

A Algorithms in Section 3

We present the algorithms for solving the subproblems of policy improvement and policy evaluation in Section 3.

Algorithm 2 Policy Improvement via SGD

- 1: **Require:** MDP $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$, current energy function f_{θ_k} , initial weights b_i , $[\theta(0)]_i$ ($i \in [m_f]$), number of iterations T , sample $\{(s_t, a_t^0)\}_{t=1}^T$
 - 2: Set stepsize $\eta \leftarrow T^{-1/2}$
 - 3: **for** $t = 0, \dots, T-1$ **do**
 - 4: $(s, a) \leftarrow (s_{t+1}, a_{t+1}^0)$
 - 5: $\theta(t+1/2) \leftarrow \theta(t) - \eta \cdot (f_{\theta(t)}(s, a) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a))) \cdot \nabla_{\theta} f_{\theta(t)}(s, a)$
 - 6: $\theta(t+1) \leftarrow \operatorname{argmin}_{\theta \in \mathcal{B}^0(R_f)} \{\|\theta - \theta(t+1/2)\|_2\}$
 - 7: **end for**
 - 8: Average over path $\bar{\theta} \leftarrow 1/T \cdot \sum_{t=0}^{T-1} \theta(t)$
 - 9: **Output:** $f_{\bar{\theta}}$
-

Algorithm 3 Policy Evaluation via TD

- 1: **Require:** MDP $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$, initial weights b_i , $[\omega(0)]_i$ ($i \in [m_Q]$), number of iterations T , sample $\{(s_t, a_t, s'_t, a'_t)\}_{t=1}^T$
 - 2: Set stepsize $\eta \leftarrow T^{-1/2}$
 - 3: **for** $t = 0, \dots, T-1$ **do**
 - 4: $(s, a, s', a') \leftarrow (s_{t+1}, a_{t+1}, s'_{t+1}, a'_{t+1})$
 - 5: $\omega(t+1/2) \leftarrow \omega(t) - \eta \cdot (Q_{\omega(t)}(s, a) - (1-\gamma) \cdot r(s, a) - \gamma Q_{\omega(t)}(s', a')) \cdot \nabla_{\omega} Q_{\omega(t)}(s, a)$
 - 6: $\omega(t+1) \leftarrow \operatorname{argmin}_{\omega \in \mathcal{B}^0(R_Q)} \{\|\omega - \omega(t+1/2)\|_2\}$
 - 7: **end for**
 - 8: Average over path $\bar{\omega} \leftarrow 1/T \cdot \sum_{t=0}^{T-1} \omega(t)$
 - 9: **Output:** $Q_{\bar{\omega}}$
-

B Supplementary Lemma in Section 3

The following lemma quantifies the policy improvement error in terms of the distance between policies, which is induced by solving (3.5).

Lemma B.1. Suppose that $\pi_{\theta_{k+1}} \propto \exp\{\tau_{k+1}^{-1} f_{\theta_{k+1}}\}$ satisfies

$$\mathbb{E}_{\tilde{\sigma}_k} [(f_{\theta_{k+1}}(s, a) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)))^2] \leq \epsilon_{k+1}.$$

We have

$$\mathbb{E}_{\tilde{\sigma}_k} [(\pi_{\theta_{k+1}}(a | s) - \hat{\pi}_{k+1}(a | s))^2] \leq \tau_{k+1}^{-2} \epsilon_{k+1} / 16,$$

where $\hat{\pi}_{k+1}$ is defined in (3.4).

Proof. Let $\tau_{k+1}^{-1} \hat{f}_{k+1} = \beta_k^{-1} Q_{\omega_k} + \tau_k^{-1} f_{\theta_k}$. Since an energy-based policy $\pi \propto \exp\{\tau^{-1} f\}$ is continuous with respect to f , by the mean value theorem, we have

$$\begin{aligned} |\pi_{\theta_{k+1}}(a | s) - \hat{\pi}_{k+1}(a | s)| &= \left| \frac{\exp\{\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, a)\}}{\sum_{a' \in \mathcal{A}} \exp\{\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, a')\}} - \frac{\exp\{\tau_{k+1}^{-1} \hat{f}_{k+1}(s, a)\}}{\sum_{a' \in \mathcal{A}} \exp\{\tau_{k+1}^{-1} \hat{f}_{k+1}(s, a')\}} \right| \\ &= \left| \frac{\partial}{\partial f(s, a)} \left(\frac{\exp\{\tau_{k+1}^{-1} \tilde{f}(s, a)\}}{\sum_{a' \in \mathcal{A}} \exp\{\tau_{k+1}^{-1} \tilde{f}(s, a')\}} \right) \right| \cdot |f_{\theta_{k+1}}(s, a) - \hat{f}_{k+1}(s, a)|, \end{aligned}$$

where \tilde{f} is a function determined by $f_{\theta_{k+1}}$ and \hat{f}_{k+1} . Furthermore, we have

$$\left| \frac{\partial}{\partial f(s, a)} \left(\frac{\exp\{\tau_{k+1}^{-1} f(s, a)\}}{\sum_{a' \in \mathcal{A}} \exp\{\tau_{k+1}^{-1} f(s, a')\}} \right) \right| = \tau_{k+1}^{-1} \cdot \pi(a | s) \cdot (1 - \pi(a | s)) \leq \tau_{k+1}^{-1} / 4.$$

Therefore, we obtain

$$\begin{aligned} & (\pi_{\theta_{k+1}}(a | s) - \hat{\pi}_{k+1}(a | s))^2 \\ & \leq \tau_{k+1}^{-2}/16 \cdot (f_{\theta_{k+1}}(s, a) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)))^2. \end{aligned} \quad (\text{B.1})$$

Taking expectation $\mathbb{E}_{\tilde{\sigma}_k}[\cdot]$ on the both sides of (B.1), we finally obtain

$$\begin{aligned} & \mathbb{E}_{\tilde{\sigma}_k}[(\pi_{\theta_{k+1}}(a | s) - \hat{\pi}_{k+1}(a | s))^2] \\ & \leq \tau_{k+1}^{-2}/16 \cdot \mathbb{E}_{\tilde{\sigma}_k}[(f_{\theta_{k+1}}(s, a) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_0}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)))^2] \leq \tau_{k+1}^{-2} \epsilon_{k+1}/16, \end{aligned}$$

which concludes the proof of Lemma B.1. \square

Lemma B.1 ensures that if the policy improvement error ϵ_{k+1} is small, then the corresponding improved policy $\pi_{\theta_{k+1}}$ is close to the ideal improved policy $\hat{\pi}_{k+1}$, which justifies solving the subproblem in (3.5) for policy improvement.

C Proof of Proposition 3.1

Proof. The subproblem of policy improvement for solving $\hat{\pi}_{k+1}$ takes the form

$$\begin{aligned} & \max_{\pi} \mathbb{E}_{\nu_k}[\langle \pi(\cdot | s), Q_{\omega_k}(s, \cdot) \rangle - \beta_k \cdot \text{KL}(\pi(\cdot | s) \parallel \pi_{\theta_k}(\cdot | s))] \\ & \text{subject to } \sum_{a \in \mathcal{A}} \pi(a | s) = 1, \text{ for any } s \in \mathcal{S}. \end{aligned}$$

The Lagrangian of the above maximization problem takes the form

$$\int_{s \in \mathcal{S}} [\langle \pi(\cdot | s), Q_{\omega_k}(s, \cdot) \rangle - \beta_k \cdot \text{KL}(\pi(\cdot | s) \parallel \pi_{\theta_k}(\cdot | s))] \nu_k(ds) + \int_{s \in \mathcal{S}} \left(\sum_{a \in \mathcal{A}} \pi(a | s) - 1 \right) \lambda(ds).$$

Plugging in $\pi_{\theta_k}(s, a) = \exp\{\tau_k^{-1} f_{\theta_k}(s, a)\} / \sum_{a' \in \mathcal{A}} \exp\{\tau_k^{-1} f_{\theta_k}(s, a')\}$, we obtain the optimality condition

$$Q_{\omega_k}(s, a) + \beta_k \tau_k^{-1} f_{\theta_k}(s, a) - \beta_k \cdot \left[\log \left(\sum_{a' \in \mathcal{A}} \exp\{\tau_k^{-1} f_{\theta_k}(s, a')\} \right) + \log \pi(a | s) + 1 \right] + \frac{\lambda(s)}{\nu_k(s)} = 0,$$

for any $a \in \mathcal{A}$ and $s \in \mathcal{S}$. Note that $\log(\sum_{a' \in \mathcal{A}} \exp\{\tau_k^{-1} f_{\theta_k}(s, a')\})$ is determined by the state s only. Hence, we have $\hat{\pi}_{k+1}(a | s) \propto \exp\{\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)\}$ for any $a \in \mathcal{A}$ and $s \in \mathcal{S}$, which concludes the proof of Proposition 3.1. \square

D Proofs for Section 4.1

The proofs in this section generalizes those of [7, 5] under a unified framework, which accounts for both SGD, and TD, which uses stochastic semi-gradient. In particular, we develop a unified global convergence analysis of a meta-algorithm with the following update,

$$\alpha(t + 1/2) \leftarrow \alpha(t) - \eta \cdot (u_{\alpha(t)}(s, a) - v(s, a) - \mu \cdot u_{\alpha(t)}(s', a')) \cdot \nabla_{\alpha} u_{\alpha(t)}(s, a), \quad (\text{D.1})$$

$$\alpha(t + 1) \leftarrow \Pi_{\mathcal{B}^0(R_u)}(\alpha(t + 1/2)) = \underset{\alpha \in \mathcal{B}^0(R_u)}{\text{argmin}} \|\alpha - \alpha(t + 1/2)\|_2, \quad (\text{D.2})$$

where $\mu \in [0, 1]$ is a constant, (s, a, s', a') is sampled from a stationary distribution ρ , and u_{α} is parametrized by the two-layer neural network $\text{NN}(\alpha; m)$ defined in (3.1). The random initialization of u_{α} is given in (3.2). We denote by $\mathbb{E}_{\text{init}}[\cdot]$ the expectation over such random initialization and $\mathbb{E}_{\rho}[\cdot]$ the expectation over (s, a) conditional on the random initialization.

