \sigma \right) \cdot \frac{1 - \sigma}{8} + 5\gamma cdM\sigma \right) \right] \\ &\geq \frac{\tau\eta^2}{2(1 - \gamma)} \left[\frac{(1 - \sigma)^2}{8} \tau - \frac{\eta}{1 - \gamma} \left(S_0\tau\sigma + \frac{(1 - \sigma)^2}{16} \tau^2 + (2 + \gamma + 5\gamma d_0) cM\sigma \right) \right] \geq 0, \end{aligned}$$

where the penultimate inequality uses $\frac{1}{2} - \sigma \leq \frac{1 - \sigma}{2}$, and the last inequality follows from the definition of ζ (cf. (48)).

- *Proving f is non-decreasing on $[\lambda^*, \infty)$.* Note that

$$\eta \leq \zeta \leq \frac{(1 - \gamma)(1 - \sigma)^2}{8S_0\sigma},$$

thus we have

$$\forall \lambda \geq \lambda^* : f'_0(\lambda) - \frac{\eta\sigma}{1 - \gamma} f'_1(\lambda) \geq (\lambda - \sigma\alpha)(\lambda - \sigma(1 + b\eta)) - \frac{\eta}{1 - \gamma} S\sigma \geq 0,$$

which indicates that $f_0 - f_1$ is non-decreasing on $[\lambda^*, \infty)$. Therefore, f is non-decreasing on $[\lambda^*, \infty)$.

B.3 Proof of Lemma 3

Note that bounding $u^{(t+1)}(s, a)$ is identical to the proof in Appendix B.1 and shall be omitted. The rest of the proof also follows closely that of Lemma 1, and we only highlight the differences due to approximation error for simplicity.

Step 2: bound $v^{(t+1)}(s, a) = \|\mathbf{T}^{(t+1)}(s, a) - \widehat{q}_\tau^{(t+1)}(s, a)\mathbf{1}_N\|_2$. Let $\mathbf{q}_\tau^{(t)} := (q_{\tau,1}^{\pi_1^{(t)}}, \dots, q_{\tau,N}^{\pi_N^{(t)}})^\top$. Similar to (104) we have

$$\begin{aligned} & \|\mathbf{T}^{(t+1)}(s, a) - \widehat{q}_\tau^{(t+1)}(s, a)\mathbf{1}_N\|_2 \\ & \leq \sigma \|\mathbf{T}^{(t)}(s, a) - \widehat{q}_\tau^{(t)}(s, a)\mathbf{1}_N\|_2 + \sigma \|\mathbf{q}_\tau^{(t+1)}(s, a) - \mathbf{q}_\tau^{(t)}(s, a)\|_2 \\ & \leq \sigma \|\mathbf{T}^{(t)}(s, a) - \widehat{q}_\tau^{(t)}(s, a)\mathbf{1}_N\|_2 + \sigma \|\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)\|_2 + 2\|\mathbf{e}\|_2. \end{aligned} \quad (140)$$

Step 3: bound $\|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty$. In the context of inexact updates, (105) writes

$$\|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty \leq \alpha \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty + (1 - \alpha) \|Q_\tau^* - \bar{Q}_\tau^{(t)}\|_\infty + (1 - \alpha) \|\bar{Q}_\tau^{(t)} - \widehat{q}_\tau^{(t)}\|_\infty.$$

For the last term, following a similar argument in (106) leads to

$$\begin{aligned} \|\bar{Q}_\tau^{(t)} - \widehat{q}_\tau^{(t)}\|_\infty &= \left\| \frac{1}{N} \sum_{n=1}^N Q_{\tau,n}^{\pi_n^{(t)}} - \frac{1}{N} \sum_{n=1}^N Q_{\tau,n}^{\bar{\pi}^{(t)}} \right\|_\infty + \left\| \frac{1}{N} \sum_{n=1}^N (Q_{\tau,n}^{\pi_n^{(t)}} - q_{\tau,n}^{\pi_n^{(t)}}) \right\|_\infty \\ &\leq M \cdot \frac{1}{N} \sum_{n=1}^N \|\log \xi_n^{(t)} - \log \bar{\xi}^{(t)}\|_\infty + \frac{1}{N} \sum_{n=1}^N e_n \\ &\leq M \|u^{(t)}\|_\infty + \|\mathbf{e}\|_\infty. \end{aligned}$$

Combining the above two inequalities, we obtain

$$\|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty \leq \alpha \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty + (1 - \alpha) \|Q_\tau^* - \bar{Q}_\tau^{(t)}\|_\infty + (1 - \alpha) (M \|u^{(t)}\|_\infty + \|\mathbf{e}\|_\infty). \quad (141)$$

Step 4: bound $\|\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)\|_2$. We remark that the bound established in (111) still holds in the inexact setting, with the same definition for $w^{(t)}$:

$$\|\mathbf{Q}_\tau^{(t+1)}(s, a) - \mathbf{Q}_\tau^{(t)}(s, a)\|_2 \leq M\sqrt{N} \|w^{(t)}\|_\infty. \quad (142)$$

To deal with the approximation error, we rewrite (113) as

$$\begin{aligned} & \|\mathbf{T}^{(t)}(s, a) - \tau \log \boldsymbol{\xi}^{(t)}(s, a) - V_\tau^*(s)\mathbf{1}_N\|_2 \\ &= \|\mathbf{T}^{(t)}(s, a) - \tau \log \boldsymbol{\xi}^{(t)}(s, a) - (Q_\tau^*(s, a) - \tau \log \pi_\tau^*(a|s))\mathbf{1}_N\|_2 \\ &\leq \|\mathbf{T}^{(t)}(s, a) - Q_\tau^*(s, a)\mathbf{1}_N\|_2 + \tau \|\log \boldsymbol{\xi}^{(t)}(s, a) - \log \pi_\tau^*(a|s)\mathbf{1}_N\|_2 \\ &\leq \|\mathbf{T}^{(t)}(s, a) - \widehat{q}_\tau^{(t)}(s, a)\mathbf{1}_N\|_2 + \|\widehat{q}_\tau^{(t)}(s, a)\mathbf{1}_N - Q_\tau^*(s, a)\mathbf{1}_N\|_2 \\ &\quad + \tau \|\log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\pi}^{(t)}(a|s)\mathbf{1}_N\|_2 + \tau \|\log \bar{\pi}^{(t)}(a|s)\mathbf{1}_N - \log \pi_\tau^*(a|s)\mathbf{1}_N\|_2 \\ &\leq \|\mathbf{T}^{(t)}(s, a) - \widehat{q}_\tau^{(t)}(s, a)\mathbf{1}_N\|_2 + \sqrt{N} |\widehat{q}_\tau^{(t)}(s, a) - Q_\tau^*(s, a)| \\ &\quad + \tau \|\log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\pi}^{(t)}(a|s)\mathbf{1}\|_2 + \tau \sqrt{N} |\log \bar{\pi}^{(t)}(a|s) - \log \pi_\tau^*(a|s)|, \end{aligned} \quad (143)$$

where the second term can be upper-bounded by

$$|\widehat{q}_\tau^{(t)}(s, a) - Q_\tau^*(s, a)| \leq \|\widehat{Q}_\tau^{(t)} - \bar{Q}_\tau^{(t)}\|_\infty + \|\bar{Q}_\tau^{(t)} - Q_\tau^*\|_\infty + \|\widehat{q}_\tau^{(t)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a)\|_\infty$$

$$\leq \|\widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}\|_\infty + \|\overline{Q}_\tau^{(t)} - Q_\tau^*\|_\infty + \|\mathbf{e}\|_\infty. \quad (144)$$

