

A. Proof of Supporting Lemmas

In this section, we present the omitted proof of the technical lemmas used in our main theorems.

A.1. Proof of Lemma 6.1

Before we prove the error bound for the local linearization, we first present some useful lemmas from recent studies of overparameterized deep neural networks. Note that in the following lemmas, $\{C_i\}_{i=1,\dots}$ are universal constants that are independent of problem parameters such as d, θ, m, L and their values can be different in different contexts. The first lemma states the uniform upper bound for the gradient of the deep neural network. Note that by definition, our parameter θ is a long vector containing the concatenation of the vectorization of all the weight matrices. Correspondingly, the gradient $\nabla_{\theta} f(\theta; \mathbf{x})$ is also a long vector.

Lemma A.1 (Lemma B.3 in Cao & Gu (2019b)). Let $\theta \in \mathbb{B}(\theta_0, \omega)$ with the radius satisfying $C_1 d^{3/2} L^{-1} m^{-3/2} \leq \omega \leq C_2 L^{-6} (\log m)^{-3/2}$. Then for all unit vectors in \mathbb{R}^d , i.e., $\mathbf{x} \in S^{d-1}$, the gradient of the neural network f defined in (4.2) is bounded as $\|\nabla_{\theta} f(\theta; \mathbf{x})\|_2 \leq C_3 \sqrt{m}$ with probability at least $1 - L^2 \exp(-C_4 m \omega^{2/3} L)$.

The second lemma provides the perturbation bound for the gradient of the neural network function. Note that the original theorem holds for any fixed d dimensional unit vector \mathbf{x} . However, due to the choice of ω and its dependency on m and d , it is easy to modify the results to hold for all $\mathbf{x} \in S^{d-1}$.

Lemma A.2 (Theorem 5 in Allen-Zhu et al. (2019b)). Let $\theta \in \mathbb{B}(\theta_0, \omega)$ with the radius satisfying

$$C_1 d^{3/2} L^{-3/2} m^{-3/2} (\log m)^{-3/2} \leq \omega \leq C_2 L^{-9/2} (\log m)^{-3}.$$

Then for all $\mathbf{x} \in S^{d-1}$, with probability at least $1 - \exp(-C_3 m \omega^{2/3} L)$ over the randomness of θ_0 , it holds that

$$\|\nabla_{\theta} f(\theta; \mathbf{x}) - \nabla_{\theta} f(\theta_0; \mathbf{x})\|_2 \leq C_4 \omega^{1/3} L^3 \sqrt{\log m} \|\nabla_{\theta} f(\theta_0; \mathbf{x})\|_2.$$

Now we are ready to bound the linearization error.

Proof of Lemma 6.1. Recall the definition of $\mathbf{g}_t(\theta_t)$ and $\mathbf{m}_t(\theta_t)$ in (4.5) and (6.2) respectively. We have

$$\begin{aligned} \|\mathbf{g}_t(\theta_t) - \mathbf{m}_t(\theta_t)\|_2 &= \|\nabla_{\theta} f(\theta_t; s_t, a_t) \Delta(s_t, a_t, s_{t+1}; \theta_t) - \nabla_{\theta} \hat{f}(\theta_t; s_t, a_t) \hat{\Delta}(s_t, a_t, s_{t+1}; \theta_t)\|_2 \\ &\leq \|(\nabla_{\theta} f(\theta_t; s_t, a_t) - \nabla_{\theta} \hat{f}(\theta_t; s_t, a_t)) \Delta(s_t, a_t, s_{t+1}; \theta_t)\|_2 \\ &\quad + \|\nabla_{\theta} \hat{f}(\theta_t; s_t, a_t) (\Delta(s_t, a_t, s_{t+1}; \theta_t) - \hat{\Delta}(s_t, a_t, s_{t+1}; \theta_t))\|_2. \end{aligned} \quad (\text{A.1})$$