Such a meta-algorithm recovers SGD for policy improvement in (3.5) when we set $\rho = \tilde{\sigma}_k$, $u_{\alpha} = f_{\theta}$, $v = \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k} + \tau_k^{-1} f_{\theta_k})$, $\mu = 0$, and $R_u = R_f$, and recovers TD for policy evaluation in (3.8) when we set $\rho = \sigma_k$, $u_{\alpha} = Q_{\omega}$, $v = (1 - \gamma) \cdot r$, $\mu = \gamma$, and $R_u = R_Q$.

To unify our analysis for SGD and TD, we assume that v in (D.1) satisfies

$$\mathbb{E}_{\rho}[(v(s, a))^2] \leq \bar{v}_1 \cdot \mathbb{E}_{\rho}[(u_{\alpha(0)}(s, a))^2] + \bar{v}_2 \cdot R_u^2 + \bar{v}_3$$

for constants $\bar{v}_1, \bar{v}_2, \bar{v}_3 \geq 0$. Also, without loss of generality, we assume that $\|(s, a)\|_2 \leq 1$ for any $s \in \mathcal{S}$ and $a \in \mathcal{A}$. In Section D.2, we set $\bar{v}_1 = 4$, $\bar{v}_2 = 4$, and $\bar{v}_3 = 0$ for SGD, and $\bar{v}_1 = 0$, $\bar{v}_2 = 0$, and $\bar{v}_3 = R_{\max}$ for TD, respectively.

For notational simplicity, we define the residual $\delta_\alpha(s, a, s', a') = u_\alpha(s, a) - v(s, a) - \mu \cdot u_\alpha(s', a')$. We denote by

$$g_{\alpha(t)}(s, a, s', a') = \delta_{\alpha(t)}(s, a, s', a') \cdot \nabla_\alpha u_{\alpha(t)}(s, a), \quad \bar{g}_{\alpha(t)} = \mathbb{E}_\rho[g_t(s, a, s', a')] \quad (\text{D.3})$$

the stochastic update vector at the t -th iteration and its population mean, respectively. For SGD, $g_{\alpha(t)}(s, a, s', a')$ corresponds to the stochastic gradient, while for TD, $g_{\alpha(t)}(s, a, s', a')$ corresponds to the stochastic semigradient.

Note that the gradient of $u_\alpha(s, a)$ with respect to α takes the form

$$\nabla_\alpha u_\alpha(s, a) = 1/\sqrt{m} \cdot (b_1 \cdot \mathbb{1}\{[\alpha]_1^\top(s, a) > 0\} \cdot (s, a)^\top, \dots, b_m \cdot \mathbb{1}\{[\alpha]_m^\top(s, a) > 0\} \cdot (s, a)^\top)^\top \in \mathbb{R}^{md}$$

almost everywhere, which yields

$$\|\nabla_\alpha u_\alpha(s, a)\|_2^2 = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{[\alpha]_i^\top(s, a) > 0\} \cdot \|(s, a)\|_2^2 \leq 1.$$

Therefore, $u_\alpha(s, a)$ is 1-Lipschitz continuous with respect to α .

In the following, we first show in Section D.1 that the overparametrization of u_α ensures that it behaves similarly as its local linearization at the random initialization $\alpha(0)$ defined in (3.2). Then in Section D.2, we establish the global convergence of the meta-algorithm defined in (D.1) and (D.2), which implies the global convergence of SGD and TD.

D.1 Local Linearization

In this section, we first define a local linearization of the two-layer neural network u_α at its random initialization and then characterize the error induced by local linearization. We define

$$u_\alpha^0(s, a) = \frac{1}{\sqrt{m}} \sum_{i=1}^m b_i \cdot \mathbb{1}\{[\alpha(0)]_i^\top(s, a) > 0\} \cdot [\alpha]_i^\top(s, a). \quad (\text{D.4})$$

The linearity of u_α^0 with respect to α yields

$$\langle \nabla_\alpha u_\alpha^0(s, a), \alpha \rangle = u_\alpha^0(s, a). \quad (\text{D.5})$$

The following lemma characterizes how far $u_{\alpha(t)}^0$ deviates from $u_{\alpha(t)}$ for $\alpha(t) \in \mathcal{B}^0(R_u)$.

Lemma D.1. For any $\alpha' \in \mathcal{B}^0(R_u)$, we have

$$\mathbb{E}_{\text{init}, \rho}[(u_{\alpha'}(s, a) - u_{\alpha'}^0(s, a))^2] = O(R_u^3 m^{-1/2}).$$

Proof. By the definition of u_α in (3.1), we have

$$|u_{\alpha'}(s, a) - u_{\alpha'}^0(s, a)| \quad (\text{D.6})$$

$$\begin{aligned} &\leq \frac{1}{\sqrt{m}} \left| \sum_{i=1}^m b_i \cdot (\mathbb{1}\{[\alpha(0)]_i^\top(s, a) > 0\} - \mathbb{1}\{[\alpha']_i^\top(s, a) > 0\}) \cdot ([\alpha(0)]_i^\top(s, a) + \|[\alpha']_i - [\alpha(0)]_i\|_2) \right| \\ &\leq \frac{1}{\sqrt{m}} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|[\alpha']_i - [\alpha(0)]_i\|_2\} \cdot (\|[\alpha(0)]_i^\top(s, a)\| + \|[\alpha']_i - [\alpha(0)]_i\|_2), \end{aligned}$$

where the second inequality follows from $|b_i| = 1$ and the fact that

$$\mathbb{1}\{[\alpha(t)]_i^\top(s, a) > 0\} \neq \mathbb{1}\{[\alpha(0)]_i^\top(s, a) > 0\}$$

implies

$$|[\alpha(0)]_i^\top(s, a)| \leq |[\alpha(t)]_i^\top(s, a) - [\alpha(0)]_i^\top(s, a)| \leq \|[\alpha(0)]_i - [\alpha(t)]_i\|_2.$$

Next, applying the inequality $\mathbb{1}\{|z| \leq y\}|z| \leq \mathbb{1}\{|z| \leq y\}y$ to the right-hand side of (D.6), we obtain

$$\begin{aligned} &|u_{\alpha'}(s, a) - u_{\alpha'}^0(s, a)| \\ &\leq \frac{2}{\sqrt{m}} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|[\alpha']_i - [\alpha(0)]_i\|_2\} \cdot \|[\alpha']_i - [\alpha(0)]_i\|_2. \end{aligned} \quad (\text{D.7})$$

Further applying the Cauchy-Schwarz inequality to (D.7) and invoking the upper bound $\|\alpha' - \alpha(0)\|_2 \leq R_u$, we obtain

$$|u_{\alpha'}(s, a) - u_{\alpha'}^0(s, a)|^2 \leq \frac{4R_u^2}{m} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|[\alpha']_i - [\alpha(0)]_i\|_2\}. \quad (\text{D.8})$$

Taking expectation on the both sides and invoking Assumption 4.4, we obtain

$$\mathbb{E}_{\text{init},\rho}[(u_{\alpha'}(s, a) - u_{\alpha'}^0(s, a))^2] \leq \frac{4cR_u^2}{m} \cdot \mathbb{E}_{\text{init}} \left[\sum_{i=1}^m \|[\alpha']_i - [\alpha(0)]_i\|_2 / \|[\alpha(0)]_i\|_2 \right]. \quad (\text{D.9})$$

By the Cauchy-Schwartz inequality, we have

$$\begin{aligned} \mathbb{E}_{\text{init}} \left[\sum_{i=1}^m \|[\alpha']_i - [\alpha(0)]_i\|_2 / \|[\alpha(0)]_i\|_2 \right] &\leq \mathbb{E}_{\text{init}} \left[\sum_{i=1}^m \|[\alpha']_i - [\alpha(0)]_i\|_2^2 \right]^{1/2} \cdot \mathbb{E}_{\text{init}} \left[\sum_{i=1}^m \|[\alpha(0)]_i\|_2^{-2} \right]^{1/2} \\ &\leq R_u \cdot \mathbb{E}_{\text{init}} \left[\sum_{i=1}^m \|[\alpha(0)]_i\|_2^{-2} \right]^{1/2}, \end{aligned}$$

where the second inequality follows from $\sum_{i=1}^m \|[\alpha']_i - [\alpha(0)]_i\|_2^2 = \|\alpha' - \alpha(0)\|_2^2 \leq R_u^2$. Therefore, we have that the right-hand side of (D.9) is $O(R_u^3 m^{-1/2})$. Thus, we obtain

$$\mathbb{E}_{\text{init},\rho}[(u_{\alpha'}(s, a) - u_{\alpha'}^0(s, a))^2] = O(R_u^3 m^{-1/2}),$$

which concludes the proof of Lemma D.1. \square

Corresponding to u_{α}^0 defined in (D.4), let $\delta_{\alpha}^0(s, a, s', a') = u_{\alpha}^0(s, a) - v(s, a) - \mu \cdot u_{\alpha}^0(s', a')$. We define the local linearization of $\bar{g}_{\alpha(t)}$, which is defined in (D.3), as

$$\bar{g}_{\alpha(t)}^0 = \mathbb{E}_{\rho}[\delta_{\alpha(t)}^0(s, a, s', a') \cdot \nabla_{\alpha} u_{\alpha(t)}^0(s, a)]. \quad (\text{D.10})$$

The following lemma characterizes the difference between $\bar{g}_{\alpha(t)}^0$ and $\bar{g}_{\alpha(t)}$.

