

## Appendix of “Provable Knowledge Transfer using Successor Feature for Deep Reinforcement Learning.”

Before moving into the technical details, we provide an overview of the structure of the appendix.

In Appendix A, we define some notations and useful lemmas to simplify the presentation and analysis. Some important notations for understanding the proof is summarized in Table 3.

In Appendix B, we provide some preliminary lemmas and proof for Theorem 1. A proof sketch is included as (i) characterization of the local convex region of the objective function in (12) and (9) (Lemma 6), (ii) Characterization of the difference between the empirical gradient in (15) and the gradient of the objective function (Lemma 7), (iii) Characterization of the relation of two consecutive iterations  $\Theta^{(t+1)}$  and  $\Theta^{(t)}$  in (62), and (iv) Mathematical induction over  $(t+1) \cdot \|\Theta^{(t)} - \Theta^*\|_2$  from  $t = 1$  to  $T$  to obtain the error bound between the learned model weights  $\Theta^{(T)}$  and the optimal  $\Theta^*$ .

In Appendix C, we provide the proof for Theorems 3 and 4. A proof sketch is included as follows: (1) Characterization of (26) by assuming knowledge of the optimal Q-function for previous tasks. (2) Characterization of the accumulated error resulting from the estimation error of the learned Q-function in previous tasks. (3) Combining the bounds from (1) and (2) leads to the error bound between (26) derived from the estimated Q-function of previous tasks and the optimal Q-function for the new tasks.

In Appendix D, we provide the proof for Theorem 2. The proof sketch is a direct application of the existing results of the convergence analysis as shown in Appendix B and the error bound between (26) derived from the estimated Q-function of previous tasks and the optimal Q-function for the new tasks as shown in Appendix C.

In Appendix E, we provide additional experiments to further support the proposed SF-DQN in Algorithm 1 and our theoretical findings.

In Appendix F, we provide the proofs for the preliminary lemmas in proving Theorems 1 and 2.

In Appendix G, we provide the proofs for the preliminary lemmas in proving Theorems 3 and 4.

In Appendix H, we provide the proof for some additional lemmas.

### A NOTATIONS AND PRELIMINARY RESULTS

**Population risk function.** We define a population risk function as

$$f_{\pi^*}(\theta) := \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \pi^*} \left\| \psi(\theta; \mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{s}' | (\mathbf{s}, \mathbf{a}), \mathbf{a}' \sim \pi^*}(\mathbf{s}') \left( \phi(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \cdot \psi(\theta^*; \mathbf{s}', \mathbf{a}') \right) \right\|_2^2. \quad (30)$$

We can see that  $\theta^*$  is the global minimal to (30) with Assumption 1. For the convenience of presentation, we simplify  $f_{\pi^*}$  as  $f$  in the supplementary materials.

Then, the gradient of (30) is

$$\begin{aligned} & \nabla f_{\pi^*}(\theta) \\ &= \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \pi^*, \mathbf{s}' | (\mathbf{s}, \mathbf{a}) \sim \mathcal{P}, \mathbf{a}' \sim \pi^*} \left( \psi(\theta; \mathbf{s}, \mathbf{a}) - \phi(\mathbf{s}, \mathbf{a}, \mathbf{s}') - \gamma \cdot \psi(\theta^*; \mathbf{s}', \mathbf{a}') \right) \cdot \nabla \psi(\theta; \mathbf{s}, \mathbf{a}). \end{aligned} \quad (31)$$

Given  $f$  is a smooth function, we have the gradient of  $f$  with respect to any  $\theta_\ell$  at the ground truth  $\theta^*$  equals to zero, namely,

$$\nabla_\ell f(\theta^*) := \nabla_{\theta_\ell} f(\theta^*) = \mathbf{0}, \quad \forall \ell \in [L]. \quad (32)$$

**Vectorized Gradient of  $\theta$  and  $w$  at iteration  $t$ .** To avoid unnecessary high-dimensional tensor analysis, the gradient with respect to  $\theta$ , denoted as  $\nabla_\theta H$  for some function  $H$ , is represented as its corresponding vectorized version,  $\nabla_{\text{vec}(\theta)} H$ .

Let  $n$  denote the dimension of  $\mathbf{W}$  defined in (1). We denote  $n_\ell$  as the dimension of the vectorized neuron weights in the  $\ell$ -th layer, namely,  $n_\ell = \dim(\text{vec}(\theta_\ell))$ .

Then, the gradient in updating  $\theta$  as

$$g^{(t)}(\theta^{(t)}; \mathcal{D}_t) = \sum_{m \in \mathcal{D}_t} \left( \psi(\theta^{(t)}; \mathbf{s}_m, a_m) - \phi(\mathbf{s}_m, a_m, \mathbf{s}'_m) - \gamma \cdot \psi(\theta^{(t)}; \mathbf{s}'_m, a'_m) \right) \cdot \nabla_{\theta} \psi(\theta^{(t)}; \mathbf{s}_m, a_m) \quad (33)$$

with  $g^{(t)}(\theta^{(t)}; \mathcal{D}_t) \in \mathbb{R}^n$ . Then, we have

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \cdot g^{(t)}(\theta^{(t)}; \mathcal{D}_t). \quad (34)$$

Similarly to (33), we define the gradient

$$l^{(t)}(\mathbf{w}^{(t)}; \mathcal{D}_t) = \sum_{m \in \mathcal{D}_t} \left( \phi(\mathbf{s}_m, a_m, \mathbf{s}'_m)^{\top} \mathbf{w}^{(t)} - r(\mathbf{s}_m, a_m, \mathbf{s}'_m) \right) \cdot \phi(\mathbf{s}_m, a_m, \mathbf{s}'_m). \quad (35)$$

In addition, without special descriptions,  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1^{\top}, \boldsymbol{\alpha}_2^{\top}, \dots, \boldsymbol{\alpha}_K^{\top}]^{\top}$  stands for any unit vector that in  $\mathbb{R}^{K_{\ell} K_{\ell-1}}$  with  $\boldsymbol{\alpha}_j \in \mathbb{R}^{K_{\ell-1}}$  ( $K_0 = d$ ). Therefore, we have

$$\begin{aligned} \|\nabla_{\ell} H\|_2 &= \max_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}^{\top} \nabla_{\ell} H\|_2 = \max_{\boldsymbol{\alpha}} \left| \sum_{j=1}^K \boldsymbol{\alpha}_j^{\top} \frac{\partial H}{\partial \mathbf{w}_{\ell,j}} \right|, \\ \|\nabla_{\ell}^2 H\|_2 &= \max_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}^{\top} \nabla_{\ell}^2 H \boldsymbol{\alpha}\|_2 = \max_{\boldsymbol{\alpha}} \left( \sum_{j=1}^K \boldsymbol{\alpha}_j^{\top} \frac{\partial H}{\partial \mathbf{w}_{\ell,j}} \right)^2. \end{aligned} \quad (36)$$

**Derivation of the gradient of deep neural networks.** We use  $h^{(\ell)}(\theta)$  to denote the input in the  $\ell$ -th layer (or the output in the  $(\ell - 1)$ -th layer) of deep neural network  $\psi(\theta)$ , and  $h^{(1)} = \mathbf{x}(\mathbf{s}, a)$ , where

$$\mathbf{h}^{(\ell)}(\theta; \mathbf{s}, a) = \sigma(\theta_{\ell-1}^{\top} \mathbf{h}^{(\ell-1)}) = \dots = \sigma\left(\theta_{\ell}^{\top} \sigma(\theta_{\ell-1}^{\top} \dots \sigma(\theta_1^{\top} \mathbf{x}(\mathbf{s}, a)))\right). \quad (37)$$

Then, we denote the dimension of  $\mathbf{h}^{(\ell)}$  as  $K_{\ell}$ . Then,  $\psi(\theta; \mathbf{s}, a)$  can be written as

$$\psi(\theta; \mathbf{s}, a) = \frac{\mathbf{1}^{\top}}{K_L} \sum_{k=1}^{K_L} \sigma(\theta_{L,k}^{\top} \mathbf{h}^{(L)}) = \frac{\mathbf{1}^{\top}}{K_L} \sigma(\theta_L^{\top} \sigma(\theta_{L-1}^{\top} \mathbf{h}^{(L-1)})), \quad (38)$$

where  $\theta_{\ell,k}$  denotes the  $k$ -th neuron weights in the  $\ell$ -th layer. Then, we define a group of functions  $\mathcal{J}_{\ell}(\theta) \in \mathbb{R}^n \rightarrow \mathbb{R}^K$  such that

$$\begin{aligned} \mathcal{J}_{\ell}(\theta) &= \begin{cases} [\mathbf{1}^{\top} \sigma'(\theta_L^{\top} \mathbf{h}^{(L)}) \theta_L^{\top} \cdot \sigma'(\theta_{L-1}^{\top} \mathbf{h}^{(L-1)}) \theta_{L-1}^{\top} \dots \sigma'(\theta_{\ell+1}^{\top} \mathbf{h}^{(\ell+1)}) \theta_{\ell+1}^{\top}]^{\top} & \text{if } \ell > 1 \\ \mathbf{1} & \text{if } \ell = 1 \end{cases} \end{aligned} \quad (39)$$

Then, the gradient of  $\psi$  can be represented as

$$\frac{\partial \psi}{\partial \theta_{\ell,k}}(\theta) = \frac{1}{K_{\ell}} \mathcal{J}_{\ell,k}(\theta) \sigma'(\theta_{\ell,k}^{\top} \mathbf{h}^{(\ell)}(\theta)) \mathbf{h}^{(\ell)}(\theta), \quad (40)$$

where  $\mathcal{J}_{\ell,k}$  stands for the  $k$ -th component of  $\mathcal{J}_{\ell}$ .

**Order-wise Analysis.** Most constant numbers will be ignored in most steps. In particular, we use  $h_1(z) \gtrsim$  (or  $\lesssim, \approx$ )  $h_2(z)$  to denote there exists some positive constant  $C$  such that  $h_1(z) \geq$  (or  $\leq, =$ )  $C \cdot h_2(z)$  when  $z \in \mathbb{R}$  is sufficiently large. In this paper, we consider the case where  $\theta_{\ell}^*$  is well-conditioned, such that its largest singular value  $\Sigma_1(\ell)$  and the condition number  $\Sigma_1(\ell)/\sigma_K(\ell)$  can be viewed as constants and will be hidden in the order-wise analysis.

#### A.1 USEFUL LEMMAS FOR MATRIX CONCENTRATION

**Lemma 1** (Weyl's inequality, (Bhatia, 2013)). *Let  $\mathbf{B} = \mathbf{A} + \mathbf{E}$  be a matrix with dimension  $m \times m$ . Let  $\lambda_i(\mathbf{B})$  and  $\lambda_i(\mathbf{A})$  be the  $i$ -th largest eigenvalues of  $\mathbf{B}$  and  $\mathbf{A}$ , respectively. Then, we have*

$$|\lambda_i(\mathbf{B}) - \lambda_i(\mathbf{A})| \leq \|\mathbf{E}\|_2, \quad \forall \quad i \in [m]. \quad (41)$$

Table 3: Notations for the proofs

$d$	Dimension of the feature mappings of the state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ .
$K$	Number of neurons in the hidden layer.
$L$	Number of hidden layers.
$T$	Number of iterations.
$\mathbf{w}_i^{(t)}$	The estimated value for reward mapping of task $i$ at $t$ -th iteration.
$\Theta_i^{(t)}$	The estimated neuron weights for the successor feature of task $i$ at $t$ -th iteration.
$\theta^{(t)}$	The value of $\Theta_1^{(t)}$ to simplify the notation in the analyses without GPI.
$g^{(t)}(\theta^{(t)}; \mathcal{D}_t)$	The pseudo-gradient function defined in (33) at point $\theta^{(t)}$ with respect to the dataset $\mathcal{D}_t$ .
$f_{\pi^*}$ or $f$	The population risk function defined in (30).
$\nabla_\ell H(\hat{\theta})$	The gradient of a function $H$ with respect to the components of $\theta_\ell$ at point $\hat{\theta}$ .
$\nabla_\ell^2 H(\hat{\theta})$	The Hessian matrix of a function $H$ with respect to the components of $\theta_\ell$ at point $\hat{\theta}$ .
$Q_i^\pi$	The Q-function of task $i$ for policy $\pi$ .
$Q_i^*$	The Q-function of task $i$ for the optimal policy $\pi^*$ .
$q^*$	A constant defined in (80), depending on task relevance $\ \mathbf{w}_i - \mathbf{w}_j\ _2$ .
$\eta_t$	The step size for updating neuron weights $\Theta_i$ for the successor feature.
$\kappa_t$	The step size for updating the parameter for the weight mapping.
$c_N$	A constant in the order of $1/\sqrt{N}$ .
$n$	The dimension of $\theta$ .
$n_\ell$	The dimension of vectorized $\theta_\ell$ .
$K_\ell$	The dimension of the input for the $\ell$ -th layer for the deep neural network. $K_0 = d$ .
$\mathcal{J}_\ell(\mathbf{W})$	A function in $\mathbb{R}^n \rightarrow \mathbb{R}^K$ , defined in (39).
$C_t$	The distribution shift between the optimal policy and behavior policy at iteration $t$ , defined in Assumption (3).
$N$	The size of the experience replay buffer.
$\phi_{\max}$	The upper bound of the transition feature.
$\rho_1$	A constant defined in (84).
$\rho_2$	The smallest eigenvalue of $\mathbb{E}\phi(s, a)\phi(s, a)^\top \in \mathbb{R}^{d \times d}$ .
$\phi_{\max}$	The upper bound of the transition feature.

**Lemma 2** ((Tropp, 2012), Theorem 1.6). *Consider a finite sequence  $\{\mathbf{Z}_k\}$  of independent, random matrices with dimensions  $d_1 \times d_2$ . Assume that such random matrix satisfies*

$$\mathbb{E}(\mathbf{Z}_k) = 0 \quad \text{and} \quad \|\mathbf{Z}_k\| \leq R \quad \text{almost surely.}$$

*Define*

$$\delta^2 := \max \left\{ \left\| \sum_k \mathbb{E}(\mathbf{Z}_k \mathbf{Z}_k^*) \right\|, \left\| \sum_k \mathbb{E}(\mathbf{Z}_k^* \mathbf{Z}_k) \right\| \right\}.$$

Then for all  $t \geq 0$ , we have

$$\text{Prob}\left\{\left\|\sum_k \mathbf{Z}_k\right\| \geq t\right\} \leq (d_1 + d_2) \exp\left(\frac{-t^2/2}{\delta^2 + Rt/3}\right).$$

**Lemma 3** (Lemma 5.2, (Vershynin, 2010)). *Let  $\mathcal{B}(0, 1) \in \{\boldsymbol{\alpha} \mid \|\boldsymbol{\alpha}\|_2 = 1, \boldsymbol{\alpha} \in \mathbb{R}^d\}$  denote a unit ball in  $\mathbb{R}^d$ . Then, a subset  $\mathcal{S}_\xi$  is called a  $\xi$ -net of  $\mathcal{B}(0, 1)$  if every point  $\mathbf{z} \in \mathcal{B}(0, 1)$  can be approximated to within  $\xi$  by some point  $\boldsymbol{\alpha} \in \mathcal{S}_\xi$ , i.e.,  $\|\mathbf{z} - \boldsymbol{\alpha}\|_2 \leq \xi$ . Then the minimal cardinality of a  $\xi$ -net  $\mathcal{S}_\xi$  satisfies*

$$|\mathcal{S}_\xi| \leq (1 + 2/\xi)^d. \quad (42)$$

**Lemma 4** (Lemma 5.3, (Vershynin, 2010)). *Let  $\mathbf{A}$  be an  $d_1 \times d_2$  matrix, and let  $\mathcal{S}_\xi(d)$  be a  $\xi$ -net of  $\mathcal{B}(0, 1)$  in  $\mathbb{R}^d$  for some  $\xi \in (0, 1)$ . Then*

$$\|\mathbf{A}\|_2 \leq (1 - \xi)^{-1} \max_{\boldsymbol{\alpha}_1 \in \mathcal{S}_\xi(d_1), \boldsymbol{\alpha}_2 \in \mathcal{S}_\xi(d_2)} |\boldsymbol{\alpha}_1^T \mathbf{A} \boldsymbol{\alpha}_2|. \quad (43)$$

**Lemma 5** (Mean Value Theorem). *Let  $\mathbf{U} \subset \mathbb{R}^{n_1}$  be open and  $\mathbf{f} : \mathbf{U} \rightarrow \mathbb{R}^{n_2}$  be continuously differentiable, and  $\mathbf{x} \in \mathbf{U}$ ,  $\mathbf{h} \in \mathbb{R}^{n_1}$  vectors such that the line segment  $\mathbf{x} + t\mathbf{h}$ ,  $0 \leq t \leq 1$  remains in  $\mathbf{U}$ . Then we have:*

$$\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) = \left(\int_0^1 \nabla \mathbf{f}(\mathbf{x} + t\mathbf{h}) dt\right) \cdot \mathbf{h},$$

where  $\nabla \mathbf{f}$  denotes the Jacobian matrix of  $\mathbf{f}$ .