Since $\hat{f}(\theta) \in \mathcal{F}_{\Theta, m}$, we have $\hat{f}(\theta) = f(\theta_0) + \langle \nabla_{\theta} f(\theta_0), \theta - \theta_0 \rangle$ and $\nabla_{\theta} \hat{f}(\theta) = \nabla_{\theta} f(\theta_0)$. Then with probability at least $1 - 2L^2 \exp(-C_1 m \omega^{2/3} L)$, we have

$$\begin{aligned} &\|(\nabla_{\theta} f(\theta_t; s_t, a_t) - \nabla_{\theta} \hat{f}(\theta_t; s_t, a_t)) \Delta(s_t, a_t, s_{t+1}; \theta_t)\|_2 \\ &= |\Delta(s_t, a_t, s_{t+1}; \theta_t)| \cdot \|(\nabla_{\theta} f(\theta_t; s_t, a_t) - \nabla_{\theta} f(\theta_0; s_t, a_t))\|_2 \\ &\leq C_2 \omega^{1/3} L^3 \sqrt{m \log m} |\Delta(s_t, a_t, s_{t+1}; \theta_t)|, \end{aligned}$$

where the inequality comes from Lemmas A.1 and A.2. By Lemma 6.4, with probability at least $1 - \delta$, it holds that

$$|\Delta(s_t, a_t, s_{t+1}; \theta_t)| = \left| f(\theta_t; s_t, a_t) - r_t - \gamma \max_{b \in \mathcal{A}} f(\theta_t; s_{t+1}, b) \right| \leq (2 + \gamma) C_3 \sqrt{\log(T/\delta)},$$

which further implies that with probability at least $1 - \delta - 2L^2 \exp(-C_1 m \omega^{2/3} L)$, we have

$$\begin{aligned} &\|(\nabla_{\theta} f(\theta_t; s_t, a_t) - \nabla_{\theta} \hat{f}(\theta_t; s_t, a_t)) \Delta(s_t, a_t, s_{t+1}; \theta_t)\|_2 \\ &\leq (2 + \gamma) C_2 C_3 \omega^{1/3} L^3 \sqrt{m \log m \log(T/\delta)}. \end{aligned}$$

For the second term in (A.1), we have

$$\|\nabla_{\theta} \hat{f}(\theta_t; s_t, a_t) (\Delta(s_t, a_t, s_{t+1}; \theta_t) - \hat{\Delta}(s_t, a_t, s_{t+1}; \theta_t))\|_2$$

$$\begin{aligned}
 &\leq \left\| \nabla_{\theta} \widehat{f}(\theta_t; s_t, a_t) (f(\theta_t; s_t, a_t) - \widehat{f}(\theta_t; s_t, a_t)) \right\|_2 \\
 &\quad + \left\| \nabla_{\theta} \widehat{f}(\theta_t; s_t, a_t) \left(\max_{b \in \mathcal{A}} f(\theta_t; s_{t+1}, b) - \max_{b \in \mathcal{A}} \widehat{f}(\theta_t; s_{t+1}, b) \right) \right\|_2 \\
 &\leq \left\| \nabla_{\theta} \widehat{f}(\theta_t; s_t, a_t) \right\|_2 \cdot |f(\theta_t; s_t, a_t) - \widehat{f}(\theta_t; s_t, a_t)| \\
 &\quad + \left\| \nabla_{\theta} \widehat{f}(\theta_t; s_t, a_t) \right\|_2 \max_{b \in \mathcal{A}} |f(\theta_t; s_{t+1}, b) - \widehat{f}(\theta_t; s_{t+1}, b)|. \tag{A.2}
 \end{aligned}$$

By Lemma 6.4, with probability at least $1 - \delta$ we have

$$|f(\theta_t; s_t, a_t) - \widehat{f}(\theta_t; s_t, a_t)| \leq \omega^{4/3} L^{11/3} \sqrt{m \log m} + C_4 \omega^2 L^4 \sqrt{m},$$

for all $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$ such that $\|\phi(s_t, a_t)\|_2 = 1$. Substituting the above result into (A.2) and applying the gradient bound in Lemma A.1, we obtain with probability at least $1 - \delta - L^2 \exp(-C_1 m \omega^{2/3} L)$ that

$$\begin{aligned}
 &\left\| \nabla_{\theta} \widehat{f}(\theta_t; s_t, a_t) (\Delta(s_t, a_t, s_{t+1}; \theta_t) - \widehat{\Delta}(s_t, a_t, s_{t+1}; \theta_t)) \right\|_2 \\
 &\leq C_5 \omega^{4/3} L^{11/3} m \sqrt{\log m} + C_6 \omega^2 L^4 m.
 \end{aligned}$$