Lemma D.2. For any $t \in [T]$, we have

$$\mathbb{E}_{\text{init}}[\|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha(t)}^0\|_2^2] = O(R_u^3 m^{-1/2}).$$

Proof. By the definition of $\bar{g}_{\alpha(t)}^0$ and $\bar{g}_{\alpha(t)}$ in (D.10) and (D.3), we have

$$\begin{aligned} \|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha(t)}^0\|_2^2 &= \|\mathbb{E}_{\rho}[\delta_{\alpha(t)}(s, a, s', a') \cdot \nabla_{\alpha} u_{\alpha(t)}(s, a) - \delta_{\alpha(t)}^0(s, a, s', a') \cdot \nabla_{\alpha} u_{\alpha(t)}^0(s, a)]\|_2^2 \\ &\leq 2 \underbrace{\mathbb{E}_{\rho}[\|\delta_{\alpha(t)}(s, a, s', a') - \delta_{\alpha(t)}^0(s, a, s', a')\|^2 \cdot \|\nabla_{\alpha} u_{\alpha(t)}(s, a)\|_2^2]}_{(i)} \\ &\quad + 2 \underbrace{\mathbb{E}_{\rho}[|\delta_{\alpha(t)}^0(s, a, s', a')| \cdot \|\nabla_{\alpha} u_{\alpha(t)}(s, a) - \nabla_{\alpha} u_{\alpha(t)}^0(s, a)\|_2^2]}_{(ii)}. \end{aligned} \quad (\text{D.11})$$

Upper Bounding (i): We have $\|\nabla_{\alpha} u_{\alpha(t)}(s, a)\|_2 \leq 1$ as $\|(s, a)\|_2 \leq 1$. Note that the difference between $\delta_{\alpha(t)}$ and $\delta_{\alpha(t)}^0$ takes the form

$$\delta_{\alpha(t)}(s, a, s', a') - \delta_{\alpha(t)}^0(s, a, s', a') = (u_{\alpha(t)}(s, a) - u_{\alpha(t)}^0(s, a)) - \mu \cdot (u_{\alpha(t)}(s', a') - u_{\alpha(t)}^0(s', a')).$$

Taking expectation on the both sides, we obtain

$$\begin{aligned} \mathbb{E}_{\text{init},\rho}[\|\delta_{\alpha(t)}(s, a, s', a') - \delta_{\alpha(t)}^0(s, a, s', a')\|^2] &\leq 2\mathbb{E}_{\text{init},\rho}[(u_{\alpha(t)}(s, a) - u_{\alpha(t)}^0(s, a))^2] + 2\mu^2 \cdot \mathbb{E}_{\text{init},\rho}[(u_{\alpha(t)}(s', a') - u_{\alpha(t)}^0(s', a'))^2] \\ &= 4\mathbb{E}_{\text{init},\rho}[(u_{\alpha(t)}(s, a) - u_{\alpha(t)}^0(s, a))^2], \end{aligned}$$

where the equality follows from $|\mu| \leq 1$ and the fact that (s, a) and (s', a') have the same marginal distribution. Thus, by Lemma D.1, we have that (i) in (D.11) is $O(R_u^3 m^{-1/2})$.

Upper Bounding (ii): First, by the Hölder's inequality, we have

$$\begin{aligned} \mathbb{E}_{\rho}[|\delta_{\alpha(t)}^0(s, a, s', a')| \cdot \|\nabla_{\alpha} u_{\alpha(t)}(s, a) - \nabla_{\alpha} u_{\alpha(t)}^0(s, a)\|_2^2] &\leq \mathbb{E}_{\rho}[\|\delta_{\alpha(t)}^0(s, a, s', a')\|^2] \cdot \mathbb{E}_{\rho}[\|\nabla_{\alpha} u_{\alpha(t)}(s, a) - \nabla_{\alpha} u_{\alpha(t)}^0(s, a)\|_2^2]. \end{aligned}$$

We use $|u_{\alpha(t)}^0(s, a) - u_{\alpha(0)}^0(s, a)| \leq \|\alpha(t) - \alpha(0)\|_2 \leq R_u$ to obtain

$$\begin{aligned} |\delta_{\alpha(t)}^0(s, a, s', a')|^2 &= (u_{\alpha(t)}^0(s, a) - v(s, a) - \mu \cdot u_{\alpha(t)}^0(s', a'))^2 \\ &\leq 3((u_{\alpha(t)}^0(s, a))^2 + (v(s, a))^2 + \mu^2 \cdot (u_{\alpha(t)}^0(s', a'))^2) \\ &\leq 3(u_{\alpha(0)}^0(s, a))^2 + 3(u_{\alpha(0)}^0(s', a'))^2 + 6R_u^2 + 3(v(s, a))^2. \end{aligned} \quad (\text{D.12})$$

Next we characterize $\|\nabla_\alpha u_{\alpha(t)}(s, a) - \nabla_\alpha u_{\alpha(0)}^0(s, a)\|_2$ in (ii). Recall that

$$\begin{aligned} \nabla_\alpha u_\alpha(s, a) &= 1/\sqrt{m} \cdot (b_1 \cdot \mathbb{1}\{[\alpha]_1^\top(s, a) > 0\} \cdot (s, a)^\top, \dots, b_m \cdot \mathbb{1}\{[\alpha]_m^\top(s, a) > 0\} \cdot (s, a)^\top)^\top, \\ \text{and} \\ \nabla_\alpha u_\alpha^0(s, a) &= 1/\sqrt{m} \cdot (b_1 \cdot \mathbb{1}\{[\alpha(0)]_1^\top(s, a) > 0\} \cdot (s, a)^\top, \dots, b_m \cdot \mathbb{1}\{[\alpha(0)]_m^\top(s, a) > 0\} \cdot (s, a)^\top)^\top. \end{aligned}$$

We have

$$\begin{aligned} \|\nabla_\alpha u_{\alpha(t)}(s, a) - \nabla_\alpha u_{\alpha(0)}^0(s, a)\|_2^2 &= \frac{1}{m} \sum_{i=1}^m (\mathbb{1}\{[\alpha(t)]_i^\top(s, a) > 0\} - \mathbb{1}\{[\alpha(0)]_i^\top(s, a) > 0\})^2 \cdot \|(s, a)\|_2^2 \\ &\leq \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|[\alpha(t)]_i - [\alpha(0)]_i\|_2\}, \end{aligned} \quad (\text{D.13})$$

where the inequality follows from the same arguments used to derive (D.6). Plugging (D.12) and (D.13) into (ii) and recalling that

$$\mathbb{E}_\rho[(v(s, a))^2] \leq \bar{v}_1 \cdot \mathbb{E}_\rho[(u_{\alpha(0)}(s, a))^2] + \bar{v}_2 \cdot R_u^2 + \bar{v}_3,$$

we find that it remains to upper bound the following two terms

$$\mathbb{E}_{\text{init}, \rho} \left[\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|[\alpha(t)]_i - [\alpha(0)]_i\|_2\} \right], \quad (\text{D.14})$$

and

$$\mathbb{E}_{\text{init}} \left[\mathbb{E}_\rho[(u_{\alpha(0)}^0(s, a))^2] \cdot \mathbb{E}_\rho \left[\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|[\alpha(t)]_i - [\alpha(0)]_i\|_2\} \right] \right]. \quad (\text{D.15})$$

