

Supplementary Materials for:

On the Convergence and Sample Complexity Analysis of Deep Q-Network with Epsilon-Greedy Exploration

The structure of the appendix mainly follows the roadmap of the proof described in Section 4.4.

In Appendix A, we define the characterizable population risk function in (31) to approximate the objective function. Also, some notations to simplify the analysis are introduced in Appendix A, and we recommend the readers to refer to Table 3 for the major notations used in the proofs.

In Appendix B, we provide the proof for Lemma 1 and Theorem 1 following the steps as (1) Characterization of the local convex region of population risk function (Lemma 2), (2) Characterization of the distance between the population risk function and the objective function (Lemma 3), (3) Characterization of the convergence of two consecutive iterations $\mathbf{W}^{(t,m+1)}$ and $\mathbf{W}^{(t,m)}$, and (4) Mathematical induction over the t and m to obtain the error bound between the convergent point $\mathbf{W}^{(T,0)}$ and the desired point \mathbf{W}^* .

In Appendix C, we provide the preliminary lemmas and the whole proof for Lemma 2, which characterizes the local convex region of the non-convex population risk function.

In Appendix D, we provide the preliminary lemmas and the whole proof for Lemma 3, which characterizes the difference of g_t and the gradient descent of defined population risk function in (31).

In Appendix E, we provide the proofs for the preliminary lemmas in proving Lemmas 2 and 3.

Before moving to the details, we provide an overview of the techniques in the proofs.

(P1.) The local convex region near \mathbf{W}^* . To characterize the local convex region, we first bound the Hessian matrix of the defined population risk function in (31) at \mathbf{W}^* . Then, we derive the changes in the Hessian matrix when the neuron weights move around the \mathbf{W}^* . Specifically, we prove that when neuron weights \mathbf{W} are not far away \mathbf{W}^* , then the Hessian matrix in this region is always positive-definite, indicating that a local convex region near \mathbf{W}^* . [90] considers the one-hidden-layer neural network, and the lower bound of the Hessian matrix only holds for Gaussian input. Instead, in this paper, we consider multi-layer cases and need to derive a lower bound for the Hessian matrix for all the layers. Instead, the input of the intermediate layer cannot be proved to be Gaussian but belong to sub-Gaussian distribution. Therefore, we built the proof for the lower bound of the Hessian matrix when the input belongs to the sub-Gaussian distribution. Compared with Gaussian input, Sub-Gaussian does not have a closed form of the probability density function. Instead of directly calculating the lower bound, we convert the problem into proving a series of functions are linearly independent over a Hilbert space (see Lemma 7 and the proof in Appendix E). Instead of directly calculating the distance of the population risk function in different points, we characterize a Gaussian variable such that the distance over the sub-Gaussian distribution can be upper bounded by the one over the Gaussian variable (see Lemma 6 and the proof in Appendix E).

(P2.) The difference between the gradient g_t and the population risk function. With the local convex region of the population risk function, we can characterize the convergence of the population risk function. With Lemma 3, we can prove that the distance between the population risk function and g_t is small enough, the behaviors of the iterations via g_t can be described by the ones in the population risk function with some additional error terms. Compared with the proof in [90], We need to address the extension from supervised learning settings to Q learning settings and the extension from the one-hidden-layer neural networks to the multi-layer neural networks. First, similar to challenges in (P1), we provide a new concentration bound to characterize the distance between the two functions for the intermediate layers (see I_1 in the proof of Lemma 3). Second, the distance between the two functions has an additional error term due to the inconsistency of the label defined in (31) and (8) (see I_2 in the proof of Lemma 3). Third, we need to develop a new concentration bound to characterize the error term caused by the distribution shift when training samples are collected by ε -greedy policy (see I_3 in the proof of Lemma 3).

(P3.) The convergence analysis of Algorithm 1. When the initialization is not far away from \mathbf{W}^* , the initialization lies in the local convex region of \mathbf{W}^* for the population risk function. When we have enough samples N and a large enough ε_t , we can guarantee that the distance between the g_t and the gradient of the population risk function is small enough such that the iterations following g_t converges to a point nearby \mathbf{W}^* as well. However, if ε_t is too large, the convergent point nearby \mathbf{W}^* can be even worse than the initial point. To avoid this issue, we have an upper bound for selecting ε_t , and the upper bound decreases as $\|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|$ decreases over t . Therefore, we build the convergence analysis of Algorithm 1.

A Definitions and Notations

In this section, we implement the details of algorithms described in Algorithm 1, and some important notations are defined to simplify the presentation of the proof.

