

## Appendix

### A. Preliminaries

#### A.1. Derivation of entropy-regularized NPG methods

This subsection establishes the equivalence between the update rules (15) and (18). Such derivations are inherently similar to the ones for the NPG update rule (without entropy regularization) (see, e.g., Agarwal et al. (2019)); we provide the proof here for pedagogical reasons.

First of all, let us follow the convention to introduce the advantage function  $A_\tau^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  of a policy  $\pi$  w.r.t. the entropy-regularized value function:

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad A_\tau^\pi(s, a) := Q_\tau^\pi(s, a) - \tau \log \pi(a|s) - V_\tau^\pi(s) \quad (\text{A.1})$$

with  $Q_\tau^\pi$  defined in (11a), which reflects the gain one can harvest by executing action  $a$  instead of following the policy  $\pi$  in state  $s$ . This advantage function plays a crucial role in the calculation of policy gradients, due to the following fundamental relation (see Appendix C.6 for the proof):

**Lemma 1** *Under softmax parameterization (7), the gradient of the regularized value function satisfies*

$$\frac{\partial V_\tau^{\pi_\theta}(\rho)}{\partial \theta(s, a)} = \frac{1}{1-\gamma} d_\rho^{\pi_\theta}(s) \cdot \pi_\theta(a|s) \cdot A_\tau^{\pi_\theta}(s, a); \quad (\text{A.2a})$$

$$\left[ (\mathcal{F}_\rho^\theta)^\dagger \nabla_\theta V_\tau^{\pi_\theta}(\rho) \right] (s, a) = \frac{1}{1-\gamma} A_\tau^{\pi_\theta}(s, a) + c(s) \quad (\text{A.2b})$$

for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , where  $c(s) := \sum_a \pi_\theta(a|s) w_{s,a}$  is some function depending only on  $s$ .

It is worth highlighting that the search direction of NPG, given in (A.2b), is invariant to the choice of  $\rho$ . With the above calculations in place, it is seen that for any  $s \in \mathcal{S}$ , the regularized NPG update rule (15) results in a policy update as follows

$$\begin{aligned} \pi^{(t+1)}(a|s) &\stackrel{\text{(i)}}{\propto} \exp(\theta^{(t+1)}(s, a)) \stackrel{\text{(ii)}}{=} \exp\left(\theta^{(t)}(s, a) + \eta \left[ (\mathcal{F}_\rho^{\theta^{(t)}})^\dagger \nabla_\theta V_\tau^{(\theta^{(t)})}(\rho) \right] (s, a)\right) \\ &\stackrel{\text{(iii)}}{\propto} \exp\left(\theta^{(t)}(s, a) + \frac{\eta}{1-\gamma} A_\tau^{(\theta^{(t)})}(s, a)\right) \\ &\stackrel{\text{(iv)}}{\propto} \pi^{(t)}(a|s) \exp\left(\frac{\eta}{1-\gamma} Q_\tau^{(\theta^{(t)})}(s, a) - \frac{\eta\tau}{1-\gamma} \log \pi^{(t)}(a|s)\right) \\ &= (\pi^{(t)}(a|s))^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta}{1-\gamma} Q_\tau^{(\theta^{(t)})}(s, a)\right). \end{aligned}$$

where we use  $A_\tau^{(\theta^{(t)})}$  to abbreviate  $A_\tau^{\pi_{\theta^{(t)}}}$ . Here, (i) uses the definition of the softmax policy, (ii) comes from the update rule (15), (iii) is a consequence of (A.2b) (since  $c(\cdot)$  does not depend on  $a$ ), whereas (iv) results from the definition (A.1) and the fact that  $V_\tau^\pi(\cdot)$  is not dependent on  $a$ . This validates the equivalence between (15) and (18).

## A.2. Basic facts about the function $\log(\|\exp(\theta)\|_1)$

In the current paper, we often encounter the function  $\log(\|\exp(\theta)\|_1) := \log(\sum_{1 \leq a \leq |\mathcal{A}|} \exp(\theta_a))$  for any vector  $\theta = [\theta_a]_{1 \leq a \leq |\mathcal{A}|} \in \mathbb{R}^{|\mathcal{A}|}$ . To facilitate analysis, we single out several basic properties concerning this function, which will be used multiple times when establishing our main results. For notational convenience, we denote by  $\pi_\theta \in \mathbb{R}^{|\mathcal{A}|}$  the softmax transform of  $\theta$  such that

$$\pi_\theta(a) = \frac{\exp(\theta_a)}{\sum_{1 \leq j \leq |\mathcal{A}|} \exp(\theta_j)}, \quad 1 \leq a \leq |\mathcal{A}|. \quad (\text{A.3})$$

By straightforward calculations, the gradient of the function  $\log(\|\exp(\theta)\|_1)$  is given by

$$\nabla_\theta \log(\|\exp(\theta)\|_1) = \frac{1}{\|\exp(\theta)\|_1} \exp(\theta) = \pi_\theta; \quad (\text{A.4})$$

*Difference of log policies.* In the analysis, we often need to control the difference of two policies, towards which the following bounds prove useful. To begin with, the mean value theorem reveals a Lipschitz continuity property (w.r.t. the  $\ell_\infty$  norm): for any  $\theta_1, \theta_2 \in \mathbb{R}^{|\mathcal{A}|}$ ,

$$\begin{aligned} |\log(\|\exp(\theta_1)\|_1) - \log(\|\exp(\theta_2)\|_1)| &= |\langle \theta_1 - \theta_2, \nabla_\theta \log(\|\exp(\theta)\|_1)|_{\theta=\theta_c} \rangle| \\ &\leq \|\theta_1 - \theta_2\|_\infty \|\nabla_\theta \log(\|\exp(\theta)\|_1)|_{\theta=\theta_c}\|_1 = \|\theta_1 - \theta_2\|_\infty, \end{aligned} \quad (\text{A.5})$$

where  $\theta_c$  is a certain convex combination of  $\theta_1$  and  $\theta_2$ , and the second line relies on (A.4). In addition, for any two vectors  $\pi_{\theta_1}$  and  $\pi_{\theta_2}$  defined w.r.t.  $\theta_1, \theta_2 \in \mathbb{R}^{|\mathcal{A}|}$  (see (A.3)), one has

$$\|\log \pi_{\theta_1} - \log \pi_{\theta_2}\|_\infty \leq 2 \|\theta_1 - \theta_2\|_\infty, \quad (\text{A.6})$$

where  $\log(\cdot)$  denotes entrywise operation. To justify (A.6), we observe from the definition (A.3) that

$$\|\log \pi_{\theta_1} - \log \pi_{\theta_2}\|_\infty \leq \|\theta_1 - \theta_2\|_\infty + \left| \log(\|\exp(\theta_1)\|_1) - \log(\|\exp(\theta_2)\|_1) \right| \leq 2 \|\theta_1 - \theta_2\|_\infty,$$

where the last inequality is a consequence of (A.5).

## B. Proof for the bandit case (Proposition 1)

We start by defining an auxiliary sequence  $\xi^{(t)} \in \mathbb{R}^{|\mathcal{A}|}$  ( $t \geq 0$ ) recursively as follows

$$\begin{aligned} \xi^{(0)} &:= \|\exp(r/\tau)\|_1 \cdot \pi^{(0)}, \\ \xi^{(t+1)}(a) &:= (\xi^{(t)}(a))^{1-\tau\eta} \exp(\eta r(a)), \quad a \in \mathcal{A}. \end{aligned}$$

When combined with (21), it is easily seen that  $\pi^{(t)}(\cdot) \propto \xi^{(t)}(\cdot)$  and, as a result,  $\pi^{(t)} = \xi^{(t)} / \|\xi^{(t)}\|_1$ .

By construction, the auxiliary sequence satisfies the following property

$$\log(\xi^{(t+1)}(a)) - r(a)/\tau = (1 - \tau\eta) \log(\xi^{(t)}(a)) + \eta r(a) - r(a)/\tau$$

$$= (1 - \tau\eta) (\log (\xi^{(t)}(a)) - r(a)/\tau),$$

thus indicating that

$$\|\log \xi^{(t)} - r/\tau\|_\infty \leq (1 - \tau\eta)^t \|\log \xi^{(0)} - r/\tau\|_\infty. \quad (\text{A.7})$$

This taken together with the optimal policy  $\pi_\tau^* = \text{softmax}(r/\tau) \propto \exp(r/\tau)$  leads to

$$\begin{aligned} \|\log \pi^{(t)} - \log \pi_\tau^*\|_\infty &\leq 2 \|\log \xi^{(t)} - r/\tau\|_\infty \leq 2(1 - \tau\eta)^t \|\log \xi^{(0)} - r/\tau\|_\infty \\ &= 2(1 - \tau\eta)^t \|\log \pi^{(0)} + (\log \|\exp(r/\tau)\|_1) \cdot \mathbf{1} - r/\tau\|_\infty \\ &= 2(1 - \tau\eta)^t \|\log \pi^{(0)} - \log \pi_\tau^*\|_\infty, \end{aligned}$$

where the first line follows from the inequality (A.6), the second line follows from the expression (A.7), whereas the last line follows from the form of  $\pi_\tau^*$ . We have thus completed the proof of Proposition 1.