## A.2 DEFINITIONS OF SUB-GAUSSIAN AND SUB-EXPONENTIAL.

**Definition 1** (Definition 5.7, (Vershynin, 2010)). *A random variable  $X$  is called a sub-Gaussian random variable if it satisfies*

$$(\mathbb{E}|X|^p)^{1/p} \leq c_1 \sqrt{p} \quad (44)$$

for all  $p \geq 1$  and some constant  $c_1 > 0$ . In addition, we have

$$\mathbb{E}e^{s(X - \mathbb{E}X)} \leq e^{c_2 \|X\|_{\psi_2}^2 s^2} \quad (45)$$

for all  $s \in \mathbb{R}$  and some constant  $c_2 > 0$ , where  $\|X\|_{\psi_2}$  is the sub-Gaussian norm of  $X$  defined as  $\|X\|_{\psi_2} = \sup_{p \geq 1} p^{-1/2} (\mathbb{E}|X|^p)^{1/p}$ .

Moreover, a random vector  $\mathbf{X} \in \mathbb{R}^d$  belongs to the sub-Gaussian distribution if one-dimensional marginal  $\boldsymbol{\alpha}^T \mathbf{X}$  is sub-Gaussian for any  $\boldsymbol{\alpha} \in \mathbb{R}^d$ , and the sub-Gaussian norm of  $\mathbf{X}$  is defined as  $\|\mathbf{X}\|_{\psi_2} = \sup_{\|\boldsymbol{\alpha}\|_2=1} \|\boldsymbol{\alpha}^T \mathbf{X}\|_{\psi_2}$ .

**Definition 2** (Definition 5.13, (Vershynin, 2010)). *A random variable  $X$  is called a sub-exponential random variable if it satisfies*

$$(\mathbb{E}|X|^p)^{1/p} \leq c_3 p \quad (46)$$

for all  $p \geq 1$  and some constant  $c_3 > 0$ . In addition, we have

$$\mathbb{E}e^{s(X - \mathbb{E}X)} \leq e^{c_4 \|X\|_{\psi_1}^2 s^2} \quad (47)$$

for  $s \leq 1/\|X\|_{\psi_1}$  and some constant  $c_4 > 0$ , where  $\|X\|_{\psi_1}$  is the sub-exponential norm of  $X$  defined as  $\|X\|_{\psi_1} = \sup_{p \geq 1} p^{-1} (\mathbb{E}|X|^p)^{1/p}$ .

## B PROOF OF THEOREM 1

**Lemma 6** (Local convexity of  $f_{\pi^*}$ ). *Given any  $\theta \in \mathbb{R}^n$ , let  $\theta$  satisfy*

$$\|\theta - \theta^*\|_2 \lesssim \frac{c_N \cdot \sigma_K}{\rho_1 \cdot K} \quad (48)$$

for some constant  $c_N \in (0, 1)$ . Then, for the  $f_{\pi^*}$  defined in (30), we have

$$\frac{(1 - c_N)\rho_1}{K^2} \preceq \nabla_\ell^2 f_{\pi^*}(\theta) \preceq \frac{7}{K}. \quad (49)$$

**Lemma 7** (Upper bound of the error gradient). *Let  $f_{\pi^*}$  be the function defined in (30). Let  $g_t$  be the function defined in (33). Then, with probability at least  $1 - q^{-K_{\ell-1}}$ , we have*

$$\begin{aligned} \left\| \nabla_{\ell} f_{\pi^*}(\theta) - g_{\ell}(\theta^{(t)}; \mathcal{D}_t) \right\|_2 &\lesssim \frac{1}{K_{\ell}} \cdot \|\theta - \theta^*\|_2 \cdot \sqrt{\frac{K_{\ell-1} \log q}{|\mathcal{D}_t|}} + \frac{\gamma}{K_{\ell}} \cdot \|\theta^{(t)} - \theta^*\|_2 \\ &\quad + \frac{R_{\max}}{1-\gamma} \cdot (1+\gamma)\tau^* \cdot \eta_{t-\tau^*} \\ &\quad + |\mathcal{A}| \cdot \frac{R_{\max}}{1-\gamma} \cdot (1 + \log_{\nu} \lambda^{-1} + \frac{1}{1-\nu}) \cdot C_t, \end{aligned} \quad (50)$$

where  $\tau^* = \min\{t \mid \lambda \nu^t \leq \eta_T\}$ , and  $\nu$  &  $\lambda$  are defined in Assumption 2.

**Lemma 8** (Convergence of  $\mathbf{w}^{(t)}$ ). *With probability at least  $1 - q^{-d}$ ,  $\mathbf{w}$  enjoys a linear convergence rate to  $\mathbf{w}^*$  as*

$$\|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 \leq \left(1 - \frac{\rho - c_N}{\phi_{\max}}\right) \cdot \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2. \quad (51)$$

*Proof of Theorem 1.* From Algorithm 1, the update of  $\theta$  can be written as

$$\begin{aligned} \theta^{(t+1)} &= \theta^{(t)} - \eta_t \cdot g^{(t)}(\theta^{(t)}; \mathcal{D}_t) \\ &= \theta^{(t)} - \eta_t \cdot \nabla f(\theta^{(t)}) + \eta_t \cdot (\nabla f(\theta^{(t)}) - g^{(t)}(\theta^{(t)}; \mathcal{D}_t)). \end{aligned} \quad (52)$$

Since  $\nabla f$  is a smooth function and  $\theta^*$  is a local (global) optimal to  $f$ , then we have

$$\begin{aligned} \nabla f(\theta^{(t)}) &= \nabla f(\theta^{(t)}) - \nabla f(\theta^*) \\ &= \int_0^1 \nabla^2 f\left(\theta^{(t)} + u \cdot (\theta^{(t)} - \theta^*)\right) du \cdot (\theta^{(t)} - \theta^*), \end{aligned} \quad (53)$$

where the last equality comes from Mean Value Theory in Lemma 5. For notational convenience, we use  $\mathbf{A}^{(t)}$  to denote the integration as

$$\mathbf{A}^{(t)} := \int_0^1 \nabla^2 f\left(\theta^{(t)} + u \cdot (\theta^{(t)} - \theta^*)\right) du. \quad (54)$$

Then, we have

$$\begin{aligned} \|\theta^{(t+1)} - \theta^*\|_2 &\leq \|\mathbf{I} - \eta_t \mathbf{A}^{(t)}\|_2 \cdot \|\theta^{(t)} - \theta^*\|_2 + \eta_t \cdot \|\nabla f(\theta^{(t)}) - g^{(t)}(\theta^{(t)}; \mathcal{D}_t)\|_2 \\ &\leq \|\mathbf{I} - \eta_t \mathbf{A}^{(t)}\|_2 \cdot \|\theta^{(t)} - \theta^*\|_2 + \eta_t \cdot \sum_{\ell=1}^L \left\| \nabla_{\ell} f(\theta^{(t)}) - g_{\ell}^{(t)}(\theta_{\ell}^{(t)}; \mathcal{D}_t) \right\|_2. \end{aligned} \quad (55)$$

From Lemma 6, we have

$$\|\mathbf{I} - \eta_t \mathbf{A}^{(t)}\|_2 \leq 1 - \eta_t \cdot \frac{(1 - c_N) \cdot \rho_1}{K^2}. \quad (56)$$

From Lemma 7, we have

$$\begin{aligned} \left\| \nabla_{\ell} f_{\pi^*}(\theta^{(t)}) - g_{\ell}(\theta^{(t)}; \mathcal{D}_t) \right\|_2 &\lesssim \frac{1}{K_{\ell}} \cdot \|\theta^{(t)} - \theta^*\|_2 \cdot \sqrt{\frac{K_{\ell-1} \log q}{|\mathcal{D}_t|}} + \frac{\gamma}{K_{\ell}} \cdot \|\theta^{(t)} - \theta^*\|_2 \\ &\quad + \frac{R_{\max}}{1-\gamma} \cdot (1+\gamma)\tau^* \cdot \eta_{t-\tau^*} \\ &\quad + |\mathcal{A}| \cdot \frac{R_{\max}}{1-\gamma} \cdot (1 + \log_{\nu} \lambda^{-1} + \frac{1}{1-\nu}) \cdot C_t. \end{aligned} \quad (57)$$

With Assumption 3, we have

$$C_t \leq C \cdot (\|\theta^{(t)} - \theta^*\|_2 + \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2).$$

When we have a sufficiently large number of samples at iteration  $t$  as

$$|\mathcal{D}_t| \gtrsim c_N^{-2} \cdot \rho_1^{-1} \cdot \left( \sum_{\ell=1}^L K_\ell \sqrt{K_{\ell-1}} \right)^2 \cdot \log q, \quad (58)$$

(55) can be simplified as

$$\|\theta^{(t+1)} - \theta^*\|_2 \leq (1 - \eta_t \cdot \xi) \cdot \|\theta^{(t)} - \theta^*\|_2 + \eta_t \cdot \Delta_t + \eta_t \cdot C^* \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2. \quad (59)$$

where

$$\begin{aligned} C^* &= |\mathcal{A}| \cdot \frac{R_{\max}}{1 - \gamma} \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1 - \nu}) \cdot C \\ \xi &= \frac{(1 - \gamma - c_N) \rho_1}{K^2} - C^* \\ \Delta_t &= \frac{R_{\max}}{1 - \gamma} \cdot (1 + \gamma) \tau^* \cdot \eta_{t - \tau^*}. \end{aligned} \quad (60)$$

Let  $\eta_t = \frac{1}{\xi \cdot (t+1)}$ , we have

$$(t+1) \cdot \|\theta^{(t+1)} - \theta^*\|_2 \leq t \cdot \|\theta^{(t)} - \theta^*\|_2 + \xi^{-1} \cdot \Delta_t + \xi^{-1} \cdot C^* \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2. \quad (61)$$

Next, we have

$$\begin{aligned} &\sum_{t=0}^{T-1} (t+1) \cdot \|\theta^{(t+1)} - \theta^*\|_2 - t \cdot \|\theta^{(t)} - \theta^*\|_2 \\ &\leq \sum_{t=0}^{T-1} \xi^{-1} \cdot (\Delta_t + C^* \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2). \end{aligned} \quad (62)$$

With the definition of  $\Delta_t$  in (60), we have

$$\begin{aligned} \sum_{t=0}^{T-1} \Delta_t &\leq \sum_{t=0}^{\tau^*} \Delta_t + \sum_{t=\tau^*}^{T-1} \lambda^{-1} \cdot \Delta_t \\ &\leq \sum_{t=0}^{\tau^*} \tau^* \cdot \frac{R_{\max}}{1 - \gamma} + \sum_{t=\tau^*}^{T-1} \frac{R_{\max} \cdot (1 + \gamma)}{1 - \gamma} \cdot \tau^* \cdot \frac{1}{T - \tau^* + 1} \\ &\lesssim \frac{R_{\max} \cdot \log^2 T}{1 - \gamma} + \frac{R_{\max} \cdot (1 + \gamma) \cdot \log^2 T}{1 - \gamma}. \end{aligned} \quad (63)$$

With Lemma 8 that  $\mathbf{w}$  enjoys a geometric decay, we have

$$\sum_{t=0}^{T-1} \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2 \lesssim \|\mathbf{w}^{(0)} - \mathbf{w}^*\|_2. \quad (64)$$

By multiplying  $1/T$  on both sides of (62), we have

$$\|\theta^{(T)} - \theta^*\|_2 \leq \frac{(2 + \gamma) \cdot R_{\max} \cdot \log^2 T + C^* \|\mathbf{w}^{(0)} - \mathbf{w}^*\|_2}{(1 - \gamma - c_N) \rho_1 K^{-2} - C^*} \cdot \frac{1}{T}. \quad (65)$$

□

## C PROOFS OF THEOREMS 3 AND 4

**Lemma 9.** Let  $Q_i^*$  be the  $Q$ -function for the optimal policy of task  $i$ , we have

$$|Q_i^* - Q_j^*| \leq \frac{1 + \gamma}{1 - \gamma} \phi_{\max} \cdot \|\mathbf{w}_i^* - \mathbf{w}_j^*\|_2. \quad (66)$$

*Proof of Theorem 3.* For any task  $j \in [n]$ , we have

$$\begin{aligned} Q_{n+1}^{\pi_{n+1}}(s, a) - Q_{n+1}^{\pi_j}(s, a) &= \max_{i \in [n]} Q_{n+1}^{\pi_i^*}(s, a) - Q_{n+1}^{\pi_j}(s, a) \\ &\geq Q_{n+1}^{\pi_j^*}(s, a) - Q_{n+1}^{\pi_j}(s, a) \\ &= (\psi_j(\Theta_j^*) - \psi_j(\Theta_j^{(T)})) \cdot \mathbf{w}_{n+1}^*. \end{aligned} \quad (67)$$

According to Theorem 1, we have

$$\|\psi_j(\Theta_j^*) - \psi_j(\Theta_j^{(T)})\|_2 \leq \frac{(2 + \gamma) \cdot R_{\max} \cdot \log^2 T + C^* \|\mathbf{w}_j^{(0)} - \mathbf{w}_j^*\|_2}{(1 - \gamma - c_N) \rho_1 K^{-2} - C^*} \cdot \frac{1}{T} := \frac{C_3}{T} \quad (68)$$

Then, we have

$$\mathcal{T}^\pi Q_{n+1}^{\pi_j}(s, a) \geq Q_{n+1}^{\pi_j}(s, a) - \gamma \cdot \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T}. \quad (69)$$

Therefore, with the contraction property of the Bellman operator  $\mathcal{T}^\pi$ , we have

$$\begin{aligned} Q_{n+1}^{\pi_{n+1}}(s, a) &= \lim_{k \rightarrow \infty} (\mathcal{T}^\pi)^k Q_{n+1}^{\pi_j}(s, a) \\ &\geq \lim_{k \rightarrow \infty} (\mathcal{T}^\pi)^{k-1} \left( Q_{n+1}^{\pi_j}(s, a) - \gamma \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \right) \\ &= \lim_{k \rightarrow \infty} (\mathcal{T}^\pi)^{k-2} \cdot \mathcal{T}^\pi \left( Q_{n+1}^{\pi_j}(s, a) - \gamma \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \right) \\ &= \lim_{k \rightarrow \infty} (\mathcal{T}^\pi)^{k-2} \cdot \left( \mathcal{T}^\pi Q_{n+1}^{\pi_j}(s, a) - \gamma^2 \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \right) \\ &= \lim_{k \rightarrow \infty} (\mathcal{T}^\pi)^{k-2} \left( Q_{n+1}^{\pi_j}(s, a) - \gamma \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} - \gamma^2 \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \right) \\ &= Q_{n+1}^{\pi_j}(s, a) - \sum_{k=1}^{\infty} \gamma^k \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \\ &= Q_{n+1}^{\pi_j}(s, a) - \frac{\gamma}{1 - \gamma} \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \\ &\geq Q_{n+1}^{\pi_j^*}(s, a) - \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} - \frac{\gamma}{1 - \gamma} \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \\ &= Q_{n+1}^{\pi_j^*}(s, a) - \frac{1}{1 - \gamma} \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{T} \end{aligned} \quad (70)$$