Combining (144), (143) and the established bounds in (112), (115), (117) leads to

$$\begin{aligned} w^{(t)}(s, a) &\leq \left(2\alpha + (1 - \alpha) \cdot \sqrt{2N}\right) \|u^{(t)}\|_\infty + \frac{1 - \alpha}{\tau} \|v^{(t)}\|_\infty \\ &\quad + \frac{1 - \alpha}{\tau} \cdot \sqrt{N} \left(\|\widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}\|_\infty + \|\overline{Q}_\tau^{(t)} - Q_\tau^*\|_\infty + \|\mathbf{e}\|_\infty \right) + \frac{1 - \alpha}{\tau} \cdot 2\sqrt{N} \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty. \end{aligned}$$

Combining the above inequality with (142) and (140) gives

$$\begin{aligned} \|v^{(t+1)}\|_\infty &\leq \sigma \left(1 + \frac{\eta M \sqrt{N}}{1 - \gamma} \right) \|v^{(t)}\|_\infty + \sigma M \sqrt{N} \left\{ \left(2\alpha + (1 - \alpha) \cdot \sqrt{2N} + \frac{1 - \alpha}{\tau} \cdot \sqrt{NM} \right) \|u^{(t)}\|_\infty \right. \\ &\quad \left. + \frac{1 - \alpha}{\tau} \cdot \sqrt{N} \left(\|\overline{Q}_\tau^{(t)} - Q_\tau^*\|_\infty + \|\mathbf{e}\|_\infty \right) + \frac{1 - \alpha}{\tau} \cdot 2\sqrt{N} \|Q_\tau^* - \tau \log \bar{\xi}^{(t)}\|_\infty \right\} + 2\sigma \sqrt{N} \|\mathbf{e}\|_\infty. \end{aligned} \quad (145)$$

Step 5: bound $\|\overline{Q}_\tau^{(t+1)} - Q_\tau^*\|_\infty$. It is straightforward to verify that (124) applies to the inexact updates as well:

$$\|Q_\tau^* - \overline{Q}_\tau^{(t+1)}\|_\infty \leq \gamma \|Q_\tau^* - \tau \log \bar{\xi}^{(t+1)}\|_\infty + \gamma \left(-\min_{s,a} \left(\overline{Q}_\tau^{(t+1)}(s, a) - \tau \log \bar{\xi}^{(t+1)}(s, a) \right) \right).$$

Plugging the above inequality into (141) and (145) establishes the bounds on $\Omega_3^{(t+1)}$ and $\Omega_2^{(t+1)}$ in (62), respectively.

Step 6: bound $-\min_{s,a} \left(\overline{Q}_\tau^{(t+1)}(s, a) - \tau \log \bar{\xi}^{(t+1)}(s, a) \right)$. We obtain the following lemma by interpreting the approximation error \mathbf{e} as part of the consensus error $\|\widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}\|_\infty$ in Lemma 9.

Lemma 10 (inexact version of Lemma 9). *Suppose $0 < \eta \leq (1 - \gamma)/\tau$. For any state-action pair $(s_0, a_0) \in \mathcal{S} \times \mathcal{A}$, one has*

$$\begin{aligned} \overline{V}_\tau^{(t+1)}(s_0) - \overline{V}_\tau^{(t)}(s_0) &\geq \frac{1}{\eta} \mathbb{E}_{s \sim d_{s_0}^{(t+1)}} \left[\alpha \text{KL}(\overline{\pi}^{(t+1)}(\cdot | s_0) \| \overline{\pi}^{(t)}(\cdot | s_0)) + \text{KL}(\overline{\pi}^{(t)}(\cdot | s_0) \| \overline{\pi}^{(t+1)}(\cdot | s_0)) \right] \\ &\quad - \frac{2}{1 - \gamma} \left(\|\widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}\|_\infty + \|\mathbf{e}\|_\infty \right), \end{aligned} \quad (146)$$

$$\overline{Q}_\tau^{(t+1)}(s_0, a_0) - \overline{Q}_\tau^{(t)}(s_0, a_0) \geq -\frac{2\gamma}{1 - \gamma} \left(\|\widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}\|_\infty + \|\mathbf{e}\|_\infty \right). \quad (147)$$

Using (147), we have

$$\begin{aligned} &\overline{Q}_\tau^{(t+1)}(s, a) - \tau \left(\alpha \log \bar{\xi}^{(t)}(s, a) + (1 - \alpha) \frac{\widehat{Q}_\tau^{(t)}(s, a)}{\tau} \right) \\ &\geq \overline{Q}_\tau^{(t)}(s, a) - \tau \left(\alpha \log \bar{\xi}^{(t)}(s, a) + (1 - \alpha) \frac{\widehat{Q}_\tau^{(t)}(s, a)}{\tau} \right) - \frac{2\gamma}{1 - \gamma} \left(\|\widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}\|_\infty + \|\mathbf{e}\|_\infty \right) \\ &\geq \alpha \left(\overline{Q}_\tau^{(t)}(s, a) - \tau \log \bar{\xi}^{(t)}(s, a) \right) - \frac{2\gamma + \eta\tau}{1 - \gamma} \|\widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}\|_\infty - \frac{2\gamma}{1 - \gamma} \|\mathbf{e}\|_\infty, \end{aligned} \quad (148)$$

which gives

$$\begin{aligned} &-\min_{s,a} \left(\overline{Q}_\tau^{(t+1)}(s, a) - \tau \log \bar{\xi}^{(t+1)}(s, a) \right) \\ &\leq -\alpha \min_{s,a} \left(\overline{Q}_\tau^{(t)}(s, a) - \tau \log \bar{\xi}^{(t)}(s, a) \right) + \frac{2\gamma + \eta\tau}{1 - \gamma} M \|u^{(t)}\|_\infty + \frac{2\gamma}{1 - \gamma} \|\mathbf{e}\|_\infty. \end{aligned} \quad (149)$$

B.4 Proof of Lemma 4

Step 1: bound $u^{(t+1)}(s, a) = \left\| \log \boldsymbol{\xi}^{(t+1)}(s, a) - \log \bar{\boldsymbol{\xi}}^{(t+1)}(s, a) \mathbf{1}_N \right\|_2$. Following the same strategy in establishing (103), we have

$$\begin{aligned} & \left\| \log \boldsymbol{\xi}^{(t+1)}(s, a) - \log \bar{\boldsymbol{\xi}}^{(t+1)}(s, a) \mathbf{1}_N \right\|_2 \\ &= \left\| \left(\mathbf{W} \log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\boldsymbol{\xi}}^{(t)}(s, a) \mathbf{1}_N \right) + \frac{\eta}{1-\gamma} \left(\mathbf{T}^{(t)}(s, a) - \widehat{\mathbf{Q}}^{(t)}(s, a) \mathbf{1}_N \right) \right\|_2 \\ &\leq \sigma \left\| \log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\boldsymbol{\xi}}^{(t)}(s, a) \mathbf{1}_N \right\|_2 + \frac{\eta}{1-\gamma} \left\| \mathbf{T}^{(t)}(s, a) - \widehat{\mathbf{Q}}^{(t)}(s, a) \mathbf{1}_N \right\|_2, \end{aligned} \quad (150)$$

or equivalently

$$\|u^{(t+1)}\|_\infty \leq \sigma \|u^{(t)}\|_\infty + \frac{\eta}{1-\gamma} \|v^{(t)}\|_\infty. \quad (151)$$

Step 2: bound $v^{(t+1)}(s, a) = \left\| \mathbf{T}^{(t+1)}(s, a) - \widehat{\mathbf{Q}}^{(t+1)}(s, a) \mathbf{1}_N \right\|_2$. In the same vein of establishing (104), we have

$$\begin{aligned} & \left\| \mathbf{T}^{(t+1)}(s, a) - \widehat{\mathbf{Q}}^{(t+1)}(s, a) \mathbf{1}_N \right\|_2 \\ &\leq \sigma \left\| \mathbf{T}^{(t)}(s, a) - \widehat{\mathbf{Q}}^{(t)}(s, a) \mathbf{1}_N \right\|_2 + \sigma \left\| \mathbf{Q}^{(t+1)}(s, a) - \mathbf{Q}^{(t)}(s, a) \right\|_2, \end{aligned} \quad (152)$$