A.1 Definition of the Empirical Risk Function and Its Corresponding Notations

Recall that the goal of Q -learning is to find the Q^* -function to minimize (6). Therefore, we have

$$Q^*(s, a) = r(s, a) + \gamma \cdot \mathbb{E}_{s'|s, a} \max_{a' \in \mathcal{A}} Q^*(s', a') \quad \text{for } (s, a) \sim \mu^*. \quad (29)$$

Since \mathbf{W}^* is the global minimal to (6), we have

$$Q(\mathbf{W}^*; s, a) = r(s, a) + \gamma \cdot \mathbb{E}_{s'|s, a} \max_{a' \in \mathcal{A}} Q(\mathbf{W}^*; s', a'). \quad (30)$$

Therefore, the population risk function is defined as

$$\begin{aligned} f(\mathbf{W}) &= \mathbb{E}_{(s, a) \sim \mu^*} [Q(\mathbf{W}; s, a) - r(s, a) - \gamma \cdot \mathbb{E}_{s'|s, a} \max_{a' \in \mathcal{A}} Q(\mathbf{W}^*; s', a')]^2 \\ &= \mathbb{E}_{(s, a) \sim \mu^*} [Q(\mathbf{W}; s, a) - Q(\mathbf{W}^*; s, a)]^2, \end{aligned} \quad (31)$$

where μ^* is the distribution of the sampled data following the optimal policy π^* .

The gradient of the (31) is

$$\begin{aligned} \nabla_{\mathbf{W}} f(\mathbf{W}) &= \mathbb{E}_{\mathbf{x} \sim \mu^*} (Q(\mathbf{W}; \mathbf{x}) - r(\mathbf{x}) - \gamma \cdot \mathbb{E}_{s' \sim p_{s, s'}^a} \max_{a' \in \mathcal{A}} Q^*(s', a')) \cdot \nabla_{\mathbf{W}} Q(\mathbf{W}; \mathbf{x}) \\ &= \mathbb{E}_{\mathbf{x} \sim \mu^*, s' \sim p_{s, s'}^a} (Q(\mathbf{W}; \mathbf{x}) - r(\mathbf{x}) - \gamma \cdot \max_{a' \in \mathcal{A}} Q(\mathbf{W}^*; s', a')) \cdot \nabla_{\mathbf{W}} Q(\mathbf{W}; \mathbf{x}). \end{aligned} \quad (32)$$

As \mathbf{W}^* is one of the ground truths to $f(\mathbf{W})$, i.e., $f(\mathbf{W}^*)$ achieves the minimum value as $f(\mathbf{W}^*) = 0 \leq f(\mathbf{W})$ for any other \mathbf{W} . Given f is a smooth function, we have the gradient of f with respect to any \mathbf{W}_ℓ at the ground truth \mathbf{W}^* equals to zero, namely,

$$\nabla_\ell f(\mathbf{W}^*) := \nabla_{\mathbf{W}_\ell} f(\mathbf{W}^*) = \mathbf{0}, \quad \forall \ell \in [L]. \quad (33)$$

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 $\boldsymbol{\alpha}_j \in \mathbb{R}_{\ell-1}^{K_0}$ ($K_0 = d$). Therefore, we have

$$\begin{aligned} \|\nabla_\ell h\|_2 &= \max_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}^\top \nabla_\ell h\|_2 = \max_{\boldsymbol{\alpha}} \left| \sum_{j=1}^K \boldsymbol{\alpha}_j^\top \frac{\partial h}{\partial \mathbf{w}_{\ell, j}} \right|, \\ \|\nabla_\ell^2 h\|_2 &= \max_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}^\top \nabla_\ell^2 h \boldsymbol{\alpha}\|_2 = \max_{\boldsymbol{\alpha}} \left(\sum_{j=1}^K \boldsymbol{\alpha}_j^\top \frac{\partial h}{\partial \mathbf{w}_{\ell, j}} \right)^2. \end{aligned} \quad (34)$$

A.2 Notations in Algorithm 1

Recall that the gradient in the t -th loop is

$$\begin{aligned} g_t(\mathbf{W}) &= \frac{1}{|\mathcal{D}_t^{(m)}|} \sum_{n \in \mathcal{D}_t^{(m)}} (Q(\mathbf{W}; \mathbf{x}_n) - y_n^{(t)}) \cdot \nabla_{\mathbf{W}} Q(\mathbf{W}; \mathbf{x}_n) \\ &= \frac{1}{N} \sum_{n=1}^N (Q(\mathbf{W}; \mathbf{x}_n) - r(\mathbf{x}_n) - \gamma \cdot \max_{a' \in \mathcal{A}} Q(\mathbf{W}^{(t-1)}; s'_n, a')) \cdot \nabla_{\mathbf{W}} Q(\mathbf{W}; \mathbf{x}_n). \end{aligned} \quad (35)$$

Then, we define $g_t^{(m)}(\mathbf{W}_\ell; \mathbf{W})$ as the components of $g_t^{(m)}(\mathbf{W})$ with respect to \mathbf{W}_ℓ . Recall that in (4) we have

$$\mathbf{W} = [\text{vec}(\mathbf{W}_1)^\top, \text{vec}(\mathbf{W}_2)^\top, \dots, \text{vec}(\mathbf{W}_L)^\top]^\top. \quad (36)$$

Then, with the definition of $g_t^{(m)}(\mathbf{W}_\ell; \mathbf{W})$, we have

$$g_t^{(m)}(\mathbf{W}) = [g_t^{(m)}(\mathbf{W}_1; \mathbf{W})^\top, g_t^{(m)}(\mathbf{W}_2; \mathbf{W})^\top, \dots, g_t^{(m)}(\mathbf{W}_L; \mathbf{W})^\top]^\top. \quad (37)$$

To simplify the analysis, the update of $\mathbf{W}^{(t,m)}$ is analyzed in the form of

$$\mathbf{W}_\ell^{(t,m+1)} = \mathbf{W}_\ell^{(t,m)} - \eta \cdot g_t^{(m)}(\mathbf{W}_\ell; \mathbf{W}^{(t,m)}) + \beta(\mathbf{W}_\ell^{(t,m)} - \mathbf{W}_\ell^{(t,m-1)}), \quad \forall \ell \in [L]. \quad (38)$$

One can see that (38) returns the same $\mathbf{W}^{(t,m+1)}$ as the gradient step at line 9 in Algorithm 1.