For any policy  $\pi_j^*$  with  $j \in [n]$ , we have

$$\begin{aligned} Q_{n+1}^{\pi_j^*}(s, a) - Q_{n+1}^{\pi_{n+1}}(s, a) &= (Q_{n+1}^{\pi_j^*}(s, a) - Q_{n+1}^{\pi_j^*}(s, a)) + (Q_{n+1}^{\pi_j^*}(s, a) - Q_{n+1}^{\pi_j}(s, a)) \\ &\leq \frac{2\gamma}{1 - \gamma} \cdot \max_{s, a} |r_{n+1}(s, a) - r_j(s, a)| + \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{(1 - \gamma)T} \\ &\leq \frac{2\gamma \cdot \phi_{\max}}{1 - \gamma} \|\mathbf{w}_{n+1} - \mathbf{w}_j\|_2 + \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{(1 - \gamma)T}. \end{aligned} \quad (71)$$

Since (71) holds for any  $j$ , we have

$$|Q_{n+1}^{\pi_j^*}(s, a) - Q_{n+1}^{\pi_{n+1}}(s, a)| \leq \frac{2\gamma \cdot \phi_{\max}}{1 - \gamma} \min_{j \in [n]} \|\mathbf{w}_{n+1} - \mathbf{w}_j\|_2 + \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{(1 - \gamma)T}. \quad (72)$$

From Lemma 9, we know that

$$Q_{n+1}^{\pi_j^*}(s, a) - Q_{n+1}^{\pi_j}(s, a) \leq \frac{1 + \gamma}{1 - \gamma} \cdot \max_{s, a} |r_{n+1}(s, a) - r_j(s, a)|. \quad (73)$$

Similar to (70) to (72), we have

$$|Q_{n+1}^*(s, a) - Q_{n+1}^{\pi_{n+1}}(s, a)| \leq \frac{(1 + \gamma) \cdot \phi_{\max}}{1 - \gamma} \min_{j \in [n]} \|\mathbf{w}_{n+1} - \mathbf{w}_j\|_2 + \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{(1 - \gamma)T}. \quad (74)$$

□

*Proof of Theorem 4.* Let  $\pi'_{n+1}$  be generalized policy with DQN via GPI. Similar to (67), we have

$$\begin{aligned} & Q_{n+1}^{\pi'_{n+1}}(s, a) - Q_{n+1}^{\pi'_j}(s, a) \\ &= \max_{i \in [n]} Q_{n+1}^{\pi_i^*}(s, a) - Q_{n+1}^{\pi'_j}(s, a) \\ &\geq Q_{n+1}^{\pi_j^*}(s, a) - Q_{n+1}^{\pi'_j}(s, a) \\ &= \psi_j(\Theta_j^*) \mathbf{w}_{n+1}^* - \psi_j(\Theta_j^{(T)}) \mathbf{w}_j^{(t)} \\ &\approx \psi_j(\Theta_j^*) \mathbf{w}_{n+1}^* - \psi_j(\Theta_j^{(T)}) \mathbf{w}_j^* \\ &= \psi_j(\Theta_j^*) \mathbf{w}_{n+1}^* - \psi_j(\Theta_j^{(T)}) \mathbf{w}_{n+1}^* + \psi_j(\Theta_j^{(T)}) \mathbf{w}_{n+1}^* - \psi_j(\Theta_j^{(T)}) \mathbf{w}_j^* \\ &\geq -\|\Theta_j^* - \Theta_j^{(T)}\| \cdot \|\mathbf{w}_{n+1}^*\|_2 - \frac{1}{1 - \gamma} \phi_{\max} \cdot \|\mathbf{w}_{n+1}^* - \mathbf{w}_j^*\|_2. \end{aligned} \quad (75)$$

Following similar steps in the proof of Theorem 3, we have

$$\begin{aligned} |Q_{n+1}^*(s, a) - Q_{n+1}^{\pi'_{n+1}}(s, a)| &\leq \frac{(1 + \gamma) \cdot \phi_{\max}}{1 - \gamma} \min_{j \in [n]} \|\mathbf{w}_{n+1} - \mathbf{w}_j\|_2 + \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{(1 - \gamma)T} \\ &\quad + \frac{1}{1 - \gamma} \phi_{\max} \cdot \min_{j \in [n]} \|\mathbf{w}_{n+1}^* - \mathbf{w}_j^*\|_2 \\ &\leq \frac{2 \cdot \phi_{\max}}{1 - \gamma} \min_{j \in [n]} \|\mathbf{w}_{n+1} - \mathbf{w}_j\|_2 + \frac{C_3 \|\mathbf{w}_{n+1}^*\|_2}{(1 - \gamma)T}. \end{aligned} \quad (76)$$

□

## D PROOF OF THEOREM 2

*Proof of Theorem 2.* For task  $i$ , let  $\pi_j$  be the policy derived from  $\psi_j(\Theta_j^{(T)}) \mathbf{w}_i^*$  with  $1 \leq j \leq i$ , where  $\Theta_j^{(T)}$  is the returned neuron weights for the successor feature of task  $j$ .

Similar to (74), we have

$$Q_i^*(s, a) - Q_i^{\pi_j}(s, a) \leq \frac{(1 + \gamma) \cdot \phi_{\max}}{1 - \gamma} \|\mathbf{w}_j - \mathbf{w}_i\|_2 + \frac{C_3 \|\mathbf{w}_i^*\|_2}{(1 - \gamma)T}. \quad (77)$$

Let  $\pi'$  be the policy derived from  $\psi_i(\Theta_i^{(t)}) \mathbf{w}_i^*$  at iteration  $t$  for task  $i$ , we have

$$Q_i^*(s, a) - Q_i^{\pi'} \leq \|\Theta_i^{(t)} - \Theta_i^*\|_2 \cdot \|\mathbf{w}_i^*\|_2. \quad (78)$$

Therefore, at iteration  $t$  for task  $i$ , we have

$$\begin{aligned} C_t &= |Q_i^*(s, a) - Q_i^{\pi_i^{(t)}}| \\ &\leq \min \left\{ \frac{(1 + \gamma) \cdot \phi_{\max}}{1 - \gamma} \min_{1 \leq j \leq i} \|\mathbf{w}_j - \mathbf{w}_i\|_2 + \frac{C_3 \|\mathbf{w}_i^*\|_2}{(1 - \gamma)T}, \|\Theta_i^{(t)} - \Theta_i^*\|_2 \cdot \|\mathbf{w}_i^*\|_2 \right\} \\ &\lesssim \min \left\{ \frac{(1 + \gamma) \cdot \phi_{\max}}{1 - \gamma} \min_{1 \leq j \leq i} \|\mathbf{w}_j - \mathbf{w}_i\|_2, \|\Theta_i^{(t)} - \Theta_i^*\|_2 \cdot \|\mathbf{w}_i^*\|_2 \right\} \quad (\text{As } T \text{ is sufficiently large}) \\ &= \min\{q_t, 1\} \cdot \|\Theta_i^{(t)} - \Theta_i^*\|_2, \end{aligned} \quad (79)$$



where

$$q_t = \frac{(1 + \gamma)R_{\max}}{1 - \gamma} \cdot \frac{\min_{1 \leq i \leq j-1} \|\mathbf{w}_i^* - \mathbf{w}_j^*\|_2}{\|\Theta_j^{(t)} - \Theta_j^*\|_2} \quad (80)$$

Following similar steps in (59) in the proof of Theorem 1, with  $C_t$  satisfying (79), we have

$$\|\theta^{(T)} - \theta^*\|_2 \leq \frac{1}{T} \sum_{t=1}^{T-1} \frac{(2 + \gamma) \cdot R_{\max} \cdot \log^2 T + C^* \|\mathbf{w}^{(0)} - \mathbf{w}^*\|_2}{(1 - \gamma - c_N)\rho_1 K^{-2} - \min\{1, q_t\} \cdot C^*} \cdot \frac{1}{T}. \quad (81)$$

□

## E ADDITIONAL NUMERICAL EXPERIMENTS

In this section we empirically validate the theoretical results obtained in the previous section, using synthetic and real-world RL benchmarks.

### E.1 SYNTHETIC DATA SETTINGS

Here, we define an MDP that contains two tasks with shared state transition dynamics. The MDP consists of a state space with  $|\mathcal{S}| = 10,000$ , an action space with  $|\mathcal{A}| = 4$ . For the first task, its successor feature is parameterized by a deep neural network with the randomly generated neuron weights  $\Theta_1^*$ , and  $\mathbf{w}_1^*$  are randomly generated as the corresponding reward mapping. We then generate  $\phi$  based on (10) with  $\psi(\Theta_1^*)$ . Since  $\phi$  is shared across all tasks, for Task 2, we randomly generate the reward mapping  $\mathbf{w}_2^*$  and then calculate  $\psi_2^*$  accordingly.

### E.2 REAL DATA: REACHER ENVIRONMENT

The reacher environment is a robotic arm manipulation task consisting of a robotic arm with two joint torque controls. The state space is continuous, and the state features consist of angular displacement and angular velocity of the two joints. The actual action space for the robot arm is continuous that consists of the torques applied to the two joints, and is discretized for 3 values (for each joint torque). Thus, the total discretized action space consists of 9 actions ( $|\mathcal{A}| = 9$ ). The discount factor used is  $\gamma = 0.9$ . Multiple tasks in this environment is defined by goal locations, and the objective of each task is to move the tip of the robotic arm towards the goal location.

The reward of each task is defined by the distance  $\delta$ , measured from the tip of the robotic arm to the corresponding goal location. Specifically, a reward of  $1 - \delta$  is given to the agent at each time step. There are 12 predefined tasks, and  $\phi$  for a given state (common to all 12 tasks) is defined by stacking the reward for each of the 12 tasks for a given state as a vector. The corresponding reward weights  $\mathbf{w}_i^*$  for  $i = 1, \dots, 12$  are defined by one hot vectors, where the  $i^{th}$  element of  $\mathbf{w}_i^*$  is 1 and other elements are 0. Thus, the inner product  $\phi^\top \mathbf{w}_i^*$  naturally recovers the reward for the  $i^{th}$  task. For running experiments with this task, we use the open source code base <https://github.com/mike-gimelfarb/deep-successor-features-for-transfer.git>.

We first provide a comparison of the performance of SF-DQN with GPI, SF-DQN without GPI, and DQN with GPI, in Figure 3a. Here we consider the average transfer performance for four tasks, after training on a source task. It can be seen that SFDQN with GPI performs better compared to its no GPI counterpart. Both of these agents perform significantly better compared to DQN with GPI. hence, this result validates our theoretical results for the performance of these three methods.

Next, we investigate the performance of the SFDQN agent when the target task reward mappings are not known and learned simultaneously with successor features. We consider varying distances from the initial target task reward mapping to the true target task reward mapping. The results are shown in Figure 3b. It can be seen that when the reward mappings are initialized far away from the true reward mappings, the convergence of the SF-DQN agent is slower compared to that is initialized closer to the true reward mappings. This aligns with our convergence analysis for the SF-DQN agent with GPI.

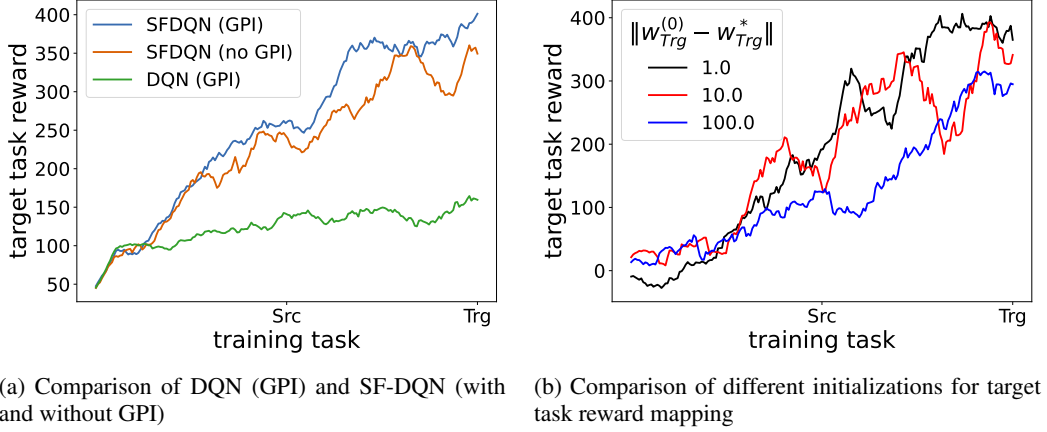
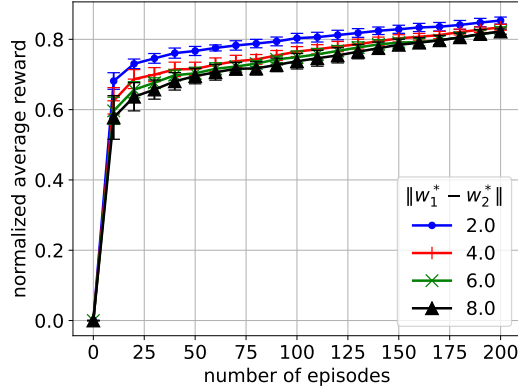


Figure 3: Experiments on Reacher environment

### E.3 ADDITIONAL EXPERIMENTS ON SYNTHETIC RL BENCHMARKS

**Effect of  $\|w_1^* - w_2^*\|$  on knowledge transfer.** We investigate the effect of the distance between  $w_1^*$  to  $w_2^*$ , on the transfer performance of the SFDQN. For this purpose, we assume SF-DQN agents have access to optimal reward mappings when training on Tasks 1 and 2. After obtaining  $\phi$  as described earlier, we initialize and train  $\Theta_2$  using  $\phi$  and  $w_2^*$ , with GPI. Reward defined by  $\phi \cdot w_2^*$  is used to obtain the average reward for Task 2. We repeat the process for different choices of  $w_2^*$ , and the results are shown in Figure 4. It can be seen that, when the task similarity is low (i.e.  $\|w_1^* - w_2^*\|$  is large), the performance of the SF-DQN agent with GPI is poor. On the other hand, when the task similarity is high, the performance becomes significantly better.

Figure 4: Effect  $\|w_1^* - w_2^*\|$  on the convergence of SF-DQN agent when training on task 2 with GPI

## F PROOF OF LEMMAS IN APPENDIX B

### F.1 PROOF OF LEMMA 6

Lemma 6 provides the lower and upper bounds for the eigenvalues of the Hessian matrix of population risk function in (30). According to Weyl’s inequality in Lemma 1, the eigenvalues of  $\nabla_\ell^2 f(\cdot)$  at any fixed point  $\theta$  can be bounded in the form of (86). Therefore, we first provide the lower and upper bounds for  $\nabla_\ell^2 f$  at the desired ground truth  $\theta^*$ . Then, the bounds for  $\nabla_\ell^2 f$  at any other point  $\theta$  is bounded through (30) by utilizing the conclusion in Lemma 10. Lemma 10 illustrates the distance between the Hessian matrix of  $f$  at  $\theta$  and  $\theta^*$ . Lemma 11 provides the lower bound of

$\mathbb{E}_{\mathbf{x}} \left( \sum_{j=1}^K \boldsymbol{\alpha}_j^\top \frac{\partial \psi}{\partial \theta_{\ell,k}}(\theta^*) \right)^2$  when  $\mathbf{x}$  belongs to sub-Gaussian distribution, which is used in proving the lower bound of the Hessian matrix in (87).