Note that the above results require that the choice of ω should satisfy all the constraints in Lemmas A.1, 6.4 and A.2, of which the intersection is

$$C_7 d^{3/2} L^{-1} m^{-3/4} \leq \omega \leq C_8 L^{-6} (\log m)^{-3}.$$

Therefore, the error of the local linearization of $\mathbf{g}_t(\theta_t)$ can be upper bounded by

$$\begin{aligned}
 |\langle \mathbf{g}_t(\theta_t) - \mathbf{m}_t(\theta_t), \theta_t - \theta^* \rangle| &\leq (2 + \gamma) C_2 C_3 \omega^{1/3} L^3 \sqrt{m \log m \log(T/\delta)} \|\theta_t - \theta^*\|_2 \\
 &\quad + (C_5 \omega^{4/3} L^{11/3} m \sqrt{\log m} + C_6 \omega^2 L^4 m) \|\theta_t - \theta^*\|_2,
 \end{aligned}$$

which holds with probability at least $1 - 2\delta - 3L^2 \exp(-C_1 m \omega^{2/3} L)$ over the randomness of the initial point. For the upper bound of the norm of \mathbf{g}_t , by Lemmas 6.4 and A.1, we have

$$\begin{aligned}
 \|\mathbf{g}_t(\theta_t)\|_2 &\leq \left\| \nabla_{\theta} f(\theta_t; s_t, a_t) \left(f(\theta_t; s_t, a_t) - r_t - \gamma \max_{b \in \mathcal{A}} f(\theta_t; s_{t+1}, b) \right) \right\|_2 \\
 &\leq (2 + \gamma) C_9 \sqrt{m \log(T/\delta)}
 \end{aligned}$$

holds with probability at least $1 - \delta - L^2 \exp(-C_1 m \omega^{2/3} L)$. \square

A.2. Proof of Lemma 6.2

Let us define $\zeta_t(\theta) = \langle \mathbf{m}_t(\theta) - \overline{\mathbf{m}}(\theta), \theta - \theta^* \rangle$, which characterizes the bias of the data. Different from the similar quantity ζ_t in Bhandari et al. (2018), our definition is based on the local linearization of f , which is essential to the analysis in our proof. It is easy to verify that $\mathbb{E}[\mathbf{m}_t(\theta)] = \overline{\mathbf{m}}(\theta)$ for any fixed and deterministic θ . However, it should be noted that $\mathbb{E}[\mathbf{m}_t(\theta_t) | \theta_t = \theta] \neq \overline{\mathbf{m}}(\theta)$ because θ_t depends on all historical states and actions $\{s_t, a_t, s_{t-1}, a_{t-1}, \dots\}$ and $\mathbf{m}_t(\cdot)$ depends on the current observation $\{s_t, a_t, s_{t+1}\}$ and thus also depends on $\{s_{t-1}, a_{t-1}, s_{t-2}, a_{t-2}, \dots\}$. Therefore, we need a careful analysis of Markov chains to decouple the dependency between θ_t and $\mathbf{m}_t(\cdot)$.

The following lemma uses data processing inequality to provide an information theoretic control of coupling.

Lemma A.3 (Control of coupling (Bhandari et al., 2018)). Consider two random variables X and Y that form the following Markov chain:

$$X \rightarrow s_t \rightarrow s_{t+\tau} \rightarrow Y,$$

where $t \in \{0, 1, 2, \dots\}$ and $\tau > 0$. Suppose Assumption 5.2 holds. Let X' and Y' be independent copies drawn from the marginal distributions of X and Y respectively, i.e., $\mathbb{P}(X' = \cdot, Y' = \cdot) = \mathbb{P}(X = \cdot) \otimes \mathbb{P}(Y = \cdot)$. Then for any bounded function $h : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$, it holds that

$$|\mathbb{E}[h(X, Y)] - \mathbb{E}[h(X', Y')]| \leq 2 \sup_{s, s'} |h(s, s')| \lambda \rho^\tau.$$