Table 3: Notations for the proofs

$g_t(\mathbf{W})$	The gradient function at point \mathbf{W} in the t -th outer loop, defined in (7).
$g_t(\mathbf{W}_\ell; \mathbf{W})$	The gradient function of $g_t(\mathbf{W})$ with respect to the components of \mathbf{W}_ℓ .
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.
\mathbf{W}^*	The desired Weights for approximating the optimal Q function.
$\mathbf{W}^{(t,m)}$	Model returned by Algorithm 1 at t -th outer loop and m -th inner loop.
f	The population risk function defined in (31).
$\nabla_{\mathbf{W}} f(\mathbf{W}^*)$	The full gradient of a function f at point \mathbf{W}^* .
$\nabla_{\ell} f(\mathbf{W}^*)$	The gradient of a function f with respect to the components of \mathbf{W}_ℓ at point \mathbf{W}^* .
$\nabla_{\ell}^2 f(\mathbf{W}^*)$	The Hessian matrix of a function f with respect to the components of \mathbf{W}_ℓ at point \mathbf{W}^* .
n	The dimension of \mathbf{W} .
n_ℓ	The dimension of vectorized \mathbf{W}_ℓ .
$\mathbf{h}^{(\ell)}(\mathbf{W})$	The input to the ℓ -th layer, defined in (39).
K_ℓ	The dimension of $\mathbf{h}^{(\ell)}$.
$\mathcal{J}_\ell(\mathbf{W})$	A function in $\mathbb{R}^n \rightarrow \mathbb{R}^K$, defined in (42).
ε_t	The value of ε in the behavior policy at t -th outer loop.
C_t	The distribution shift between the optimal policy and behavior policy at iteration t .
N	The size of the experience replay buffer.
R_{\max}	The upper bound of the reward.

A.3 Notations for the Deep Neural Networks.

Let n denote the dimension of \mathbf{W} defined in (4). We denote n_ℓ as the dimension of the vectorized neuron weights in the ℓ -th layer, namely, $n_\ell = \dim(\text{vec}(\mathbf{W}_\ell))$.

Then, let $\mathbf{h}^{(\ell)}(\mathbf{W})$ denote the input in the ℓ -th layer (or the output in the $(\ell-1)$ -th layer) with respect to the neuron weights as \mathbf{W} , and $\mathbf{h}^{(1)} = (s, a)$, where

$$\mathbf{h}^{(\ell)}(\mathbf{W}) = \phi(\mathbf{W}_{\ell-1}^\top \mathbf{h}^{(\ell-1)}) = \dots = \phi\left(\mathbf{W}_\ell^\top \phi\left(\mathbf{W}_{\ell-1} \dots \phi\left(\mathbf{W}_1^\top \mathbf{x}\right)\right)\right). \quad (39)$$

$\mathbf{h}^{(\ell)}(\mathbf{W})$ may be shortened as $\mathbf{h}^{(\ell)}$ when the neuron weights are clear from the contexts. Then, we denote the dimension of $\mathbf{h}^{(\ell)}$ as K_ℓ , where

$$K_\ell = \begin{cases} K, & \text{if } \ell > 1 \\ d, & \text{if } \ell = 1. \end{cases} \quad (40)$$

Then, $Q(\mathbf{W}; \mathbf{s}, a)$ can be written as

$$Q(\mathbf{W}; \mathbf{s}, a) = \frac{\mathbf{1}^\top}{K} \phi(\mathbf{w}_{L,k}^\top \mathbf{h}^{(L)}) = \frac{\mathbf{1}^\top}{K} \phi(\mathbf{W}_L^\top \phi(\mathbf{W}_{L-1}^\top \mathbf{h}^{(L-1)})), \quad (41)$$

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

$$\mathcal{J}_\ell(\mathbf{W}) = \begin{cases} [\mathbf{1}^\top \phi'(\mathbf{W}_L^\top \mathbf{h}^{(L)}) \mathbf{W}_L^\top \cdot \phi'(\mathbf{W}_{L-1}^\top \mathbf{h}^{(L-1)}) \mathbf{W}_{L-1}^\top \cdots \phi'(\mathbf{W}_{\ell+1}^\top \mathbf{h}^{(\ell+1)}) \mathbf{W}_{\ell+1}^\top]^\top & \text{if } \ell > 1 \\ \mathbf{1} & \text{if } \ell = 1. \end{cases} \quad (42)$$

Then, the gradient of Q can be represented as

$$\frac{\partial Q}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) = \frac{1}{K} \mathcal{J}_{\ell,k}(\mathbf{W}) \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^{(\ell)}(\mathbf{W})) \mathbf{h}^{(\ell)}(\mathbf{W}), \quad (43)$$

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

A.4 Notations for Order-wise Analysis

Without loss of generality, we consider the case that $d \gg K$. If $K \gg d$, we can always switch the order of K and d in the proof. Let $\sigma_i(L)$ denote the i -th largest singular value of \mathbf{W}_L^* . In this paper, we consider the case that \mathbf{W}_L^* is well-conditioned and bounded, i.e., $\sigma_1(L)$ and $\sigma_1(L)/\sigma_K(L)$ can be viewed as the constant and will be ignored in the analysis. In addition, some 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.