**Lemma 10.** *Let  $f(\theta)$  be the population risk function defined in (30). If  $\theta$  is close to  $\theta^*$  such that*

$$\|\theta - \theta^*\|_2 \lesssim \frac{\rho_1}{K} \quad (82)$$

*we have*

$$\|\nabla_\ell^2 f(\theta) - \nabla_\ell^2 f(\theta^*)\|_2 \lesssim \frac{1}{K} \cdot \|\theta - \theta^*\|_2. \quad (83)$$

**Lemma 11.** *Suppose the following assumptions hold:*

1.  $\{\theta_j\}_{j=1}^K \in \mathbb{R}^{K_\ell}$  are linear independent,
2. Let  $p(\mathbf{h}) : \mathbb{R}^{K_\ell} \rightarrow [0, 1]$  be the probability density for  $\mathbf{h}$  such that  $\mathbb{E}_{\mathbf{h}} \|\mathbf{h}\|_2^2 \leq +\infty$ .

*Let  $\boldsymbol{\alpha} \in \mathbb{R}^{K_\ell K_{\ell-1}}$  be the unit vector defined in (36), we have*

$$\rho_1 := \min_{\|\boldsymbol{\alpha}\|_2=1} \int_{\mathcal{R}} \left( \sum_{j=1}^K \boldsymbol{\alpha}^\top \mathbf{h} \phi'(\theta_{\ell,j}^\top \mathbf{h}) \right)^2 p_H(\mathbf{h}) \cdot d\mathbf{h} > 0, \quad (84)$$

*where  $\mathcal{R} \subset \mathbb{R}^{K_\ell}$  with  $\int_{\mathcal{R}} f_H(\mathbf{h}) > 0$ . Moreover, if further assuming  $\mathbf{h}$  belongs to Gaussian distribution, we have  $\rho_1 > 0.091$ .*

**Lemma 12.** *Let  $\mathbf{h}^{(\ell)}(\theta)$  be the function defined in (37). When  $\theta$  is sufficiently close to  $\theta^*$ , i.e.,  $\|\theta - \theta^*\|_2$  is smaller than some positive constant  $c < 1$ , we have*

$$\begin{aligned} \|\mathbf{h}^{(\ell)}(\theta)\|_2 &\lesssim \|\mathbf{x}\|_2, \\ \|\mathbf{h}^{(\ell)}(\theta) - \mathbf{h}^{(\ell)}(\theta^*)\|_2 &\lesssim \|\theta - \theta^*\|_2 \cdot \|\mathbf{x}\|_2. \end{aligned} \quad (85)$$

*Proof of Lemma 6.* Let  $\lambda_{\max}(\theta)$  and  $\lambda_{\min}(\theta)$  denote the largest and smallest eigenvalues of  $\nabla_\ell^2 f(\theta)$  at  $\theta$ , respectively. Then, from Lemma 1, we have

$$\begin{aligned} \lambda_{\max}(\theta) &\leq \lambda_{\max}(\theta^*) + \|\nabla_\ell^2 f(\theta) - \nabla_\ell^2 f(\theta^*)\|_2, \\ \lambda_{\min}(\theta) &\geq \lambda_{\min}(\theta^*) - \|\nabla_\ell^2 f(\theta) - \nabla_\ell^2 f(\theta^*)\|_2. \end{aligned} \quad (86)$$

Then, we provide the lower bound of the Hessian matrix of the population function at  $\theta^*$ . Let  $\mathcal{P}$  be the distribution for  $\mathbf{h}^{(\ell)}(\theta)$  when  $\mathbf{x} \sim \mu^*$  with probability density function denoted as  $p_H$ . For any  $\boldsymbol{\alpha} \in \mathbb{R}^{K_\ell K}$  with  $\|\boldsymbol{\alpha}\|_2 = 1$ , we have

$$\begin{aligned} &\min_{\|\boldsymbol{\alpha}\|_2=1} \boldsymbol{\alpha}^\top \nabla_\ell^2 f(\theta^*) \boldsymbol{\alpha} \\ &= \frac{1}{K^2} \min_{\|\boldsymbol{\alpha}\|_2=1} \mathbb{E}_{\mathbf{h} \sim \mathcal{P}} \left( \sum_{j=1}^K \boldsymbol{\alpha}_j^\top \mathbf{h}^{(\ell)} \mathcal{J}_{\ell,k} \phi'(\theta_{\ell,j}^\top \mathbf{h}^{(\ell)}) \right)^2 \\ &= \frac{1}{K^2} \min_{\|\boldsymbol{\alpha}\|_2=1} \int_{\mathbb{R}^{K_\ell-1}} \left( \sum_{j=1}^K \boldsymbol{\alpha}_j^\top \mathbf{h}^{(\ell)} \mathcal{J}_{\ell,k} \phi'(\theta_{\ell,j}^\top \mathbf{h}^{(\ell)}) \right)^2 p_H(\mathbf{h}^{(\ell)}) \cdot d\mathbf{h}^{(\ell)} \\ &= \frac{1}{K^2} \min_{\|\boldsymbol{\alpha}\|_2=1} \int_{\{\mathbf{h}^{(\ell)} | \mathcal{J}_{\ell,k} \neq 0\}} \left( \sum_{j=1}^K \boldsymbol{\alpha}_j^\top \mathbf{h}^{(\ell)} \phi'(\theta_{\ell,j}^\top \mathbf{h}^{(\ell)}) \right)^2 p_H(\mathbf{h}^{(\ell)}) \cdot d\mathbf{h}^{(\ell)} \\ &\gtrsim \frac{\rho_1}{K^2}, \end{aligned} \quad (87)$$

where the last inequality comes from Lemma 11, and Lemma 11 holds since  $\mathbf{h}^{(\ell)}$  belongs to sub-Gaussian distribution and  $\theta_\ell$  is full rank.

Next, the upper bound of  $\nabla_\ell^2 f$  can be bounded as

$$\begin{aligned}
& \max_{\|\alpha\|_2=1} \alpha^\top \nabla_\ell^2 f(\theta^*) \alpha \\
&= \frac{1}{K^2} \max_{\|\alpha\|_2=1} \mathbb{E}_{\mathbf{x}} \left( \sum_{j=1}^K \alpha_j^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\theta_{\ell,j}^{*\top} \mathbf{h}^{(\ell)}) \right)^2 \\
&= \frac{1}{K^2} \max_{\|\alpha\|_2=1} \mathbb{E}_{\mathbf{x}} \sum_{j_1=1}^K \sum_{j_2=1}^K \alpha_{j_1}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\theta_{\ell,j_1}^{*\top} \mathbf{h}^{(\ell)}) \cdot \alpha_{j_2}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\theta_{\ell,j_2}^{*\top} \mathbf{h}^{(\ell)}) \\
&= \frac{1}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \mathbb{E}_{\mathbf{x}} \alpha_{j_1}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\theta_{\ell,j_1}^{*\top} \mathbf{h}^{(\ell)}) \cdot \alpha_{j_2}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\theta_{\ell,j_2}^{*\top} \mathbf{h}^{(\ell)}) \\
&\leq \frac{1}{K^2} \max_{\|\alpha\|_2=1} \sum_{j_1=1}^K \sum_{j_2=1}^K \left[ \mathbb{E}_{\mathbf{x}} (\alpha_{j_1}^\top \mathbf{h}^{(\ell)})^4 \cdot \mathbb{E}_{\mathbf{x}} (\phi'(\theta_{\ell,j_1}^{*\top} \mathbf{h}^{(\ell)}))^4 \cdot \mathbb{E}_{\mathbf{x}} (\alpha_{j_2}^\top \mathbf{h}^{(\ell)})^4 \cdot \mathbb{E}_{\mathbf{x}} (\phi'(\theta_{\ell,j_2}^{*\top} \mathbf{h}^{(\ell)}))^4 \right]^{1/4} \\
&\leq \frac{1}{K^2} \max_{\|\alpha\|_2=1} \sum_{j_1=1}^K \sum_{j_2=1}^K \left[ \mathbb{E}_{\mathbf{x}} (\alpha_{j_1}^\top \mathbf{x})^4 \cdot \mathbb{E}_{\mathbf{x}} (\alpha_{j_2}^\top \mathbf{x})^4 \right]^{1/4} \\
&\leq \frac{3}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \|\alpha_{j_1}\|_2 \cdot \|\alpha_{j_2}\|_2 \leq \frac{6}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \frac{1}{2} (\|\alpha_{j_1}\|_2^2 + \|\alpha_{j_2}\|_2^2) \\
&= \frac{6}{K}.
\end{aligned} \tag{88}$$

Therefore, we have

$$\lambda_{\max}(\theta^*) = \max_{\|\alpha\|_2=1} \alpha^\top \nabla_\ell^2 f(\theta^*; p) \alpha \leq \frac{6}{K}. \tag{89}$$

Then, given (82), we have

$$\|\theta - \theta^*\|_2 \lesssim \frac{2\rho_1}{K}. \tag{90}$$

Combining (90) and Lemma 10, we have

$$\|\nabla_\ell^2 f(\theta) - \nabla_\ell^2 f(\theta^*)\|_2 \lesssim \frac{\rho_1}{K^2}. \tag{91}$$

Therefore, from (91) and (86), we have

$$\begin{aligned}
\lambda_{\max}(\theta) &\leq \lambda_{\max}(\theta^*) + \|\nabla_\ell^2 f(\theta) - \nabla_\ell^2 f(\theta^*)\|_2 \leq \frac{6}{K} + \frac{\rho_1}{2K^2} \leq \frac{7}{K}, \\
\lambda_{\min}(\theta) &\geq \lambda_{\min}(\theta^*) - \|\nabla_\ell^2 f(\theta) - \nabla_\ell^2 f(\theta^*)\|_2 \geq \frac{\rho_1}{K^2} - \frac{\rho_1}{2K^2} = \frac{\rho_1}{2K^2},
\end{aligned} \tag{92}$$

which completes the proof.  $\square$

## F.2 PROOF OF LEMMA 7

The error bound between  $\|\nabla_\ell f - g_t\|_2$  is divided into bounding  $I_1$ ,  $I_2$ ,  $I_3$ , and  $I_4$  as shown in (98).  $I_1$  represents the deviation of the gradient of  $\mathcal{D}_t$  to their expectation, which can be bounded through concentration inequality.  $I_2$  is derived from the distribution shift between the trajectory and its stationary distribution, which can be bounded with assumption 2.  $I_3$  come from the data distribution shift between the behavior policy and optimal policy.  $I_4$  comes from the inconsistency of the "noisy" label and the "ground truth" label in the population risk function (30). To ensure a smooth presentation, we will defer the proof of  $I_1 - I_4$  until we have completed the main proof of Lemma 7.

*Proof of Lemma 7.* From (33), we know that

$$\begin{aligned}
& g^{(t)}(\theta_{\ell,k}^{(t)}; \mathcal{X}_m) \\
&= \sum_{m \in \mathcal{D}_t} (\psi(\theta^{(t)}; \mathbf{s}_m, a_m) - y_m^{(t)}) \cdot \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \\
&= \sum_{m \in \mathcal{D}_t} \left( \psi(\theta^{(t)}; \mathbf{s}_m, a_m) - \phi(\theta^*; \mathbf{s}_m, a_m) - \gamma \cdot \psi(\mathbf{s}'_m, a'_m; \theta^{(t)}) \right) \cdot \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \\
&= \sum_{m \in \mathcal{D}_t} \left( \psi(\theta^{(t,n)}; \mathbf{s}_m, a_m) - \psi(\theta^*; \mathbf{s}_m, a_m) + \gamma \cdot \max_{a'} \psi(\mathbf{s}'_m, a'; \theta^*) \right. \\
&\quad \left. - \gamma \cdot \psi(\mathbf{s}'_m, a'_m; \theta^{(t)}) \right) \cdot \frac{\partial \psi(\theta^{(t,n)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \\
&= \sum_{m \in \mathcal{D}_t} \left( \psi(\theta^{(t)}; \mathbf{s}_m, a_m) - \psi(\theta^*; \mathbf{s}_m, a_m) \right) \cdot \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \\
&\quad + \gamma \cdot \left( \max_{a'} \psi(\mathbf{s}'_m, a'; \theta^*) - \psi(\mathbf{s}'_m, a'_m; \theta^{(t)}) \right) \cdot \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \\
&:= \sum_{m \in \mathcal{D}_t} b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) + \Delta b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m),
\end{aligned} \tag{93}$$

where we have

$$b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) = \left( \psi(\theta^{(t)}; \mathbf{s}_m, a_m) - \psi(\theta^*; \mathbf{s}_m, a_m) \right) \cdot \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \tag{94}$$

and

$$\Delta b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) = \left( \max_{a'} \psi(\theta^*; \mathbf{s}'_m, a') - \psi(\theta^{(t-1)}; \mathbf{s}'_m, a'_m) \right) \cdot \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}}. \tag{95}$$

Then, let us define  $\bar{b}_{\ell,k}^{(t)}$  as

$$\bar{b}_{\ell,k}^{(t)}(\theta; \mathcal{X}) = \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t} \left( \psi(\theta; \mathbf{s}, a) - \psi(\theta^*; \mathbf{s}, a) \right) \cdot \nabla_{\theta} \psi(\theta; \mathbf{s}, a). \tag{96}$$

From (30), we know that

$$\frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta^{(t)}) = \mathbb{E}_{(\mathbf{s}, a) \sim \mu^*} \left( \phi(\theta^{(t)}; \mathbf{s}, a) - \phi(\theta^*; \mathbf{s}, a) \right) \cdot \frac{\partial \phi(\theta^{(t)}; \mathbf{s}, a)}{\partial \theta_{\ell,k}}. \tag{97}$$

Then, from (93) and (97), we have

$$\begin{aligned}
& g^{(t)}(\theta_{\ell,k}^{(t)}; \mathcal{X}_m) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta^{(t)}; \mathcal{X}_m) \\
&= \sum_{m \in \mathcal{D}_t} b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) + \Delta b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta^{(t)}; \mathcal{X}_m) \\
&= \left[ b_{\ell,k}^{(t)}(\theta_{\ell,k}^{(t)}; \mathcal{X}_m) - \mathbb{E}_{\mathcal{X}_m \sim \mathcal{D}_t} b_{\ell,k}^{(t)}(\theta_{\ell,k}^{(t)}; \mathcal{X}_m) \right] + \left[ \mathbb{E}_{\mathcal{X}_m \sim \mathcal{D}_t} b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) - \bar{b}_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) \right] \\
&\quad + \left[ \bar{b}_{\ell,k}^{(t)}(\theta^{(t)}) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta^{(t)}) \right] + \mathbb{E}_{\mathcal{X}_m \sim \mathcal{D}_t} \Delta b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) \\
&:= \mathbf{I}_1 + \mathbf{I}_2 + \mathbf{I}_3 + \mathbf{I}_4.
\end{aligned} \tag{98}$$

Therefore, we have

$$\left\| g^{(t)}(\theta_{\ell,k}^{(t)}; \mathcal{X}_m) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta^{(t)}) \right\|_2 \leq \|\mathbf{I}_1\|_2 + \|\mathbf{I}_2\|_2 + \|\mathbf{I}_3\|_2 + \|\mathbf{I}_4\|_2. \tag{99}$$

Next, we first provide the bound for  $\|\mathbf{I}_1\|_2$ ,  $\|\mathbf{I}_2\|_2$ ,  $\|\mathbf{I}_3\|_2$ , and  $\|\mathbf{I}_4\|_2$  as

$$\begin{aligned}\|\mathbf{I}_1\|_2 &\leq \frac{1}{K_\ell} \cdot \|\theta - \theta^*\|_2 \cdot \sqrt{\frac{d \log q}{|\mathcal{D}_t|}}, \\ \|\mathbf{I}_2\|_2 &\leq \frac{R_{\max}}{1-\gamma} \cdot (1+\gamma)\tau^* \cdot \eta_{t-\tau^*}, \\ \|\mathbf{I}_3\|_2 &\leq |\mathcal{A}| \cdot \frac{R_{\max}}{1-\gamma} \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1-\nu}) \cdot C_t, \\ \|\mathbf{I}_4\|_2 &\leq \frac{\gamma}{K_\ell} \cdot \|\theta^{(t)} - \theta^*\|_2,\end{aligned}\tag{100}$$

where  $|\mathcal{A}|$  is the size of action space. The details for the derivation of  $I_1$ - $I_4$  can be found after the proof.