We already show in the proof of Lemma D.1 that (D.14) is $O(R_u m^{-1/2})$. We characterize (D.15) in the following. For the random initialization of $u_\alpha(s, a)$ in (3.2), we have

$$\mathbb{E}_\rho[(u_{\alpha(0)}^0(s, a))^2] = \frac{1}{m} \cdot \mathbb{E}_\rho \left[\sum_{i=1}^m \sigma([\alpha(0)]_i^\top(s, a))^2 + \sum_{1 \leq i \neq j \leq m} b_i b_j \cdot \sigma([\alpha(0)]_i^\top(s, a)) \cdot \sigma([\alpha(0)]_j^\top(s, a)) \right],$$

plugging which into (D.15) gives

$$\begin{aligned} &\mathbb{E}_{\text{init}} \left[\mathbb{E}_\rho[(u_{\alpha(0)}^0(s, a))^2] \cdot \mathbb{E}_\rho \left[\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|\alpha(t) - \alpha(0)\|_2\} \right] \right] \\ &\leq \mathbb{E}_{\text{init}} \left[\frac{1}{m} \cdot \mathbb{E}_\rho \left[\sum_{i=1}^m \sigma([\alpha(0)]_i^\top(s, a))^2 + \sum_{1 \leq i \neq j \leq m} b_i b_j \cdot \sigma([\alpha(0)]_i^\top(s, a)) \cdot \sigma([\alpha(0)]_j^\top(s, a)) \right] \right. \\ &\quad \cdot \frac{c}{m} \cdot \left(\sum_{i=1}^m \|[\alpha(t)]_i - [\alpha(0)]_i\|_2^2 \right)^{1/2} \cdot \left(\sum_{i=1}^m \frac{1}{\|[\alpha(0)]_i\|_2^2} \right)^{1/2} \Big], \end{aligned}$$

where we use the same arguments applied to (D.8) in the proof of Lemma D.1. Note that b_i, b_j are independent of $\alpha(0)$, $\mathbb{E}_{\text{init}}[b_i b_j] = 0$, and $\sum_{i=1}^m \|[\alpha(t)]_i - [\alpha(0)]_i\|_2^2 = \|\alpha(t) - \alpha(0)\|_2^2 \leq R_u^2$. We further obtain

$$\begin{aligned} &\mathbb{E}_{\text{init}} \left[\mathbb{E}_\rho[(u_{\alpha(0)}^0(s, a))^2] \cdot \mathbb{E}_\rho \left[\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{[\alpha(0)]_i^\top(s, a) \leq \|[\alpha(t)]_i - [\alpha(0)]_i\|_2\} \right] \right] \\ &\leq \frac{c R_u}{m^2} \cdot \mathbb{E}_{\text{init}} \left[\mathbb{E}_\rho \left[\sum_{i=1}^m \sigma([\alpha(0)]_i^\top(s, a))^2 \right] \cdot \left(\sum_{i=1}^m \frac{1}{\|[\alpha(0)]_i\|_2^2} \right)^{1/2} \right] \\ &\leq \frac{c R_u}{m^2} \cdot \mathbb{E}_{\text{init}} \left[\left(\sum_{i=1}^m \|[\alpha(0)]_i\|_2^2 \right) \cdot \left(\sum_{i=1}^m \frac{1}{\|[\alpha(0)]_i\|_2^2} \right)^{1/2} \right]. \end{aligned}$$

Finally, by the Cauchy-Schwarz inequality, we have

$$\begin{aligned} &\mathbb{E}_{\text{init}} \left[\left(\sum_{i=1}^m \|[\alpha(0)]_i\|_2^2 \right) \cdot \left(\sum_{i=1}^m \frac{1}{\|[\alpha(0)]_i\|_2^2} \right)^{1/2} \right] \\ &\leq \mathbb{E}_{\text{init}} \left[\left(\sum_{i=1}^m \|[\alpha(0)]_i\|_2^2 \right)^2 \right]^{1/2} \cdot \mathbb{E}_{\text{init}} \left[\sum_{i=1}^m \frac{1}{\|[\alpha(0)]_i\|_2^2} \right]^{1/2}, \end{aligned}$$

whose right-hand side is $O(m^{3/2})$. Thus, we obtain that (D.15) is $O(R_u m^{-1/2})$ and (ii) in (D.11) is $O(R_u^3 m^{-1/2})$, which concludes the proof of Lemma D.2. \square

D.2 Global Convergence

In this section, we establish the global convergence of the meta-algorithm defined in (D.1) and (D.2). We first present the following lemma for characterizing the variance of the stochastic update vector $g_{\alpha(t)}(s, a, s', a')$ defined in (D.3), which later allows us to focus on tracking its mean in the global convergence analysis.

Lemma D.3 (Variance of the Stochastic Update Vector). There exists a constant $\xi_g^2 = O(R_u^2)$ independent of t , such that for any $t \leq T$, it holds that

$$\mathbb{E}_{\text{init}, \rho}[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha(t)}\|_2^2] \leq \xi_g^2.$$

Proof. Since we have

$$\begin{aligned} \mathbb{E}_{\text{init}, \rho}[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha(t)}\|_2^2] &= \mathbb{E}_{\text{init}}[\mathbb{E}_{\rho}[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha(t)}\|_2^2]] \\ &\leq \mathbb{E}_{\text{init}}[\mathbb{E}_{\rho}[\|g_{\alpha(t)}(s, a, s', a')\|_2^2]] = \mathbb{E}_{\text{init}, \rho}[\|g_{\alpha(t)}(s, a, s', a')\|_2^2], \end{aligned}$$

it suffices to prove that $\mathbb{E}[\|g_{\alpha(t)}(s, a, s', a')\|_2^2] = O(R_u^2)$. By the definition of $\mathbb{E}_{\rho}[\|g_{\alpha(t)}(s, a, s', a')\|_2^2]$ in (D.3), using $\|\nabla_{\alpha(t)} u_{\alpha(t)}(s, a)\|_2^2 \leq 1$, we obtain

$$\begin{aligned} \mathbb{E}_{\rho}[\|g_{\alpha(t)}(s, a, s', a')\|_2^2] &= \mathbb{E}_{\rho}[\|\delta_{\alpha(t)}(s, a, s', s') \cdot \nabla_{\alpha} u_{\alpha(t)}(s, a)\|_2^2] \\ &\leq \mathbb{E}_{\rho}[\|\delta_{\alpha(t)}(s, a, s', s')\|^2]. \end{aligned} \tag{D.16}$$

Then, by similar arguments used in the derivation of (D.12), we obtain

$$\begin{aligned} \mathbb{E}_{\text{init}, \rho}[\|\delta_{\alpha(t)}(s, a, s', s')\|^2] &\leq 6\mathbb{E}_{\text{init}, \rho}[(u_{\alpha(0)}(s, a))^2] + 6R_u^2 + 3\mathbb{E}_{\text{init}, \rho}[(v(s, a))^2] \\ &\leq (6 + 3\bar{v}_1) \cdot \mathbb{E}_{\text{init}, \rho}[(u_{\alpha(0)}(s, a))^2] + (6 + \bar{v}_2)R_u^2 + 3\bar{v}_3^2. \end{aligned} \tag{D.17}$$

Note that by $\|(s, a)\|_2 \leq 1$, we have

$$\mathbb{E}_{\text{init}, \rho}[(u_{\alpha(0)}(s, a))^2] = \mathbb{E}_{z \sim \mathcal{N}(0, I_d/d), \rho}[\sigma(z^\top(s, a))^2] \leq \mathbb{E}_{z \sim \mathcal{N}(0, I_d/d)}[\|z\|_2^2] = 1,$$

which together with (D.16) and (D.17) implies $\mathbb{E}_{\text{init}, \rho}[\|g_{\alpha(t)}(s, a, s', a')\|_2^2] = O(R_u^2)$. Thus, we complete the proof of Lemma D.3. \square

Before presenting the global convergence result of the meta-algorithm defined in (D.1), we first define $u_{\alpha^*}^0$, which later become the exact learning target of the meta-algorithm defined in (D.1) and (D.2). In specific, we define the approximate stationary point as $\alpha^* \in \mathcal{B}^0(R_u)$ such that

$$\alpha^* = \Pi_{\mathcal{B}^0(R_u)}(\alpha^* - \eta \cdot \bar{g}_{\alpha^*}^0), \tag{D.18}$$

which is equivalent to the condition

$$\langle \bar{g}_{\alpha^*}^0, \alpha - \alpha^* \rangle \geq 0, \text{ for any } \alpha \in \mathcal{B}^0(R_u). \tag{D.19}$$

Then we establish the uniqueness and existence of $u_{\alpha^*}^0$ with α^* defined in D.18. We first define the operator

$$\mathcal{T}u(s, a) = \mathbb{E}[v(s, a) + \mu \cdot u(s', a') \mid s' \sim \mathcal{P}(\cdot \mid s, a), a \sim \pi(\cdot \mid s')]. \tag{D.20}$$

Then using the definition of \mathcal{T} in (D.20) and plugging the definition of $\bar{g}_{\alpha^*}^0$ in (D.4) into (D.19), we obtain

$$\langle u_{\alpha^*}^0 - \mathcal{T}u_{\alpha^*}^0, u_{\alpha^*}^0 - u_{\alpha^*}^0 \rangle_{\rho} \geq 0, \text{ for any } u_{\alpha^*}^0 \in \mathcal{F}_{B, m},$$

which is equivalent to $u_{\alpha^*}^0 = \Pi_{\mathcal{F}_{B, m}} \mathcal{T}u_{\alpha^*}^0$. Here the projection $\Pi_{\mathcal{F}_{B, m}}$ is defined with respect to the ℓ_2 -distance under measure ρ . Finally, as we have the following contraction inequality

$$\begin{aligned} &\mathbb{E}_{\rho}[(\Pi_{\mathcal{F}_{B, m}} \mathcal{T}u_{\alpha}^0(s, a) - \Pi_{\mathcal{F}_{B, m}} \mathcal{T}u_{\alpha'}^0(s, a))^2] \\ &\leq \mathbb{E}_{\rho}[(\mathcal{T}u_{\alpha}^0(s, a) - \mathcal{T}u_{\alpha'}^0(s, a))^2] \\ &= \mu^2 \cdot \mathbb{E}_{\rho}[(\mathbb{E}[u_{\alpha}^0(s', a') \mid s' \sim \mathcal{P}(\cdot \mid s, a), a' \sim \pi(\cdot \mid s')]) - \mathbb{E}[u_{\alpha'}^0(s', a') \mid s' \sim \mathcal{P}(\cdot \mid s, a), a' \sim \pi(\cdot \mid s')])^2] \\ &\leq \mu^2 \cdot \mathbb{E}_{\rho}[(u_{\alpha}^0(s, a) - u_{\alpha'}^0(s, a))^2], \end{aligned}$$

we know that such fixed-point solution $u_{\alpha^*}^0$ uniquely exists.