The term $\left\| \mathbf{Q}^{(t+1)}(s, a) - \mathbf{Q}^{(t)}(s, a) \right\|_2$ can be bounded in a similar way in (111):

$$\left\| \mathbf{Q}^{(t+1)}(s, a) - \mathbf{Q}^{(t)}(s, a) \right\|_2 \leq \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \sqrt{N} \|w_0^{(t)}\|_\infty, \quad (153)$$

where the coefficient $\frac{(1+\gamma)\gamma}{(1-\gamma)^2}$ comes from M in Lemma 8 when $\tau = 0$, and $w_0^{(t)} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ is defined as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad w_0^{(t)}(s, a) := \left\| \log \boldsymbol{\xi}^{(t+1)}(s, a) - \log \boldsymbol{\xi}^{(t)}(s, a) - \frac{\eta}{1-\gamma} V^*(s) \mathbf{1}_N \right\|_2. \quad (154)$$

It remains to bound $\|w_0^{(t)}\|_\infty$. Towards this end, we rewrite (112) as

$$\begin{aligned} & w_0^{(t)}(s, a) \\ &= \left\| \mathbf{W} \log \boldsymbol{\xi}^{(t)}(s, a) + \frac{\eta}{1-\gamma} \mathbf{T}^{(t)}(s, a) - \log \boldsymbol{\xi}^{(t)}(s, a) - \frac{\eta}{1-\gamma} V^*(s) \mathbf{1}_N \right\|_2 \\ &= \left\| (\mathbf{W} - \mathbf{I}) \left(\log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\boldsymbol{\xi}}^{(t)}(s, a) \mathbf{1}_N \right) + \frac{\eta}{1-\gamma} \left(\mathbf{T}^{(t)}(s, a) - V^*(s) \mathbf{1}_N \right) \right\|_2 \\ &\leq 2 \left\| \log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\boldsymbol{\xi}}^{(t)}(s, a) \mathbf{1}_N \right\|_2 + \frac{\eta}{1-\gamma} \left\| \mathbf{T}^{(t)}(s, a) - V^*(s) \mathbf{1}_N \right\|_2 \\ &\leq 2 \left\| \log \boldsymbol{\xi}^{(t)}(s, a) - \log \bar{\boldsymbol{\xi}}^{(t)}(s, a) \mathbf{1}_N \right\|_2 + \frac{\eta}{1-\gamma} \left\| \mathbf{T}^{(t)}(s, a) - \widehat{\mathbf{Q}}^{(t)}(s, a) \mathbf{1}_N \right\|_2 + \frac{\eta}{1-\gamma} \cdot \sqrt{N} \left| \widehat{\mathbf{Q}}^{(t)}(s, a) - V^*(s) \right|. \end{aligned} \quad (155)$$

Note that it holds for all $(s, a) \in \mathcal{S} \times \mathcal{A}$:

$$\left| \widehat{\mathbf{Q}}^{(t)}(s, a) - V^*(s) \right| \leq \frac{1}{1-\gamma}$$

since $\widehat{\mathbf{Q}}^{(t)}(s, a)$ and $V^*(s)$ are both in $[0, 1/(1-\gamma)]$. This along with (155) gives

$$w_0^{(t)}(s, a) \leq 2 \|u^{(t)}\|_\infty + \frac{\eta}{1-\gamma} \|v^{(t)}\|_\infty + \frac{\eta \sqrt{N}}{(1-\gamma)^2}.$$

Combining the above inequality with (153) and (152), we arrive at

$$\|v^{(t+1)}\|_\infty \leq \sigma \left(1 + \frac{(1+\gamma)\gamma\sqrt{N}\eta}{(1-\gamma)^3} \right) \|v^{(t)}\|_\infty + \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \sqrt{N} \sigma \left\{ 2 \|u^{(t)}\|_\infty + \frac{\eta}{(1-\gamma)^2} \cdot \sqrt{N} \right\}. \quad (156)$$

Step 3: establish the descent equation. The following lemma characterizes the improvement in $\phi^{(t)}(\eta)$ for every iteration of Algorithm 1, with the proof postponed to Appendix C.4.

Lemma 11 (Performance improvement of exact FedNPG). *For all starting state distribution $\rho \in \Delta(\mathcal{S})$, we have the iterates of FedNPG satisfy*

$$\phi^{(t+1)}(\eta) \leq \phi^{(t)}(\eta) + \frac{2\eta}{(1-\gamma)^2} \|\widehat{Q}^{(t)} - \overline{Q}^{(t)}\|_\infty - \eta \left(V^*(\rho) - \overline{V}^{(t)}(\rho) \right), \quad (157)$$

where

$$\phi^{(t)}(\eta) := \mathbb{E}_{s \sim d_{\rho^*}} \left[\text{KL}(\pi^*(\cdot|s) \parallel \overline{\pi}^{(t)}(\cdot|s)) \right] - \frac{\eta}{1-\gamma} \overline{V}^{(t)}(d_{\rho^*}), \quad \forall t \geq 0. \quad (158)$$

It remains to control the term $\|\overline{Q}^{(t)} - \widehat{Q}^{(t)}\|_\infty$. Similar to (106), for all $t \geq 0$, we have

$$\begin{aligned} \|\overline{Q}^{(t)} - \widehat{Q}^{(t)}\|_\infty &= \left\| \frac{1}{N} \sum_{n=1}^N Q_n^{\pi_n^{(t)}} - \frac{1}{N} \sum_{n=1}^N Q_n^{\overline{\pi}^{(t)}} \right\|_\infty \\ &\stackrel{(a)}{\leq} \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \cdot \frac{1}{N} \sum_{n=1}^N \|\log \xi_n^{(t)} - \log \overline{\xi}^{(t)}\|_\infty \\ &\stackrel{(b)}{\leq} \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \|u^{(t)}\|_\infty, \end{aligned} \quad (159)$$

where (a) invokes Lemma 8 with $\tau = 0$ and (b) stems from the definition of $u^{(t)}$. This along with (157) gives

$$\phi^{(t+1)}(\eta) \leq \phi^{(t)}(\eta) + \frac{2(1+\gamma)\gamma}{(1-\gamma)^4} \eta \|u^{(t)}\|_\infty - \eta \left(V^*(\rho) - \overline{V}^{(t)}(\rho) \right).$$

B.5 Proof of Lemma 5

The bound on $u^{(t+1)}(s, a)$ is already established in Step 1 in Appendix B.1 and shall be omitted. As usual we only highlight the key differences with the proof of Lemma 4 due to approximation error.