Let  $\boldsymbol{\alpha} \in \mathbb{R}^{Kd}$  and  $\boldsymbol{\alpha}_j \in \mathbb{R}^d$  with  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1^T, \boldsymbol{\alpha}_2^T, \dots, \boldsymbol{\alpha}_K^T]^T$ , with probability at least  $1 - q^{-d}$ , we have

$$\begin{aligned}\|g^{(t)}(\theta_\ell; \theta) - \nabla_\ell f_{\pi^*}(\theta)\|_2^2 &= \left| \boldsymbol{\alpha}^T (g^{(t)}(\theta) - \nabla f_{\pi^*}(\theta)) \right|^2 \\ &\leq \sum_{k=1}^K \left| \boldsymbol{\alpha}_k^T (g^{(t)}(\theta_{\ell,k}; \theta) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta)) \right|^2 \\ &\leq \sum_{k=1}^K \left\| g^{(t)}(\theta_{\ell,k}; \theta) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta) \right\|_2^2 \cdot \|\boldsymbol{\alpha}_k\|_2^2 \\ &\leq \max_k \left\| g^{(t)}(\theta_{\ell,k}; \theta) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta) \right\|_2^2.\end{aligned}\tag{101}$$

In conclusion, we have

$$\begin{aligned}&\|g^{(t)}(\theta_\ell; \theta) - \nabla_\ell f_{\pi^*}(\theta)\|_2 \\ &\leq \max_k \left\| g^{(t)}(\theta_{\ell,k}; \theta) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta) \right\|_2 \\ &\leq \max_k \|\mathbf{I}_1(k)\|_2 + \|\mathbf{I}_2(k)\|_2 + \|\mathbf{I}_3(k)\|_2 + \|\mathbf{I}_4(k)\|_2 \\ &\leq \frac{1}{K_\ell} \cdot \|\theta - \theta^*\|_2 \cdot \sqrt{\frac{d \log q}{|\mathcal{D}_t|}} + \frac{R_{\max}}{1-\gamma} \cdot (1+\gamma)\tau^* \cdot \eta_{t-\tau^*} \\ &\quad + |\mathcal{A}| \cdot \frac{R_{\max}}{1-\gamma} \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1-\nu}) \cdot C_t + \frac{\gamma}{K_\ell} \cdot \|\theta^{(t)} - \theta^*\|_2,\end{aligned}\tag{102}$$

where  $\tau^* = \min\{t \mid \lambda \nu^t \leq \eta_T\}$  □

### F.2.1 PROOF OF UPPER BOUND OF $I_1$

*Proof.* We define a random variable

$$Z^{(\ell)}(k) = (\psi(\theta; \mathbf{s}, a) - \psi(\theta^*; \mathbf{s}, a)) \cdot \mathcal{J}_{\ell,k} \cdot \boldsymbol{\alpha}^T \mathbf{h}^{(\ell)}(\theta)$$

with  $(\mathbf{s}, a) \sim \mathcal{D}_t$  and

$$Z_m^{(\ell)}(k) = (Q(\mathbf{x}_m; \theta) - Q(\mathbf{x}_m; \theta^*)) \cdot \mathcal{J}_{\ell,k} \cdot \boldsymbol{\alpha}^T \mathbf{h}_n^{(\ell)}(\theta)$$

as the realization of  $Z^{(\ell)}$  for  $m \in \mathcal{D}_t$ , where  $\boldsymbol{\alpha}$  is any fixed unit vector.

According to the definition of  $I_1$  in (98), we can rewrite  $\mathbf{I}_1$  as

$$\mathbf{I}_1 = \frac{1}{K_\ell} \left[ \sum_{m \in \mathcal{D}_t} Z_m^{(\ell)}(k) - \mathbb{E}_{(\mathbf{s}, a) \sim \mathcal{D}_t} Z^{(\ell)}(k) \right].\tag{103}$$

Then, for any  $p \in \mathbb{N}^+$ , we have

$$\begin{aligned}
(\mathbb{E}|Z^{(\ell)}|^p)^{1/p} &= \left( \mathbb{E}_{\mathcal{X} \sim \mathcal{D}_t} |\psi(\theta; \mathbf{s}, a) - \psi(\theta^*; \mathbf{s}, a)|^p \cdot |\mathcal{J}_{\ell,k} \sigma'(\mathbf{w}_{\ell,k}^\top \mathbf{x})| \cdot |\boldsymbol{\alpha}^T \mathbf{h}^{(\ell)}|^p \right)^{1/p} \\
&\leq \left( \mathbb{E}_{\mathcal{X} \sim \mathcal{D}_t} |\psi(\theta; \mathbf{s}, a) - \psi(\theta^*; \mathbf{s}, a)|^p \cdot |\boldsymbol{\alpha}^T \mathbf{h}^{(\ell)}|^p \right)^{1/p} \\
&\leq \left( \mathbb{E}_{\mathcal{X} \sim \mathcal{D}_t} \|\theta - \theta^*\|_2 \cdot \|\mathbf{x}(\mathbf{s}, a)\|_2 \right)^p \cdot |\boldsymbol{\alpha}^T \mathbf{x}(\mathbf{s}, a)|^p \right)^{1/p} \\
&\lesssim \|\theta - \theta^*\|_2 \cdot p.
\end{aligned} \tag{104}$$

From Definition 2, we know that  $Z^{(\ell)}$  belongs to sub-exponential distribution with  $\|Z^{(\ell)}\|_{\psi_1} \lesssim \|\theta - \theta^*\|_2$ . Therefore, by Chernoff inequality, for any  $s \in \mathbb{R}$ , we have

$$\mathbb{P} \left\{ \left| \frac{1}{|\mathcal{D}_t|} \sum_{m \in \mathcal{D}_t} Z_m^{(\ell)}(k) - \mathbb{E} Z^{(\ell)}(k) \right| < t \right\} \leq 1 - \frac{e^{-(\|\theta - \theta^*\|_2)^2 \cdot |\mathcal{D}_t| \cdot s^2}}{e^{|\mathcal{D}_t| \cdot st}}. \tag{105}$$

Let  $t = \|\theta - \theta^*\|_2 \sqrt{\frac{d \log q}{N}}$  and  $s = \frac{2}{\|\theta - \theta^*\|_2} \cdot t$  for some large constant  $q > 0$ . Then, with probability at least  $1 - q^{-d}$ , we have

$$\left| \frac{1}{|\mathcal{D}_t|} \sum_{m \in \mathcal{D}_t} Z_m^{(\ell)}(k) - \mathbb{E} Z^{(\ell)}(k) \right| \lesssim \|\theta - \theta^*\|_2 \cdot \sqrt{\frac{d \log q}{|\mathcal{D}_t|}}. \tag{106}$$

From Lemma 4 and (103), with probability at least  $1 - |\mathcal{S}_{\frac{1}{2}}(d)| \cdot q^{-d}$ , we have

$$\|\mathbf{I}_1\|_2 \leq 2 \cdot \frac{1}{K_\ell} \left| \frac{1}{|\mathcal{D}_t|} \sum_{m \in \mathcal{D}_t} Z_m^{(\ell)} - \mathbb{E} Z^{(\ell)} \right| \lesssim \frac{1}{K_\ell} \|\theta - \theta^*\|_2 \cdot \sqrt{\frac{d \log q}{|\mathcal{D}_t|}}. \tag{107}$$

From Lemma 3, we know that  $|\mathcal{S}_{\frac{1}{2}}(d)| \leq 5^d$ . Therefore, the probability for (107) holds is at least  $1 - \left(\frac{q}{5}\right)^{-d}$ . Because  $q \gg 5$ , we denote the probability as  $1 - q^{-d}$  for convenience.  $\square$

### F.2.2 PROOF OF UPPER BOUND OF $I_2$

*Proof.*  $I_2$  is the bias of the data because the data  $(\mathbf{s}, a)$  at iteration  $t$  depends on the neural network parameters  $\theta^{(t)}$ . Recall the definition of  $b_{\ell,k}^{(t)}$  and  $\bar{b}_{\ell,k}^{(t)}$ , we define

$$\Delta_t = b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m) - \bar{b}_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m). \tag{108}$$

It is easy to verify that

$$\begin{aligned}
\|b_{\ell,k}^{(t)}(\theta; \mathcal{X}_m) - b_{\ell,k}^{(t)}(\tilde{\theta}; \mathcal{X}_m)\|_2 &\leq (1 + \gamma) \cdot \|\theta - \tilde{\theta}\|_2, \\
\|\bar{b}_{\ell,k}^{(t)}(\theta; \mathcal{X}_m) - \bar{b}_{\ell,k}^{(t)}(\tilde{\theta}; \mathcal{X}_m)\|_2 &\leq (1 + \gamma) \cdot \|\theta - \tilde{\theta}\|_2, \\
\text{and} \quad \|b_{\ell,k}^{(t)}\| &\lesssim \frac{R_{\max}}{1 - \gamma}.
\end{aligned} \tag{109}$$

Then, we have

$$\Delta_t(\theta) - \Delta_t(\tilde{\theta}) \lesssim (1 + \gamma) \cdot \|\theta - \tilde{\theta}\|_2. \tag{110}$$

Therefore, we have

$$\Delta_t(\theta^{(t)}) \leq \Delta_t(\theta^{(t-\tau)}) + \frac{1 + \gamma}{1 - \gamma} \cdot R_{\max} \cdot \sum_{i=t-\tau}^{t-1} \eta_i. \tag{111}$$

Then, we need to bound  $\delta_t(\theta^{(t-\tau)})$ .

Let us define the observed tuple  $O_t(\mathbf{s}, a, \mathbf{s}')$  as the collection of the state, action, and the next state at the  $t$ -th iteration. Note that

$$\theta^{(t-\tau)} \longrightarrow \mathbf{s}_{t-\tau} \longrightarrow \mathbf{s}_t \longrightarrow O_t \tag{112}$$

forms a Markov chain introduced by the policy  $\pi_t$ .

Let  $\tilde{\theta}^{(t-\tau,0)}$  and  $\tilde{O}_t$  be independently drawn from the marginal distributions of  $\theta^{(t-\tau,0)}$  and  $O_t$ , respectively.

With Lemma 9 in Bhandari et al. (2018), we have

$$\mathbb{E} \Delta_t(\theta^{(t-\tau)}, O_t) - \mathbb{E} \Delta_t(\tilde{\theta}^{(t-\tau)}, \tilde{O}_t) \lesssim 2 \sup_{\theta, O} |\Delta_t(\theta, O)| \cdot \lambda \cdot \nu^\tau. \quad (113)$$

By definition, we have  $\mathbb{E} \Delta_m(\tilde{\theta}^{(t-\tau)}, \tilde{O}_t) = 0$  and

$$|\Delta_t(\theta, O)| \leq \frac{2 R_{\max}}{1 - \gamma}. \quad (114)$$

Therefore, we have

$$\begin{aligned} \mathbb{E} \Delta_t(\theta^{(t)}) &\leq \mathbb{E} \Delta_t(\theta^{(t-\tau)}) + \frac{1 + \gamma}{1 - \gamma} \cdot R_{\max} \cdot \sum_{i=t-\tau}^{t-1} \eta_i \\ &\leq \frac{R_{\max}}{1 - \gamma} \left( \lambda \cdot \nu^\tau + (1 + \gamma) \cdot \tau \cdot \eta_{t-\tau} \right), \end{aligned} \quad (115)$$

where the last inequality comes from the fact that the step size  $\eta_m$  is non-increasing.

Choose  $\tau^* = \min \{t = 0, 1, 2, \dots \mid \lambda \nu^\tau \leq \eta_T\}$ . When  $t \leq \tau^*$ , we choose  $\tau = t$  and have

$$\mathbb{E} \Delta_t(\theta^{(t)}) \leq \frac{R_{\max}}{1 - \gamma} \cdot \tau^* \cdot \eta_0. \quad (116)$$

When  $n > \tau^*$ , we can choose  $\tau = \tau^*$  and obtain

$$\mathbb{E} \Delta_t(\theta^{(t)}) \leq \frac{R_{\max}}{1 - \gamma} \cdot (1 + \gamma) \tau^* \cdot \eta_{t-\tau^*}. \quad (117)$$

Combining (116) and (117), we have

$$|I_2| \leq \frac{R_{\max}}{1 - \gamma} \cdot (1 + \gamma) \tau^* \cdot \eta_{\max\{0, t-\tau^*\}}, \quad (118)$$

where  $\tau^* = \min\{t \mid \lambda \nu^t \leq \eta_T\}$ . □

### F.2.3 PROOF OF BOUND OF $I_3$

*Proof.* We have

$$\begin{aligned} I_3 &= \bar{b}_{\ell,k}^{(t)}(\theta^{(t)}) - \frac{\partial f_{\pi^*}}{\partial \theta_{\ell,k}}(\theta^{(t)}) \\ &= \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t} \left( \psi(\theta; \mathbf{s}, a) - \psi(\theta^*; \mathbf{s}, a) \right) \cdot \frac{\partial \psi(\theta; \mathbf{s}, a)}{\partial \theta_{\ell,k}} \\ &\quad - \mathbb{E}_{(\mathbf{s}, a) \sim \mu^*} \left( \psi(\theta; \mathbf{s}, a) - \psi(\theta^*; \mathbf{s}, a) \right) \cdot \frac{\partial \psi(\theta; \mathbf{s}, a)}{\partial \theta_{\ell,k}} \\ &= \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t} \left( \psi(\theta; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \mathbb{E}_{\mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \max_{a'} \psi(\theta^*; \mathbf{s}', a') \right) \cdot \frac{\partial \psi(\theta; \mathbf{s}, a)}{\partial \theta_{\ell,k}} \\ &\quad - \mathbb{E}_{(\mathbf{s}, a) \sim \mu^*} \left( \psi(\theta; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \mathbb{E}_{\mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \max_{a'} \psi(\theta^*; \mathbf{s}', a') \right) \cdot \frac{\partial \psi(\theta; \mathbf{s}, a)}{\partial \theta_{\ell,k}} \\ &= \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t, \mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \left( \psi(\theta; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \max_{a'} \psi(\theta^*; \mathbf{s}', a') \right) \cdot \frac{\partial \psi(\theta; \mathbf{s}, a)}{\partial \theta_{\ell,k}} \\ &\quad - \mathbb{E}_{(\mathbf{s}, a) \sim \mu^*, \mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \left( \psi(\theta; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \max_{a'} \psi(\theta^*; \mathbf{s}', a') \right) \cdot \frac{\partial \psi(\theta; \mathbf{s}, a)}{\partial \theta_{\ell,k}} \end{aligned} \quad (119)$$



Then, we have

$$\begin{aligned}
& \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} (\mu^*(d\mathbf{s}, da) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mu_t(d\mathbf{s}, da) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \\
&= \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} (\mathcal{P}^*(d\mathbf{s}) \pi^*(da|\mathbf{s}) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mathcal{P}_t(d\mathbf{s}) \pi_t(da|\mathbf{s}) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \\
&\leq \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} (\mathcal{P}^*(d\mathbf{s}) - \mathcal{P}_t(d\mathbf{s})) \pi^*(da|\mathbf{s}) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) \right| \\
&\quad + \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} \mathcal{P}_t(d\mathbf{s}) (\pi_t(da|\mathbf{s}) - \pi^*(da|\mathbf{s})) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) \right|.
\end{aligned} \tag{120}$$

From Theorem 3.1 in Mitrophanov (2005), we know that

$$\begin{aligned}
& \left| \int_{(\mathbf{s},a)} (\mathcal{P}^*(d\mathbf{s}) - \mathcal{P}_t(d\mathbf{s})) \right| \leq |\mathcal{A}| (\log_\nu \lambda^{-1} + \frac{1}{1-\nu}) C_t \\
& \text{and} \quad \left\| \pi_t(da|\mathbf{s}) - \pi^*(da|\mathbf{s}) \right\| \leq C_t.
\end{aligned} \tag{121}$$