Proof of Lemma 6.2. The proof of this lemma is adapted from (Bhandari et al., 2018), where the result was originally proved for linear function approximation of temporal difference learning. We first show that $\zeta_t(\boldsymbol{\theta})$ is Lipschitz. For any $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{B}(\boldsymbol{\theta}_0, \omega)$, we have

$$\begin{aligned}\zeta_t(\boldsymbol{\theta}) - \zeta_t(\boldsymbol{\theta}') &= \langle \mathbf{m}_t(\boldsymbol{\theta}) - \bar{\mathbf{m}}(\boldsymbol{\theta}), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle - \langle \mathbf{m}_t(\boldsymbol{\theta}') - \bar{\mathbf{m}}(\boldsymbol{\theta}'), \boldsymbol{\theta}' - \boldsymbol{\theta}^* \rangle \\ &= \langle \mathbf{m}_t(\boldsymbol{\theta}) - \bar{\mathbf{m}}(\boldsymbol{\theta}) - (\mathbf{m}_t(\boldsymbol{\theta}') - \bar{\mathbf{m}}(\boldsymbol{\theta}')), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle \\ &\quad + \langle \mathbf{m}_t(\boldsymbol{\theta}') - \bar{\mathbf{m}}(\boldsymbol{\theta}'), \boldsymbol{\theta} - \boldsymbol{\theta}' \rangle,\end{aligned}$$

which directly implies

$$\begin{aligned}|\zeta_t(\boldsymbol{\theta}) - \zeta_t(\boldsymbol{\theta}')| &\leq \|\mathbf{m}_t(\boldsymbol{\theta}) - \mathbf{m}_t(\boldsymbol{\theta}')\|_2 \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 + \|\bar{\mathbf{m}}(\boldsymbol{\theta}) - \bar{\mathbf{m}}(\boldsymbol{\theta}')\|_2 \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \\ &\quad + \|\mathbf{m}_t(\boldsymbol{\theta}') - \bar{\mathbf{m}}(\boldsymbol{\theta}')\|_2 \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2.\end{aligned}$$

By the definition of \mathbf{m}_t , we have

$$\begin{aligned}\|\mathbf{m}_t(\boldsymbol{\theta}) - \mathbf{m}_t(\boldsymbol{\theta}')\|_2 &= \left\| \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0) \left((f(\boldsymbol{\theta}; s, a) - f(\boldsymbol{\theta}'; s, a)) - \gamma \left(\max_{b \in \mathcal{A}} f(\boldsymbol{\theta}; s', b) - \max_{b \in \mathcal{A}} f(\boldsymbol{\theta}'; s', b) \right) \right) \right\|_2 \\ &\leq (1 + \gamma) C_3^2 m \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2,\end{aligned}$$

which holds with probability at least $1 - L^2 \exp(-C_4 m \omega^{2/3} L)$, where we used the fact that the neural network function is Lipschitz with parameter $C_3 \sqrt{m}$ by Lemma A.1. Similar bound can also be established for $\|\bar{\mathbf{m}}_t(\boldsymbol{\theta}) - \bar{\mathbf{m}}_t(\boldsymbol{\theta}')\|$ in the same way. Note that for $\boldsymbol{\theta} \in \mathbb{B}(\boldsymbol{\theta}_0, \omega)$ with ω and m satisfying the conditions in Lemma 6.1, we have by the definition in (6.2) that

$$\begin{aligned}\|\mathbf{m}_t(\boldsymbol{\theta})\|_2 &\leq \left(|\widehat{f}(\boldsymbol{\theta}; s, a)| + r(s, a) + \gamma \max_b \widehat{f}(\boldsymbol{\theta}; s', b) \right) \|\nabla_{\boldsymbol{\theta}} \widehat{f}(\boldsymbol{\theta})\|_2 \\ &\leq (2 + \gamma) (|f(\boldsymbol{\theta}_0)| + \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)\|_2 \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2) \|\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)\|_2 \\ &\leq (2 + \gamma) C_3 (C_8 \sqrt{m} \sqrt{\log(T/\delta)} + C_3 m \omega).\end{aligned}\tag{A.3}$$

The same bound can be established for $\|\bar{\mathbf{m}}_t\|$ in a similar way. Therefore, we have $|\zeta_t(\boldsymbol{\theta}) - \zeta_t(\boldsymbol{\theta}')| \leq \ell_{m,L} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2$, where $\ell_{m,L}$ is defined as

$$\ell_{m,L} = 2(1 + \gamma) C_3^2 m \omega + (2 + \gamma) C_3 (C_8 \sqrt{m} \sqrt{\log(T/\delta)} + C_3 m \omega).$$