Step 1: bound $v^{(t+1)}(s, a) = \|\mathbf{T}^{(t+1)}(s, a) - \widehat{q}^{(t+1)}(s, a)\mathbf{1}_N\|_2$. Let $\mathbf{q}^{(t)} := (q_1^{\pi_1^{(t)}}, \dots, q_N^{\pi_N^{(t)}})^\top$. From (87), we have

$$\begin{aligned} &\left\| \mathbf{T}^{(t+1)}(s, a) - \widehat{q}^{(t+1)}(s, a)\mathbf{1}_N \right\|_2 \\ &= \left\| \mathbf{W} \left(\mathbf{T}^{(t)}(s, a) + \mathbf{q}^{(t+1)}(s, a) - \mathbf{q}^{(t)}(s, a) \right) - \widehat{q}^{(t+1)}(s, a)\mathbf{1}_N \right\|_2 \\ &= \left\| \left(\mathbf{W}\mathbf{T}^{(t)}(s, a) - \widehat{q}^{(t)}(s, a)\mathbf{1}_N \right) + \mathbf{W} \left(\mathbf{q}^{(t+1)}(s, a) - \mathbf{q}^{(t)}(s, a) \right) + \left(\widehat{q}^{(t)}(s, a) - \widehat{q}^{(t+1)}(s, a) \right) \mathbf{1}_N \right\|_2 \\ &\leq \sigma \|\mathbf{T}^{(t)}(s, a) - \widehat{q}^{(t)}(s, a)\mathbf{1}_N\|_2 + \sigma \left\| \left(\mathbf{q}^{(t+1)}(s, a) - \mathbf{q}^{(t)}(s, a) \right) + \left(\widehat{q}^{(t)}(s, a) - \widehat{q}^{(t+1)}(s, a) \right) \mathbf{1}_N \right\|_2 \\ &\leq \sigma \|\mathbf{T}^{(t)}(s, a) - \widehat{q}^{(t)}(s, a)\mathbf{1}_N\|_2 + \sigma \|\mathbf{q}^{(t+1)}(s, a) - \mathbf{q}^{(t)}(s, a)\|_2 \\ &\leq \sigma \|\mathbf{T}^{(t)}(s, a) - \widehat{q}^{(t)}(s, a)\mathbf{1}_N\|_2 + \sigma \|\mathbf{Q}^{(t+1)}(s, a) - \mathbf{Q}^{(t)}(s, a)\|_2 + 2\sigma\sqrt{N} \|\mathbf{e}\|_\infty. \end{aligned} \quad (160)$$

Note that (153) still holds for inexact FedNPG:

$$\left\| \mathbf{Q}^{(t+1)}(s, a) - \mathbf{Q}^{(t)}(s, a) \right\|_2 \leq \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \sqrt{N} \|w_0^{(t)}\|_\infty, \quad (161)$$

where $w_0^{(t)}$ is defined in (154). We rewrite (155), the bound on $w_0^{(t)}(s, a)$, as

$$w_0^{(t)}(s, a) \leq 2 \|\log \boldsymbol{\xi}^{(t)}(s, a) - \log \overline{\boldsymbol{\xi}}^{(t)}(s, a)\mathbf{1}_N\|_2$$

$$+ \frac{\eta}{1-\gamma} \|\mathbf{T}^{(t)}(s, a) - \widehat{q}^{(t)}(s, a) \mathbf{1}_N\|_2 + \frac{\eta}{1-\gamma} \cdot \sqrt{N} |\widehat{q}^{(t)}(s, a) - V^*(s)|. \quad (162)$$

With the following bound

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad |\widehat{q}^{(t)}(s, a) - V^*(s)| \leq \|\widehat{q}^{(t)} - \overline{Q}^{(t)}\|_\infty + \frac{1}{1-\gamma}$$

in mind, we write (155) as

$$w_0^{(t)}(s, a) \leq 2\|u^{(t)}\|_\infty + \frac{\eta}{1-\gamma} \|v^{(t)}\|_\infty + \frac{\eta}{1-\gamma} \cdot \sqrt{N} \left(\|\widehat{q}^{(t)} - \overline{q}^{(t)}\|_\infty + \frac{1}{1-\gamma} \right).$$

Putting all pieces together, we obtain

$$\begin{aligned} \|v^{(t+1)}\|_\infty &\leq \sigma \left(1 + \frac{(1+\gamma)\gamma\sqrt{N}\eta}{(1-\gamma)^3} \right) \|v^{(t)}\|_\infty \\ &\quad + \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \sqrt{N} \sigma \left\{ \left(2 + \frac{(1+\gamma)\gamma\sqrt{N}\eta}{(1-\gamma)^3} \right) \|u^{(t)}\|_\infty + \frac{\eta\sqrt{N}}{(1-\gamma)^2} + \frac{\eta\sqrt{N}}{1-\gamma} \|\mathbf{e}\|_\infty \right\} \\ &\quad + 2\sigma\sqrt{N} \|\mathbf{e}\|_\infty. \end{aligned} \quad (163)$$

Step 2: establish the descent equation. Note that Lemma 11 directly applies by replacing $\widehat{Q}^{(t)}$ with $\widehat{q}^{(t)}$:

$$\phi^{(t+1)}(\eta) \leq \phi^{(t)}(\eta) + \frac{2\eta}{(1-\gamma)^2} \|\widehat{q}^{(t)} - \overline{Q}^{(t)}\|_\infty - \eta \left(V^*(\rho) - \overline{V}^{(t)}(\rho) \right).$$

To bound the middle term, for all $t \geq 0$, we have

$$\begin{aligned} \|\overline{Q}^{(t)} - \widehat{q}^{(t)}\|_\infty &= \left\| \frac{1}{N} \sum_{n=1}^N Q_n^{\pi_n^{(t)}} - \frac{1}{N} \sum_{n=1}^N Q_n^{\overline{\pi}^{(t)}} \right\|_\infty \\ &\leq \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \cdot \frac{1}{N} \sum_{n=1}^N \left\| \log \xi_n^{(t)} - \log \overline{\xi}^{(t)} \right\|_\infty + \frac{1}{N} \left\| \sum_{n=0}^N \left(q_n^{\pi_n^{(t)}} - Q_n^{\pi_n^{(t)}} \right) \right\|_\infty + \frac{1}{N} \sum_{n=1}^N e_n \\ &\leq \frac{(1+\gamma)\gamma}{(1-\gamma)^2} \|u^{(t)}\|_\infty + \|\mathbf{e}\|_\infty. \end{aligned} \quad (164)$$

Hence, (93) is established by combining the above two inequalities.

C Proof of auxiliary lemmas

C.1 Proof of Lemma 6

The first claim is easily verified as $\log \xi_n^{(t)}(s, \cdot)$ always deviate from $\log \pi_n^{(t)}(\cdot|s)$ by a global constant shift, as long as it holds for $t = 0$:

$$\begin{aligned} \log \xi_n^{(t+1)}(s, \cdot) &= \alpha \sum_{n'=1}^N [W]_{n,n'} \log \xi_{n'}^{(t)}(s, \cdot) + (1-\alpha) T_n^{(t)}(s, \cdot) / \tau \\ &= \alpha \sum_{n'=1}^N [W]_{n,n'} \left(\log \pi_{n'}^{(t)}(s, \cdot) + c_{n'}^{(t)}(s) \mathbf{1}_{|\mathcal{A}|} \right) + (1-\alpha) T_n^{(t)}(s, \cdot) / \tau \\ &= \alpha \sum_{n'=1}^N [W]_{n,n'} \log \pi_{n'}^{(t)}(s, \cdot) + (1-\alpha) T_n^{(t)}(s, \cdot) / \tau - z_n^{(t)}(s) \mathbf{1}_{|\mathcal{A}|} + c_n^{(t+1)}(s) \mathbf{1}_{|\mathcal{A}|} \\ &= \log \pi_n^{(t+1)}(\cdot|s) + c_n^{(t+1)}(s) \mathbf{1}_{|\mathcal{A}|}, \end{aligned}$$

where $z_n^{(t)}$ is the normalization term (cf. line 5, Algorithm 2) and $\{c_n^{(t)}(s)\}$ are some constants. To prove the second claim, $\forall t \geq 0, \forall (s, a) \in \mathcal{S} \times \mathcal{A}$, let

$$\bar{T}^{(t)}(s, a) := \frac{1}{N} \mathbf{1}^\top \mathbf{T}^{(t)}(s, a). \quad (165)$$

Taking inner product with $\frac{1}{N} \mathbf{1}$ for both sides of (25) and using the double stochasticity property of \mathbf{W} , we get

$$\bar{T}^{(t+1)}(s, a) = \bar{T}^{(t)}(s, a) + \widehat{Q}_\tau^{(t+1)}(s, a) - \widehat{Q}_\tau^{(t)}(s, a). \quad (166)$$