Therefore, the bound of  $\mathbf{I}_3$  can be found as

$$\begin{aligned}
\|\mathbf{I}_3\|_2 &\leq \frac{R_{\max}}{1-\gamma} \cdot |\mathcal{A}| \cdot C_t \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1-\nu}) \\
&= |\mathcal{A}| \cdot \frac{R_{\max}}{1-\gamma} \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1-\nu}) \cdot C_t.
\end{aligned} \tag{122}$$

□

#### F.2.4 PROOF OF BOUND OF $I_4$

*Proof.* We have

$$\begin{aligned}
\|\mathbf{I}_4\| &= \|\Delta b_{\ell,k}^{(t)}(\theta^{(t)}; \mathcal{X}_m)\|_2 \\
&= \max_{\mathbf{s},a} \gamma \cdot \left( \max_{a'} \psi(\mathbf{s}'_m, a'; \theta^*) - \psi(\mathbf{s}'_m, a'_m; \theta^{(t)}) \right) \cdot \left\| \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \right\|_2 \\
&\leq \max_{\mathbf{s},a} \gamma \cdot \left( \max_{a'} \psi(\mathbf{s}'_m, a'; \theta^*) - \max_{a'} \psi(\mathbf{s}'_m, a'; \theta^{(t)}) \right) \cdot \left\| \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \right\|_2 \\
&\leq \gamma \cdot \max_{\mathbf{s},a,a'} \left| \psi(\mathbf{s}'_m, a'; \theta^*) - \psi(\mathbf{s}'_m, a'; \theta^{(t)}) \right| \cdot \left\| \frac{\partial \psi(\theta^{(t)}; \mathcal{X}_m)}{\partial \theta_{\ell,k}} \right\|_2 \\
&\lesssim \gamma \cdot \|\theta^{(t)} - \theta^*\|_2 \cdot \frac{1}{K_\ell} \\
&\leq \frac{\gamma}{K_\ell} \|\theta^{(t)} - \theta^*\|_2.
\end{aligned} \tag{123}$$

□

#### F.3 PROOF OF LEAMMA 8

*Proof of Lemma 8.* From the update rule of  $\mathbf{w}$  in Algorithm 1, we have

$$\begin{aligned}
\mathbf{w}^{(t+1)} - \mathbf{w}^* &= \mathbf{w}^{(t)} - \mathbf{w}^* - \kappa_t \cdot \sum_{m \in \mathcal{D}_t} (\phi_m^\top \mathbf{w}^{(t)} - r_m) \cdot \phi_m \\
&= \mathbf{w}^{(t)} - \mathbf{w}^* - \kappa_t \cdot \sum_{m \in \mathcal{D}_t} (\phi_m^\top \mathbf{w}^{(t)} - \phi_m^\top \mathbf{w}^*) \cdot \phi_m \\
&= \left( \mathbf{I} - \kappa_t \sum_{m \in \mathcal{D}_m} \phi_m^\top \phi_m \right) \cdot (\mathbf{w}^{(t)} - \mathbf{w}^*).
\end{aligned} \tag{124}$$

For any unit vector  $\alpha \in \dim(\mathbf{w})$ , we have

$$\begin{aligned} |\alpha^\top \mathbb{E}_{\mathcal{D}_t} \phi^\top \phi \alpha| &\leq \max_{\|\phi\|_2} |\alpha^\top \phi|^2 \leq \phi_{\max}^2, \\ |\alpha^\top \mathbb{E}_{\mathcal{D}_t} \phi^\top \phi \alpha| &\geq |\alpha^\top \phi_{\min}|^2 \geq 0. \end{aligned} \quad (125)$$

Also, it is easy to verify that  $|\alpha^\top \mathbb{E}_{\mathcal{D}_t} \phi^\top \phi \alpha| = 0$  if only and if  $\phi_m$  are all parallel to each other. As  $\phi_m$  does not parallel to each other, let  $\rho_2 > 0$  denote the minimal eigenvalue of  $\mathbb{E}_{\mathcal{D}_t} \phi^\top \phi$ .

Given  $\phi$  is bounded,  $\phi$  belongs to the sub-Gaussian distribution. Similar to (106), with Chebyshev's inequality, we have

$$\left\| \sum_{m \in \mathcal{D}_m} \phi_m^\top \phi_m - \mathbb{E}_{\mathcal{D}_t} \phi^\top \phi \right\|_2 \leq \sqrt{\frac{d \log q}{|\mathcal{D}_t|}} \quad (126)$$

with probability at least  $1 - d^{-q}$ . Let  $N \geq c_N^{-2} d \log q$ , according to Lemma 1, we have

$$\lambda_{\min} \left( \sum_{m \in \mathcal{D}_m} \phi_m^\top \phi_m \right) \leq \lambda_{\min}(\mathbb{E}_{\mathcal{D}_t} \phi^\top \phi) - c_N \leq \rho_2 - c_N. \quad (127)$$

When we choose  $\kappa_t = \frac{1}{\phi_{\max}}$ , we have

$$\|\mathbf{w}^{(t+1)} - \mathbf{w}^*\|_2 \leq \left(1 - \frac{\rho_2 - c_N}{\phi_{\max}}\right) \cdot \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_2. \quad (128)$$

□

## G PROOF OF LEMMAS IN APPENDIX C

**Lemma 13.** *We have*

$$|Q_i^{\pi_i^*}(s, a) - Q_i^{\pi_j^*}(s, a)| \leq \frac{2\gamma}{1-\gamma} \cdot \max_{s,a} |r_i(s, a) - r_j(s, a)|. \quad (129)$$

*Proof.*  $|Q_i^{\pi_i^*}(s, a) - Q_i^{\pi_j^*}(s, a)|$  can be upper bounded as

$$\begin{aligned} &|Q_i^{\pi_i^*}(s, a) - Q_i^{\pi_j^*}(s, a)| \\ &= \left| r_i + \gamma \cdot \sum_{s'} p_{s,s'}^a Q_i^{\pi_i^*}(s', \pi_i^*(s')) - \left( r_i + \gamma \cdot \sum_{s'} p_{s,s'}^a Q_i^{\pi_j^*}(s', \pi_j^*(s')) \right) \right| \\ &= \gamma \cdot \left| \sum_{s'} p_{s,s'}^a Q_i^{\pi_i^*}(s', \pi_i^*(s')) - \sum_{s'} p_{s,s'}^a Q_i^{\pi_j^*}(s', \pi_j^*(s')) \right| \\ &\leq \gamma \cdot \sum_{s'} p_{s,s'}^a \cdot \left| Q_i^{\pi_i^*}(s', \pi_i^*(s')) - Q_i^{\pi_j^*}(s', \pi_j^*(s')) \right| \\ &\leq \gamma \cdot \sum_{s'} p_{s,s'}^a \cdot \left[ \left| Q_i^{\pi_i^*}(s', \pi_i^*(s')) - Q_j^{\pi_j^*}(s', \pi_j^*(s')) \right| + \left| Q_j^{\pi_j^*}(s', \pi_j^*(s')) - Q_i^{\pi_j^*}(s', \pi_j^*(s')) \right| \right] \\ &= \gamma \cdot \sum_{s'} p_{s,s'}^a \cdot \left[ \left| \max_{a'} Q_i^{\pi_i^*}(s', a') - \max_{a'} Q_j^{\pi_j^*}(s', a') \right| + \left| Q_j^{\pi_j^*}(s', \pi_j^*(s')) - Q_i^{\pi_j^*}(s', \pi_j^*(s')) \right| \right] \\ &\leq \gamma \cdot \sum_{s'} p_{s,s'}^a \cdot \left[ \max_{a'} \left| Q_i^{\pi_i^*}(s', a') - Q_j^{\pi_j^*}(s', a') \right| + \left| Q_j^{\pi_j^*}(s', \pi_j^*(s')) - Q_i^{\pi_j^*}(s', \pi_j^*(s')) \right| \right] \\ &\leq \gamma \cdot \sum_{s'} p_{s,s'}^a \cdot \left[ \max_{s',a'} \left| Q_i^{\pi_i^*}(s', a') - Q_j^{\pi_j^*}(s', a') \right| + \max_{s'} \left| Q_j^{\pi_j^*}(s', \pi_j^*(s')) - Q_i^{\pi_j^*}(s', \pi_j^*(s')) \right| \right] \end{aligned} \quad (130)$$

Let

$$I_5 = \max_{s,a} \left| Q_i^{\pi_i^*}(s, a) - Q_j^{\pi_j^*}(s, a) \right|$$

and

$$I_6 = \max_{\mathbf{s}, a} \left| Q_j^{\pi_j^*}(\mathbf{s}, a) - Q_i^{\pi_i^*}(\mathbf{s}, a) \right| \geq \max_{\mathbf{s}} \left| Q_j^{\pi_j^*}(\mathbf{s}, \pi_j^*(\mathbf{s})) - Q_i^{\pi_i^*}(\mathbf{s}, \pi_j^*(\mathbf{s})) \right|.$$

Then, we have

$$\begin{aligned} I_5 &= \max_{\mathbf{s}, a} \left| r_i + \gamma \cdot \sum_{\mathbf{s}'} p_{\mathbf{s}, \mathbf{s}'}^a \max_{a'} Q_i^{\pi_i^*}(\mathbf{s}', a') - r_j - \gamma \cdot \sum_{\mathbf{s}'} p_{\mathbf{s}, \mathbf{s}'}^a \max_{a'} Q_j^{\pi_j^*}(\mathbf{s}', a') \right| \\ &\leq \max_{\mathbf{s}, a} |r_i(\mathbf{s}, a) - r_j(\mathbf{s}, a)| + \gamma \max_{\mathbf{s}, a} \sum_{\mathbf{s}'} p_{\mathbf{s}, \mathbf{s}'}^a \cdot \max_{a'} |Q_i^{\pi_i^*}(\mathbf{s}', a') - Q_j^{\pi_j^*}(\mathbf{s}', a')| \\ &\leq \max_{\mathbf{s}, a} |r_i(\mathbf{s}, a) - r_j(\mathbf{s}, a)| + \gamma \cdot I_5. \end{aligned} \quad (131)$$

Therefore, we have

$$I_5 \leq \frac{1}{1 - \gamma} \max_{\mathbf{s}, a} |r_i(\mathbf{s}, a) - r_j(\mathbf{s}, a)|. \quad (132)$$

Similar to (131), we have

$$\begin{aligned} I_6 &\leq \max_{\mathbf{s}, a} |r_i(\mathbf{s}, a) - r_j(\mathbf{s}, a)| + \gamma \max_{\mathbf{s}, a} \sum_{\mathbf{s}'} p_{\mathbf{s}, \mathbf{s}'}^a \cdot |Q_j^{\pi_j^*}(\mathbf{s}', \pi_j^*(\mathbf{s}')) - Q_i^{\pi_i^*}(\mathbf{s}', \pi_j^*(\mathbf{s}'))| \\ &\leq \max_{\mathbf{s}, a} |r_i(\mathbf{s}, a) - r_j(\mathbf{s}, a)| + \gamma \cdot I_6. \end{aligned} \quad (133)$$

Therefore, we have

$$I_6 \leq \frac{1}{1 - \gamma} \max_{\mathbf{s}, a} |r_i(\mathbf{s}, a) - r_j(\mathbf{s}, a)|. \quad (134)$$

Therefore, we have

$$|Q_i^{\pi_i^*}(\mathbf{s}, a) - Q_j^{\pi_j^*}(\mathbf{s}, a)| \leq \gamma(I_5 + I_6) \leq \frac{2\gamma}{1 - \gamma} \cdot \max_{\mathbf{s}, a} |r_i(\mathbf{s}, a) - r_j(\mathbf{s}, a)|. \quad (135)$$

□

## H ADDITIONAL PROOF OF THE LEMMAS

### H.1 PROOF OF LEMMA 10

The distance of the second order derivatives of the population risk function  $f(\cdot)$  at point  $\theta$  and  $\theta^*$  can be converted into bounding  $\mathbf{P}_1$ ,  $\mathbf{P}_2$ , which are defined in (137). The major idea in proving  $\mathbf{P}_1$  is to connect the error bound to the angle between  $\theta$  and  $\theta^*$  given  $\mathbf{h}^{(\ell)}$  belongs to the sub-Gaussian distribution.

*Proof of Lemma 10.* From the definition of  $f$  in (30), we have

$$\begin{aligned} \frac{\partial^2 f}{\partial \theta_{\ell, j_1} \partial \theta_{\ell, j_2}}(\theta^*) &= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \mathcal{J}_{\ell, k} \sigma'(\theta_{j_1}^{*T} \mathbf{h}) \cdot \mathcal{J}_{\ell, k} \sigma'(\theta_{j_2}^{*T} \mathbf{h}) \cdot \mathbf{h}^* \mathbf{h}^{*T}, \\ \text{and } \frac{\partial^2 f}{\partial \theta_{\ell, j_1} \partial \theta_{\ell, j_2}}(\theta) &= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \mathcal{J}_{\ell, k}^* \sigma'(\theta_{\ell, j_1}^\top \mathbf{h}) \cdot \mathcal{J}_{\ell, k}^* \sigma'(\theta_{\ell, j_2}^\top \mathbf{h}) \cdot \mathbf{h} \mathbf{h}^\top, \end{aligned} \quad (136)$$

where  $\mathbf{h} = \mathbf{h}^{(\ell)}(\theta)$  and  $\mathbf{h}^* = \mathbf{h}^{(\ell)}(\theta^*)$ .

Then, we have

$$\begin{aligned} &\frac{\partial^2 f}{\partial \theta_{\ell, j_1} \partial \theta_{\ell, j_2}}(\theta^*) - \frac{\partial^2 f}{\partial \theta_{\ell, j_1} \partial \theta_{\ell, j_2}}(\theta) \\ &= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} [\mathcal{J}_{\ell, k}^* \sigma'(\theta_{\ell, j_1}^\top \mathbf{h}^*) \mathcal{J}_{\ell, k}^* \sigma'(\theta_{\ell, j_2}^\top \mathbf{h}^*) \mathbf{h}^* \mathbf{h}^{*T} \\ &\quad - \mathcal{J}_{\ell, k} \sigma'(\theta_{\ell, j_1}^\top \mathbf{h}) \mathcal{J}_{\ell, k} \sigma'(\theta_{\ell, j_2}^\top \mathbf{h}) \mathbf{h} \mathbf{h}^\top] \\ &= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} [\mathcal{J}_{\ell, k}^* \sigma'(\theta_{\ell, j_1}^\top \mathbf{h}^*) (\mathcal{J}_{\ell, k}^* \sigma'(\theta_{\ell, j_2}^\top \mathbf{h}^*) \mathbf{h}^* \mathbf{h}^{*T} - \mathcal{J}_{\ell, k} \sigma'(\theta_{\ell, j_2}^\top \mathbf{h}) \mathbf{h} \mathbf{h}^\top) \\ &\quad + \mathcal{J}_{\ell, k} \sigma'(\theta_{\ell, j_2}^\top \mathbf{h}) (\mathcal{J}_{\ell, k}^* \sigma'(\theta_{\ell, j_1}^\top \mathbf{h}^*) \mathbf{h}^* \mathbf{h}^{*T} - \mathcal{J}_{\ell, k} \sigma'(\theta_{\ell, j_1}^\top \mathbf{h}) \mathbf{h} \mathbf{h}^\top)] \\ &:= \frac{1}{K^2} (\mathbf{P}_1 + \mathbf{P}_2). \end{aligned} \quad (137)$$

For any  $\mathbf{a} \in \mathbb{R}^{K_\ell}$  with  $\|\mathbf{a}\|_2 = 1$ , we have

$$\mathbf{a}^\top \mathbf{P}_1 \mathbf{a} = \mathbb{E}_{\mathbf{x}} \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_1}^{*T} \mathbf{h}^*) \left( \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k} \sigma'(\theta_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right). \quad (138)$$