B Proof of Lemma 1 and Theorem 1

The main idea in proving Theorem 1 is to characterize the gradient descent term by the *Mean Value Theorem* (MVT) in Lemma 4 as shown in (47) and (48). The MVT is not directly applied in g_t because it is not smooth. However, the population risk functions defined in (31), which are the expectations over random variables, are smooth. Lemma 2 characterizes the bounds of the Hessian matrix defined in (49). Lemma 3 characterizes the bounds of gradient differences between the population risk function defined in (31) and g_t in (7) as shown in (60). Furthermore, according to Lemma 3, we know that the distance $\|\nabla_\ell f(\mathbf{W}) - \nabla_\ell f(\mathbf{W}^*)\|_2$ is upper bounded in the order of $\|\mathbf{W} - \mathbf{W}^*\|_2$ as shown in (60). Then, we can establish the connection between $\|\mathbf{W}^{(t,m+1)} - \mathbf{W}^*\|_2$ and $\|\mathbf{W}^{(t,m)} - \mathbf{W}^*\|_2$ as shown in (59). Then, by mathematical induction over m , one can characterize the iteration of $\{\|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2\}_{t=1}^T$ as shown in (65), which completes the proof of Lemma 1. Finally, selecting ε_t based on (68) for all $t \in [T]$, we derive the error bound of $\|\mathbf{W}^{(T,0)} - \mathbf{W}^*\|_2$ by mathematical induction over t , which completes the proof of Theorem 1.

Lemma 2. Given any $\mathbf{W} \in \mathbb{R}^n$, let \mathbf{W} satisfy

$$\|\mathbf{W} - \mathbf{W}^*\|_2 \lesssim \frac{\rho \cdot c_I \cdot \sigma_K}{K} \quad (44)$$

for some constant $c_I \in (0, 1)$. Then, for the f defined in (31), we have

$$\frac{(1 - c_I)\rho}{K^2} \preceq \nabla_\ell^2 f(\mathbf{W}) \preceq \frac{7}{K}. \quad (45)$$

Lemma 3. Let f be the function defined in (31). Let g_t be the function defined in (7). Then, we have

$$\begin{aligned} \|\nabla_\ell f(\mathbf{W}) - g_t(\mathbf{W}_\ell; \mathbf{W})\|_2 &\lesssim \frac{2 - \varepsilon_t}{K} \sqrt{\frac{K_\ell \cdot \log q}{N}} \cdot \|\mathbf{W} - \mathbf{W}^*\|_2 \\ &+ \frac{(1 - \varepsilon_t/2) \cdot \gamma}{K} \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2 \\ &+ C_d \cdot (C_t + (1 - C_t)\varepsilon) \cdot \frac{R_{\max}}{1 - \gamma}. \end{aligned} \quad (46)$$

with probability at least $1 - q^{-K_\ell}$.

Lemma 4 (Mean Value Theorem). *Let $\mathbf{U} \subset \mathbb{R}^{d_1}$ be open and $\mathbf{f} : \mathbf{U} \rightarrow \mathbb{R}^{d_2}$ be continuously differentiable, and $\mathbf{x} \in \mathbf{U}$, $\mathbf{h} \in \mathbb{R}^{d_1}$ 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} .

Proof of Theorem 1. Let \mathbf{W}_ℓ denote the neuron weights in the ℓ -th layer. From Algorithm 1 and (38), in the s -th iteration and t -th episode, we have

$$\begin{aligned} \mathbf{W}_\ell^{(t,m+1)} &= \mathbf{W}_\ell^{(t,m)} - \eta g_t^{(m)}(\mathbf{W}_\ell; \mathbf{W}^{(t,m)}) + \beta(\mathbf{W}_\ell^{(t,m)} - \mathbf{W}_\ell^{(t,m-1)}) \\ &= \mathbf{W}_\ell^{(t,m)} - \eta \nabla_\ell f(\mathbf{W}^{(t,m)}) + \beta(\mathbf{W}_\ell^{(t,m)} - \mathbf{W}_\ell^{(t,m-1)}) \\ &\quad + \eta \cdot (\nabla_\ell f(\mathbf{W}^{(t,m)}) - g_t^{(m)}(\mathbf{W}_\ell; \mathbf{W}^{(t,m)})). \end{aligned} \quad (47)$$