By the choice of $\mathbf{T}^{(0)}$ (line 2 of Algorithm 2), we have $\bar{T}^{(0)} = \widehat{Q}_\tau^{(0)}$ and hence by induction

$$\forall t \geq 0 : \quad \bar{T}^{(t)} = \widehat{Q}_\tau^{(t)}. \quad (167)$$

This implies

$$\begin{aligned} \log \bar{\xi}^{(t+1)}(s, a) - \alpha \log \bar{\xi}^{(t)}(s, a) &= (1 - \alpha) \widehat{Q}_\tau^{(t)}(s, a) / \tau \\ &= (1 - \alpha) \bar{T}^{(t)}(s, a) / \tau \\ &= \frac{1}{N} \mathbf{1}^\top \log \boldsymbol{\xi}^{(t+1)}(s, a) - \alpha \frac{1}{N} \mathbf{1}^\top \log \boldsymbol{\xi}^{(t)}(s, a). \end{aligned}$$

Therefore, to prove (98), it suffices to verify the claim for $t = 0$:

$$\begin{aligned} \frac{1}{N} \mathbf{1}^\top \log \boldsymbol{\xi}^{(0)}(s, a) &= \log \|\exp(Q_\tau^*(s, \cdot) / \tau)\|_1 + \frac{1}{N} \mathbf{1}^\top \log \boldsymbol{\pi}^{(0)}(a|s) - \log \left\| \exp \left(\frac{1}{N} \sum_{n=1}^N \log \pi_n^{(0)}(\cdot|s) \right) \right\|_1 \\ &= \log \|\exp(Q_\tau^*(s, \cdot) / \tau)\|_1 + \log \bar{\pi}^{(0)}(a|s) = \log \bar{\xi}^{(0)}(s, a). \end{aligned}$$

By taking logarithm over both sides of the definition of $\bar{\pi}^{(t+1)}$ (cf. (24)), we get

$$\log \bar{\pi}^{(t+1)}(a|s) = \alpha \log \bar{\pi}^{(t)}(a|s) + (1 - \alpha) \widehat{Q}_\tau^{(t)}(s, a) / \tau - z^{(t)}(s) \quad (168)$$

for some constant $z^{(t)}(s)$, which deviate from the update rule of $\log \bar{\xi}^{(t+1)}$ by a global constant shift and hence verifies (99).

C.2 Proof of Lemma 8

For notational simplicity, we let $Q_\tau^{\theta'}$ and Q_τ^θ denote $Q_\tau^{\pi_{\theta'}}$ and $Q_\tau^{\pi_\theta}$, respectively. From (7a) we immediately know that to bound $\|Q_\tau^{\theta'} - Q_\tau^\theta\|_\infty$, it suffices to control $|V_\tau^\theta(s) - V_\tau^{\theta'}(s)|$ for each $s \in \mathcal{S}$. By (4) we have

$$|V_\tau^\theta(s) - V_\tau^{\theta'}(s)| \leq |V^\theta(s) - V^{\theta'}(s)| + \tau |\mathcal{H}(s, \pi_\theta) - \mathcal{H}(s, \pi_{\theta'})|, \quad (169)$$

so in the following we bound both terms in the RHS of (169).

Step 1: bounding $|\mathcal{H}(s, \pi_\theta) - \mathcal{H}(s, \pi_{\theta'})|$. We first bound $|\mathcal{H}(s, \pi_\theta) - \mathcal{H}(s, \pi_{\theta'})|$ using the idea in the proof of Lemma 14 in Mei et al. (2020). We let

$$\theta_t = \theta + t(\theta' - \theta), \quad \forall t \in \mathbb{R}, \quad (170)$$

and let $h_t \in \mathbb{R}^{|\mathcal{S}|}$ be

$$\forall s \in \mathcal{S} : \quad h_t(s) := - \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) \log \pi_{\theta_t}(a|s). \quad (171)$$

Note that $\|h_t\|_\infty \leq \log |\mathcal{A}|$. We also denote $H_t : \mathcal{S} \rightarrow \mathbb{R}^{|\mathcal{A}| \times |\mathcal{A}|}$ by:

$$\forall s \in \mathcal{S} : \quad H_t(s) := \frac{\partial \pi_{\theta_t}(\cdot|s)}{\partial \theta} \Big|_{\theta=\theta_t} = \text{diag}\{\pi_{\theta_t}(\cdot|s)\} - \pi_{\theta_t}(\cdot|s) \pi_{\theta_t}(\cdot|s)^\top, \quad (172)$$

then we have

$$\begin{aligned}
\forall s \in \mathcal{S} : \quad \left| \frac{dh_t(s)}{dt} \right| &= \left| \left\langle \frac{\partial h_t(s)}{\partial \theta_t(\cdot|s)}, \theta'(s, \cdot) - \theta(s, \cdot) \right\rangle \right| \\
&= |\langle H_t(s) \log \pi_{\theta_t}(\cdot|s), \theta'(s, \cdot) - \theta(s, \cdot) \rangle| \\
&\leq \|H_t(s) \log \pi_{\theta_t}(\cdot|s)\|_1 \|\theta'(s, \cdot) - \theta(s, \cdot)\|_\infty,
\end{aligned} \tag{173}$$

where $\frac{\partial h_t(s)}{\partial \theta_t(\cdot|s)}$ stands for $\frac{\partial h_t(s)}{\partial \theta(\cdot|s)}|_{\theta=\theta_t}$. The first term in (173) is further upper bounded as

$$\begin{aligned}
\|H_t(s) \log \pi_{\theta_t}(\cdot|s)\|_1 &= \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) |\log \pi_{\theta_t}(a|s) - \pi_{\theta_t}(\cdot|s)^\top \log \pi_{\theta_t}(\cdot|s)| \\
&\leq \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) (|\log \pi_{\theta_t}(a|s)| + |\pi_{\theta_t}(\cdot|s)^\top \log \pi_{\theta_t}(\cdot|s)|) \\
&= -2 \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a, s) \log \pi_{\theta_t}(a|s) \leq 2 \log |\mathcal{A}|.
\end{aligned}$$

By Lagrange mean value theorem, there exists $t \in (0, 1)$ such that

$$|h_1(s) - h_0(s)| = \left| \frac{dh_t(s)}{dt} \right| \leq 2 \log |\mathcal{A}| \|\theta'(s, \cdot) - \theta(s, \cdot)\|_\infty,$$

where the inequality follows from (173) and the above inequality. Combining (5) with the above inequality, we arrive at

$$|\mathcal{H}(s, \pi_\theta) - \mathcal{H}(s, \pi_{\theta'})| \leq \frac{2 \log |\mathcal{A}|}{1 - \gamma} \|\theta' - \theta\|_\infty. \tag{174}$$

Step 2: bounding $|V^\theta(s) - V^{\theta'}(s)|$. Similar to the previous proof, we bound $|V^\theta(s) - V^{\theta'}(s)|$ by bounding $\left| \frac{dV^{\theta_t}}{dt}(s) \right|$. By Bellman's consistency equation, the value function of π_{θ_t} is given by

$$V^{\theta_t}(s) = \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) r(s, a) + \gamma \sum_a \pi_{\theta_t}(a|s) \sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) V^{\theta_t}(s'),$$

which can be represented in a matrix-vector form as

$$V^{\theta_t}(s) = e_s^\top \mathbf{M}_t r_t, \tag{175}$$

where $e_s \in \mathbb{R}^{|\mathcal{S}|}$ is a one-hot vector whose s -th entry is 1,

$$\mathbf{M}_t := (\mathbf{I} - \gamma \mathbf{P}_t)^{-1}, \tag{176}$$

with $\mathbf{P}_t \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ denoting the induced state transition matrix by π_{θ_t}

$$\mathbf{P}_t(s, s') = \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) \mathcal{P}(s'|s, a), \tag{177}$$

and $r_t \in \mathbb{R}^{|\mathcal{S}|}$ is given by

$$\forall s \in \mathcal{S} : \quad r_t(s) := \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) r(s, a). \tag{178}$$

Taking derivative w.r.t. t in (175), we obtain (Petersen and Pedersen, 2008)

$$\frac{dV^{\theta_t}(s)}{dt} = \gamma \cdot e_s^\top \mathbf{M}_t \frac{d\mathbf{P}_t}{dt} \mathbf{M}_t r_t + e_s^\top \mathbf{M}_t \frac{dr_t}{dt}. \tag{179}$$

We now calculate each term respectively.