Now, with a well-defined learning target $u_{\alpha^*}^0$, we are ready to prove the global convergence of the meta-algorithm defined in (D.1) and (D.2) with two-layer neural network approximation.

Theorem D.4. Suppose that we run $T \geq 64/(1 - \mu)^2$ iterations of the meta-algorithm defined in (D.1) and (D.2). Setting the stepsize $\eta = T^{-1/2}$, we have

$$\mathbb{E}_{\text{init}, \rho}[(u_{\bar{\alpha}}(s, a) - u_{\alpha^*}^0(s, a))^2] = O(R_u^2 T^{-1/2} + R_u^{5/2} m^{-1/4} + R_u^3 m^{-1/2}),$$

where $\bar{\alpha} = 1/T \cdot \sum_{t=0}^{T-1} \alpha(t)$ and α^* is the approximate stationary point defined in (D.18).

Proof. The proof of the theorem consists of two parts. We first analyze the progress of each step. Then based on such one-step analysis, we establish the error bound of the approximation via two-layer neural network u_α .

One-Step Analysis: For any $t < T$, using the stationarity condition in (D.18) and the convexity of $\mathcal{B}^0(R_u)$, we obtain

$$\begin{aligned} & \mathbb{E}_\rho[\|\alpha(t+1) - \alpha^*\|_2^2 | \alpha(t)] \\ &= \mathbb{E}_\rho[\|\Pi_{\mathcal{B}^0(R_u)}(\alpha(t) - \eta \cdot g_{\alpha(t)}(s, a, s', a')) - \Pi_{\mathcal{B}^0(R_u)}(\alpha^* - \eta \bar{g}_{\alpha^*}^0)\|_2^2 | \alpha(t)] \\ &\leq \mathbb{E}_\rho[\|(\alpha(t) - \alpha^*) - \eta \cdot (g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha^*}^0)\|_2^2 | \alpha(t)] \\ &= \|\alpha(t) - \alpha^*\|_2^2 - 2\eta \cdot \langle \bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle + \eta^2 \cdot \mathbb{E}_\rho[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha^*}^0\|_2^2 | \alpha(t)]. \end{aligned} \quad (\text{D.21})$$

In the following, we upper bound the last two terms in (D.21). First, to upper bound $\mathbb{E}_\rho[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha^*}^0\|_2^2 | \alpha(t)]$, by the Cauchy-Schwarz inequality we have

$$\begin{aligned} & \mathbb{E}_\rho[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha^*}^0\|_2^2 | \alpha(t)] \\ &\leq 2\mathbb{E}_\rho[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha(t)}\|_2^2 | \alpha(t)] + 2\|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0\|_2^2 \\ &\leq 2\mathbb{E}_\rho[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha(t)}\|_2^2 | \alpha(t)] + 4\|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0\|_2^2 + 4\|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0\|_2^2, \end{aligned} \quad (\text{D.22})$$

where the total expectation on the first two terms on the right-hand side are characterized in Lemmas D.3 and D.2, respectively. To characterize $\|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0\|_2^2$, again using $\|(s, a)\|_2 \leq 1$, we have

$$\begin{aligned} \|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0\|_2^2 &= \mathbb{E}_\rho[(\delta_{\alpha(t)}(s, a, s', a') - \delta_{\alpha^*}(s, a, s', a'))^2 \cdot \|\nabla_\alpha u_{\alpha(t)}^0(s, a)\|_2^2] \\ &\leq \mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a)) - \mu \cdot (u_{\alpha(t)}^0(s', a') - u_{\alpha^*}^0(s', a'))]^2. \end{aligned} \quad (\text{D.23})$$

For the right-hand side of (D.23), we use the Cauchy-Schwarz inequality on the interaction term and obtain

$$\begin{aligned} & \mathbb{E}_\rho[(u_{\alpha(t)}^0(s', a') - u_{\alpha^*}^0(s', a')) \cdot (u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))] \\ &\leq \mathbb{E}_\rho[(u_{\alpha(t)}^0(s', a') - u_{\alpha^*}^0(s', a'))^2]^{1/2} \cdot \mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2]^{1/2} \\ &= \mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2], \end{aligned} \quad (\text{D.24})$$

where in the last line we use the fact that (s, a) and (s', a') have the same marginal distribution. Thus, we obtain

$$\|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0\|_2^2 \leq 4\mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2]. \quad (\text{D.25})$$

Next, to upper bound $\langle \bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle$, we use the Hölder's inequality to obtain

$$\begin{aligned} \langle \bar{g}_{\alpha(t)} - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle &= \langle \bar{g}_{\alpha(t)} - \bar{g}_{\alpha(t)}^0, \alpha(t) - \alpha^* \rangle + \langle \bar{g}_{\alpha(t)}^0 - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle \\ &\geq -\|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha(t)}^0\|_2 \cdot \|\alpha(t) - \alpha^*\|_2 + \langle \bar{g}_{\alpha(t)}^0 - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle \\ &\geq -R_u \|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha(t)}^0\|_2 + \langle \bar{g}_{\alpha(t)}^0 - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle, \end{aligned} \quad (\text{D.26})$$

where the second inequality follows from $\|\alpha(t) - \alpha^*\|_2 \leq R_u$. For the term $\langle \bar{g}_{\alpha(t)}^0 - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle$ on the right-hand side of (D.26), we have

$$\begin{aligned} & \langle \bar{g}_{\alpha(t)}^0 - \bar{g}_{\alpha^*}^0, \alpha(t) - \alpha^* \rangle \\ &= \mathbb{E}_\rho\left[\left((u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a)) - \mu \cdot (u_{\alpha(t)}^0(s', a') - u_{\alpha^*}^0(s', a'))\right) \cdot \langle \nabla_\alpha u_{\alpha(t)}^0(s, a), \alpha(t) - \alpha^* \rangle\right] \\ &= \mathbb{E}_\rho\left[\left((u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a)) - \mu \cdot (u_{\alpha(t)}^0(s', a') - u_{\alpha^*}^0(s', a'))\right) \cdot (u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))\right] \\ &\geq \mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2] - \mu \cdot \mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2] \\ &\geq (1 - \mu) \cdot \mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2], \end{aligned} \quad (\text{D.27})$$

where the second equality and the first inequality follow from (D.5) and (D.24), respectively.

Therefore, combining (D.21) with (D.22), (E.4), (D.26), and (D.27), we obtain

$$\begin{aligned} & \mathbb{E}_\rho[\|\alpha(t+1) - \alpha^*\|_2^2 | \alpha(t)] \\ &\leq \|\alpha(t) - \alpha^*\|_2^2 - (2\eta(1 - \gamma) - 8\eta^2) \cdot \mathbb{E}_\rho[(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2 | \alpha(t)] \\ &\quad + 2\eta^2 \|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha(t)}^0\|_2^2 + 2\eta R_u \|\bar{g}_{\alpha(t)} - \bar{g}_{\alpha(t)}^0\|_2 + \eta^2 \cdot \mathbb{E}_\rho[\|g_{\alpha(t)}(s, a, s', a') - \bar{g}_{\alpha(t)}\|_2^2 | \alpha(t)]. \end{aligned} \quad (\text{D.28})$$

Error Bound: Rearranging (D.28), we obtain

$$\begin{aligned} & \mathbb{E}_\rho[(u_{\alpha(t)}(s, a) - u_{\alpha^*}^0(s, a))^2 | \alpha(t)] \\ & \leq \mathbb{E}_\rho[2(u_{\alpha(t)}(s, a) - u_{\alpha(t)}^0(s, a))^2 + 2(u_{\alpha(t)}^0(s, a) - u_{\alpha^*}^0(s, a))^2 | \alpha(t)] \\ & \leq (\eta(1 - \gamma) - 4\eta^2)^{-1} \cdot (\|\alpha(t) - \alpha^*\|_2^2 - \mathbb{E}_\rho[\|\alpha(t+1) - \alpha^*\|_2^2 | \alpha(t)] + \xi_\alpha^2 \eta^2) \\ & \quad + O(R_u^{5/2} m^{-1/4} + R_u^3 m^{-1/2}). \end{aligned} \quad (\text{D.29})$$

Taking total expectation on both sides of (D.29) and telescoping for $t + 1 \in [T]$, we further obtain

$$\begin{aligned} \mathbb{E}_{\text{init}, \rho}[(u_{\alpha}(s, a) - u_{\alpha^*}^0(s, a))^2] & \leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{\text{init}, \rho}[(u_{\alpha(t)}(s, a) - u_{\alpha^*}^0(s, a))^2] \\ & \leq T^{-1} \cdot (\eta(1 - \gamma) - 4\eta^2)^{-1} \cdot (\mathbb{E}_{\text{init}}[\|\alpha(0) - \alpha^*\|_2^2] + T\xi_\alpha^2 \eta^2) \\ & \quad + O(R_u^{5/2} m^{-1/4} + R_u^3 m^{-1/2}). \end{aligned} \quad (\text{D.30})$$