## C. Proof for key lemmas

### C.1. Proof of Lemma 1

To begin with, the regularized NPG update rule (see (18) in Algorithm 1) indicates that

$$\log \pi^{(t+1)}(a|s) = \left(1 - \frac{\eta\tau}{1-\gamma}\right) \log \pi^{(t)}(a|s) + \frac{\eta}{1-\gamma} Q_\tau^{(t)}(s, a) - \log Z^{(t)}(s), \quad (\text{A.8})$$

where  $Z^{(t)}$  is some quantity depending only on the state  $s$  (but not the action  $a$ ). Rearranging terms gives

$$-\tau \log \pi^{(t)}(a|s) + Q_\tau^{(t)}(s, a) = \frac{1-\gamma}{\eta} (\log \pi^{(t+1)}(a|s) - \log \pi^{(t)}(a|s)) + \frac{1-\gamma}{\eta} \log Z^{(t)}(s). \quad (\text{A.9})$$

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

$$\begin{aligned} V_\tau^{(t)}(s_0) &= \mathbb{E}_{a_0 \sim \pi^{(t)}(\cdot|s_0)} [-\tau \log \pi^{(t)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0)] \\ &= \mathbb{E}_{a_0 \sim \pi^{(t)}(\cdot|s_0)} \left[ \frac{1-\gamma}{\eta} \log Z^{(t)}(s_0) \right] + \mathbb{E}_{a_0 \sim \pi^{(t)}(\cdot|s_0)} \left[ \frac{1-\gamma}{\eta} (\log \pi^{(t+1)}(a_0|s_0) - \log \pi^{(t)}(a_0|s_0)) \right] \\ &= \frac{1-\gamma}{\eta} \log Z^{(t)}(s_0) - \frac{1-\gamma}{\eta} \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \pi^{(t+1)}(\cdot|s_0)) \\ &= \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ \frac{1-\gamma}{\eta} \log Z^{(t)}(s_0) \right] - \frac{1-\gamma}{\eta} \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \pi^{(t+1)}(\cdot|s_0)), \end{aligned} \quad (\text{A.10})$$

where the first identity makes use of the definitions (8) and (11a), the second line follows from (A.9), the third line relies on the definition of the KL divergence, and the last line follows since  $Z^{(t)}(s)$  does not depend on  $a$ . Invoking (A.9) again to rewrite  $\log Z^{(t)}(s_0)$  appearing in the first term of (A.10), we reach

$$V_\tau^{(t)}(s_0) = \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ -\tau \log \pi^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0) + \left( \tau - \frac{1-\gamma}{\eta} \right) (\log \pi^{(t+1)}(a_0|s_0) - \log \pi^{(t)}(a_0|s_0)) \right]$$

$$\begin{aligned}
& - \frac{1-\gamma}{\eta} \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \pi^{(t+1)}(\cdot|s_0)) \\
= & \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ -\tau \log \pi^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0) \right] + \left( \tau - \frac{1-\gamma}{\eta} \right) \text{KL}(\pi^{(t+1)}(\cdot|s_0) \parallel \pi^{(t)}(\cdot|s_0)) \\
& - \frac{1-\gamma}{\eta} \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \pi^{(t+1)}(\cdot|s_0)) \\
= & \mathbb{E}_{\substack{a_0 \sim \pi^{(t+1)}(\cdot|s_0), \\ s_1 \sim P(\cdot|s_0, a_0)}} \left[ -\tau \log \pi^{(t+1)}(a_0|s_0) + r(s_0, a_0) + \gamma V_\tau^{(t)}(s_1) \right] \\
& - \left( \frac{1-\gamma}{\eta} - \tau \right) \text{KL}(\pi^{(t+1)}(\cdot|s_0) \parallel \pi^{(t)}(\cdot|s_0)) - \frac{1-\gamma}{\eta} \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \pi^{(t+1)}(\cdot|s_0)), \tag{A.11}
\end{aligned}$$

where the second line uses the definition of the KL divergence, and the third line expands  $Q_\tau^{(t)}$  using the definition (11a).

To finish up, applying the above relation (A.11) recursively to expand  $V_\tau^{(t)}(s_i)$  ( $i \geq 1$ ), we arrive at

$$\begin{aligned}
V_\tau^{(t)}(s_0) = & \mathbb{E}_{\substack{a_i \sim \pi^{(t+1)}(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[ \sum_{i=0}^{\infty} \gamma^i \{ r(s_i, a_i) - \tau \log \pi^{(t+1)}(a_i|s_i) \} \right. \\
& \left. - \sum_{i=0}^{\infty} \gamma^i \left\{ \left( \frac{1-\gamma}{\eta} - \tau \right) \text{KL}(\pi^{(t+1)}(\cdot|s_i) \parallel \pi^{(t)}(\cdot|s_i)) + \frac{1-\gamma}{\eta} \text{KL}(\pi^{(t)}(\cdot|s_i) \parallel \pi^{(t+1)}(\cdot|s_i)) \right\} \right] \\
= & V_\tau^{(t+1)}(s_0) - \mathbb{E}_{s \sim d_{s_0}^{(t+1)}} \left[ \left( \frac{1}{\eta} - \frac{\tau}{1-\gamma} \right) \text{KL}(\pi^{(t+1)}(\cdot|s) \parallel \pi^{(t)}(\cdot|s)) + \frac{1}{\eta} \text{KL}(\pi^{(t)}(\cdot|s) \parallel \pi^{(t+1)}(\cdot|s)) \right], \tag{A.12}
\end{aligned}$$

where the second line follows since the regularized value function  $V_\tau^{(t+1)}$  can be viewed as the value function of  $\pi^{(t+1)}$  with adjusted rewards  $r_\tau^{(t+1)}(s, a) := r(s, a) - \tau \log \pi^{(t+1)}(a|s)$ . Averaging the initial state  $s_0$  over the distribution  $\rho$  concludes the proof.

## C.2. Proof of Lemma 2

In the sequel, we prove each claim in Lemma 2 in order.

*Proof of Eqn. (36).* Jensen's inequality tells us that: for any  $s \in \mathcal{S}$  one has

$$\begin{aligned}
\mathbb{E}_{a \sim \pi(\cdot|s)} \left[ Q(s, a) - \tau \log \pi(a|s) \right] &= \tau \sum_a \pi(a|s) \log \left( \frac{\exp(Q(s, a)/\tau)}{\pi(a|s)} \right) \\
&\leq \tau \log \left( \sum_a \pi(a|s) \frac{\exp(Q(s, a)/\tau)}{\pi(a|s)} \right) \\
&= \tau \log \left( \sum_a \exp(Q(s, a)/\tau) \right) = \tau \log (\|\exp(Q(s, \cdot)/\tau)\|_1), \tag{A.13}
\end{aligned}$$

where in the second line, equality is attained if  $\pi(\cdot|s) \propto \exp(Q(s, \cdot)/\tau)$ . This immediately gives rise to

$$\mathcal{T}_\tau(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \max_{\pi(\cdot|s') \in \Delta(\mathcal{A})} \mathbb{E}_{a' \sim \pi(\cdot|s')} \left[ Q(s', a') - \tau \log \pi(a'|s') \right] \right]$$

$$= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \tau \log \left( \left\| \exp \left( Q(s', \cdot) / \tau \right) \right\|_1 \right) \right].$$

*Proof of Eqn. (37).* Recall the characterization of  $\pi_\tau^*$  and  $V_\tau^*$  established in Nachum et al. (2017):

$$\pi_\tau^*(a|s) = \exp \left( \frac{Q_\tau^*(s, a) - V_\tau^*(s)}{\tau} \right), \quad (\text{A.14a})$$

$$V_\tau^*(s) = \tau \log \left( \left\| \exp \left( Q_\tau^*(s, \cdot) / \tau \right) \right\|_1 \right). \quad (\text{A.14b})$$

Substitution into the expression (36) tells us that for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned} \mathcal{T}_\tau(Q_\tau^*)(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \tau \log \left( \left\| \exp \left( Q_\tau^*(s', \cdot) / \tau \right) \right\|_1 \right) \right] \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V_\tau^*(s') \right] \\ &= Q_\tau^*(s, a), \end{aligned}$$

where the second line results from (A.14b), and the last line follows from the definition of the soft Q-function.

*Proof of Eqn. (38).* Invoking again the expression (36), we can demonstrate that for any  $Q_1$  and  $Q_2$ ,

$$\begin{aligned} & \left| \mathcal{T}_\tau(Q_1)(s, a) - \mathcal{T}_\tau(Q_2)(s, a) \right| \\ &= \left| \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \tau \log \left( \left\| \exp \left( Q_1(s', \cdot) / \tau \right) \right\|_1 \right) \right] - \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \tau \log \left( \left\| \exp \left( Q_2(s', \cdot) / \tau \right) \right\|_1 \right) \right] \right| \\ &= \gamma \tau \left| \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \log \left( \left\| \exp \left( Q_1(s', \cdot) / \tau \right) \right\|_1 \right) - \log \left( \left\| \exp \left( Q_2(s', \cdot) / \tau \right) \right\|_1 \right) \right] \right| \\ &\leq \gamma \tau \|Q_1 / \tau - Q_2 / \tau\|_\infty \\ &= \gamma \|Q_1 - Q_2\|_\infty \end{aligned}$$

holds for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , where the inequality follows from the Lipschitz property (A.5).

### C.3. Proof of Lemma 3

For any state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we observe that

$$\begin{aligned} & Q_\tau^*(s, a) - Q_\tau^{(t+1)}(s, a) \\ &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V_\tau^*(s') \right] - \left( r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ V_\tau^{(t+1)}(s') \right] \right) \\ &= \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \tau \log \left( \left\| \exp \left( \frac{Q_\tau^*(s', \cdot)}{\tau} \right) \right\|_1 \right) \right] - \gamma \mathbb{E}_{\substack{s' \sim P(\cdot|s, a), \\ a' \sim \pi^{(t+1)}(\cdot|s')}} \left[ Q_\tau^{(t+1)}(s', a') - \tau \log \pi^{(t+1)}(a'|s') \right], \end{aligned} \quad (\text{A.15})$$

where the first step invokes the definition (11a) of  $Q_\tau$ , and the second step is due to the expression (A.14b) of  $V_\tau^*$ . To continue, recall that  $\pi^{(t)}$  is related to  $\xi^{(t)}$  as

$$\forall s \in \mathcal{S}: \quad \pi^{(t)}(\cdot|s) = \frac{1}{\|\xi^{(t)}(s, \cdot)\|_1} \xi^{(t)}(s, \cdot) \quad (\text{A.16})$$

which can be seen by comparing (42) with (18). This in turn leads to

$$\begin{aligned} \log \pi^{(t+1)}(a|s) &= \log \xi^{(t+1)}(s, a) - \log (\|\xi^{(t+1)}(s, \cdot)\|_1) \\ &= \alpha \log \xi^{(t)}(s, a) + (1 - \alpha) \frac{Q_\tau^{(t)}(s, a)}{\tau} - \log (\|\xi^{(t+1)}(s, \cdot)\|_1), \end{aligned} \quad (\text{A.17})$$

where the second line comes from (42b). By plugging (A.17) into (A.15) we obtain

$$\begin{aligned} Q_\tau^*(s, a) - Q_\tau^{(t+1)}(s, a) &= \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \tau \log (\|\exp(Q_\tau^*(s', \cdot)/\tau)\|_1) - \tau \log (\|\xi^{(t+1)}(s', \cdot)\|_1) \right] \\ &\quad - \gamma \mathbb{E}_{\substack{s' \sim P(\cdot|s, a), \\ a' \sim \pi^{(t+1)}(\cdot|s')}} \left[ \underbrace{Q_\tau^{(t+1)}(s', a') - \tau \left( \alpha \log \xi^{(t)}(s', a') + (1 - \alpha) \frac{Q_\tau^{(t)}(s', a')}{\tau} \right)}_{= \log \xi^{(t+1)}(s', a')} \right] \end{aligned} \quad (\text{A.18})$$

for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . In the sequel, we bound each term on the right-hand side of (A.18) separately.