Applying the above inequality recursively, for all $\tau = 0, \dots, t$, we have

$$\begin{aligned}\zeta_t(\boldsymbol{\theta}_t) &\leq \zeta_t(\boldsymbol{\theta}_{t-\tau}) + \ell_{m,L} \sum_{i=t-\tau}^{t-1} \|\boldsymbol{\theta}_{i+1} - \boldsymbol{\theta}_i\|_2 \\ &\leq \zeta_t(\boldsymbol{\theta}_{t-\tau}) + (2 + \gamma) C_3 (C_8 \sqrt{m} \sqrt{\log(T/\delta)} + C_3 m \omega) \ell_{m,L} \sum_{i=t-\tau}^{t-1} \eta_i.\end{aligned}\tag{A.4}$$

Next, we need to bound $\zeta_t(\boldsymbol{\theta}_{t-\tau})$. Define the observed tuple $O_t = (s_t, a_t, s_{t+1})$ as the collection of the current state and action and the next state. Note that $\boldsymbol{\theta}_{t-\tau} \rightarrow s_{t-\tau} \rightarrow s_t \rightarrow O_t$ forms a Markov chain induced by the target policy π . Recall that $\mathbf{m}_t(\cdot)$ depends on the observation O_t . Let's rewrite $\mathbf{m}(\boldsymbol{\theta}, O_t) = \mathbf{m}_t(\boldsymbol{\theta})$. Similarly, we can rewrite $\zeta_t(\boldsymbol{\theta})$ as $\zeta(\boldsymbol{\theta}, O_t)$. Let $\boldsymbol{\theta}'_{t-\tau}$ and O'_t be independently drawn from the marginal distributions of $\boldsymbol{\theta}_{t-\tau}$ and O_t respectively. Applying Lemma A.3 yields

$$\mathbb{E}[\zeta(\boldsymbol{\theta}_{t-\tau}, O_t)] - \mathbb{E}[\zeta(\boldsymbol{\theta}'_{t-\tau}, O'_t)] \leq 2 \sup_{\boldsymbol{\theta}, O} |\zeta(\boldsymbol{\theta}, O)| \lambda \rho^\tau,$$

where we used the uniform mixing result in Assumption 5.2. By definition $\boldsymbol{\theta}'_{t-\tau}$ and O'_t are independent, which implies $\mathbb{E}[\mathbf{m}(\boldsymbol{\theta}'_{t-\tau}, O'_t) | \boldsymbol{\theta}'_{t-\tau}] = \bar{\mathbf{m}}(\boldsymbol{\theta}'_{t-\tau})$ and

$$\mathbb{E}[\zeta(\boldsymbol{\theta}'_{t-\tau}, O'_t)] = \mathbb{E}[\mathbb{E}[\langle \mathbf{m}(\boldsymbol{\theta}'_{t-\tau}, O'_t) - \bar{\mathbf{m}}(\boldsymbol{\theta}'_{t-\tau}), \boldsymbol{\theta}'_{t-\tau} - \boldsymbol{\theta}^* \rangle | \boldsymbol{\theta}'_{t-\tau}]] = 0.$$

Moreover, by the definition of $\zeta(\cdot, \cdot)$, we have

$$|\zeta(\boldsymbol{\theta}, O) \leq \|\mathbf{m}_t(\boldsymbol{\theta}) - \bar{\mathbf{m}}(\boldsymbol{\theta})\|_2 \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \leq 2(2 + \gamma)C_3(C_8\sqrt{m}\sqrt{\log(T/\delta)} + C_3m\omega)\omega,$$

where the second inequality is due to (A.3) and that $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2 \leq \omega$. Therefore, for any $\tau = 0, \dots, t$, we have

$$\begin{aligned} \mathbb{E}[\zeta_t(\boldsymbol{\theta}_t)] &\leq \mathbb{E}\zeta_t(\boldsymbol{\theta}_{t-\tau}) + (2 + \gamma)C_3(C_8\sqrt{m}\sqrt{\log(T/\delta)} + C_3m\omega)\ell_{m,L} \sum_{i=t-\tau}^{t-1} \eta_i \\ &\leq (2 + \gamma)C_3(C_8\sqrt{m}\sqrt{\log(T/\delta)} + C_3m\omega)(\omega\lambda\rho^\tau + \ell_{m,L}\tau\eta_{t-\tau}). \end{aligned} \quad (\text{A.5})$$