- For the first term, it follows that

$$\begin{aligned}
\left| \gamma \cdot e_s^\top \mathbf{M}_t \frac{d\mathbf{P}_t}{dt} \mathbf{M}_t r_t \right| &\leq \gamma \left\| \mathbf{M}_t \frac{d\mathbf{P}_t}{dt} \mathbf{M}_t r_t \right\|_\infty \\
&\leq \frac{\gamma}{1-\gamma} \left\| \frac{d\mathbf{P}_t}{dt} \mathbf{M}_t r_t \right\|_\infty \\
&\leq \frac{2\gamma}{1-\gamma} \|\mathbf{M}_t r_t\|_\infty \|\theta' - \theta\|_\infty \\
&\leq \frac{2\gamma}{(1-\gamma)^2} \|r_t\|_\infty \|\theta' - \theta\|_\infty \\
&\leq \frac{2\gamma}{(1-\gamma)^2} \|\theta' - \theta\|_\infty .
\end{aligned} \tag{180}$$

where the second and fourth lines use the fact $\|\mathbf{M}_t\|_1 \leq 1/(1-\gamma)$ (Li et al., 2023b, Lemma 10), and the last line follow from

$$\|r_t\|_\infty = \max_{s \in \mathcal{S}} \left| \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) r(s, a) \right| \leq 1.$$

We defer the proof of (180) to the end of proof.

- For the second term, it follows that

$$\left| e_s^\top \mathbf{M}_t \frac{dr_t}{dt} \right| \leq \frac{1}{1-\gamma} \left\| \frac{dr_t}{dt} \right\|_\infty \leq \frac{1}{1-\gamma} \|\theta' - \theta\|_\infty . \tag{182}$$

where the first inequality follows again from $\|\mathbf{M}_t\|_1 \leq 1/(1-\gamma)$, and the second inequality follows from

$$\begin{aligned}
\left\| \frac{dr_t}{dt} \right\|_\infty &= \max_{s \in \mathcal{S}} \left| \frac{dr_t(s)}{dt} \right| = \max_{s \in \mathcal{S}} \left| \left\langle \frac{\partial \pi_{\theta_t}(\cdot|s)^\top r(s, \cdot)}{\partial \theta_t(s, \cdot)}, \theta'(s, \cdot) - \theta(s, \cdot) \right\rangle \right| \\
&\leq \max_{s \in \mathcal{S}} \left\| \frac{\partial \pi_{\theta_t}(\cdot|s)^\top r(s, \cdot)}{\partial \theta_t(s, \cdot)} \right\|_1 \|\theta'(s, \cdot) - \theta(s, \cdot)\|_\infty \\
&= \max_{s \in \mathcal{S}} \left(\sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) |r(s, a) - \pi_{\theta_t}(\cdot|s)^\top r(s, \cdot)| \right) \|\theta'(s, \cdot) - \theta(s, \cdot)\|_\infty \\
&\leq \max_{s \in \mathcal{S}} \underbrace{\max_{a \in \mathcal{A}} |r(s, a) - \pi_{\theta_t}(\cdot|s)^\top r(s, \cdot)|}_{\leq 1 \text{ since } r(s, a) \in [0, 1]} \|\theta'(s, \cdot) - \theta(s, \cdot)\|_\infty \\
&\leq \max_{s \in \mathcal{S}} \|\theta'(s, \cdot) - \theta(s, \cdot)\|_\infty = \|\theta' - \theta\|_\infty .
\end{aligned}$$

Plugging the above two inequalities into (179) and using Lagrange mean value theorem, we have

$$|V^\theta(s) - V^{\theta'}(s)| \leq \frac{1+\gamma}{(1-\gamma)^2} \|\theta' - \theta\|_\infty . \tag{183}$$

Step 3: sum up. Combining (183), (174) and (169), we have

$$\forall s \in \mathcal{S} : |V_\tau^\theta(s) - V_\tau^{\theta'}(s)| \leq \frac{1+\gamma+2\tau(1-\gamma)\log|\mathcal{A}|}{(1-\gamma)^2} \|\log \pi - \log \pi'\|_\infty . \tag{184}$$

Combining (184) and (7a), (107) immediately follows.

Proof of (180). For any vector $x \in \mathbb{R}^{|\mathcal{S}|}$, we have

$$\left[\frac{d\mathbf{P}_t}{dt} x \right]_s = \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} \frac{d\pi_{\theta_t}(a|s)}{dt} \mathcal{P}(s'|s, a) x(s'),$$

from which we can bound the l_∞ norm as

$$\begin{aligned} \left\| \frac{d\mathbf{P}_t}{dt} x \right\|_\infty &\leq \max_s \sum_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) \left| \frac{d\pi_{\theta_t}(a|s)}{dt} \right| \|x\|_\infty \\ &= \max_s \sum_{a \in \mathcal{A}} \left| \frac{d\pi_{\theta_t}(a|s)}{dt} \right| \|x\|_\infty \\ &\leq 2 \|\theta' - \theta\|_\infty \|x\|_\infty \end{aligned}$$

as desired, where the last line follows from the following fact:

$$\begin{aligned} \sum_{a \in \mathcal{A}} \left| \frac{d\pi_{\theta_t}(a|s)}{dt} \right| &= \sum_{a \in \mathcal{A}} \left| \left\langle \frac{\partial \pi_{\theta_t}(a|s)}{\partial \theta_t}, \theta' - \theta \right\rangle \right| \\ &= \sum_{a \in \mathcal{A}} \left| \left\langle \frac{\partial \pi_{\theta_t}(a|s)}{\partial \theta_t(s, \cdot)}, \theta'(s, \cdot) - \theta(s, \cdot) \right\rangle \right| \\ &= \sum_{a \in \mathcal{A}} \pi_{\theta_t}(a|s) |(\theta'(s, a) - \theta(s, a)) - \pi_{\theta_t}(\cdot|s)^\top (\theta'(s, \cdot) - \theta(s, \cdot))| \\ &\leq \max_a |\theta'(s, a) - \theta(s, a)| + |\pi_{\theta_t}(\cdot|s)^\top (\theta'(s, \cdot) - \theta(s, \cdot))| \\ &\leq 2 \|\theta' - \theta\|_\infty. \end{aligned}$$

C.3 Proof of Lemma 9

To simplify the notation, we denote

$$\delta^{(t)} := \widehat{Q}_\tau^{(t)} - \overline{Q}_\tau^{(t)}. \quad (185)$$

We first rearrange the terms of (168) and obtain

$$-\tau \log \overline{\pi}^{(t)}(a|s) + \left(\overline{Q}_\tau^{(t)}(s, a) + \delta^{(t)}(s, a) \right) = \frac{1-\gamma}{\eta} \left(\log \overline{\pi}^{(t+1)}(a|s) - \log \overline{\pi}^{(t)}(a|s) \right) + \frac{1-\gamma}{\eta} z^{(t)}(s). \quad (186)$$