From (31), we can see that \mathbf{W}^* is the global optimal to f because $f(\mathbf{W}^*)$ achieves the minimum value as 0. Therefore, we have $\nabla_\ell f_t(\mathbf{W}^*) = \mathbf{0}$. Since $\nabla_\ell f$ is a smooth function \mathbf{W}^* , from the *Mean Value Theorem* in Lemma 4, we have

$$\begin{aligned} \nabla_\ell f(\mathbf{W}^{(t,m)}) &= \nabla_\ell f(\mathbf{W}^{(t,m)}) - \nabla_\ell f(\mathbf{W}^*) \\ &= \int_0^1 \nabla_\ell^2 f(\mathbf{W}^{(t,m)} + u \cdot (\mathbf{W}^{(t,m)} - \mathbf{W}^*)) du \cdot (\mathbf{W}_\ell^{(t,m)} - \mathbf{W}_\ell^*). \end{aligned} \quad (48)$$

For notational convenience, we use \mathbf{H} to denote the integration as

$$\mathbf{H} := \int_0^1 \nabla_\ell^2 f(\mathbf{W}^{(t,m)} + u \cdot (\mathbf{W}^{(t,m)} - \mathbf{W}^*)) du. \quad (49)$$

Then, we have

$$\begin{aligned} \begin{bmatrix} \mathbf{W}^{(t,m+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t,m)} - \mathbf{W}^* \end{bmatrix} &= \begin{bmatrix} \mathbf{I} - \eta \mathbf{H} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{W}^{(t,m)} - \mathbf{W}^* \\ \mathbf{W}^{(t,m-1)} - \mathbf{W}^* \end{bmatrix} \\ &\quad + \eta \begin{bmatrix} \nabla_\ell f(\mathbf{W}^{(t,m)}) - g_t^{(m)}(\mathbf{W}_\ell; \mathbf{W}^{(t,m)}) \\ \mathbf{0} \end{bmatrix}. \end{aligned} \quad (50)$$

Let $\mathbf{H} = \mathbf{S} \mathbf{\Lambda} \mathbf{S}^\top$ be the eigen-decomposition of \mathbf{H} . Then, we define

$$\mathbf{A}(\beta) := \begin{bmatrix} \mathbf{S}^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^\top \end{bmatrix} \mathbf{A}(\beta) \begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \eta \mathbf{\Lambda} + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}. \quad (51)$$

Since $\begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{S}^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$, we know $\mathbf{A}(\beta)$ and $\begin{bmatrix} \mathbf{I} - \eta \mathbf{\Lambda} + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}$ share the same eigenvalues. Let $\lambda_i^{(\mathbf{A})}$ be the i -th eigenvalue of $\mathbf{H}_i^{(\ell)}$, then the corresponding i -th eigenvalue of (51), denoted by $\lambda_i^{(\mathbf{A})}$, satisfies

$$(\lambda_i^{(\mathbf{A})}(\beta))^2 - (1 - \eta \lambda_i^{(\mathbf{A})} + \beta) \lambda_i^{(\mathbf{A})}(\beta) + \beta = 0. \quad (52)$$

By simple calculation, we have

$$|\lambda_i^{(\mathbf{A})}(\beta)| = \begin{cases} \sqrt{\beta}, & \text{if } \beta \geq (1 - \sqrt{\eta \lambda_i^{(\mathbf{A})}})^2, \\ \frac{1}{2} \left| (1 - \eta \lambda_i^{(\mathbf{A})} + \beta) + \sqrt{(1 - \eta \lambda_i^{(\mathbf{A})} + \beta)^2 - 4\beta} \right|, & \text{otherwise.} \end{cases} \quad (53)$$

Specifically, we have

$$\lambda_i^{(\mathbf{A})}(0) > \lambda_i^{(\mathbf{A})}(\beta), \quad \text{for } \forall \beta \in (0, (1 - \eta \lambda_i^{(\mathbf{A})})^2), \quad (54)$$

and $\lambda_i^{(\mathbf{A})}$ achieves the minimum $\lambda_i^{(\mathbf{A})\star} = \left|1 - \sqrt{\eta\lambda_i^{(\mathbf{A})}}\right|$ when $\beta^\star = \left(1 - \sqrt{\eta\lambda_i^{(\mathbf{A})}}\right)^2$. From Lemma 2, for any $\mathbf{a} \in \mathbb{R}^d$ with $\|\mathbf{a}\|_2 = 1$, we have

$$\begin{aligned} \mathbf{a}^\top \nabla_\ell f(\mathbf{W}^{(t,m)}) \mathbf{a} &= \int_0^1 \mathbf{a}^\top \nabla_\ell^2 f\left(\mathbf{W}^{(t,m)} + u \cdot (\mathbf{W}^{(t,m)} - \mathbf{W}^\star)\right) \mathbf{a} \cdot du \\ &\leq \int_0^1 \lambda_{\max} \|\mathbf{a}\|_2^2 du = \lambda_{\max}, \\ \mathbf{a}^\top \nabla_\ell f(\mathbf{W}^{(t,m)}) \mathbf{a} &= \int_0^1 \mathbf{a}^\top \nabla_\ell^2 f\left(\mathbf{W}^{(t,m)} + u \cdot (\mathbf{W}^{(t,m)} - \mathbf{W}^\star)\right) \mathbf{a} \cdot du \\ &\geq \int_0^1 \lambda_{\min} \|\mathbf{a}\|_2^2 du = \lambda_{\min}, \end{aligned} \tag{55}$$

where $\lambda_{\max} \approx \frac{1}{K}$, and $\lambda_{\min} \approx \frac{\rho}{K^2}$. Therefore, we have

$$\lambda_{\min}^{(\mathbf{A})} \approx \frac{(1 - c_I)\rho}{K^2}, \quad \text{and} \quad \lambda_{\max}^{(\mathbf{A})} \approx \frac{1}{K}. \tag{56}$$