Define τ^* as the mixing time of the Markov chain that satisfies

$$\tau^* = \min\{t = 0, 1, 2, \dots \mid \lambda\rho^t \leq \eta_T\}.$$

When $t \leq \tau^*$, we choose $\tau = t$ in (A.5). Note that η_t is nondecreasing. We obtain

$$\begin{aligned} \mathbb{E}[\zeta_t(\boldsymbol{\theta}_t)] &\leq \mathbb{E}[\zeta_t(\boldsymbol{\theta}_0)] + 2(2 + \gamma)C_3(C_8\sqrt{m}\sqrt{\log(T/\delta)} + C_3m\omega)\ell_{m,L}\tau^*\eta_0 \\ &= 2(2 + \gamma)C_3(C_8\sqrt{m}\sqrt{\log(T/\delta)} + C_3m\omega)\ell_{m,L}\tau^*\eta_0, \end{aligned}$$

where we used the fact that the initial point $\boldsymbol{\theta}_0$ is independent of $\{s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0\}$ and thus independent of $\zeta_t(\cdot)$. When $t > \tau^*$, we can choose $\tau = \tau^*$ in (A.5) and obtain

$$\begin{aligned} \mathbb{E}[\zeta_t(\boldsymbol{\theta}_t)] &\leq (2 + \gamma)C_3(C_8\sqrt{m}\sqrt{\log(T/\delta)} + C_3m\omega)(\omega\eta_T + \ell_{m,L}\tau^*\eta_{t-\tau^*}) \\ &\leq \tilde{C}(m \log(T/\delta) + m^2\omega^2)\tau^*\eta_{t-\tau^*}, \end{aligned}$$

where $\tilde{C} > 0$ is a universal constant, which completes the proof. \square

A.3. Proof of Lemma 6.3

Proof of Lemma 6.3. To simplify the notation, we use \mathbb{E}_π to denote $\mathbb{E}_{\mu, \pi, \mathcal{P}}$, namely, the expectation over $s \in \mu, a \sim \pi(\cdot|s)$ and $s' \sim \mathcal{P}(\cdot|s, a)$, in the rest of the proof. By the definition of $\bar{\mathbf{m}}$ in (6.2), we have

$$\begin{aligned} &\langle \bar{\mathbf{m}}(\boldsymbol{\theta}) - \bar{\mathbf{m}}(\boldsymbol{\theta}^*), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle \\ &= \mathbb{E}_{\mu, \pi, \mathcal{P}} [(\hat{\Delta}(s, a, s'; \boldsymbol{\theta}) - \hat{\Delta}(s, a, s'; \boldsymbol{\theta}^*)) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle] \\ &= \mathbb{E}_{\mu, \pi, \mathcal{P}} [(\hat{f}(\boldsymbol{\theta}; s, a) - \hat{f}(\boldsymbol{\theta}^*; s, a)) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle] \\ &\quad - \gamma \mathbb{E}_{\mu, \pi, \mathcal{P}} \left[\left(\max_{b \in \mathcal{A}} \hat{f}(\boldsymbol{\theta}; s', b) - \max_{b \in \mathcal{A}} \hat{f}(\boldsymbol{\theta}^*; s', b) \right) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle \right], \end{aligned}$$

where in the first equation we used the fact that $\nabla_{\boldsymbol{\theta}} \hat{f}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)$ for all $\boldsymbol{\theta} \in \Theta$ and $\hat{f} \in \mathcal{F}_{\Theta, m}$. Further by the property of the local linearization of f at $\boldsymbol{\theta}_0$, we have

$$\hat{f}(\boldsymbol{\theta}; s, a) - \hat{f}(\boldsymbol{\theta}^*; s, a) = \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle, \quad (\text{A.6})$$

which further implies

$$\begin{aligned} &\mathbb{E}_{\mu, \pi, \mathcal{P}} [(\hat{f}(\boldsymbol{\theta}; s, a) - \hat{f}(\boldsymbol{\theta}^*; s, a)) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle | \boldsymbol{\theta}_0] \\ &= (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \mathbb{E} [\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a) \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a)^\top | \boldsymbol{\theta}_0] (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \\ &= m \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_\pi}^2. \end{aligned}$$