- In view of the property (A.5), the first term on the right-hand side of (A.18) can be bounded by

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

- Regarding the second term, the monotonicity (33) of the soft Q-function allows us to derive

$$\begin{aligned} &Q_\tau^{(t+1)}(s, a) - \tau \left( \alpha \log \xi^{(t)}(s, a) + (1 - \alpha) \frac{Q_\tau^{(t)}(s, a)}{\tau} \right) \\ &\geq Q_\tau^{(t)}(s, a) - \tau \left( \alpha \log \xi^{(t)}(s, a) + (1 - \alpha) \frac{Q_\tau^{(t)}(s, a)}{\tau} \right) \\ &= \alpha \left( Q_\tau^{(t)}(s, a) - \tau \log \xi^{(t)}(s, a) \right) \\ &\stackrel{(i)}{=} \alpha \left( \alpha \left( Q_\tau^{(t-1)}(s, a) - \tau \log \xi^{(t-1)}(s, a) \right) + Q_\tau^{(t)}(s, a) - Q_\tau^{(t-1)}(s, a) \right) \\ &\stackrel{(ii)}{\geq} \alpha^2 \left( Q_\tau^{(t-1)}(s, a) - \tau \log \xi^{(t-1)}(s, a) \right) \\ &\stackrel{(iii)}{\geq} \alpha^{t+1} \left( Q_\tau^{(0)}(s, a) - \tau \log \xi^{(0)}(s, a) \right) \\ &\stackrel{(iv)}{\geq} -\alpha^{t+1} \|Q_\tau^{(0)} - \tau \log \xi^{(0)}\|_\infty \end{aligned}$$

for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . Here, (i) follows by construction (42b), (ii) invokes the monotonicity property (33) (so that  $Q_\tau^{(t)} \geq Q_\tau^{(t-1)}$ ), and (iii) follows by repeating the arguments (i) and (ii) recursively.

Combining the preceding two bounds with the expression (A.18), we conclude that

$$0 \leq Q_\tau^*(s, a) - Q_\tau^{(t+1)}(s, a) \leq \gamma \|Q_\tau^* - \tau \log \xi^{(t+1)}\|_\infty + \gamma \alpha^{t+1} \|Q_\tau^{(0)} - \tau \log \xi^{(0)}\|_\infty \quad (\text{A.19})$$

for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , thus concluding the proof.

#### C.4. Proof of Lemma 4

Recall that, in this scenario, the policies are updated using inexact policy evaluation via (26), namely,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad \pi^{(t+1)}(a|s) = \frac{(\pi^{(t)}(a|s))^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta}{1-\gamma} \widehat{Q}_\tau^{(t)}(s, a)\right)}{\widehat{Z}^{(t)}(s)}, \quad (\text{A.20})$$

where  $\widehat{Z}^{(t)}(s) := \sum_{a'} \pi^{(t)}(a'|s)^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta}{1-\gamma} \widehat{Q}_\tau^{(t)}(s, a')\right)$ . To facilitate analysis, we further introduce another auxiliary policy sequence  $\{\check{\pi}^{(t)}\}$ , which corresponds to the policy update as if we had access to exact soft Q-function of  $\pi^{(t)}$  in the  $t$ -th iteration; this is defined as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad \check{\pi}^{(t+1)}(a|s) = \frac{(\pi^{(t)}(a|s))^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta}{1-\gamma} Q_\tau^{(t)}(s, a)\right)}{Z^{(t)}(s)}, \quad (\text{A.21})$$

where we abuse the notation by letting  $Z^{(t)}(s) := \sum_{a'} \pi^{(t)}(a'|s)^{1-\frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta}{1-\gamma} Q_\tau^{(t)}(s, a')\right)$ . It is worth emphasizing that  $\check{\pi}^{(t+1)}$  is produced on the basis of  $\pi^{(t)}$  as opposed to  $\check{\pi}^{(t)}$ ; it should be viewed as a one-step perfect update from a given policy  $\pi^{(t)}$ .

We first make note of the following fact: for any step size  $0 < \eta \leq (1-\gamma)/\tau$ , it follows from (A.6) — together with the construction (A.20) and (A.21) — that

$$\begin{aligned} & \left\| \log \pi^{(t+1)} - \log \check{\pi}^{(t+1)} \right\|_\infty \\ & \leq 2 \left\| \log \left( \pi^{(t)}(a|s)^{1-\eta\tau/(1-\gamma)} \exp\left(\frac{\eta}{1-\gamma} \widehat{Q}_\tau^{(t)}(s, a)\right) \right) - \log \left( \pi^{(t)}(a|s)^{1-\eta\tau/(1-\gamma)} \exp\left(\frac{\eta}{1-\gamma} Q_\tau^{(t)}(s, a)\right) \right) \right\|_\infty \\ & = \frac{2\eta}{1-\gamma} \left\| \widehat{Q}_\tau^{(t)} - Q_\tau^{(t)} \right\|_\infty. \end{aligned} \quad (\text{A.22})$$

Next, let us recall the inequality (A.10) in the proof of Lemma 1 under exact policy evaluation  $\check{\pi}^{(t+1)}(\cdot|s)$ ; when applied to the current setting, it essentially indicates that

$$\begin{aligned} V_\tau^{(t)}(s_0) &= \mathbb{E}_{a_0 \sim \check{\pi}^{(t+1)}(\cdot|s_0)} \left[ \frac{1-\gamma}{\eta} \log Z^{(t)}(s_0) \right] - \frac{1-\gamma}{\eta} \text{KL} \left( \pi^{(t)}(\cdot|s_0) \parallel \check{\pi}^{(t+1)}(\cdot|s_0) \right) \\ &= \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ \frac{1-\gamma}{\eta} \log Z^{(t)}(s_0) \right] - \frac{1-\gamma}{\eta} \text{KL} \left( \pi^{(t)}(\cdot|s_0) \parallel \check{\pi}^{(t+1)}(\cdot|s_0) \right), \end{aligned} \quad (\text{A.23})$$

where the last step follows since the quantity  $Z^{(t)}(s)$  does not depend on  $a$  at all. In order to control the first term of (A.23), we invoke the definition of  $\check{\pi}^{(t+1)}(\cdot|s)$  to show that

$$\begin{aligned} & \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ \frac{1-\gamma}{\eta} \log Z^{(t)}(s_0) \right] \\ & \stackrel{(i)}{=} \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ -\tau \log \check{\pi}^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0) + \left( \tau - \frac{1-\gamma}{\eta} \right) (\log \check{\pi}^{(t+1)}(a_0|s_0) - \log \pi^{(t)}(a_0|s_0)) \right] \\ & = \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ -\tau \log \pi^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0) \right] + \left( \tau - \frac{1-\gamma}{\eta} \right) \text{KL} \left( \pi^{(t+1)}(\cdot|s_0) \parallel \pi^{(t)}(\cdot|s_0) \right) \end{aligned}$$

$$\begin{aligned}
& -\frac{1-\gamma}{\eta} \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} [\log \tilde{\pi}^{(t+1)}(a_0|s_0) - \log \pi^{(t)}(a_0|s_0)] \\
\leq & \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} [-\tau \log \pi^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0)] + \left( \tau - \frac{1-\gamma}{\eta} \right) \text{KL}(\pi^{(t+1)}(\cdot|s_0) \parallel \pi^{(t)}(\cdot|s_0)) \\
& + 2 \|\widehat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty, \tag{A.24}
\end{aligned}$$

where the final step results from (A.22). Putting the above bound together with (A.23) guarantees that

$$\begin{aligned}
V_\tau^{(t)}(s_0) & \leq \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} [-\tau \log \pi^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0)] - \frac{1-\gamma}{\eta} \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \tilde{\pi}^{(t+1)}(\cdot|s_0)) \\
& \quad - \left( \frac{1-\gamma}{\eta} - \tau \right) \text{KL}(\pi^{(t+1)}(\cdot|s_0) \parallel \pi^{(t)}(\cdot|s_0)) + 2 \|\widehat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty \\
& \leq \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} [-\tau \log \pi^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0)] + 2 \|\widehat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty \\
& \leq \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ -\tau \log \pi^{(t+1)}(a_0|s_0) + r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P(\cdot|s_0, a_0)} [V_\tau^{(t)}(s_1)] \right] + 2 \|\widehat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty.
\end{aligned}$$

where the last identity makes use of the relation  $Q_\tau^{(t)}(s_0, a_0) = r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P(\cdot|s_0, a_0)} [V_\tau^{(t)}(s_1)]$ . Invoking the above inequality recursively as in the expression (A.12) (see Lemma 1), we can expand it to establish

$$V_\tau^{(t)}(s_0) \leq V_\tau^{(t+1)}(s_0) + 2 \|\widehat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty \sum_{i=0}^{\infty} \gamma^i = V_\tau^{(t+1)}(s_0) + \frac{2}{1-\gamma} \|\widehat{Q}_\tau^{(t)} - Q_\tau^{(t)}\|_\infty.$$