Let $T \geq 64/(1 - \mu)^2$ and $\eta = T^{-1/2}$, it holds that $T^{-1/2} \cdot (\eta(1 - \gamma) - 4\eta^2)^{-1} \leq 16(1 - \gamma)^{-1/2}$ and $T\eta^2 \leq 1$, which together with (D.30) implies

$$\begin{aligned} & \mathbb{E}_{\text{init}, \rho}[(u_{\alpha(t)}(s, a) - u_{\alpha^*}^0(s, a))^2 | \alpha(t)] \\ & \leq \frac{16}{(1 - \mu)^2 \sqrt{T}} \cdot (\mathbb{E}_{\text{init}}[\|\alpha(0) - \alpha^*\|_2^2] + \xi_\alpha^2) + O(R_u^{5/2} m^{-1/4} + R_u^3 m^{-1/2}) \\ & \leq \frac{16(R_\alpha^2 + \xi_\alpha^2)}{(1 - \mu)^2 \sqrt{T}} + O(R_u^{5/2} m^{-1/4} + R_u^3 m^{-1/2}) = O(R_u^2 T^{-1/2} + R_u^{5/2} m^{-1/4} + R_u^3 m^{-1/2}), \end{aligned}$$

where in the second inequality we use $\|\alpha(0) - \alpha^*\|_2 \leq R_u$ and in the equality we use Lemma D.3. Thus, we conclude the proof of Theorem D.4. \square

Following the definition of u_α^0 in (D.4), we define the local linearization of Q_ω at the initialization as

$$Q_\omega^0(s, a) = \frac{1}{\sqrt{m_Q}} \sum_{i=1}^{m_Q} b_i \cdot \mathbf{1}\{[\omega(0)]_i^\top(s, a) \geq 0\} \cdot [\omega]_i^\top(s, a).$$

Similarly, for f_θ we define

$$f_\theta^0(s, a) = \frac{1}{\sqrt{m_f}} \sum_{i=1}^{m_f} b_i \cdot \mathbf{1}\{[\theta(0)]_i^\top(s, a) \geq 0\} \cdot [\theta]_i^\top(s, a).$$

In the sequel, we show that Theorem D.4 implies both Theorems 4.5 and 4.6.

To obtain Theorem 4.5, we set $\rho = \tilde{\sigma}_k$, $u_\alpha = f_\theta$, $v = \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k} + \tau_k^{-1} f_{\theta_k})$, $\mu = 0$, and $R_u = R_f$. Using τ_{k+1} , τ_k , and β_k specified in Algorithm 1, we have

$$\begin{aligned} \mathbb{E}_{\tilde{\sigma}_k}[(v(s, a))^2] & \leq 2\tau_{k+1}^2 \cdot (\beta_k^{-2} \cdot \mathbb{E}_{\tilde{\sigma}_k}[(Q_{\omega_k}(s, a))^2] + \tau_k^{-2} \cdot \mathbb{E}_{\tilde{\sigma}_k}[(f_{\theta_k}(s, a))^2]) \\ & \leq 4\mathbb{E}_{\tilde{\sigma}_k}[(f_{\theta(0)}(s, a))^2] + 4R_f^2, \end{aligned}$$

where in the second inequality we use $\tau_{k+1}^2 \beta_k^{-2} + \tau_{k+1}^2 \tau_k^{-2} \leq 1$ and the fact that $(Q_{\omega_k}(s, a))^2 \leq 2(Q_{\omega(0)}(s, a))^2 + 2R_Q^2$ and $(f_{\theta_k}(s, a))^2 \leq 2(f_{\theta(0)}(s, a))^2 + 2R_f^2$, which is a consequence of the 1-Lipschitz continuity of the neural network with respect to the weights. Also note that $Q_{\omega(0)}(s, a) = f_{\theta(0)}(s, a)$ due to the fact that Q_{ω_k} and f_{θ_k} share the same initialization. Thus, we have $\bar{v}_1 = 4$, $\bar{v}_2 = 4$, and $\bar{v}_3 = 0$. Moreover, by $f_{\theta^*}^0 = \Pi_{\mathcal{F}_{R_f, m_f}} \mathcal{T} f_{\theta^*}^0 = \Pi_{\mathcal{F}_{R_f, m}}(\tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k} + \tau_k^{-1} f_{\theta_k}))$, we have

$$f_{\theta^*}^0 = \operatorname{argmin}_{f \in \mathcal{F}_{R_f, m_f}} \left\{ \|f - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k} + \tau_k^{-1} f_{\theta_k})\|_{2, \tilde{\sigma}_k} \right\},$$

which together with the fact that $\tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}^0(s, a) + \tau_k^{-1} f_{\theta_k}^0(s, a)) \in \mathcal{F}_{R_f, m_f}$ implies

$$\begin{aligned} & \mathbb{E}_{\text{init}, \tilde{\sigma}_k}[(f_{\theta^*}^0(s, a) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}^0(s, a) + \tau_k^{-1} f_{\theta_k}^0(s, a)))^2] \\ & \leq \mathbb{E}_{\text{init}, \tilde{\sigma}_k}[(\tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}^0(s, a) + \tau_k^{-1} f_{\theta_k}^0(s, a)) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)))^2] \\ & \leq \tau_{k+1}^2 \beta_k^{-2} \cdot \mathbb{E}_{\text{init}, \tilde{\sigma}_k}[(Q_{\omega_k}^0(s, a) - Q_{\omega_k}(s, a))^2] + \tau_{k+1}^2 \tau_k^{-2} \cdot \mathbb{E}_{\text{init}, \tilde{\sigma}_k}[(f_{\theta_k}^0(s, a) - f_{\theta_k}(s, a))^2] \\ & = O(R_f^3 m_f^{-1/2}). \end{aligned} \quad (\text{D.31})$$

Finally, plugging (D.31) into Theorem D.4 for f_θ , we obtain

$$\begin{aligned} & \mathbb{E}_{\text{init}, \tilde{\sigma}_k} \left[\left(f_{\bar{\theta}}(s, a) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)) \right)^2 \right] \\ & \leq 2\mathbb{E}_{\text{init}, \tilde{\sigma}_k} \left[(f_{\bar{\theta}}(s, a) - f_{\theta^*}^0(s, a))^2 \right] + 2\mathbb{E}_{\text{init}, \tilde{\sigma}_k} \left[\left(f_{\theta^*}^0(s, a) - \tau_{k+1} \cdot (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)) \right)^2 \right] \\ & = O(R_f^2 T^{-1/2} + R_f^{5/2} m_f^{-1/4} + R_f^3 m_f^{-1/2}) \end{aligned}$$

which gives Theorem 4.5.

To obtain Theorem 4.6, we set $\rho = \sigma_k$, $u_\alpha = Q_\omega$, $v = (1 - \gamma) \cdot r$, $\mu = \gamma$ and $R_u = R_Q$. Correspondingly, we have $\bar{v}_1 = 0$, $\bar{v}_2 = 0$, $\bar{v}_3 = R_{\max}^2$ and $u_{\alpha^*}^0 = Q_{\omega^*}^0$. Moreover, by the definition of the operator \mathcal{T} in (D.20), we have $\mathcal{T} = \mathcal{T}^{\pi_{\theta_k}}$, which implies $Q^{\pi_{\theta_k}} = \mathcal{T}Q^{\pi_{\theta_k}}$. Meanwhile, by Assumption 4.3, we have $Q^{\pi_{\theta_k}} \in \mathcal{F}_{R_Q, m_Q}$, which implies $Q^{\pi_{\theta_k}} = \Pi_{\mathcal{F}_{R_Q, m_Q}} Q^{\pi_{\theta_k}} = \Pi_{\mathcal{F}_{R_Q, m_Q}} \mathcal{T}Q^{\pi_{\theta_k}}$. Since we already show that $Q_{\omega^*}^0$ is the unique solution to the equation $Q = \Pi_{\mathcal{F}_{R_Q, m_Q}} \mathcal{T}Q$, we obtain $Q_{\alpha^*}^0 = Q^{\pi_{\theta_k}}$. Therefore, we can substitute $Q_{\alpha^*}^0$ with $Q^{\pi_{\theta_k}}$ in Theorem D.4 to obtain Theorem 4.6.