This in turn allows us to express $\overline{V}_\tau^{(t)}(s_0)$ for any $s_0 \in \mathcal{S}$ as follows

$$\begin{aligned} \overline{V}_\tau^{(t)}(s_0) &= \mathbb{E}_{a_0 \sim \overline{\pi}^{(t)}(\cdot|s_0)} \left[-\tau \log \overline{\pi}^{(t)}(a_0|s_0) + \overline{Q}_\tau^{(t)}(s_0, a_0) \right] \\ &= \mathbb{E}_{a_0 \sim \overline{\pi}^{(t)}(\cdot|s_0)} \left[\frac{1-\gamma}{\eta} z^{(t)}(s_0) \right] + \mathbb{E}_{a_0 \sim \overline{\pi}^{(t)}(\cdot|s_0)} \left[\frac{1-\gamma}{\eta} \left(\log \overline{\pi}^{(t+1)}(a_0|s_0) - \log \overline{\pi}^{(t)}(a_0|s_0) \right) - \delta^{(t)}(s_0, a_0) \right] \\ &= \frac{1-\gamma}{\eta} z^{(t)}(s_0) - \frac{1-\gamma}{\eta} \text{KL}(\overline{\pi}^{(t)}(\cdot|s_0) \parallel \overline{\pi}^{(t+1)}(\cdot|s_0)) - \mathbb{E}_{a_0 \sim \overline{\pi}^{(t)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right] \\ &= \mathbb{E}_{a_0 \sim \overline{\pi}^{(t+1)}(\cdot|s_0)} \left[\frac{1-\gamma}{\eta} z^{(t)}(s_0) \right] - \frac{1-\gamma}{\eta} \text{KL}(\overline{\pi}^{(t)}(\cdot|s_0) \parallel \overline{\pi}^{(t+1)}(\cdot|s_0)) - \mathbb{E}_{a_0 \sim \overline{\pi}^{(t)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right], \end{aligned} \quad (187)$$

where the first identity makes use of (7b), the second line follows from (186). Invoking (7b) again to rewrite the $z(s_0)$ appearing in the first term of (187), we reach

$$\begin{aligned} \overline{V}_\tau^{(t)}(s_0) &= \mathbb{E}_{a_0 \sim \overline{\pi}^{(t+1)}(\cdot|s_0)} \left[-\tau \log \overline{\pi}^{(t+1)}(a_0|s_0) + \overline{Q}_\tau^{(t)}(s_0, a_0) + \left(\tau - \frac{1-\gamma}{\eta} \right) \left(\log \overline{\pi}^{(t+1)}(a_0|s_0) - \log \overline{\pi}^{(t)}(a_0|s_0) \right) \right] \end{aligned}$$

$$\begin{aligned}
& -\frac{1-\gamma}{\eta} \text{KL}(\bar{\pi}^{(t)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)) - \mathbb{E}_{a_0 \sim \bar{\pi}^{(t)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right] + \mathbb{E}_{a_0 \sim \bar{\pi}^{(t+1)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right] \\
= & \mathbb{E}_{\substack{a_0 \sim \bar{\pi}^{(t+1)}(\cdot|s_0), \\ s_1 \sim P(\cdot|s_0, a_0)}} \left[-\tau \log \bar{\pi}^{(t+1)}(a_0|s_0) + r(s_0, a_0) + \gamma \bar{V}_\tau^{(t)}(s_0) \right] \\
& - \left(\frac{1-\gamma}{\eta} - \tau \right) \text{KL}(\bar{\pi}^{(t+1)}(\cdot|s_0) \parallel \bar{\pi}^{(t)}(\cdot|s_0)) - \frac{1-\gamma}{\eta} \text{KL}(\bar{\pi}^{(t)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)) \\
& - \mathbb{E}_{a_0 \sim \bar{\pi}^{(t)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right] + \mathbb{E}_{a_0 \sim \bar{\pi}^{(t+1)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right]. \tag{188}
\end{aligned}$$

Note that for any $(s_0, a_0) \in \mathcal{S} \times \mathcal{A}$, we have

$$\begin{aligned}
& - \mathbb{E}_{a_0 \sim \bar{\pi}^{(t)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right] + \mathbb{E}_{a_0 \sim \bar{\pi}^{(t+1)}(\cdot|s_0)} \left[\delta^{(t)}(s_0, a_0) \right] \\
= & \sum_{a_0 \in \mathcal{A}} \left(\bar{\pi}^{(t+1)}(a_0|s_0) - \bar{\pi}^{(t)}(a_0|s_0) \right) \delta^{(t)}(s_0, a_0) \\
\leq & \|\bar{\pi}^{(t+1)}(\cdot|s_0) - \bar{\pi}^{(t)}(\cdot|s_0)\|_1 \|\delta^{(t)}\|_\infty \leq 2\|\delta^{(t)}\|_\infty. \tag{189}
\end{aligned}$$

To finish up, applying (188) recursively to expand $\bar{V}_\tau^{(t)}(s_i)$, $i \geq 1$ and making use of (189), we arrive at

$$\begin{aligned}
& \bar{V}_\tau^{(t)}(s_0) \\
\leq & \sum_{i=1}^{\infty} \gamma^i \cdot 2 \|\delta^{(t)}\|_\infty + \mathbb{E}_{\substack{a_i \sim \bar{\pi}^{(t+1)}(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[\sum_{i=1}^{\infty} \gamma^i \left\{ r(s_i, a_i) - \tau \log \bar{\pi}^{(t+1)}(a_i|s_i) \right\} \right. \\
& \left. - \sum_{i=1}^{\infty} \gamma^i \left\{ \left(\frac{1-\gamma}{\eta} - \tau \right) \text{KL}(\bar{\pi}^{(t+1)}(\cdot|s_i) \parallel \bar{\pi}^{(t)}(\cdot|s_i)) + \frac{1-\gamma}{\eta} \text{KL}(\bar{\pi}^{(t)}(\cdot|s_i) \parallel \bar{\pi}^{(t+1)}(\cdot|s_i)) \right\} \right] \\
= & \frac{2}{1-\gamma} \|\delta^{(t)}\|_\infty + \bar{V}_\tau^{(t+1)}(s_0) \\
& - \mathbb{E}_{s \sim d_{s_0}^{\bar{\pi}^{(t+1)}}} \left[\left(\frac{1}{\eta} - \frac{\tau}{1-\gamma} \right) \text{KL}(\bar{\pi}^{(t+1)}(\cdot|s_i) \parallel \bar{\pi}^{(t)}(\cdot|s_i)) + \frac{1}{\eta} \text{KL}(\bar{\pi}^{(t)}(\cdot|s_i) \parallel \bar{\pi}^{(t+1)}(\cdot|s_i)) \right], \tag{190}
\end{aligned}$$

where the third line follows since $\bar{V}_\tau^{(t+1)}$ can be viewed as the value function of $\bar{\pi}^{(t+1)}$ with adjusted rewards $\bar{r}^{(t+1)}(s, a) := r(s, a) - \tau \log \bar{\pi}^{(t+1)}(s|a)$. And (125) follows immediately from the above inequality (190). By (7a) we can easily see that (126) is a consequence of (125).

C.4 Proof of Lemma 11

We first introduce the famous performance difference lemma which will be used in our proof.

Lemma 12 (Performance difference lemma). *For all policies π, π' and state s_0 , we have*

$$V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[A^{\pi'}(s, a) \right]. \tag{191}$$

The proof of Lemma 12 can be found in, for example, Appendix A of Agarwal et al. (2021).