Then, we have

$$\begin{aligned} |\mathbf{a}^\top \mathbf{P}_1 \mathbf{a}| &= \left| \mathbb{E}_{\mathbf{x}} \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_1}^{*T} \mathbf{h}^*) \left( \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k} \sigma'(\theta_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right) \right| \\ &\leq \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k} \sigma'(\theta_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right| \\ &\leq \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h})^2 \right| \\ &\quad + \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k}^* \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h})^2 - \mathcal{J}_{\ell,k} \sigma'(\theta_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right| \\ &\quad + \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k} \sigma'(\theta_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 - \mathcal{J}_{\ell,k} \sigma'(\theta_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right| \\ &\lesssim \|\theta - \theta^*\|_2 + \|\theta - \theta^*\|_2 + \mathbb{E}_{\mathbf{x}} \left| (\sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}) - \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*)) \cdot (\mathbf{a}^\top \mathbf{h})^2 \right| \\ &\lesssim \|\theta - \theta^*\|_2 + \mathbb{E}_{\mathbf{x}} \left| (\sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}) - \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}^*)) \cdot (\mathbf{a}^\top \mathbf{h})^2 \right|. \end{aligned} \quad (139)$$

Utilizing the Gram-Schmidt process, we can demonstrate the existence of a set of normalized orthonormal vectors denoted as  $\mathcal{B} = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{a}_4^\perp, \dots, \mathbf{a}_d^\perp\} \in \mathbb{R}^d$ . This set forms an orthogonal and normalized basis for  $\mathbb{R}^d$ , wherein the subspace spanned by  $\mathbf{a}, \mathbf{b}, \mathbf{c}$  includes  $\mathbf{a}, \theta_{\ell,j_2}$ , and  $\theta_{\ell,j_2}^*$ . Then, for any  $\mathbf{x} \in \mathbb{R}^d$ , we have a unique  $\mathbf{z} = [z_1, z_2, \dots, z_d]^\top$  such that

$$\mathbf{h} = z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c} + \dots + z_d \mathbf{a}_d^\perp.$$

Because (i)  $\mathbf{a}, \theta_{\ell,j_2}$ , and  $\theta_{\ell,j_2}^*$  belongs to the subspace spanned by vectors  $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$  and (ii)  $\mathbf{a}_4^\perp, \dots, \mathbf{a}_d^\perp$  are orthogonal to  $\mathbf{a}, \mathbf{b}$ , and  $\mathbf{c}$ . Then, we know that

$$\begin{aligned} \theta_{\ell,j_2}^{*T} \mathbf{h} &= \theta_{\ell,j_2}^{*T} (z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c} + \dots + z_d \mathbf{a}_d^\perp) \\ &= z_1 \theta_{\ell,j_2}^{*T} \mathbf{a} + z_2 \theta_{\ell,j_2}^{*T} \mathbf{b} + z_3 \theta_{\ell,j_2}^{*T} \mathbf{c} + \dots + z_d \theta_{\ell,j_2}^{*T} \mathbf{a}_d^\perp \\ &= z_1 \theta_{\ell,j_2}^{*T} \mathbf{a} + z_2 \theta_{\ell,j_2}^{*T} \mathbf{b} + z_3 \theta_{\ell,j_2}^{*T} \mathbf{c} + 0 \\ &= \theta_{\ell,j_2}^{*T} (z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c}) \\ &:= \theta_{\ell,j_2}^{*T} \tilde{\mathbf{h}}. \end{aligned} \quad (140)$$

where  $\tilde{\mathbf{h}} = z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c}$ . Similar to (140), we have  $\theta_{\ell,j_2}^\top \mathbf{h} = \theta_{\ell,j_2}^\top \tilde{\mathbf{h}}$  and  $\mathbf{a}^\top \mathbf{h} = \mathbf{a}^\top \tilde{\mathbf{h}}$ .

Then, we define  $I_4$  as

$$\begin{aligned} I_4 &:= \mathbb{E}_{\mathbf{h}} \left| (\sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h}) - \sigma'(\theta_{\ell,j_2}^\top \mathbf{h})) \cdot (\mathbf{a}^\top \mathbf{h}) \right| \\ &= \int_{\mathcal{R}_{\mathbf{h}}} |\sigma'(\theta_{\ell,j_2}^\top \mathbf{h}) - \sigma'(\theta_{\ell,j_2}^{*T} \mathbf{h})| \cdot |\mathbf{a}^\top \mathbf{h}|^2 \cdot f_H(\mathbf{h}) d\mathbf{h} \\ &= \int_{\mathcal{R}_{\mathbf{z}}} |\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{h}}) - \sigma'(\theta_{\ell,j_2}^{*T} \tilde{\mathbf{h}})| \cdot |\mathbf{a}^\top \tilde{\mathbf{h}}|^2 \cdot f_Z(\mathbf{z}) \cdot |\mathbf{J}_{\mathbf{h}}(\mathbf{z})| d\mathbf{z} \end{aligned} \quad (141)$$

where  $|\mathbf{J}_{\mathbf{h}}(\mathbf{z})|$  is the determinant of the Jacobian matrix  $\frac{\partial \mathbf{h}}{\partial \mathbf{z}}$ . Since  $\mathbf{z}$  is a representation of  $\mathbf{h}$  based on an orthogonal and normalized basis, we have  $|\mathbf{J}_{\mathbf{h}}(\mathbf{z})| = 1$ . According to (140),  $I_4$  can be rewritten as

$$\begin{aligned} I_4 &= \int_{\mathcal{R}_{\mathbf{z}}} |\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{h}}) - \sigma'(\theta_{\ell,j_2}^{*T} \tilde{\mathbf{h}})| \cdot |\mathbf{a}^\top \tilde{\mathbf{h}}|^2 \cdot f_Z(\mathbf{z}) d\mathbf{z} \\ &= \int_{\mathcal{R}_{\mathbf{z}}} |\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{h}}) - \sigma'(\theta_{\ell,j_2}^{*T} \tilde{\mathbf{h}})| \cdot |\mathbf{a}^\top \tilde{\mathbf{h}}|^2 \cdot f_Z(z_1, z_2, z_3) dz_1 dz_2 dz_3 \end{aligned} \quad (142)$$

where in the last equality we abuse  $f_Z(z_1, z_2, z_3)$  to represent the probability density function of  $(z_1, z_2, z_3)$  defined in region  $\mathcal{R}_{\mathbf{z}}$ .

Next, we show that  $\mathbf{z}$  is rotational invariant over  $\mathcal{R}_z$ . Let  $\mathbf{R} = [\mathbf{a} \ \mathbf{b} \ \mathbf{c} \ \cdots \ \mathbf{a}_d^\perp]$ , we have  $\mathbf{h} = \mathbf{R}\mathbf{z}$ . For any  $\mathbf{z}^{(1)}$  and  $\mathbf{z}^{(2)}$  with  $\|\mathbf{z}^{(1)}\|_2 = \|\mathbf{z}^{(2)}\|_2$ . We define  $\mathbf{h}^{(1)} = \mathbf{R}\mathbf{z}^{(1)}$  and  $\mathbf{h}^{(2)} = \mathbf{R}\mathbf{z}^{(2)}$ . Since  $\mathbf{x}$  is rotational invariant and  $\|\mathbf{h}^{(1)}\|_2 = \|\mathbf{h}^{(2)}\|_2 = \|\mathbf{z}^{(1)}\|_2 = \|\mathbf{z}^{(2)}\|_2$ , then we know  $\mathbf{h}^{(1)}$  and  $\mathbf{h}^{(2)}$  has the same distribution density. Then,  $\mathbf{z}^{(1)}$  and  $\mathbf{z}^{(2)}$  has the same distribution density as well. Therefore,  $\mathbf{z}$  is rotational invariant over  $\mathcal{R}_z$ .

Then, we consider spherical coordinates with  $z_1 = R \cos \sigma_1$ ,  $z_2 = R \sin \sigma_1 \sin \sigma_2$ ,  $z_3 = R \sin \sigma_1 \cos \sigma_2$ . Hence, we have

$$I_4 = \int |\sigma'(\theta_{\ell, j_2}^\top \tilde{\mathbf{h}}) - \sigma'(\theta_{\ell, j_2}^{\star \top} \tilde{\mathbf{h}})| \cdot |R \cos \sigma_1|^2 \cdot f_Z(R, \sigma_1, \sigma_2) \cdot R^2 \sin \sigma_1 \cdot dR d\sigma_1 d\sigma_2. \quad (143)$$

Since  $\mathbf{z}$  is rotational invariant, we have that

$$f_Z(R, \sigma_1, \sigma_2) = f_Z(R). \quad (144)$$

Then, we have

$$\begin{aligned} I_4 &= \int |\sigma'(\theta_{\ell, j_2}^\top (\tilde{\mathbf{h}}/R)) - \sigma'(\theta_{\ell, j_2}^{\star \top} (\tilde{\mathbf{h}}/R))| \cdot |R \cos \sigma_1|^2 \cdot f_Z(R) R^2 \sin \sigma_1 dR d\sigma_1 d\sigma_2 \\ &= \int_0^\infty R^4 f_Z(R) dR \int_0^{\psi_1(R)} \int_0^{\psi_2(R)} |\cos \sigma_1|^2 \cdot \sin \sigma_1 \\ &\quad \cdot |\sigma'(\theta_{\ell, j_2}^\top (\tilde{\mathbf{h}}/R)) - \sigma'(\theta_{\ell, j_2}^{\star \top} (\tilde{\mathbf{h}}/R))| d\sigma_1 d\sigma_2 \\ &\leq \int_0^\infty R^4 f_Z(R) dR \int_0^\pi \int_0^{2\pi} \sin \sigma_1 \cdot |\sigma'(\theta_{\ell, j_2}^\top \bar{\mathbf{x}}) - \sigma'(\theta_{\ell, j_2}^{\star \top} \bar{\mathbf{x}})| d\sigma_1 d\sigma_2, \end{aligned} \quad (145)$$

where the first equality holds because  $\sigma'(\theta_{\ell, j_2}^\top \mathbf{h})$  only depends on the direction of  $\mathbf{h}$ , and  $\bar{\mathbf{x}} := \mathbf{h}/R = (\cos \sigma_1, \sin \sigma_1 \sin \sigma_2, \sin \sigma_1 \cos \sigma_2)$  in the last inequality.

Because  $\mathbf{z}$  belongs to the sub-Gaussian distribution, we have  $F_z(R) \geq 1 - 2e^{-\frac{R^2}{\sigma^2}}$  for some constant  $\sigma > 0$ . Then, the integration of  $R$  can be represented as

$$\begin{aligned} \int_0^\infty R^4 f_Z(R) dR &= \int_0^\infty R^4 d(1 - F_z(R)) \\ &\leq \int_0^\infty 4R^3 (1 - F_z(R)) dR \\ &\leq \int_0^\infty 8R^3 e^{-\frac{R^2}{\sigma^2}} dR \\ &\leq \frac{32}{\sqrt{2\pi}} \sigma \int_0^\infty R^2 e^{-\frac{R^2}{\sigma^2}} dR \\ &= 32\sigma^2 \int_0^\infty R^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{R^2}{\sigma^2}} dR, \end{aligned} \quad (146)$$

where the last inequality comes from the calculation that

$$\begin{aligned} \int_0^\infty 2R^2 e^{-\frac{R^2}{\sigma^2}} dR &= \sqrt{2\pi}\sigma^3, \\ \int_0^\infty 2R^3 e^{-\frac{R^2}{\sigma^2}} dR &= 4\sigma^4. \end{aligned} \quad (147)$$

Then, we define  $\tilde{\mathbf{x}} \in \mathbb{R}^{K_\ell}$  belongs to Gaussian distribution as  $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ . Therefore, we have

$$\begin{aligned} I_4 &\leq 32\sigma^2 \cdot \int_0^\infty R^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{R^2}{\sigma^2}} dR \int_0^\pi \int_0^{2\pi} \sin \sigma_1 \cdot |\sigma'(\theta_{\ell, j_2}^\top \bar{\mathbf{x}}) - \sigma'(\theta_{\ell, j_2}^{\star \top} \bar{\mathbf{x}})| d\sigma_1 d\sigma_2 \\ &= 32\sigma^2 \cdot \mathbb{E}_{z_1, z_2, z_3} |\sigma'(\theta_{\ell, j_2}^\top \tilde{\mathbf{x}}) - \sigma'(\theta_{\ell, j_2}^{\star \top} \tilde{\mathbf{x}})| \\ &\approx \mathbb{E}_{\tilde{\mathbf{x}}} |\sigma'(\theta_{\ell, j_2}^\top \tilde{\mathbf{x}}) - \sigma'(\theta_{\ell, j_2}^{\star \top} \tilde{\mathbf{x}})|, \end{aligned} \quad (148)$$

where  $\tilde{\mathbf{x}}$  belongs to Gaussian distribution.

Therefore, the inequality bound over a sub-Gaussian distribution is bounded by the one over a Gaussian distribution. In the following contexts, we provide the upper bound of  $\mathbb{E}_{\tilde{\mathbf{x}}} |\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \sigma'(\theta_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})|$ .

Define a set  $\mathcal{A}_1 = \{\mathbf{x} | (\theta_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})(\theta_{\ell,j_2}^\top \tilde{\mathbf{x}}) < 0\}$ . If  $\tilde{\mathbf{x}} \in \mathcal{A}_1$ , then  $\theta_{\ell,j_2}^{*\top} \tilde{\mathbf{x}}$  and  $\theta_{\ell,j_2}^\top \tilde{\mathbf{x}}$  have different signs, which means the value of  $\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{x}})$  and  $\sigma'(\theta_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})$  are different. This is equivalent to say that

$$|\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \sigma'(\theta_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})| = \begin{cases} 1, & \text{if } \tilde{\mathbf{x}} \in \mathcal{A}_1 \\ 0, & \text{if } \tilde{\mathbf{x}} \in \mathcal{A}_1^c \end{cases}. \quad (149)$$

Moreover, if  $\tilde{\mathbf{x}} \in \mathcal{A}_1$ , then we have

$$|\theta_{\ell,j_2}^{*T} \tilde{\mathbf{x}}| \leq |\theta_{\ell,j_2}^{*T} \tilde{\mathbf{x}} - \theta_{\ell,j_2}^\top \tilde{\mathbf{x}}| \leq \|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2 \cdot \|\tilde{\mathbf{x}}\|_2. \quad (150)$$

Let us define a set  $\mathcal{A}_2$  such that

$$\begin{aligned} \mathcal{A}_2 &= \left\{ \tilde{\mathbf{x}} \mid \frac{|\theta_{\ell,j_2}^{*T} \tilde{\mathbf{x}}|}{\|\theta_{\ell,j_2}^*\|_2 \|\tilde{\mathbf{x}}\|_2} \leq \frac{\|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2}{\|\theta_{\ell,j_2}^*\|_2} \right\} \\ &= \left\{ \tilde{\mathbf{x}} \mid \left| \cos \theta_{\tilde{\mathbf{x}}, \theta_{\ell,j_2}^*} \right| \leq \frac{\|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2}{\|\theta_{\ell,j_2}^*\|_2} \right\}. \end{aligned} \quad (151)$$

Hence, we have that

$$\begin{aligned} \mathbb{E}_{\tilde{\mathbf{x}}} |\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \sigma'(\theta_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})|^2 &= \mathbb{E}_{\tilde{\mathbf{x}}} |\sigma'(\theta_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \sigma'(\theta_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})| \\ &= \text{Prob}(\tilde{\mathbf{x}} \in \mathcal{A}_1) \\ &\leq \text{Prob}(\tilde{\mathbf{x}} \in \mathcal{A}_2). \end{aligned} \quad (152)$$

Since  $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \|\mathbf{a}\|_2^2 \mathbf{I})$ ,  $\theta_{\tilde{\mathbf{x}}, \theta_{\ell,j_2}^*}$  belongs to the uniform distribution on  $[-\pi, \pi]$ , we have

$$\begin{aligned} \text{Prob}(\tilde{\mathbf{x}} \in \mathcal{A}_2) &= \frac{\pi - \arccos \frac{\|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2}{\|\theta_{\ell,j_2}^*\|_2}}{\pi} \leq \frac{1}{\pi} \tan(\pi - \arccos \frac{\|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2}{\|\theta_{\ell,j_2}^*\|_2}) \\ &= \frac{1}{\pi} \cot(\arccos \frac{\|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2}{\|\theta_{\ell,j_2}^*\|_2}) \\ &\leq \frac{2}{\pi} \frac{\|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2}{\|\theta_{\ell,j_2}^*\|_2} \\ &\leq \|\theta_{\ell,j_2}^* - \theta_{\ell,j_2}\|_2 \end{aligned} \quad (153)$$

Hence, (145) and (153) suggest that

$$\begin{aligned} I_4 &\lesssim \|\theta_i - \theta_i^*\|_2 \cdot \|\mathbf{a}\|_2^2, \\ \text{and } \|\mathbf{P}_1\|_2 &\leq \|\theta - \theta^*\|_2 + I_4 \lesssim \|\theta - \theta^*\|_2, \end{aligned} \quad (154)$$

The same bound that is shown in (154) holds for  $\mathbf{P}_2$  as well.