### C.5. Proof of Lemma 5

First of all, we follow the definition (8) of the entropy-regularized value function to deduce that

$$\begin{aligned}
V_\tau^*(\rho) - V_\tau^{(t)}(\rho) & = \mathbb{E}_{\substack{s_0 \sim \rho, a_i \sim \pi_\tau^*(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[ \sum_{i=0}^{\infty} \gamma^i (r(s_i, a_i) - \tau \log \pi_\tau^*(a_i|s_i)) \right] - V_\tau^{(t)}(\rho) \\
& = \mathbb{E}_{\substack{s_0 \sim \rho, a_i \sim \pi_\tau^*(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[ \sum_{i=0}^{\infty} \gamma^i (r(s_i, a_i) - \tau \log \pi_\tau^*(a_i|s_i) + V_\tau^{(t)}(s_i) - V_\tau^{(t)}(s_i)) \right] - V_\tau^{(t)}(\rho) \\
& = \mathbb{E}_{\substack{s_0 \sim \rho, a_i \sim \pi_\tau^*(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[ V_\tau^{(t)}(s_0) + \sum_{i=0}^{\infty} \gamma^i (r(s_i, a_i) - \tau \log \pi_\tau^*(a_i|s_i) + \gamma V_\tau^{(t)}(s_{i+1}) - V_\tau^{(t)}(s_i)) \right] - V_\tau^{(t)}(\rho) \\
& \stackrel{(i)}{=} \mathbb{E}_{\substack{s_0 \sim \rho, a_i \sim \pi_\tau^*(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[ \sum_{i=0}^{\infty} \gamma^i (r(s_i, a_i) - \tau \log \pi_\tau^*(a_i|s_i) + \gamma V_\tau^{(t)}(s_{i+1}) - V_\tau^{(t)}(s_i)) \right] \\
& \stackrel{(ii)}{=} \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^*} \left[ \sum_a \pi_\tau^*(a|s) \left( r(s, a) - \tau \log \pi_\tau^*(a|s) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V_\tau^{(t)}(s')] - V_\tau^{(t)}(s) \right) \right] \\
& \stackrel{(iii)}{=} \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\rho^*} \left[ \sum_a \pi_\tau^*(a|s) (Q_\tau^{(t)}(s, a) - \tau \log \pi_\tau^*(a|s)) - V_\tau^{(t)}(s) \right]. \tag{A.25}
\end{aligned}$$

Here, (i) is due to the definition  $V_\tau^{(t)}(\rho) = \mathbb{E}_{s_0 \sim \rho} [V_\tau^{(t)}(s_0)]$ , (ii) follows by aggregating terms corresponding to the same state-action pair and the definition of  $d_\rho^{\pi_\tau^*}$  (cf. (5)), whereas (iii) results from the definition (11a) of the regularized Q-function.

To continue, we shall attempt to control each part of (A.25) separately. To begin with, observe that the first part of (A.25) can be bounded by Jensen's inequality, namely,

$$\begin{aligned} \sum_a \pi_\tau^*(a|s) (Q_\tau^{(t)}(s, a) - \tau \log \pi_\tau^*(a|s)) &= \tau \sum_a \pi_\tau^*(a|s) \log \left( \frac{\exp(Q_\tau^{(t)}(s, a)/\tau)}{\pi_\tau^*(a|s)} \right) \\ &\leq \tau \log \left( \sum_a \pi_\tau^*(a|s) \frac{\exp(Q_\tau^{(t)}(s, a)/\tau)}{\pi_\tau^*(a|s)} \right) \\ &= \tau \log \left( \sum_a \exp(Q_\tau^{(t)}(s, a)/\tau) \right). \end{aligned} \quad (\text{A.26})$$

With regards to the second part of (A.25), it is seen from the definition of  $\pi^{(t+1)}$  (cf. (17)) that

$$Q_\tau^{(t)}(s, a) = \tau \log \pi_\tau^{(t+1)}(a|s) + \tau \log \left( \sum_a \exp(Q_\tau^{(t)}(s, a)/\tau) \right), \quad (\text{A.27})$$

thus allowing one to derive

$$\begin{aligned} V_\tau^{(t)}(s) &= \sum_a \pi_\tau^{(t)}(a|s) (Q_\tau^{(t)}(s, a) - \tau \log \pi_\tau^{(t)}(a|s)) \\ &\stackrel{(i)}{=} \tau \sum_a \pi_\tau^{(t)}(a|s) \left\{ \log \pi_\tau^{(t+1)}(a|s) + \log \left( \sum_a \exp(Q_\tau^{(t)}(s, a)/\tau) \right) - \log \pi_\tau^{(t)}(a|s) \right\} \\ &= \tau \log \left( \sum_a \exp(Q_\tau^{(t)}(s, a)/\tau) \right) + \tau \sum_a \pi_\tau^{(t)}(a|s) (\log \pi_\tau^{(t+1)}(a|s) - \log \pi_\tau^{(t)}(a|s)) \\ &= \tau \log \left( \sum_a \exp(Q_\tau^{(t)}(s, a)/\tau) \right) - \tau \text{KL} \left( \pi_\tau^{(t)}(a|s) \parallel \pi_\tau^{(t+1)}(\cdot|s) \right), \end{aligned} \quad (\text{A.28})$$

where (i) relies on the identity (A.27). Substituting the inequalities (A.26) and (A.28) into the expression (A.25), we can demonstrate with a little algebra that

$$V_\tau^*(\rho) - V_\tau^{(t)}(\rho) \leq \frac{1}{\eta} \mathbb{E}_{s \sim d_\rho^{\pi_\tau^*}} \left[ \text{KL} \left( \pi^{(t)}(\cdot|s) \parallel \pi^{(t+1)}(\cdot|s) \right) \right].$$

## C.6. Proof of Lemma 1

The results of this lemma, or some similar versions, have appeared in prior work (e.g. Mei et al. (2020, Lemma 10) and Agarwal et al. (2020, Lemma 5.6)). We include the proof here primarily for the sake of self-completeness.

*Proof of Eqn. (A.2a).* The policy gradient of the unregularized value function  $V^{\pi_\theta}(s_0)$  is well-known as the policy gradient theorem (Sutton et al. 2000). Here, we deal with a slightly different variant – an entropy-regularized value function  $V_\tau^{\pi_\theta}(s_0)$  in the expression (2) with the softmax policy parameterization in (7). Invoking the Bellman equation and recognizing that  $V_\tau^{\pi_\theta}(s_0)$  can be viewed as an unregularized value function with instantaneous rewards  $r(s, a) - \tau \log \pi_\theta(a|s)$  for any  $(s, a)$ , we obtain

$$\begin{aligned} \nabla_\theta V_\tau^{\pi_\theta}(s_0) &= \nabla_\theta \left[ \sum_{a_0} \pi_\theta(a_0|s_0) \left( r(s_0, a_0) - \tau \log \pi_\theta(a_0|s_0) + \gamma \mathbb{E}_{s' \sim P(\cdot|s_0, a_0)} [V_\tau^{\pi_\theta}(s')] \right) \right] \\ &\stackrel{(i)}{=} \nabla_\theta \left[ \sum_{a_0} \pi_\theta(a_0|s_0) \left( Q_\tau^{\pi_\theta}(s_0, a_0) - \tau \log \pi_\theta(a_0|s_0) \right) \right] \\ &= \sum_{a_0} (\nabla_\theta \pi_\theta(a_0|s_0)) \left( Q_\tau^{\pi_\theta}(s_0, a_0) - \tau \log \pi_\theta(a_0|s_0) \right) + \sum_{a_0} \pi_\theta(a_0|s_0) \nabla_\theta \left( Q_\tau^{\pi_\theta}(s_0, a_0) - \tau \log \pi_\theta(a_0|s_0) \right) \\ &\stackrel{(ii)}{=} \sum_{a_0} \left( \pi_\theta(a_0|s_0) \nabla_\theta \log \pi_\theta(a_0|s_0) \right) \left( Q_\tau^{\pi_\theta}(s_0, a_0) - \tau \log \pi_\theta(a_0|s_0) \right) \\ &\quad + \sum_{a_0} \pi_\theta(a_0|s_0) \nabla_\theta \left( r(s_0, a_0) + \gamma \sum_{s_1} P(s_1|s_0, a_0) V_\tau^{\pi_\theta}(s_1) - \tau \log \pi_\theta(a_0|s_0) \right), \end{aligned}$$

where (i) relies on the definition (11a) of  $Q_\tau^{\pi_\theta}$ , and (ii) makes use of the identity

$$\nabla_\theta \pi_\theta(a_0|s_0) = \pi_\theta(a_0|s_0) \nabla_\theta \log \pi_\theta(a_0|s_0)$$

as well as the definition (11a) of  $Q_\tau^{\pi_\theta}$ . Given that

$$\sum_{a_0} \pi_\theta(a_0|s_0) \nabla_\theta \log \pi_\theta(a_0|s_0) = \sum_{a_0} \nabla_\theta \pi_\theta(a_0|s_0) = \nabla_\theta \left( \sum_{a_0} \pi_\theta(a_0|s_0) \right) = \nabla_\theta 1 = 0 \quad (\text{A.29})$$

and that  $r(s, a)$  is independent of  $\theta$ , one can continue the above derivative to reach

$$\begin{aligned} \nabla_\theta V_\tau^{\pi_\theta}(s_0) &= \sum_{a_0} \left( \pi_\theta(a_0|s_0) \nabla_\theta \log \pi_\theta(a_0|s_0) \right) \left( Q_\tau^{\pi_\theta}(s_0, a_0) - \tau \log \pi_\theta(a_0|s_0) \right) \\ &\quad + \gamma \sum_{a_0} \pi_\theta(a_0|s_0) \sum_{s_1} P(s_1|s_0, a_0) \nabla_\theta V_\tau^{\pi_\theta}(s_1) \\ &= \mathbb{E}_{\substack{a_i \sim \pi_\theta(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[ (\nabla_\theta \log \pi_\theta(a_0|s_0)) \left( Q_\tau^{\pi_\theta}(s_0, a_0) - \tau \log \pi_\theta(a_0|s_0) \right) + \gamma \nabla_\theta V_\tau^{\pi_\theta}(s_1) \right]. \end{aligned}$$