Thus, when $\eta \leq \frac{1}{2\lambda_{\max}^{(\mathbf{A})}} \lesssim K$, $\|\mathbf{A}(\beta^\star)\|_2$ can be bounded by

$$\|\mathbf{A}(\beta^\star)\|_2 = 1 - \sqrt{\eta \cdot \lambda_{\min}^{(\mathbf{A})}} \leq 1 - \sqrt{\frac{(1 - c_I)\eta\rho}{K^2}}. \tag{57}$$

Therefore, we have

$$\begin{aligned} \|\mathbf{W}_\ell^{(t,m+1)} - \mathbf{W}_\ell^\star\|_2 &\leq \left(1 - \sqrt{\frac{(1 - c_I)\eta\rho}{K^2}}\right) \cdot \|\mathbf{W}_\ell^{(t,m)} - \mathbf{W}_\ell^\star\|_2 \\ &\quad + \eta \cdot \|\nabla_\ell f(\mathbf{W}^{(t,m)}) - g_t^{(m)}(\mathbf{W}^{(t,m)})\|_2 \\ &\lesssim \left(1 - (1 - \frac{c_I}{2})\sqrt{\frac{\eta\rho}{K^2}}\right) \cdot \|\mathbf{W}_\ell^{(t,m)} - \mathbf{W}_\ell^\star\|_2 \\ &\quad + \eta \cdot \|\nabla_\ell f(\mathbf{W}^{(t,m)}) - g_t^{(m)}(\mathbf{W}^{(t,m)})\|_2. \end{aligned} \tag{58}$$

Take the sum of (58) from $\ell = 1$ to $\ell = L$, we have

$$\begin{aligned} \|\mathbf{W}^{(t,m+1)} - \mathbf{W}^\star\|_2 &\leq \left(1 - (1 - \frac{c_I}{2})\sqrt{\frac{\eta\rho}{K^2}}\right) \cdot \|\mathbf{W}^{(t,m)} - \mathbf{W}^\star\|_2 \\ &\quad + \eta \cdot \sum_{\ell}^L \|\nabla_\ell f(\mathbf{W}^{(t,m)}) - g_t^{(m)}(\mathbf{W}^{(t,m)})\|_2. \end{aligned} \tag{59}$$

From Lemma 3, we have

$$\begin{aligned} \left\|\nabla_\ell f(\mathbf{W}^{(t,m)}) - g_t^{(m)}(\mathbf{W}_\ell; \mathbf{W}^{(t,m)})\right\|_2 &\lesssim \frac{2 - \varepsilon_t}{K} \sqrt{\frac{K_\ell \log q}{N_t}} \cdot \|\mathbf{W}^{(t,m)} - \mathbf{W}^\star\|_2 \\ &\quad + \frac{(1 - \varepsilon_t/2)\gamma}{K} \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^\star\|_2 \\ &\quad + C_d \cdot (C_t + (1 - C_t)\varepsilon) \cdot \frac{R_{\max}}{1 - \gamma}. \end{aligned} \tag{60}$$

For some small constant $c_N \geq 0$, let

$$\eta \cdot \frac{1}{K} \sqrt{\frac{K_\ell \log q}{N_t}} \leq \frac{c_N}{L} \sqrt{\frac{\eta\rho}{K^2}}, \tag{61}$$

which requires

$$\begin{aligned} N_t &\gtrsim c_N^{-2} \cdot \rho^{-1} \cdot \eta^{-1} \cdot L^2 \cdot \max_{\ell} K_\ell \cdot \log q \\ &= c_N^{-2} \cdot \rho^{-1} \cdot L \cdot d \cdot \log q. \end{aligned} \tag{62}$$

Then, the sample complexity

$$N = \sum_{t=1}^T N_t \gtrsim c_N^{-2} \cdot \rho^{-1} \cdot L \cdot d \cdot \log q \cdot T. \quad (63)$$

Therefore, we have

$$\begin{aligned} \|\mathbf{W}^{(t,m+1)} - \mathbf{W}^*\|_2 &\leq \left(1 - (1 - (2 - \varepsilon_t)c_N - \frac{c_I}{2})\sqrt{\frac{\rho}{TK^2}} \right) \cdot \|\mathbf{W}^{(t,m)} - \mathbf{W}^*\|_2 \\ &\quad + \sqrt{\eta} \cdot \frac{(1 - \varepsilon_t/2)\gamma}{K} \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2 \\ &\quad + \eta \cdot C_d \cdot (C_t + (1 - C_t)\varepsilon) \cdot \frac{R_{\max}}{1 - \gamma}. \end{aligned} \quad (64)$$