Therefore, we have

$$\begin{aligned} \|\nabla_\ell^2 f(\theta^*) - \nabla_\ell^2 f(\theta)\|_2 &= \max_{\|\boldsymbol{\alpha}\|_2 \leq 1} \left| \boldsymbol{\alpha}^\top (\nabla_\ell^2 f(\theta^*) - \nabla_\ell^2 f(\theta)) \boldsymbol{\alpha} \right| \\ &\leq \frac{1}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \|\mathbf{P}_1 + \mathbf{P}_2\|_2 \cdot \|\boldsymbol{\alpha}_{j_1}\|_2 \cdot \|\boldsymbol{\alpha}_{j_2}\|_2 \\ &\lesssim \frac{1}{K^2} \cdot \sum_{j_1=1}^K \sum_{j_2=1}^K \|\theta - \theta^*\|_2 \cdot \|\boldsymbol{\alpha}_{j_1}\|_2 \|\boldsymbol{\alpha}_{j_2}\|_2 \\ &\lesssim \frac{1}{K^2} \cdot \sum_{j_1=1}^K \sum_{j_2=1}^K \|\theta - \theta^*\|_2 \cdot \left( \frac{\|\boldsymbol{\alpha}_{j_1}\|_2^2 + \|\boldsymbol{\alpha}_{j_2}\|_2^2}{2} \right) \\ &\lesssim \frac{1}{K} \cdot \|\theta^* - \theta\|_2, \end{aligned} \quad (155)$$

where  $\boldsymbol{\alpha} \in \mathbb{R}^{Kd}$  and  $\boldsymbol{\alpha}_j \in \mathbb{R}^{K\ell}$  with  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1^\top, \boldsymbol{\alpha}_2^\top, \dots, \boldsymbol{\alpha}_K^\top]^\top$ .  $\square$

## H.2 PROOF OF LEMMA 11

We aim to prove that  $\int_{\mathcal{R}} \left( \sum_{j=1}^K \alpha^\top \mathbf{h} \sigma'(\theta_{\ell,j}^\top \mathbf{h}) \right)^2 p_H(\mathbf{h}) \cdot d\mathbf{h}$  is strictly greater than zero for any  $\alpha$ . Therefore, the  $\rho_1$  in (6) is strictly greater than zero. The proof is inspired by Theorem 3.1 in (Du et al., 2019). It is obviously that  $(\sum_{j=1}^K \alpha^\top \mathbf{h} \sigma'(\theta_{\ell,j}^\top \mathbf{h}))^2$  is greater or equal to zero. Given  $(\sum_{j=1}^K \alpha^\top \mathbf{h} \sigma'(\theta_{\ell,j}^\top \mathbf{h}))^2$  is continuous, we only need to show that  $\alpha$  such that  $\sum_{j=1}^K \alpha^\top \mathbf{h} \sigma'(\theta_{\ell,j}^\top \mathbf{h}) \neq 0$  for any  $\alpha$ , namely,  $\{\mathbf{h} \sigma'(\theta_{\ell,j}^\top \mathbf{h})\}_{j=1}^K$  are linear independent.

*Proof of Lemma 11.* Let  $\mathcal{H}$  be a Hilbert space on  $\mathbb{R}^{K_\ell}$ , and the inner product of  $\mathcal{H}$  is defined as

$$\langle f, g \rangle = \int_{\mathcal{R}} f(\mathbf{h})^\top g(\mathbf{h}) f_H(\mathbf{h}) \cdot d\mathbf{h}, \quad \forall f, g \in \mathcal{H}, \quad (156)$$

where the Lebesgue measure of  $\mathcal{R}$  over  $\mathbb{R}^{K_\ell}$  is non-zero. Instead of directly proving  $\int_{\mathcal{R}} \left( \sum_{k=1}^K \alpha^\top \mathbf{h} \sigma'(\theta_k^\top \mathbf{h}) \right)^2 f_H(\mathbf{h}) \cdot d\mathbf{h} > 0$  for any  $\alpha$ , we note that it is sufficient to prove that  $\{\mathbf{h} \sigma'(\theta_k^\top \mathbf{h})\}_{k \in [K]}$  are linear independent over the Hilbert space  $\mathcal{H}$ . Namely, if  $\{\mathbf{h} \sigma'(\theta_k^\top \mathbf{h})\}_{k \in [K]}$  are linear independent, we have

$$\alpha^\top \mathbf{h} \sigma'(\theta_k^\top \mathbf{h}) \neq 0 \quad \text{almost everywhere.} \quad (157)$$

Therefore, we can know that  $\int_{\mathcal{R}} \left( \sum_{j=1}^K \alpha^\top \mathbf{h} \sigma'(\theta_{\ell,j}^\top \mathbf{h}) \right)^2 p_H(\mathbf{h}) \cdot d\mathbf{h}$  is strictly greater than zero.

Next, we provide the whole proof for that  $\{\mathbf{h} \sigma'(\theta_k^\top \mathbf{h})\}_{k \in [K]}$  are linear independent over the Hilbert space  $\mathcal{H}$ .

We define a group of functions  $\{\psi_j(\mathbf{h})\}_{j=1}^K$ , where  $\psi_j(\mathbf{h}) = \mathbf{h} \sigma'(\theta_j^\top \mathbf{h})$ . From the assumption in Lemma 11, we can justify that  $\mathbb{E}_{\mathbf{h} \sim \mathcal{D}} |\psi_j(\mathbf{h})|^2 \leq \mathbb{E}_{\mathbf{h} \sim \mathcal{D}} |\mathbf{h}|^2 < \infty$ .

Let  $\mathcal{X}_i = \{\mathbf{h} \mid \theta_i^\top \mathbf{h} = 0\}$  for any  $i \in [K]$ . For any fixed  $k$ , we can justify that  $\mathcal{X}_k$  cannot be covered by other sets  $\{\mathcal{X}_j\}_{j \neq k}$  as long as  $\theta_k$  does not parallel to any other weights  $\theta_j$  with  $j \neq k$ . Namely,  $\mathcal{X}_k \not\subset \cup_{j \neq k} \mathcal{X}_j$ . The idea of proving the claim above is that the intersection of  $\mathcal{X}_j$  and  $\mathcal{X}_k$  is only a hyperplane in  $\mathcal{X}_k$ . The union of finite many hyperplanes is not even a measurable space and thus cannot cover the original space. Formally, we provide the formal proof for this claim as follows.

Let  $\lambda$  be the Lebesgue measure on  $\mathcal{X}_k$ , then  $\lambda(\mathcal{X}_k) > 0$ . When  $\theta_j$  does not parallel to  $\theta_k$ ,  $\mathcal{X}_k \cap \mathcal{X}_j$  is only a hyperplane in  $\mathcal{X}_k$  for  $j \neq k$ . Hence, we have  $\lambda(\mathcal{X}_j \cap \mathcal{X}_k) = 0$ . Next, we have

$$\lambda(\mathcal{X}_k \cap (\cup_{j \neq k} \mathcal{X}_j)) \leq \sum_{j \neq k} \lambda(\mathcal{X}_k \cap \mathcal{X}_j) = 0. \quad (158)$$

Therefore, we have

$$\lambda(\mathcal{X}_k / (\cup_{j \neq k} \mathcal{X}_j)) = \lambda(\mathcal{X}_k) - \lambda(\mathcal{X}_k \cap (\cup_{j \neq k} \mathcal{X}_j)) = \lambda(\mathcal{X}_k) > 0. \quad (159)$$

Therefore, we have  $\mathcal{X}_k / (\cup_{j \neq k} \mathcal{X}_j)$  is not empty, which means that  $\mathcal{X}_k \not\subset \cup_{j \neq k} \mathcal{X}_j$ .

Next, Since  $\mathcal{X}_k / (\cup_{j \neq k} \mathcal{X}_j)$  is not an empty set, there exists a point  $\mathbf{z}_k \in \mathcal{X}_k / (\cup_{j \neq k} \mathcal{X}_j)$  and  $r_0 > 0$  such that

$$\mathcal{B}(\mathbf{z}_k, r) \cap \mathcal{D}_j = \emptyset \quad \text{with} \quad \forall r \leq r_0 \text{ and } j \neq k, \quad (160)$$

where  $\mathcal{B}(\mathbf{z}_k, r)$  stands for a ball centered at  $\mathbf{z}_k$  with a radius of  $r$ . Then, we divide  $\mathcal{B}(\mathbf{z}_k, r)$  into two disjoint subsets such that

$$\begin{aligned} \mathcal{B}_r^+ &= \mathcal{B}(\mathbf{z}_k, r) \cap \{\mathbf{h} \mid \theta_k^\top \mathbf{h} > 0\}, \\ \mathcal{B}_r^- &= \mathcal{B}(\mathbf{z}_k, r) \cap \{\mathbf{h} \mid \theta_k^\top \mathbf{h} < 0\}. \end{aligned} \quad (161)$$

Because  $\mathbf{z}_k$  is a boundary point of  $\{\mathbf{h} \mid \theta_k^\top \mathbf{h} = 0\}$ , both  $\mathcal{B}_r^+$  and  $\mathcal{B}_r^-$  are non-empty.

Note that  $\psi_j(\mathbf{h})$  is continuous at any point except for the ones in  $\mathcal{X}_j$ . Then, for any  $j \neq k$ , we know that  $\sigma_j(\theta_k^\top \mathbf{h})$  is continuous at point  $\mathbf{z}_k$  since  $\mathbf{z}_k \notin \mathcal{X}_j$ . Hence, it is easy to verify that

$$\lim_{r \rightarrow 0^+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \psi_k(\mathbf{h}) d\mathbf{h} = \lim_{r \rightarrow 0^-} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \psi_k(\mathbf{h}) d\mathbf{h} = \psi_k(\mathbf{z}_k). \quad (162)$$

While for  $\psi_k$ , we know that  $\psi_k(\mathbf{h}) \equiv 0$  for  $\mathbf{h} \in \mathcal{B}_r^-$ , (ii)  $\psi_k(\mathbf{h}) = \mathbf{h}$  for  $\mathbf{h} \in \mathcal{B}_r^+$ . Hence, it is easy to verify that

$$\begin{aligned} \lim_{r \rightarrow 0_+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \psi_k(\mathbf{h}) d\mathbf{h} &= \mathbf{z}_k \\ \lim_{r \rightarrow 0_-} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \psi_k(\mathbf{h}) d\mathbf{h} &= 0. \end{aligned} \quad (163)$$

Now let us proof that  $\{\psi_j\}_{j=1}^K$  are linear independent by contradiction. Suppose  $\{\psi_j\}_{j=1}^K$  are linear dependent, we have

$$\sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) \equiv 0, \quad \forall \mathbf{h}. \quad (164)$$

Then, we have

$$\begin{aligned} \lim_{r \rightarrow 0_+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} &= 0 \\ \lim_{r \rightarrow 0_-} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} &= 0 \end{aligned} \quad (165)$$

Then, we have

$$\begin{aligned} 0 &= \lim_{r \rightarrow 0_+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} - \lim_{r \rightarrow 0_+} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} \\ &= \alpha_k \mathbf{z}_k \end{aligned} \quad (166)$$

where the last equality comes from (162) and (163).

Note that  $\mathbf{z}_k$  cannot be  $\mathbf{0}$  because  $\mathbf{z}_k \notin \mathcal{X}_j$ . Therefore, we have  $\alpha_k = 0$ . Similarly to (166), we can obtain that  $\alpha_j = 0$  by define  $\mathbf{z}_j$  following the definition of  $\mathbf{z}_k$  for any  $j \in [K]$ . Then, we know that (164) holds if and only if  $\alpha = \mathbf{0}$ , which contradicts the assumption that  $\{\psi_j\}_{j=1}^K$  are linear dependent.

In conclusion, we know that  $\{\psi_j\}_{j=1}^K$  are linear independent, and  $\int_{\mathcal{R}} \left( \sum_{j=1}^K \alpha^\top \mathbf{h} \sigma'(\theta_{\ell,j}^\top \mathbf{h}) \right)^2 p_H(\mathbf{h}) \cdot d\mathbf{h}$  is strictly greater than zero.  $\square$

### H.3 PROOF OF LEMMA 12

*Proof of Lemma 12.* From the definition of (37), we have

$$\begin{aligned} &\|\mathbf{h}^{(\ell)}(\theta) - \mathbf{h}^{(\ell)}(\theta^*)\|_2 \\ &= \|\sigma(\theta_{\ell-1}^\top \mathbf{h}^{(\ell-1)}(\theta)) - \sigma(\theta_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\theta^*))\|_2 \\ &= \|\sigma(\theta_{\ell-1}^\top \mathbf{h}^{(\ell-1)}(\theta)) - \sigma(\theta_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\theta)) + \sigma(\theta_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\theta)) - \sigma(\theta_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\theta^*))\|_2 \\ &\leq \|\sigma(\theta_{\ell-1}^\top \mathbf{h}^{(\ell-1)}(\theta)) - \sigma(\theta_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\theta))\|_2 + \|\sigma(\theta_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\theta)) - \sigma(\theta_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\theta^*))\|_2 \\ &\leq \|\theta_{\ell-1} - \theta_{\ell-1}^*\|_2 \cdot \|\mathbf{h}^{(\ell-1)}(\theta)\|_2 + \|\mathbf{h}^{(\ell-1)}(\theta) - \mathbf{h}^{(\ell-1)}(\theta^*)\|_2. \end{aligned} \quad (167)$$

With the assumption in the Lemma 12 such that  $\theta$  is close enough to  $\theta^*$ , we have

$$\|\theta_i\|_2 \leq \|\theta_i^*\|_2 + \|\theta_i - \theta_i^*\|_2 \lesssim 1. \quad (168)$$

Therefore, we have

$$\|\mathbf{h}^{(i)}(\theta)\|_2 \leq \|\theta_i\|_2 \cdots \|\theta_1\|_2 \cdot \|\mathbf{x}\|_2 \lesssim \|\mathbf{x}\|_2. \quad (169)$$



Then, we have

$$\begin{aligned}
& \|\mathbf{h}^{(\ell)}(\theta) - \mathbf{h}^{(\ell)}(\theta^*)\|_2 \\
& \leq \|\theta_{\ell-1} - \theta_{\ell-1}^*\|_2 \cdot \|\mathbf{x}\|_2 + \|\mathbf{h}^{(\ell-1)}(\theta) - \mathbf{h}^{(\ell-1)}(\theta^*)\|_2 \\
& \leq \sum_{i=1}^{\ell-1} \|\theta_i - \theta_i^*\|_2 \cdot \|\mathbf{x}\|_2 + \|\mathbf{h}^{(1)}(\theta) - \mathbf{h}^{(1)}(\theta^*)\|_2 \\
& = \sum_{i=1}^{\ell-1} \|\theta_i - \theta_i^*\|_2 \cdot \|\mathbf{x}\|_2 + \|\mathbf{x} - \mathbf{x}\|_2 \\
& = \sum_{i=1}^{\ell-1} \|\theta_i - \theta_i^*\|_2 \cdot \|\mathbf{h}^{(i-1)}(\theta)\|_2 \\
& \leq \|\theta - \theta^*\|_2 \cdot \|\mathbf{x}\|_2,
\end{aligned} \tag{170}$$

which completes the proof.  $\square$