Repeating the above calculations recursively, we arrive at

$$\begin{aligned} \nabla_\theta V_\tau^{\pi_\theta}(s_0) &= \mathbb{E}_{\substack{a_i \sim \pi_\theta(\cdot|s_i), \\ s_{i+1} \sim P(\cdot|s_i, a_i), \forall i \geq 0}} \left[ \sum_{t=0}^{\infty} \gamma^t (\nabla_\theta \log \pi_\theta(a_t|s_t)) \left( Q_\tau^{\pi_\theta}(s_t, a_t) - \tau \log \pi_\theta(a_t|s_t) \right) \right] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[ (\nabla_\theta \log \pi_\theta(a|s)) \left( Q_\tau^{\pi_\theta}(s, a) - \tau \log \pi_\theta(a|s) \right) \right] \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[ (\nabla_\theta \log \pi_\theta(a|s)) \left( A_\tau^{\pi_\theta}(s, a) + V_\tau^{\pi_\theta}(s) \right) \right] \\
 &= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[ (\nabla_\theta \log \pi_\theta(a|s)) A_\tau^{\pi_\theta}(s, a) \right], \tag{A.30}
 \end{aligned}$$

where the second line follows by aggregating the terms corresponding to the same state-action pair, and the third line invokes the definition (A.1) of  $A_\tau^{\pi_\theta}$ . To see why the last line holds, invoke (A.29) to reach

$$\begin{aligned}
 \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[ V_\tau^{\pi_\theta}(s) \nabla_\theta \log \pi_\theta(a|s) \right] &= \sum_a V_\tau^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) \\
 &= V_\tau^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) = 0.
 \end{aligned}$$

Further, it is easily seen that under the softmax parametrization in (7),

$$\frac{\partial \log \pi_\theta(a'|s')}{\partial \theta(s, a)} = \mathbf{1}[s' = s] (\mathbf{1}[a' = a] - \pi_\theta(a|s)) \tag{A.31}$$

for any  $(s, a), (s', a') \in \mathcal{S} \times \mathcal{A}$ . Combining with (A.30), it further implies that

$$\begin{aligned}
 \frac{\partial V_\tau^{\pi_\theta}(s_0)}{\partial \theta(s, a)} &= \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \frac{\partial \log \pi_\theta(a'|s')}{\partial \theta(s, a)} A_\tau^{\pi_\theta}(s', a') \right] \\
 &= \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \left( \mathbf{1}[s' = s] (\mathbf{1}[a' = a] - \pi_\theta(a|s)) \right) A_\tau^{\pi_\theta}(s', a') \right] \\
 &\stackrel{(i)}{=} \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \mathbf{1}[(s', a') = (s, a)] A_\tau^{\pi_\theta}(s', a') \right] \\
 &= \frac{1}{1-\gamma} d_{s_0}^{\pi_\theta}(s) \pi_\theta(a|s) A_\tau^{\pi_\theta}(s, a).
 \end{aligned}$$

where (i) follows from  $\mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} A_\tau^{\pi_\theta}(s', a') = \sum_{a'} \pi_\theta(a'|s') A_\tau^{\pi_\theta}(s', a') = 0$  due to the definition (A.1). The proof regarding  $V_\tau^{\pi_\theta}(\rho)$  can be obtained by averaging the initial state  $s_0$  over the distribution  $\rho$ .

*Proof of Eqn. (A.2b).* In order to establish (A.2b), a crucial observation is that  $w_\theta := (\mathcal{F}_\rho^\theta)^\dagger \nabla_\theta V_\tau^{\pi_\theta}(\rho)$  is exactly the solution to the following least-squares problem

$$\text{minimize}_{w \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}} \left\| \mathcal{F}_\rho^\theta w - \nabla_\theta V_\tau^{\pi_\theta}(\rho) \right\|_2^2. \tag{A.32}$$

From the definition (14) of the Fisher information matrix, we have

$$\mathcal{F}_\rho^\theta w = \mathbb{E}_{s \sim d_\rho^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[ (\nabla_\theta \log \pi_\theta(a|s)) (\nabla_\theta \log \pi_\theta(a|s))^\top w \right].$$

for any fixed vector  $w = [w_{s,a}]_{(s,a) \in \mathcal{S} \times \mathcal{A}}$ . As a result, for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$  one has

$$(\mathcal{F}_\rho^\theta w)_{s,a} = \mathbb{E}_{s' \sim d_\rho^{\pi_\theta}} \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \frac{\partial \log \pi_\theta(a'|s')}{\partial \theta(s, a)} \left( \sum_{\tilde{s}, \tilde{a}} \frac{\partial \log \pi_\theta(a'|s')}{\partial \theta(\tilde{s}, \tilde{a})} w_{\tilde{s}, \tilde{a}} \right) \right]$$

$$\begin{aligned}
&\stackrel{(i)}{=} \mathbb{E}_{s' \sim d_\rho^{\pi_\theta}} \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \mathbf{1}[s' = s] \left( \mathbf{1}[a' = a] - \pi_\theta(a|s) \right) \left( \sum_{\tilde{s}, \tilde{a}} \mathbf{1}[\tilde{s} = s'] \left( \mathbf{1}[\tilde{a} = a'] - \pi_\theta(\tilde{a}|\tilde{s}) \right) w_{\tilde{s}, \tilde{a}} \right) \right] \\
&= \mathbb{E}_{s' \sim d_\rho^{\pi_\theta}} \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \mathbf{1}[s' = s] \left( \mathbf{1}[a' = a] - \pi_\theta(a|s) \right) \left( w_{s', a'} - \sum_{\tilde{a}} \pi_\theta(\tilde{a}|s') w_{s', \tilde{a}} \right) \right] \\
&= d_\rho^{\pi_\theta}(s) \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \left( \mathbf{1}[a' = a] - \pi_\theta(a|s) \right) (w_{s, a'} - c(s)) \right] \\
&= d_\rho^{\pi_\theta}(s) \mathbb{E}_{a' \sim \pi_\theta(\cdot|s')} \left[ \mathbf{1}[a' = a] w_{s, a'} - \pi_\theta(a|s) w_{s, a'} - \mathbf{1}[a' = a] c(s) + \pi_\theta(a|s) c(s) \right] \\
&= d_\rho^{\pi_\theta}(s) \left[ \pi_\theta(a|s) w_{s, a} - \pi_\theta(a|s) c(s) - \pi_\theta(a|s) c(s) + \pi_\theta(a|s) c(s) \right] \\
&= d_\rho^{\pi_\theta}(s) \pi_\theta(a|s) [w_{s, a} - c(s)],
\end{aligned}$$

where (i) makes use of the derivative calculation (A.31), and we define  $c(s) := \sum_a \pi_\theta(a|s) w_{s, a}$ . Consequently, the objective function of (A.32) can be written as

$$\begin{aligned}
\|\mathcal{F}_\rho^\theta w - \nabla_\theta V_\tau^{\pi_\theta}(\rho)\|_2^2 &= \sum_{s, a} \left( d_\rho^{\pi_\theta}(s) \pi_\theta(a|s) [w_{s, a} - c(s)] - \frac{1}{1-\gamma} d_\rho^{\pi_\theta}(s) \pi_\theta(a|s) A_\tau^{\pi_\theta}(s, a) \right)^2 \\
&= \sum_{s, a} \left( d_\rho^{\pi_\theta}(s) \pi_\theta(a|s) \left( w_{s, a} - c(s) - \frac{1}{1-\gamma} A_\tau^{\pi_\theta}(s, a) \right) \right)^2,
\end{aligned}$$

which is minimized by choosing  $w_{s, a} = \frac{1}{1-\gamma} A_\tau^{\pi_\theta}(s, a) + c(s)$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . This concludes the proof.

#### D. Convergence guarantees for CPI-style policy updates

Employing the SPI update as the improved policy, we arrive at the following CPI-style update

$$\pi^{(t+1)} = (1 - \beta)\pi^{(t)} + \beta\bar{\pi}^{(t+1)}. \quad (\text{A.33a})$$

Here,  $\bar{\pi}^{(t+1)}$  corresponds to a one-step SPI update from  $\pi^{(t)}$ , namely,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad \bar{\pi}^{(t+1)}(a|s) = \frac{1}{\bar{Z}^{(t)}(s)} \exp(Q_\tau^{(t)}(s, a)/\tau), \quad (\text{A.33b})$$

where we denote

$$\bar{Z}^{(t)}(s) = \sum_{a \in \mathcal{A}} \exp(Q_\tau^{(t)}(s, a)/\tau) \quad \text{and} \quad Q_\tau^{(t)} = Q_\tau^{\pi^{(t)}}$$

as usual. Here,  $\beta \in (0, 1]$  is a parameter that controls the ‘‘conservatism’’ of the updates. We characterize the convergence rate of this update rule (A.33) in the following theorem.

**Theorem 1 (Linear convergence of CPI-style updates)** *For any  $0 < \beta \leq 1$ , the update rule (A.33) satisfies*

$$V_\tau^*(\rho) - V_\tau^{(t)}(\rho) \leq \left\| \frac{\rho}{\mu_\tau^*} \right\|_\infty (1 - \beta(1 - \gamma))^t (V_\tau^*(\mu_\tau^*) - V_\tau^{(0)}(\mu_\tau^*)), \quad \forall t \geq 0, \quad (\text{A.34})$$

where  $\mu_\tau^*$  is the stationary distribution defined in (29).