E Proofs for Section 4.2

Proof of Lemma 4.7. We first have

$$\pi_{k+1}(a | s) = \exp\{\beta_k^{-1} Q^{\pi_{\theta_k}}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)\} / Z_{k+1}(s),$$

and

$$\pi_{\theta_{k+1}}(a | s) = \exp\{\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, a)\} / Z_{\theta_{k+1}}(s).$$

Here $Z_{k+1}(s), Z_{\theta_{k+1}}(s) \in \mathbb{R}$ are normalization factors, which are defined as

$$\begin{aligned} Z_{k+1}(s) &= \sum_{a' \in \mathcal{A}} \exp\{\beta_k^{-1} Q^{\pi_{\theta_k}}(s, a') + \tau_k^{-1} f_{\theta_k}(s, a')\}, \\ Z_{\theta_{k+1}}(s) &= \sum_{a' \in \mathcal{A}} \exp\{\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, a')\}, \end{aligned} \tag{E.1}$$

respectively. Thus, we reformulate the inner product in (4.5) as

$$\begin{aligned} & \langle \log \pi_{k+1}(\cdot | s) - \log \pi_{\theta_{k+1}}(\cdot | s), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\ &= \langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - (\beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot) + \tau_k^{-1} f_{\theta_k}(s, \cdot)), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle, \end{aligned} \tag{E.2}$$

where we use the fact that

$$\begin{aligned} & \langle \log Z_{k+1}(s) - \log Z_{\theta_{k+1}}(s), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\ &= (\log Z_{k+1}(s) - \log Z_{\theta_{k+1}}(s)) \sum_{a' \in \mathcal{A}} (\pi^*(a' | s) - \pi_{\theta_k}(a' | s)) = 0. \end{aligned}$$

Thus, it remains to upper bound the right-hand side of (E.2). We first decompose it to two terms, namely the error from learning the Q-function and the error from fitting the improved policy, that is,

$$\begin{aligned} & \langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - (\beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot) + \tau_k^{-1} f_{\theta_k}(s, \cdot)), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\ &= \underbrace{\langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - (\beta_k^{-1} Q_{\omega_k}(s, \cdot) + \tau_k^{-1} f_{\theta_k}(s, \cdot)), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle}_{(i)} \\ & \quad + \underbrace{\langle \beta_k^{-1} Q_{\omega_k}(s, \cdot) - \beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle}_{(ii)}. \end{aligned} \tag{E.3}$$

Upper Bounding (i): We have

$$\begin{aligned} & \langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - (\beta_k^{-1} Q_{\omega_k}(s, \cdot) + \tau_k^{-1} f_{\theta_k}(s, \cdot)), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\ &= \left\langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - (\beta_k^{-1} Q_{\omega_k}(s, \cdot) + \tau_k^{-1} f_{\theta_k}(s, \cdot)), \pi_0(\cdot | s) \cdot \left(\frac{\pi^*(\cdot | s)}{\pi_0(\cdot | s)} - \frac{\pi_{\theta_k}(\cdot | s)}{\pi_0(\cdot | s)} \right) \right\rangle. \end{aligned} \tag{E.4}$$

Taking expectation with respect to $s \sim \nu^*$ on the both sides of (E.4) and using the Cauchy-Schwarz inequality, we obtain

$$\begin{aligned}
& |\mathbb{E}_{\nu^*} [\langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - (\beta_k^{-1} Q_{\omega_k}(s, \cdot) + \tau_k^{-1} f_{\theta_k}(s, \cdot)), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle]| \\
&= \left| \int_{\mathcal{S}} \left\langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - (\beta_k^{-1} Q_{\omega_k}(s, \cdot) + \tau_k^{-1} f_{\theta_k}(s, \cdot)), \pi_0(\cdot | s) \cdot \nu_k(s) \cdot \left(\frac{\pi^*(\cdot | s)}{\pi_0(\cdot | s)} - \frac{\pi_{\theta_k}(\cdot | s)}{\pi_0(\cdot | s)} \right) \right\rangle \cdot \frac{\nu^*(s)}{\nu_k(s)} ds \right| \\
&= \left| \int_{\mathcal{S} \times \mathcal{A}} (\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, a) - (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a))) \cdot \left(\frac{\sigma^*(a | s)}{\tilde{\sigma}_k(a | s)} - \frac{\pi_{\theta_k}(a | s) \cdot \nu^*(s)}{\tilde{\sigma}_k(a | s)} \right) d\tilde{\sigma}_k(s, a) \right| \\
&\leq \mathbb{E}_{\tilde{\sigma}_k} [(\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, a) - (\beta_k^{-1} Q_{\omega_k}(s, a) + \tau_k^{-1} f_{\theta_k}(s, a)))^2]^{1/2} \cdot \mathbb{E}_{\tilde{\sigma}_k} \left[\left| \frac{d\sigma^*}{d\tilde{\sigma}_k} - \frac{d(\pi_{\theta_k} \nu^*)}{d\tilde{\sigma}_k} \right|^2 \right]^{1/2} \\
&\leq \tau_{k+1}^{-1} \epsilon_{k+1} \cdot \phi_k^*, \tag{E.5}
\end{aligned}$$

where in the last inequality we use the error bound in (4.3) and the definition of ϕ_k^* in (4.2).

Upper Bounding (ii): By the Cauchy-Schwarz inequality, we have

$$\begin{aligned}
& |\mathbb{E}_{\nu^*} [\langle \beta_k^{-1} Q_{\omega_k}(s, \cdot) - \beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle]| \\
&= \left| \int_{\mathcal{S} \times \mathcal{A}} (\beta_k^{-1} Q_{\omega_k}(s, a) - \beta_k^{-1} Q^{\pi_{\theta_k}}(s, a)) \cdot \left(\frac{\pi^*(a | s)}{\pi_{\theta_k}(a | s)} - \frac{\pi_{\theta_k}(a | s)}{\pi_{\theta_k}(a | s)} \right) \cdot \frac{\nu^*(s)}{\nu_k(s)} d\sigma_k(s, a) \right| \\
&\leq \mathbb{E}_{\sigma_k} [(\beta_k^{-1} Q_{\omega_k}(s, a) - \beta_k^{-1} Q^{\pi_{\theta_k}}(s, a))^2]^{1/2} \cdot \mathbb{E}_{\sigma_k} \left[\left| \frac{d\sigma^*}{d\sigma_k} - \frac{d\nu^*}{d\nu_k} \right|^2 \right]^{1/2} \\
&\leq \beta_k^{-1} \epsilon'_k \cdot \psi_k^*, \tag{E.6}
\end{aligned}$$

where in the last inequality we use the error bound in (4.4) and the definition of ψ_k^* in (4.2). Finally, combining (E.2), (E.3), (E.5), and (E.6), we have

$$\begin{aligned}
& |\mathbb{E}_{\nu^*} [\langle \log \pi_{\theta_{k+1}}(\cdot | s) - \log \pi_{k+1}(\cdot | s), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle]| \\
&\leq \tau_{k+1}^{-1} \epsilon_{k+1} \cdot \phi_k^* + \beta_k^{-1} \epsilon'_k \cdot \psi_k^*,
\end{aligned}$$

which concludes the proof of Lemma 4.7. \square

Proof of Lemma 4.8. By the triangle inequality, we have

$$\begin{aligned}
& \|\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - \tau_k^{-1} f_{\theta_k}(s, \cdot)\|_{\infty}^2 \\
&\leq 2\|\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - \tau_k^{-1} f_{\theta_k}(s, \cdot) - \beta_k^{-1} Q_{\omega_k}(s, \cdot)\|_{\infty}^2 + 2\|\beta_k^{-1} Q_{\omega_k}(s, \cdot)\|_{\infty}^2. \tag{E.7}
\end{aligned}$$

For the first term on the right-hand side of (E.7), we have

$$\mathbb{E}_{\nu^*} [\|\tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - \tau_k^{-1} f_{\theta_k}(s, \cdot) - \beta_k^{-1} Q_{\omega_k}(s, \cdot)\|_{\infty}^2] \leq |\mathcal{A}| \cdot \tau_{k+1}^{-2} \epsilon_{k+1}^2. \tag{E.8}$$

For the second term on the right-hand side of (E.7), we have

$$\mathbb{E}_{\nu^*} [\|\beta_k^{-1} Q_{\omega_k}(s, \cdot)\|_{\infty}^2] \leq \beta_k^{-2} \cdot \mathbb{E}_{\nu^*} \left[\max_{a \in \mathcal{A}} 2(Q_{\omega_0}(s, a))^2 + 2R_f^2 \right] = \beta_k^{-2} M, \tag{E.9}$$

where we use the 1-Lipschitz continuity of Q_{ω} in ω and the constraint $\|\omega_k - \omega_0\|_2 \leq R_{\omega}$. Then, taking expectation with respect to $s \sim \nu^*$ on the both sides of (E.7) and plugging in (E.8) and (E.9), we finish the proof of Lemma 4.8. \square

F Proof of Corollary 4.10

Proof. By Theorems 4.5 and 4.6, we have $\epsilon_{k+1} = O(R_f^2 T^{-1/2} + R_f^{5/2} m_f^{-1/4} + R_f^3 m_f^{-1/2})$ and $\epsilon'_k = O(R_Q^2 T^{-1/2} + R_Q^{5/2} m_Q^{-1/4} + R_Q^3 m_Q^{-1/2})$, which gives

$$\begin{aligned}
\tau_{k+1}^{-1} \epsilon_{k+1} \cdot \phi_{k+1}^* &= O(k K^{-1/2} \cdot \phi_k^* \cdot (R_f^2 T^{-1/2} + R_f^{5/2} m_f^{-1/4})), \\
|\mathcal{A}| \cdot \tau_{k+1}^{-2} \epsilon_{k+1}^2 &= O(k^2 K^{-1} \cdot |\mathcal{A}| \cdot (R_f^2 T^{-1/2} + R_f^{5/2} m_f^{-1/4})^2), \\
\beta_k^{-1} \epsilon'_k \cdot \psi_k^* &= O(K^{-1/2} \cdot \psi_k^* \cdot (R_Q^2 T^{-1/2} + R_Q^{5/2} m_Q^{-1/4})),
\end{aligned}$$

when $m_f = \Omega(R_f^2)$ and $m_Q = \Omega(R_Q^2)$.