By mathematical induction, when $M = \log \gamma^{-1}$ and $\eta = 1/T = 1/\Theta(N)$, we have

$$\begin{aligned} &\|\mathbf{W}^{(t,M)} - \mathbf{W}^*\|_2 \\ &\lesssim \sqrt{\frac{K^2}{N}} \cdot C_d \cdot (C_t + (1 - C_t)\varepsilon_t) \cdot \frac{R_{\max}}{1 - \gamma} + (1 - \varepsilon_t/2)\gamma \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2 \\ &\leq \frac{c_N \cdot C_d \cdot (C_t + (1 - C_t) \cdot \varepsilon_t)}{K} \cdot \frac{R_{\max}}{1 - \gamma} + (1 - \varepsilon_t/2)\gamma \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2 \\ &\leq \frac{c_N \cdot C_d \cdot (C_{\max} + (1 - C_{\max}) \cdot \varepsilon_t)}{K} \cdot \frac{R_{\max}}{1 - \gamma} + (1 - \varepsilon_t/2)\gamma \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2. \end{aligned} \quad (65)$$

From Algorithm 1, we know that $\mathbf{W}^{(t+1,0)} = \mathbf{W}^{(t,M)}$. To guarantee that iteration converge to the ground truth \mathbf{W}^* , namely, $\|\mathbf{W}^{(t+1,0)} - \mathbf{W}^*\|_2 < \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2$, we need

$$\varepsilon_t \leq \frac{(1 - \gamma)^2 \cdot K \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2}{(1 - C_t) \cdot c_N \cdot C_d \cdot R_{\max}} - \frac{C_t}{1 - C_t}. \quad (66)$$

To guarantee that $\varepsilon_T \geq 0$, then we have

$$\|\mathbf{W}^{(T,0)} - \mathbf{W}^*\|_F \gtrsim \frac{C_T \cdot c_N \cdot C_d \cdot R_{\max}}{(1 - \gamma)^2 \cdot K}. \quad (67)$$

Specifically, let

$$\varepsilon_t = \frac{c_\varepsilon \cdot K \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2}{(1 - C_t) \cdot c_N \cdot C_d \cdot R_{\max}} - \frac{C_t}{1 - C_t}, \quad (68)$$

we have

$$\begin{aligned} \|\mathbf{W}^{(t+1,0)} - \mathbf{W}^*\|_2 &\lesssim \gamma + c_\varepsilon(1 - \gamma) \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2, \\ \text{and } \|\mathbf{W}^{(T,0)} - \mathbf{W}^*\|_2 &\lesssim [\gamma + c_\varepsilon(1 - \gamma)]^T \cdot \|\mathbf{W}^{(0,0)} - \mathbf{W}^*\|_2, \end{aligned} \quad (69)$$

which completes the proof. \square

C Proof of Lemma 2

Lemma 2 provides the lower and upper bounds for the eigenvalues of the Hessian matrix of population risk function in (31). According to Weyl's inequality in Lemma 5, the eigenvalues of $\nabla_\ell^2 f(\cdot)$ at any fixed point \mathbf{W} can be bounded in the form of (75). Therefore, we first provide the lower and upper bounds for $\nabla_\ell^2 f$ at the desired ground truth \mathbf{W}^* . Then, the bounds for $\nabla_\ell^2 f$ at any other point \mathbf{W} is bounded through (31) by utilizing the conclusion in Lemma 6. Lemma 6 illustrates the distance between the Hessian matrix of f at \mathbf{W} and \mathbf{W}^* . Lemma 7 provides the lower bound of $\mathbb{E}_{\mathbf{x}} \left(\sum_{j=1}^K \boldsymbol{\alpha}_j^\top \frac{\partial Q}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}^*) \right)^2$ when \mathbf{x} belongs to sub-Gaussian distribution, which is used in proving the lower bound of the Hessian matrix in (76).

Lemma 5 (Weyl's inequality, [5]). 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

$$|\lambda_i(\mathbf{B}) - \lambda_i(\mathbf{A})| \leq \|\mathbf{E}\|_2, \quad \forall i \in [m]. \quad (70)$$

Lemma 6. Let $f(\mathbf{W})$ be the population risk function defined in (31). If \mathbf{W} is close to \mathbf{W}^* such that

$$\|\mathbf{W} - \mathbf{W}^*\|_2 \lesssim \frac{\rho}{K} \quad (71)$$

we have

$$\|\nabla_\ell^2 f(\mathbf{W}) - \nabla_\ell^2 f(\mathbf{W}^*)\|_2 \lesssim \frac{1}{K} \cdot \|\mathbf{W} - \mathbf{W}^*\|_2. \quad (72)$$

Lemma 7. Suppose the following assumptions hold:

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

Let $\boldsymbol{\alpha} \in \mathbb{R}^{K_1 K_2}$ be the unit vector defined in (34), we have

$$\rho := \min_{\|\boldsymbol{\alpha}\|_2=1} \int_{\mathcal{R}} \left(\sum_{j=1}^K \boldsymbol{\alpha}^\top \mathbf{h} \phi'(\mathbf{w}_{\ell,j}^\top \mathbf{h}) \right)^2 p_H(\mathbf{h}) \cdot d\mathbf{h} > 0, \quad (73)$$

where $\mathcal{R} \subset \mathbb{R}^{K_\ell}$ with $\int_{\mathcal{R}} f_H(\mathbf{h}) > 0$. Moreover, if further assuming \mathcal{P} is Gaussian distribution and $\mathcal{R} = \mathbb{R}^{K_\ell}$, we have $\rho > 0.091$.