For all $t \geq 0$, we define the advantage function $\bar{A}^{(t)}$ as:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : \bar{A}^{(t)}(s, a) := \bar{Q}^{(t)}(s, a) - \bar{V}^{(t)}(s). \tag{192}$$

Then for Alg. 1, the update rule of $\bar{\pi}$ (Eq. (168)) can be written as

$$\log \bar{\pi}^{(t+1)}(a|s) = \log \bar{\pi}^{(t)}(a|s) + \frac{\eta}{1-\gamma} \left(\bar{A}^{(t)}(s, a) + \delta^{(t)}(s, a) \right) - \log \hat{z}^{(t)}(s), \tag{193}$$

where $\delta^{(t)}$ is defined in (185) and

$$\begin{aligned}
\log \widehat{z}^{(t)}(s) &= \log \sum_{a' \in \mathcal{A}} \bar{\pi}^{(t)}(a'|s) \exp \left\{ \frac{\eta}{1-\gamma} \left(\bar{A}^{(t)}(s, a') + \delta^{(t)}(s, a') \right) \right\} \\
&\geq \sum_{a' \in \mathcal{A}} \bar{\pi}^{(t)}(a'|s) \log \exp \left\{ \frac{\eta}{1-\gamma} \left(\bar{A}^{(t)}(s, a') + \delta^{(t)}(s, a') \right) \right\} \\
&= \frac{\eta}{1-\gamma} \sum_{a' \in \mathcal{A}} \bar{\pi}^{(t)}(a'|s) \left(\bar{A}^{(t)}(s, a') + \delta^{(t)}(s, a') \right) \\
&= \frac{\eta}{1-\gamma} \sum_{a' \in \mathcal{A}} \bar{\pi}^{(t)}(a'|s) \delta^{(t)}(s, a') \geq -\frac{\eta}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty}, \tag{194}
\end{aligned}$$

where the first inequality follows by Jensen's inequality on the concave function $\log x$ and the last equality uses $\sum_{a' \in \mathcal{A}} \bar{\pi}^{(t)}(a'|s) \bar{A}^{(t)}(s, a') = 0$.

For all starting state distribution μ , we use $d^{(t+1)}$ as shorthand for $d_{\mu}^{\bar{\pi}^{(t+1)}}$, the performance difference lemma (Lemma 12) implies:

$$\begin{aligned}
&\bar{V}^{(t+1)}(\mu) - \bar{V}^{(t)}(\mu) \\
&= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{(t+1)}} \sum_{a \in \mathcal{A}} \bar{\pi}^{(t+1)}(a|s) \left(\bar{A}^{(t)}(s, a) + \delta^{(t)}(s, a) \right) - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{(t+1)}} \mathbb{E}_{a \sim \bar{\pi}^{(t+1)}(\cdot|s)} \left[\delta^{(t)}(s, a) \right] \\
&= \frac{1}{\eta} \mathbb{E}_{s \sim d^{(t+1)}} \sum_{a \in \mathcal{A}} \bar{\pi}^{(t+1)}(a|s) \log \frac{\bar{\pi}^{(t+1)}(a|s) \widehat{z}^{(t)}(s)}{\bar{\pi}^{(t)}(a|s)} - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{(t+1)}} \mathbb{E}_{a \sim \bar{\pi}^{(t+1)}(\cdot|s)} \left[\delta^{(t)}(s, a) \right] \\
&= \frac{1}{\eta} \mathbb{E}_{s \sim d^{(t+1)}} \text{KL}(\bar{\pi}^{(t+1)}(\cdot|s) \| \bar{\pi}^{(t)}(\cdot|s)) + \frac{1}{\eta} \mathbb{E}_{s \sim d^{(t+1)}} \log \widehat{z}^{(t)}(s) - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{(t+1)}} \mathbb{E}_{a \sim \bar{\pi}^{(t+1)}(\cdot|s)} \left[\delta^{(t)}(s, a) \right] \\
&\geq \frac{1}{\eta} \mathbb{E}_{s \sim d^{(t+1)}} \left(\log \widehat{z}^{(t)}(s) + \frac{\eta}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty} \right) - \frac{2}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty},
\end{aligned}$$

from which we can see that

$$\bar{V}^{(t+1)}(\mu) - \bar{V}^{(t)}(\mu) \geq -\frac{2}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty}, \tag{195}$$

where we use (194), and that

$$\bar{V}^{(t+1)}(\mu) - \bar{V}^{(t)}(\mu) \geq \frac{1-\gamma}{\eta} \mathbb{E}_{s \sim \mu} \left(\log \widehat{z}^{(t)}(s) + \frac{\eta}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty} \right) - \frac{2}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty}, \tag{196}$$

which follows from $d^{(t+1)} = d_{\mu}^{\bar{\pi}^{(t+1)}} \geq (1-\gamma)\mu$ and the fact that $\log \widehat{z}^{(t)}(s) + \frac{\eta}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty} \geq 0$ (by (194)).

For any fixed ρ , we use d^{\star} as shorthand for $d_{\rho}^{\pi^{\star}}$. By the performance difference lemma (Lemma 12),

$$\begin{aligned}
&V^{\star}(\rho) - \bar{V}^{(t)}(\rho) \\
&= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\star}} \sum_{a \in \mathcal{A}} \pi^{\star}(a|s) \left(\bar{A}^{(t)}(s, a) + \delta^{(t)}(s, a) \right) - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\star}} \mathbb{E}_{a \sim \pi^{\star}(\cdot|s)} \left[\delta^{(t)}(s, a) \right] \\
&= \frac{1}{\eta} \mathbb{E}_{s \sim d^{\star}} \sum_{a \in \mathcal{A}} \pi^{\star}(a|s) \log \frac{\bar{\pi}^{(t+1)}(a|s) \widehat{z}^{(t)}(s)}{\bar{\pi}^{(t)}(a|s)} - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\star}} \mathbb{E}_{a \sim \pi^{\star}(\cdot|s)} \left[\delta^{(t)}(s, a) \right] \\
&= \frac{1}{\eta} \mathbb{E}_{s \sim d^{\star}} \left(\text{KL}(\pi^{\star}(\cdot|s) \| \bar{\pi}^{(t)}(\cdot|s)) - \text{KL}(\pi^{\star}(\cdot|s) \| \bar{\pi}^{(t+1)}(\cdot|s)) + \log \widehat{z}^{(t)}(s) \right) - \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\star}} \mathbb{E}_{a \sim \pi^{\star}(\cdot|s)} \left[\delta^{(t)}(s, a) \right] \\
&\leq \frac{1}{\eta} \mathbb{E}_{s \sim d^{\star}} \left(\text{KL}(\pi^{\star}(\cdot|s) \| \bar{\pi}^{(t)}(\cdot|s)) - \text{KL}(\pi^{\star}(\cdot|s) \| \bar{\pi}^{(t+1)}(\cdot|s)) + \left(\log \widehat{z}^{(t)}(s) + \frac{\eta}{1-\gamma} \left\| \delta^{(t)} \right\|_{\infty} \right) \right), \tag{197}
\end{aligned}$$

where we use (193) in the second equality.

By applying (196) with $\mu = d^*$ as the initial state distribution, we have

$$\frac{1}{\eta} \mathbb{E}_{s \sim \mu} \left(\log \widehat{z}^{(t)}(s) + \frac{\eta}{1-\gamma} \|\delta^{(t)}\|_{\infty} \right) \leq \frac{1}{1-\gamma} \left(\bar{V}^{(t+1)}(d^*) - \bar{V}^{(t)}(d^*) \right) + \frac{2}{(1-\gamma)^2} \|\delta^{(t)}\|_{\infty}.$$

Plugging the above equation into (197), we obtain

$$\begin{aligned} V^*(\rho) - \bar{V}^{(t)}(\rho) &\leq \frac{1}{\eta} \mathbb{E}_{s \sim d^*} \left(\text{KL}(\pi^*(\cdot|s) \parallel \bar{\pi}^{(t)}(\cdot|s)) - \text{KL}(\pi^*(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) \right) \\ &\quad + \frac{1}{1-\gamma} \left(\bar{V}^{(t+1)}(d^*) - \bar{V}^{(t)}(d^*) \right) + \frac{2}{(1-\gamma)^2} \|\delta^{(t)}\|_{\infty}, \end{aligned}$$

which gives Lemma 11.