Next, setting $m_f = R_f^{10} \cdot \Omega(K^6 \cdot \phi_k^{*4} + K^4 \cdot |\mathcal{A}|^2)$, $m_Q = \Omega(K^2 R_Q^{10} \cdot \psi_k^{*4})$ and $T = \Omega(K^3 R_f^4 \cdot \phi_k^{*2} + K R_Q^4 \cdot \psi_k^{*2})$, we further have

$$\varepsilon_k = \tau_{k+1}^{-1} \epsilon_{k+1} \cdot \phi_k^* + \beta_k^{-1} \epsilon'_k \cdot \psi_k^* = O(K^{-1}). \quad (\text{F.1})$$

Meanwhile, setting $m_f = \Omega(K^4 R_f^{10} \cdot |\mathcal{A}|^2)$ and $T = \Omega(K^2 R_f^4 \cdot |\mathcal{A}|)$, we have

$$\varepsilon'_k = |\mathcal{A}| \cdot \tau_{k+1}^{-2} \epsilon_{k+1}^2 = O(K^{-1}). \quad (\text{F.2})$$

Summing up (F.1) and (F.2) for $k+1 \in [K]$ and plugging it into Theorem 4.9, we obtain

$$\min_{0 \leq k \leq K} \{\mathcal{L}(\pi^*) - \mathcal{L}(\pi_{\theta_k})\} \leq \frac{\beta^2 \log |\mathcal{A}| + M + O(1)}{(1-\gamma)\beta \cdot \sqrt{K}},$$

which completes the proof of Corollary 4.10. \square

G Proofs of Section 5

Proof of Lemma 5.1. The proof follows that of Lemma 6.1 in [24]. By the definition of $V^\pi(s)$ in (2.1), we have

$$\begin{aligned} \mathbb{E}_{\nu^*}[V^{\pi^*}(s)] &= \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{E}_{a_t \sim \pi^*(\cdot | s_t), s_t \sim (\mathcal{P}^{\pi^*})^t \nu^*} [(1-\gamma) \cdot r(s_t, a_t)] \\ &= \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{E}_{a_t \sim \pi^*(\cdot | s_t), s_t \sim (\mathcal{P}^{\pi^*})^t \nu^*} [(1-\gamma) \cdot r(s_t, a_t) + V^\pi(s_t) - V^\pi(s_t)] \\ &= \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{E}_{s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t), a_t \sim \pi^*(\cdot | s_t), s_t \sim (\mathcal{P}^{\pi^*})^t \nu^*} [(1-\gamma) \cdot r(s_t, a_t) + \gamma \cdot V^\pi(s_{t+1}) - V^\pi(s_t)] \\ &\quad + \mathbb{E}_{\nu^*}[V^\pi(s)], \end{aligned} \quad (\text{G.1})$$

where the third inequality is obtained by taking $\mathbb{E}_{\nu^*}[V^\pi(s_0)] = \mathbb{E}_{\nu^*}[V^\pi(s)]$ out and, correspondingly, delaying $V^\pi(s_t)$ by one time step to $V^\pi(s_{t+1})$ in each term of the summation. Note that for the advantage function, by definition of the action-value function, we have

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s) = (1-\gamma) \cdot r(s, a) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} [V^\pi(s')] - V^\pi(s),$$

which together with (G.1) implies

$$\begin{aligned} \mathbb{E}_{\nu^*}[V^{\pi^*}(s)] &= \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{E}_{a_t \sim \pi^*(\cdot | s_t), s_t \sim (\mathcal{P}^{\pi^*})^t \nu^*} [A^\pi(s_t, a_t)] + \mathbb{E}_{\nu^*}[V^\pi(s)] \\ &= (1-\gamma)^{-1} \cdot \mathbb{E}_{\sigma^*}[A^\pi(s, a)] + \mathbb{E}_{\nu^*}[V^\pi(s)]. \end{aligned} \quad (\text{G.2})$$

Here the second equality follows from $(\mathcal{P}^{\pi^*})^t \nu^* = \nu^*$ for any $t \geq 0$ and $\sigma^* = \pi^* \nu^*$. Finally, note that for any given $s \in \mathcal{S}$,

$$\begin{aligned} \mathbb{E}_{\pi^*}[A^\pi(s, a)] &= \mathbb{E}_{\pi^*}[Q^\pi(s, a) - V^\pi(s)] = \langle Q^\pi(s, \cdot), \pi^*(\cdot | s) \rangle - \langle Q^\pi(s, \cdot), \pi(\cdot | s) \rangle \\ &= \langle Q^\pi(s, \cdot), \pi^*(\cdot | s) - \pi(\cdot | s) \rangle. \end{aligned} \quad (\text{G.3})$$

Plugging (G.3) into (G.2) and recalling the definition of $\mathcal{L}(\pi)$ in (4.6), we finish the proof of Lemma 5.1. \square

Proof of Lemma 5.2. First, we have

$$\begin{aligned} &\text{KL}(\pi^*(\cdot | s) \| \pi_{\theta_k}(\cdot | s)) - \text{KL}(\pi^*(\cdot | s) \| \pi_{\theta_{k+1}}(\cdot | s)) \\ &= \langle \log(\pi_{\theta_{k+1}}(\cdot | s) / \pi_{\theta_k}(\cdot | s)), \pi^*(\cdot | s) \rangle \\ &= \langle \log(\pi_{\theta_{k+1}}(\cdot | s) / \pi_{\theta_k}(\cdot | s)), \pi^*(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle + \text{KL}(\pi_{\theta_{k+1}}(\cdot | s) \| \pi_{\theta_k}(\cdot | s)) \\ &= \langle \log(\pi_{\theta_{k+1}}(\cdot | s) / \pi_{\theta_k}(\cdot | s)) - \beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\ &\quad + \beta_k^{-1} \cdot \langle Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle + \text{KL}(\pi_{\theta_{k+1}}(\cdot | s) \| \pi_{\theta_k}(\cdot | s)) \\ &\quad + \langle \log(\pi_{\theta_{k+1}}(\cdot | s) / \pi_{\theta_k}(\cdot | s)), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle. \end{aligned} \quad (\text{G.4})$$

Recall that $\pi_{k+1} \propto \exp\{\tau_k^{-1} f_{\theta_k} + \beta_k^{-1} Q^{\pi_{\theta_k}}\}$ and $Z_{k+1}(s)$ and $Z_{\theta_k}(s)$ are defined in (E.1). Also recall that we have $\langle \log Z_{\theta_k}(s), \pi(\cdot | s) - \pi'(\cdot | s) \rangle = \langle \log Z_k(s), \pi(\cdot | s) - \pi'(\cdot | s) \rangle = 0$ for all k ,

π , and π' , which implies that, on the right-hand-side of (G.4),

$$\begin{aligned}
& \langle \log \pi_{\theta_k}(\cdot | s) + \beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\
&= \langle \tau_k^{-1} f_{\theta_k}(s, \cdot) + \beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle - \langle \log Z_{\theta_k}(s), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\
&= \langle \tau_k^{-1} f_{\theta_k}(s, \cdot) + \beta_k^{-1} Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle - \langle \log Z_{k+1}(s), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\
&= \langle \log \pi_{k+1}(\cdot | s), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle, \tag{G.5}
\end{aligned}$$

and

$$\begin{aligned}
& \langle \log(\pi_{\theta_{k+1}}(\cdot | s) / \pi_{\theta_k}(\cdot | s)), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle \\
&= \langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - \tau_k^{-1} f_{\theta_k}(s, \cdot), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle \\
&\quad - \langle \log Z_{\theta_{k+1}}(s), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle + \langle \log Z_{\theta_k}(s), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle \\
&= \langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - \tau_k^{-1} f_{\theta_k}(s, \cdot), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle. \tag{G.6}
\end{aligned}$$

Plugging (G.5) and (G.6) into (G.4), we obtain

$$\begin{aligned}
& \text{KL}(\pi^*(\cdot | s) \| \pi_{\theta_k}(\cdot | s)) - \text{KL}(\pi^*(\cdot | s) \| \pi_{\theta_{k+1}}(\cdot | s)) \tag{G.7} \\
&= \langle \log(\pi_{\theta_{k+1}}(\cdot | s) / \pi_{\theta_k}(\cdot | s)), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle + \beta_k^{-1} \cdot \langle Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\
&\quad + \langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - \tau_k^{-1} f_{\theta_k}(s, \cdot), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle + \text{KL}(\pi_{\theta_{k+1}}(\cdot | s) \| \pi_{\theta_k}(\cdot | s)) \\
&\geq \langle \log(\pi_{\theta_{k+1}}(\cdot | s) / \pi_{\theta_k}(\cdot | s)), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle + \beta_k^{-1} \cdot \langle Q^{\pi_{\theta_k}}(s, \cdot), \pi^*(\cdot | s) - \pi_{\theta_k}(\cdot | s) \rangle \\
&\quad + \langle \tau_{k+1}^{-1} f_{\theta_{k+1}}(s, \cdot) - \tau_k^{-1} f_{\theta_k}(s, \cdot), \pi_{\theta_k}(\cdot | s) - \pi_{\theta_{k+1}}(\cdot | s) \rangle + 1/2 \cdot \|\pi_{\theta_{k+1}}(\cdot | s) - \pi_{\theta_k}(\cdot | s)\|_1^2,
\end{aligned}$$

where in the last inequality we use the Pinsker's inequality. Rearranging the terms in (G.7), we finish the proof of Lemma 5.2. \square