According to Theorem 1, it takes the CPI-style policy update (A.33) at most

$$\frac{1}{\beta(1-\gamma)} \log \left( \left\| \frac{\rho}{\mu_\tau^*} \right\|_\infty \frac{V_\tau^*(\mu_\tau^*) - V_\tau^{(0)}(\mu_\tau^*)}{\epsilon} \right)$$

iterations to reach  $V_\tau^*(\rho) - V_\tau^{(t)}(\rho) \leq \epsilon$ . As it turns out, the CPI-style update rule can be analyzed using our framework through the following performance improvement lemma, which is an adaptation of Lemma 1. In what follows, we use  $\bar{Q}_\tau^{(t+1)}$  and  $\bar{V}_\tau^{(t+1)}$  to abbreviate  $Q_\tau^{\bar{\pi}^{(t+1)}}$  and  $V_\tau^{\bar{\pi}^{(t+1)}}$ , respectively.

**Lemma 2 (Performance improvement of CPI-style updates)** *Consider the policy update rule (A.33a) with any  $\beta \in (0, 1]$ . For any distribution  $\rho$ , one has*

$$V_\tau^{(t+1)}(\rho) - V_\tau^{(t)}(\rho) \geq \frac{\beta\tau}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{(t+1)}} [\text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s))].$$

See Appendix D.1.

Combining the above result with Lemma 5 and following a similar approach to (60) give

$$\begin{aligned} V_\tau^*(\rho) - V_\tau^{(t+1)}(\rho) &= V_\tau^*(\rho) - V_\tau^{(t)}(\rho) + (V_\tau^{(t)}(\rho) - V_\tau^{(t+1)}(\rho)) \\ &\stackrel{(i)}{\leq} V_\tau^*(\rho) - V_\tau^{(t)}(\rho) - \frac{\beta\tau}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{(t+1)}} [\text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s))] \\ &\stackrel{(ii)}{\leq} V_\tau^*(\rho) - V_\tau^{(t)}(\rho) - \frac{\beta\tau}{1-\gamma} \left\| \frac{d_\rho^{\pi_\tau^*}}{d_\rho^{(t+1)}} \right\|_\infty^{-1} \mathbb{E}_{s \sim d_\rho^{\pi_\tau^*}} [\text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s))] \\ &\stackrel{(iii)}{\leq} V_\tau^*(\rho) - V_\tau^{(t)}(\rho) - \beta \left\| \frac{d_\rho^{\pi_\tau^*}}{d_\rho^{(t+1)}} \right\|_\infty^{-1} (V_\tau^*(\rho) - V_\tau^{(t)}(\rho)) \\ &= \left( 1 - \beta \left\| \frac{d_\rho^{\pi_\tau^*}}{d_\rho^{(t+1)}} \right\|_\infty^{-1} \right) (V_\tau^*(\rho) - V_\tau^{(t)}(\rho)). \end{aligned} \tag{A.35}$$

Here, (i) arises from Lemma 2, (ii) employs the pre-factor  $\|d_\rho^{\pi_\tau^*}/d_\rho^{(t+1)}\|_\infty^{-1}$  to accommodate the change of distributions, whereas (iii) follows from Lemma 5 and the constraint that  $0 \leq \eta \leq \frac{1-\gamma}{\tau}$ .

By taking  $\rho$  to be the stationary distribution  $\mu_\tau^*$  (cf. (29)), one has

$$\begin{aligned} V_\tau^*(\mu_\tau^*) - V_\tau^{(t+1)}(\mu_\tau^*) &\leq \left( 1 - \beta \left\| \frac{d_{\mu_\tau^*}^{\pi_\tau^*}}{d_{\mu_\tau^*}^{(t+1)}} \right\|_\infty^{-1} \right) (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)) \\ &\leq \left( 1 - \beta \left\| \frac{\mu_\tau^*}{(1-\gamma)\mu_\tau^*} \right\|_\infty^{-1} \right) (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)) \\ &= (1 - \beta(1-\gamma)) (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)), \end{aligned}$$

where we have used  $d_{\mu_\tau^*}^{\pi_\tau^*} = \mu_\tau^*$  (cf. (29)) and  $d_{\mu_\tau^*}^{(t+1)} \geq (1-\gamma)\mu_\tau^*$  in the second step. This immediately concludes the proof.

### D.1. Proof of Lemma 2

First of all, we claim that

$$V_\tau^{(t+1)}(\rho) - V_\tau^{(t)}(\rho) = \frac{\tau}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{(t+1)}} \left[ \text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) - \text{KL}(\pi^{(t+1)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) \right], \quad (\text{A.36})$$

which we shall establish momentarily. Since the KL divergence  $\text{KL}(\pi(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s))$  is convex in  $\pi(\cdot|s)$  (Cover 1999), the update rule (A.33a) together with Jensen's inequality necessarily implies that

$$\begin{aligned} \text{KL}(\pi^{(t+1)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) &\leq \beta \text{KL}(\bar{\pi}^{(t+1)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) + (1-\beta) \text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) \\ &= (1-\beta) \text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)). \end{aligned}$$

Substituting the above inequality into (A.36) allows us to conclude that

$$V_\tau^{(t+1)}(\rho) - V_\tau^{(t)}(\rho) \geq \frac{\beta\tau}{1-\gamma} \mathbb{E}_{s \sim d_\rho^{(t+1)}} \left[ \text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) \right].$$

The rest of this proof is then dedicated to establishing the claim (A.36), which is similar to the proof of Lemma 1. To begin with, we express  $V_\tau^{(t)}(s_0)$  as follows

$$\begin{aligned} V_\tau^{(t)}(s_0) &= \mathbb{E}_{a_0 \sim \pi^{(t)}(\cdot|s_0)} \left[ -\tau \log \pi^{(t)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0) \right] \\ &= \mathbb{E}_{a_0 \sim \pi^{(t)}(\cdot|s_0)} \left[ -\tau \log \pi^{(t)}(a_0|s_0) + \tau \log \bar{\pi}^{(t+1)}(a_0|s_0) \right] + \tau \log \bar{Z}^{(t)}(s_0) \\ &= \tau \log \bar{Z}^{(t)}(s_0) - \tau \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)) \\ &= \tau \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ \log \bar{Z}^{(t)}(s_0) \right] - \tau \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)), \end{aligned}$$

where the first line makes use of the definitions (8) and (11a), the second line follows from (A.33), the third line uses the definition of the KL divergence, and the last line follows since  $\bar{Z}^{(t)}(s_0)$  does not depend on  $a$ . To continue, we subtract and add  $\tau \text{KL}(\pi^{(t+1)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0))$  to obtain

$$\begin{aligned} V_\tau^{(t)}(s_0) &= \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ \tau \log \bar{Z}^{(t)}(s_0) - \tau \log \pi^{(t+1)}(a_0|s_0) + \tau \log \bar{\pi}^{(t+1)}(a_0|s_0) \right] \\ &\quad + \tau \text{KL}(\pi^{(t+1)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)) - \tau \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)) \\ &= \mathbb{E}_{a_0 \sim \pi^{(t+1)}(\cdot|s_0)} \left[ -\tau \log \pi^{(t+1)}(a_0|s_0) + Q_\tau^{(t)}(s_0, a_0) \right] + \tau \text{KL}(\pi^{(t+1)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)) \\ &\quad - \tau \text{KL}(\pi^{(t)}(\cdot|s_0) \parallel \bar{\pi}^{(t+1)}(\cdot|s_0)) \\ &= V_\tau^{(t+1)}(s_0) + \frac{\tau}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^{(t+1)}} \left[ \text{KL}(\pi^{(t+1)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) - \text{KL}(\pi^{(t)}(\cdot|s) \parallel \bar{\pi}^{(t+1)}(\cdot|s)) \right]. \end{aligned}$$

Here, the first step relies on the definition of KL divergence, the second step comes from (A.33), while the last step is obtained by using the relation  $Q_\tau^{(t)}(s_0, a_0) = r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P(\cdot|s_0, a_0)} [V_\tau^{(t)}(s_1)]$  and then invoking the above equality recursively as in the expression (A.12) (see Lemma 1). Averaging the equality over the initial state distribution  $s_0 \sim \rho$  thus establishes the claim (A.36).

## E. Proof for approximate entropy-regularized NPG (Theorem 2)

In this section, we complete the proofs of Theorem 2 in Section 4.3, which consists of (i) establishing the linear system in (58) and (ii) extracting the convergence rate from (58).

*Step 1: establishing the linear system (58).* In what follows, we shall justify the linear system relation by checking each row separately.

(1) **Bounding**  $\|Q_\tau^* - \tau \log \widehat{\xi}^{(t+1)}\|_\infty$ . From the construction (57b) of  $\widehat{\xi}^{(t+1)}$ , we have

$$Q_\tau^* - \tau \log \widehat{\xi}^{(t+1)} = \alpha(Q_\tau^* - \tau \log \widehat{\xi}^{(t)}) + (1 - \alpha)(Q_\tau^* - Q_\tau^{(t)}) + (1 - \alpha)(Q_\tau^{(t)} - \widehat{Q}_\tau^{(t)}).$$

Taken together with the triangle inequality and the assumption  $\|Q_\tau^{(t)} - \widehat{Q}_\tau^{(t)}\|_\infty \leq \delta$ , this gives

$$\|Q_\tau^* - \tau \log \widehat{\xi}^{(t+1)}\|_\infty \leq \alpha \|Q_\tau^* - \tau \log \widehat{\xi}^{(t)}\|_\infty + (1 - \alpha) \|Q_\tau^* - Q_\tau^{(t)}\|_\infty + (1 - \alpha) \delta. \quad (\text{A.37})$$