where Σ_π is defined in Assumption 5.3. For the other term, we define $b_{\max}^\boldsymbol{\theta} = \operatorname{argmax}_{b \in \mathcal{A}} \hat{f}(\boldsymbol{\theta}; s', b)$ and $b_{\max}^{\boldsymbol{\theta}^*} = \operatorname{argmax}_{b \in \mathcal{A}} \hat{f}(\boldsymbol{\theta}^*; s', b)$. Then we have

$$\mathbb{E}_{\mu, \pi, \mathcal{P}} \left[\left(\max_{b \in \mathcal{A}} \hat{f}(\boldsymbol{\theta}; s', b) - \max_{b \in \mathcal{A}} \hat{f}(\boldsymbol{\theta}^*; s', b) \right) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle \right]$$

$$= \mathbb{E}_{\mu, \pi, \mathcal{P}} [(\widehat{f}(\boldsymbol{\theta}; s', b_{\max}^{\boldsymbol{\theta}}) - \widehat{f}(\boldsymbol{\theta}^*; s', b_{\max}^{\boldsymbol{\theta}^*})) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle]. \quad (\text{A.7})$$

Applying Cauchy-Schwarz inequality, we have

$$\begin{aligned} & \mathbb{E}_{\mu, \pi, \mathcal{P}} [(\widehat{f}(\boldsymbol{\theta}; s', b_{\max}^{\boldsymbol{\theta}}) - \widehat{f}(\boldsymbol{\theta}^*; s', b_{\max}^{\boldsymbol{\theta}^*})) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle] \\ & \leq \sqrt{\mathbb{E}_{\mu, \pi, \mathcal{P}} [(\max_b |(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s', b)|)^2]} \sqrt{\mathbb{E}_{\mu, \pi, \mathcal{P}} [(\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*))^2]} \\ & = m \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}^*(\boldsymbol{\theta} - \boldsymbol{\theta}^*)} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}}, \end{aligned}$$

where we used the fact that $\Sigma_{\pi}^*(\boldsymbol{\theta} - \boldsymbol{\theta}^*) = 1/m \mathbb{E}_{\mu, \pi, \mathcal{P}} [\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, \widetilde{b}_{\max}) \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, \widetilde{b}_{\max})^\top]$ and $\widetilde{b}_{\max} = \operatorname{argmax}_{b \in \mathcal{A}} |\langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, b), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle|$ according to (5.6). Substituting the above results into (A.7), we obtain

$$\begin{aligned} & \mathbb{E}_{\mu, \pi, \mathcal{P}} [(\max_{b \in \mathcal{A}} \widehat{f}(\boldsymbol{\theta}; s', b) - \max_{b \in \mathcal{A}} \widehat{f}(\boldsymbol{\theta}^*; s', b)) \langle \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0; s, a), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle] \\ & \leq m \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}^*(\boldsymbol{\theta} - \boldsymbol{\theta}^*)} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}}, \end{aligned}$$

which immediately implies

$$\begin{aligned} \langle \overline{\mathbf{m}}(\boldsymbol{\theta}) - \overline{\mathbf{m}}(\boldsymbol{\theta}^*), \boldsymbol{\theta} - \boldsymbol{\theta}^* \rangle & \geq m \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}} \cdot \left(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}} - \gamma \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}^*(\boldsymbol{\theta} - \boldsymbol{\theta}^*)} \right) \\ & = m \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}} \cdot \frac{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}}^2 - \gamma^2 \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}^*(\boldsymbol{\theta} - \boldsymbol{\theta}^*)}^2}{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}} + \gamma \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}^*(\boldsymbol{\theta} - \boldsymbol{\theta}^*)}} \\ & \geq m(1 - \alpha^{-1/2}) \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\Sigma_{\pi}}^2 \\ & = (1 - \alpha^{-1/2}) \mathbb{E}_{\mu, \pi, \mathcal{P}} [(\widehat{f}(\boldsymbol{\theta}) - \widehat{f}(\boldsymbol{\theta}^*))^2 | \boldsymbol{\theta}_0], \end{aligned}$$

where the second inequality is due to Assumption 5.3 and the last equation is due to (A.6) and the definition of Σ_{π} in (5.5). \square