(2) **Bounding**  $-\min_{s,a} (Q_\tau^{(t+1)}(s,a) - \tau \log \widehat{\xi}^{(t+1)}(s,a))$ . Invoking the definition (57b) of  $\widehat{\xi}^{(t+1)}$  again implies that for any  $(s,a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned} & - \left( Q_\tau^{(t+1)}(s,a) - \tau \log \widehat{\xi}^{(t+1)}(s,a) \right) \\ &= - \left( Q_\tau^{(t+1)}(s,a) - \tau \left( \alpha \log \widehat{\xi}^{(t)}(s,a) + (1 - \alpha) \widehat{Q}_\tau^{(t)}(s,a) / \tau \right) \right) \\ &= -\alpha \left( Q_\tau^{(t)}(s,a) - \tau \log \widehat{\xi}^{(t)}(s,a) \right) + (1 - \alpha) \left( \widehat{Q}_\tau^{(t)}(s,a) - Q_\tau^{(t)}(s,a) \right) + \left( Q_\tau^{(t)}(s,a) - Q_\tau^{(t+1)}(s,a) \right) \\ &\leq -\alpha \left( Q_\tau^{(t)}(s,a) - \tau \log \widehat{\xi}^{(t)}(s,a) \right) + (1 - \alpha) \delta + \frac{2\gamma\delta}{1 - \gamma}, \end{aligned}$$

where the last inequality follows from  $\|Q_\tau^{(t)} - \widehat{Q}_\tau^{(t)}\|_\infty \leq \delta$  and (56). Taking the maximum over  $(s,a) \in \mathcal{S} \times \mathcal{A}$  on both sides and using the definition  $\alpha = 1 - \frac{\eta\tau}{1 - \gamma}$  yield

$$-\min_{s,a} \left( Q_\tau^{(t+1)}(s,a) - \tau \log \widehat{\xi}^{(t+1)}(s,a) \right) \leq -\alpha \min_{s,a} \left( Q_\tau^{(t)}(s,a) - \tau \log \widehat{\xi}^{(t)}(s,a) \right) + (1 - \alpha) \delta \left( 1 + \frac{2\gamma}{\eta\tau} \right). \quad (\text{A.38})$$

(3) **Bounding**  $\|Q_\tau^* - Q_\tau^{(t+1)}\|_\infty$ . Following the same arguments as for (A.18), we obtain

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

where the last line follows from (A.5). By plugging (A.37) and (A.38) into the above inequality, we arrive at the claimed bound regarding this term.

*Step 2: deducing convergence guarantees from the linear system (58).* We start by pinning down the eigenvalues and eigenvectors of the matrix  $B$ . Specifically, the three eigenvalues can be calculated as

$$\lambda_1 = \alpha + \gamma(1 - \alpha) = 1 - \eta\tau, \quad \lambda_2 = \alpha \quad \text{and} \quad \lambda_3 = 0, \quad (\text{A.39})$$

whose corresponding eigenvectors are given respectively by

$$v_1 = \begin{bmatrix} \gamma \\ 1 \\ 0 \end{bmatrix}, \quad v_2 = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}, \quad \text{and} \quad v_3 = \begin{bmatrix} \alpha \\ \alpha - 1 \\ 0 \end{bmatrix}. \quad (\text{A.40})$$

With some elementary computation, one can show that  $z_0$  and  $b$  introduced in (59) can be related to the eigenvectors of  $B$  in the following way:

$$\begin{aligned} z_0 &\leq \begin{bmatrix} \|Q_\tau^* - Q_\tau^{(0)}\|_\infty \\ \|Q_\tau^* - \tau \log \widehat{\xi}^{(0)}\|_\infty \\ \|Q_\tau^{(0)} - \tau \log \widehat{\xi}^{(0)}\|_\infty \end{bmatrix} \\ &= \frac{1}{1 - \eta\tau} \left[ (1 - \alpha) \|Q_\tau^* - Q_\tau^{(0)}\|_\infty + \alpha \left( \|Q_\tau^* - \tau \log \widehat{\xi}^{(0)}\|_\infty + \|Q_\tau^{(0)} - \tau \log \widehat{\xi}^{(0)}\|_\infty \right) \right] v_1 \\ &\quad + \|Q_\tau^{(0)} - \tau \log \widehat{\xi}^{(0)}\|_\infty v_2 + c_z v_3 \\ &\leq \frac{1}{1 - \eta\tau} \left( \|Q_\tau^* - Q_\tau^{(0)}\|_\infty + 2\alpha\tau \|\log \pi_\tau^* - \log \pi^{(0)}\|_\infty \right) v_1 + \|Q_\tau^{(0)} - \tau \log \widehat{\xi}^{(0)}\|_\infty v_2 + c_z v_3, \end{aligned} \quad (\text{A.41})$$

where  $c_z$  is some scalar whose value is immaterial since the eigenvalue corresponding to  $v_3$  is  $\lambda_3 = 0$ , and the last line follows from the same reasoning for (52). Another useful identity is:

$$b = (1 - \alpha)\delta \begin{bmatrix} \gamma \left( 2 + \frac{2\gamma}{\eta\tau} \right) \\ 1 \\ 1 + \frac{2\gamma}{\eta\tau} \end{bmatrix} = (1 - \alpha)\delta \left[ \left( 2 + \frac{2\gamma}{\eta\tau} \right) v_1 + \left( 1 + \frac{2\gamma}{\eta\tau} \right) v_2 \right]. \quad (\text{A.42})$$

With these preparations in place, we can now invoke the recursion relationship (58) and the non-negativity of  $B$  to obtain

$$\begin{aligned} z_{t+1} &\leq B^{t+1} z_0 + \sum_{s=0}^t B^{t-s} b \\ &\leq B^{t+1} \left[ \frac{1}{1 - \eta\tau} \left( \|Q_\tau^* - Q_\tau^{(0)}\|_\infty + 2\alpha\tau \|\log \pi_\tau^* - \log \pi^{(0)}\|_\infty \right) v_1 + \|Q_\tau^{(0)} - \tau \log \widehat{\xi}^{(0)}\|_\infty v_2 + c_z v_3 \right] \\ &\quad + (1 - \alpha)\delta \sum_{s=0}^t B^{t-s} \left[ \left( 2 + \frac{2\gamma}{\eta\tau} \right) v_1 + \left( 1 + \frac{2\gamma}{\eta\tau} \right) v_2 \right] \\ &= \left[ \lambda_1^t \left( \|Q_\tau^* - Q_\tau^{(0)}\|_\infty + 2\alpha\tau \|\log \pi_\tau^* - \log \pi^{(0)}\|_\infty \right) + (1 - \alpha)\delta \left( 2 + \frac{2\gamma}{\eta\tau} \right) \frac{1 - \lambda_1^{t+1}}{1 - \lambda_1} \right] v_1 \\ &\quad + \left[ \lambda_2^{t+1} \|Q_\tau^{(0)} - \tau \log \widehat{\xi}^{(0)}\|_\infty + (1 - \alpha)\delta \left( 1 + \frac{2\gamma}{\eta\tau} \right) \frac{1 - \lambda_2^{t+1}}{1 - \lambda_2} \right] v_2, \end{aligned}$$

where the eigenvalues and eigenvectors of  $B$  are given in (A.39) and (A.40), respectively, and the second inequality relies on (A.41) and (A.42). Note that we are only interested in the first two

entries of the vector  $z_t$ . Since the first two entries of the eigenvector  $v_2$  are non-positive, we can safely drop the term involving  $v_2$  in the above inequality to obtain

$$\begin{aligned}
 & \left[ \begin{array}{l} \|Q_\tau^* - Q_\tau^{(t+1)}\|_\infty \\ \|Q_\tau^* - \tau \log \widehat{\xi}^{(t+1)}\|_\infty \end{array} \right] \\
 & \leq \left\{ \lambda_1^t (\|Q_\tau^* - Q_\tau^{(0)}\|_\infty + 2\alpha\tau \|\log \pi_\tau^* - \log \pi^{(0)}\|_\infty) + (1-\alpha)\delta \left( 2 + \frac{2\gamma}{\eta\tau} \frac{1-\lambda_1^{t+1}}{1-\lambda_1} \right) \right\} \begin{bmatrix} \gamma \\ 1 \end{bmatrix} \\
 & \leq \left\{ (1-\eta\tau)^t \left( \|Q_\tau^* - Q_\tau^{(0)}\|_\infty + 2 \left( 1 - \frac{\eta\tau}{1-\gamma} \right) \tau \|\log \pi_\tau^* - \log \pi^{(0)}\|_\infty \right) + \frac{2\delta}{1-\gamma} \left( 1 + \frac{\gamma}{\eta\tau} \right) \right\} \begin{bmatrix} \gamma \\ 1 \end{bmatrix}.
 \end{aligned} \tag{A.43}$$

When it comes to the log policies, we recall again the fact that  $\pi^{(t)}$  is related to  $\widehat{\xi}^{(t)}$  as

$$\forall s \in \mathcal{S}: \quad \pi^{(t)}(\cdot|s) = \frac{1}{\|\widehat{\xi}^{(t)}(s, \cdot)\|_1} \widehat{\xi}^{(t)}(s, \cdot). \tag{A.44}$$

Invoking the elementary property (A.6), we reach

$$\|\log \pi_\tau^* - \log \pi^{(t+1)}\|_\infty \leq 2 \|Q_\tau^*/\tau - \log \widehat{\xi}^{(t+1)}\|_\infty.$$

This together with the bound on  $\|Q_\tau^* - \tau \log \widehat{\xi}^{(t+1)}\|_\infty$  in (A.43) establishes our claim for  $\|\log \pi_\tau^* - \log \pi^{(t+1)}\|_\infty$ .

## F. Proof for local quadratic convergence (Theorem 3)

Assuming that the policy  $\pi^{(t)}$  obeys Condition (30), we can control the difference of the corresponding discounted state visitation probabilities in terms of the sub-optimality gap w.r.t. the log policy. This is stated in the following lemma, whose proof is deferred to Section F.1.

**Lemma 3** *Consider any policy  $\pi$  satisfying  $\|\log \pi - \log \pi_\tau^*\|_\infty \leq 1$ . It follows that*

$$\left\| 1 - \frac{d_\rho^{\pi_\tau^*}}{d_\rho^\pi} \right\|_\infty \leq 2 \left( \frac{1}{1-\gamma} \left\| \frac{d_\rho^{\pi_\tau^*}}{\rho} \right\|_\infty - 1 \right) \|\log \pi - \log \pi_\tau^*\|_\infty.$$

*In particular, by taking  $\rho = \mu_\tau^*$  one has*

$$\left\| 1 - \frac{\mu_\tau^*}{d_{\mu_\tau^*}^\pi} \right\|_\infty = \left\| 1 - \frac{d_{\mu_\tau^*}^{\pi_\tau^*}}{d_{\mu_\tau^*}^\pi} \right\|_\infty \leq \frac{2\gamma}{1-\gamma} \|\log \pi - \log \pi_\tau^*\|_\infty.$$

First, by virtue of the SPI update rule (17) and the inequality (A.6), it is guaranteed that

$$\|\log \pi_\tau^* - \log \pi^{(t+1)}\|_\infty \leq \frac{2}{\tau} \|Q_\tau^* - Q_\tau^{(t)}\|_\infty \leq \frac{2\gamma}{\tau} \|V_\tau^* - V_\tau^{(t)}\|_\infty \leq \frac{2\gamma}{\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)), \tag{A.45}$$

where the last inequality comes from a change of distributions argument. Armed with Lemma 3 and the inequality (A.45), we arrive at

$$\left\| 1 - \frac{d_{\mu_\tau^*}^{\pi_\tau^*}}{d_{\mu_\tau^*}^{(t+1)}} \right\|_\infty \leq \frac{2\gamma}{1-\gamma} \|\log \pi^{(t+1)} - \log \pi_\tau^*\|_\infty \leq \frac{4\gamma^2}{(1-\gamma)\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)). \quad (\text{A.46})$$

Substitution into (60) gives

$$\begin{aligned} V_\tau^*(\mu_\tau^*) - V_\tau^{(t+1)}(\mu_\tau^*) &\leq \left( 1 - \left\| \frac{d_{\mu_\tau^*}^{\pi_\tau^*}}{d_{\mu_\tau^*}^{(t+1)}} \right\|_\infty^{-1} \right) (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)) \\ &= \left( 1 - \left\| 1 + \frac{d_{\mu_\tau^*}^{\pi_\tau^*} - d_{\mu_\tau^*}^{(t+1)}}{d_{\mu_\tau^*}^{(t+1)}} \right\|_\infty^{-1} \right) (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)) \\ &\leq \left( 1 - \frac{1}{1 + \frac{4\gamma^2}{(1-\gamma)\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*))} \right) (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)) \\ &= \frac{\frac{4\gamma^2}{(1-\gamma)\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*))^2}{1 + \frac{4\gamma^2}{(1-\gamma)\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*))} \leq \frac{4\gamma^2}{(1-\gamma)\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*))^2, \end{aligned}$$

where the second inequality makes use of the bound (A.46). This in turn reveals that

$$\frac{4\gamma^2}{(1-\gamma)\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t+1)}(\mu_\tau^*)) \leq \left( \frac{4\gamma^2}{(1-\gamma)\tau} \left\| \frac{1}{\mu_\tau^*} \right\|_\infty (V_\tau^*(\mu_\tau^*) - V_\tau^{(t)}(\mu_\tau^*)) \right)^2,$$

which leads to our claimed result by a standard change of distributions.

### F.1. Proof of Lemma 3

For any policy  $\pi$ , denote by  $P_\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$  the state transition matrix induced by  $\pi$  as follows

$$\forall s, s' \in \mathcal{S}: \quad [P_\pi]_{s,s'} := \mathbb{E}_{a \sim \pi(\cdot|s)} [P(s'|s, a)]. \quad (\text{A.47})$$

For any policy  $\pi$  satisfying  $\|\log \pi - \log \pi_\tau^*\|_\infty \leq 1$ , we develop an upper bound on  $\left| [P_\pi - P_{\pi_\tau^*}]_{s,s'} \right|$  as follows

$$\begin{aligned} \left| [P_\pi - P_{\pi_\tau^*}]_{s,s'} \right| &= \left| \sum_a P(s'|s, a) (\pi(a|s) - \pi_\tau^*(a|s)) \right| \leq \sum_a P(s'|s, a) \pi_\tau^*(a|s) \left| \frac{\pi(a|s)}{\pi_\tau^*(a|s)} - 1 \right| \\ &\stackrel{(i)}{\leq} (e-1) \sum_a P(s'|s, a) \pi_\tau^*(a|s) \left| \log \pi(a|s) - \log \pi_\tau^*(a|s) \right| \\ &\leq \|\log \pi - \log \pi_\tau^*\|_\infty (e-1) \sum_a P(s'|s, a) \pi_\tau^*(a|s) \\ &\leq 2 [P_{\pi_\tau^*}]_{s,s'} \|\log \pi - \log \pi_\tau^*\|_\infty, \end{aligned}$$

where (i) uses the assumption  $\|\log \pi^* - \log \pi\|_\infty \leq 1$  together with the elementary inequality  $|x| \leq (e-1) |\log(1+x)|$  when  $-1 < x \leq e-1$ . With the preceding bound in mind, we can demonstrate that

$$\left| (d_\rho^{\pi_\tau^*})^\top (P_\pi - P_{\pi_\tau^*}) \right| \leq 2 \|\log \pi - \log \pi_\tau^*\|_\infty (d_\rho^{\pi_\tau^*})^\top P_{\pi_\tau^*}. \quad (\text{A.48})$$

Here and throughout, we overload the notation  $|z|$  for any vector  $z \in \mathbb{R}^{|\mathcal{S}|}$  to denote  $[|z_i|]_{1 \leq i \leq |\mathcal{S}|}$ .

In addition, the definitions of  $d_\rho^\pi$  and  $d_\rho^{\pi^*}$  admit the following matrix-vector representation:

$$(d_\rho^\pi)^\top = (1 - \gamma)\rho^\top (I - \gamma P_\pi)^{-1}, \quad (\text{A.49})$$

$$(d_\rho^{\pi^*})^\top = (1 - \gamma)\rho^\top (I - \gamma P_{\pi^*})^{-1}, \quad (\text{A.50})$$

thus allowing one to derive

$$\begin{aligned} (d_\rho^\pi - d_\rho^{\pi^*})^\top &= (1 - \gamma)\rho^\top (I - \gamma P_{\pi^*})^{-1} [(I - \gamma P_{\pi^*}) - (I - \gamma P_\pi)] (I - \gamma P_\pi)^{-1} \\ &= \gamma (d_\rho^{\pi^*})^\top (P_\pi - P_{\pi^*}) (I - \gamma P_\pi)^{-1}. \end{aligned}$$

This together with the non-negativity of the matrix  $(I - \gamma P_\pi)^{-1}$  (Li et al. 2020, Lemma 7) enables the following bound

$$\begin{aligned} \left| (d_\rho^\pi - d_\rho^{\pi^*})^\top \right| &\leq \gamma \left| (d_\rho^{\pi^*})^\top (P_\pi - P_{\pi^*}) \right| (I - \gamma P_\pi)^{-1} \\ &\leq 2 \|\log \pi - \log \pi_\tau^*\|_\infty \gamma (d_\rho^{\pi^*})^\top P_{\pi_\tau^*} (I - \gamma P_\pi)^{-1}, \end{aligned} \quad (\text{A.51})$$

where the last inequality results from (A.48).

Furthermore, we make the observation that

$$\begin{aligned} \gamma (d_\rho^{\pi^*})^\top P_{\pi_\tau^*} &= (1 - \gamma)\gamma \rho^\top (I - \gamma P_{\pi_\tau^*})^{-1} P_{\pi_\tau^*} = (1 - \gamma)\gamma \rho^\top \left[ \sum_{i=0}^{\infty} (\gamma P_{\pi_\tau^*})^i \right] P_{\pi_\tau^*} \\ &= (1 - \gamma)\rho^\top \left[ \sum_{i=1}^{\infty} (\gamma P_{\pi_\tau^*})^i \right] \\ &= (1 - \gamma)\rho^\top \left[ (I - \gamma P_{\pi_\tau^*})^{-1} - I \right] = (d_\rho^{\pi^*})^\top - (1 - \gamma)\rho^\top \leq \left( \left\| \frac{d_\rho^{\pi^*}}{\rho} \right\|_\infty - (1 - \gamma) \right) \rho^\top, \end{aligned}$$

where the last line comes from a change of distributions argument. Combining this bound with (A.51) gives

$$\begin{aligned} \left| (d_\rho^\pi - d_\rho^{\pi^*})^\top \right| &\leq 2 \|\log \pi - \log \pi_\tau^*\|_\infty \left( \frac{1}{1 - \gamma} \left\| \frac{d_\rho^{\pi^*}}{\rho} \right\|_\infty - 1 \right) (1 - \gamma)\rho^\top (I - \gamma P_\pi)^{-1} \\ &= 2 \|\log \pi - \log \pi_\tau^*\|_\infty \left( \frac{1}{1 - \gamma} \left\| \frac{d_\rho^{\pi^*}}{\rho} \right\|_\infty - 1 \right) (d_\rho^\pi)^\top, \end{aligned}$$

where the last line arises from the expression (A.49). As a result, we establish the claimed bound

$$\left\| 1 - \frac{d_\rho^{\pi^*}}{d_\rho^\pi} \right\|_\infty \leq 2 \|\log \pi - \log \pi_\tau^*\|_\infty \left( \frac{1}{1 - \gamma} \left\| \frac{d_\rho^{\pi^*}}{\rho} \right\|_\infty - 1 \right).$$

## References

- Agarwal, A., Jiang, N., and Kakade, S. M. (2019). Reinforcement learning: Theory and algorithms. Technical report.
- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. (2020). Optimality and approximation with policy gradient methods in Markov decision processes. In *Conference on Learning Theory*, pages 64–66. PMLR.
- Cover, T. M. (1999). *Elements of information theory*. John Wiley & Sons.
- Li, G., Wei, Y., Chi, Y., Gu, Y., and Chen, Y. (2020). Sample complexity of asynchronous Q-learning: Sharper analysis and variance reduction. *arXiv preprint arXiv:2006.03041*.
- Mei, J., Xiao, C., Szepesvari, C., and Schuurmans, D. (2020). On the global convergence rates of softmax policy gradient methods. *arXiv preprint arXiv:2005.06392*.
- 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.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063.