Lemma 8. Let $\mathbf{h}^{(\ell)}(\mathbf{W})$ be the function defined in (39). When \mathbf{W} is sufficiently close to \mathbf{W}^* , i.e., $\|\mathbf{W} - \mathbf{W}^*\|_2$ is smaller than some positive constant $c < 1$, we have

$$\begin{aligned} \|\mathbf{h}^{(\ell)}(\mathbf{W})\|_2 &\lesssim \|\mathbf{x}\|_2, \\ \|\mathbf{h}^{(\ell)}(\mathbf{W}) - \mathbf{h}^{(\ell)}(\mathbf{W}^*)\|_2 &\lesssim \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \|\mathbf{x}\|_2. \end{aligned} \quad (74)$$

Proof of Lemma 2. Let $\lambda_{\max}(\mathbf{W})$ and $\lambda_{\min}(\mathbf{W})$ denote the largest and smallest eigenvalues of $\nabla_\ell^2 f(\mathbf{W})$ at a point \mathbf{W} , respectively. Then, from Lemma 5, we have

$$\begin{aligned} \lambda_{\max}(\mathbf{W}) &\leq \lambda_{\max}(\mathbf{W}^*) + \|\nabla_\ell^2 f(\mathbf{W}) - \nabla_\ell^2 f(\mathbf{W}^*)\|_2, \\ \lambda_{\min}(\mathbf{W}) &\geq \lambda_{\min}(\mathbf{W}^*) - \|\nabla_\ell^2 f(\mathbf{W}) - \nabla_\ell^2 f(\mathbf{W}^*)\|_2. \end{aligned} \quad (75)$$

Then, we provide the lower bound of the Hessian matrix of the population function at \mathbf{W}^* . Let \mathcal{P} be the distribution for $\mathbf{h}^{(\ell)}(\mathbf{W})$ when $\mathbf{x} \sim \mu_t$ with probability density function denoted as p_H . For any $\boldsymbol{\alpha} \in \mathbb{R}^{K_\ell K}$ defined in (34) with $\|\boldsymbol{\alpha}\|_2 = 1$, we have

$$\begin{aligned} &\min_{\|\boldsymbol{\alpha}\|_2=1} \boldsymbol{\alpha}^\top \nabla_\ell^2 f(\mathbf{W}^*) \boldsymbol{\alpha} \\ &= \frac{1}{K^2} \min_{\|\boldsymbol{\alpha}\|_2=1} \mathbb{E}_{\mathbf{h} \sim \mathcal{P}} \left(\sum_{j=1}^K \boldsymbol{\alpha}_j^\top \mathbf{h}^{(\ell)} \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j}^{*\top} \mathbf{h}^{(\ell)}) \right)^2 \\ &= \frac{1}{K^2} \min_{\|\boldsymbol{\alpha}\|_2=1} \int_{\mathbb{R}^{K_\ell-1}} \left(\sum_{j=1}^K \boldsymbol{\alpha}_j^\top \mathbf{h}^{(\ell)} \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j}^{*\top} \mathbf{h}^{(\ell)}) \right)^2 p_H(\mathbf{h}^{(\ell)}) \cdot d\mathbf{h}^{(\ell)} \\ &= \frac{1}{K^2} \min_{\|\boldsymbol{\alpha}\|_2=1} \int_{\{\mathbf{h}^{(\ell)} | \mathcal{J}_{\ell,k} \neq 0\}} \left(\sum_{j=1}^K \boldsymbol{\alpha}_j^\top \mathbf{h}^{(\ell)} \phi'(\mathbf{w}_{\ell,j}^{*\top} \mathbf{h}^{(\ell)}) \right)^2 p_H(\mathbf{h}^{(\ell)}) \cdot d\mathbf{h}^{(\ell)} \\ &\gtrsim \frac{\rho}{K^2}, \end{aligned} \quad (76)$$

where the last inequality comes from Lemma 7, and Lemma 7 holds since $\mathbf{h}^{(\ell)}$ belongs to sub-Gaussian distribution and \mathbf{W}_ℓ is full rank.

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

$$\begin{aligned}
& \max_{\|\alpha\|_2=1} \alpha^\top \nabla_\ell^2 f(\mathbf{W}^*) \alpha \\
&= \frac{1}{K^2} \max_{\|\alpha\|_2=1} \mathbb{E}_{\mathbf{x}} \left(\sum_{j=1}^K \alpha_j^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j}^{*\top} \mathbf{h}^{(\ell)}) \right)^2 \\
&= \frac{1}{K^2} \max_{\|\alpha\|_2=1} \mathbb{E}_{\mathbf{x}} \sum_{j_1=1}^K \sum_{j_2=1}^K \alpha_{j_1}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_1}^{*\top} \mathbf{h}^{(\ell)}) \cdot \alpha_{j_2}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \mathbf{h}^{(\ell)}) \\
&= \frac{1}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \mathbb{E}_{\mathbf{x}} \alpha_{j_1}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_1}^{*\top} \mathbf{h}^{(\ell)}) \cdot \alpha_{j_2}^\top \mathbf{h}^{(\ell)} \cdot \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \mathbf{h}^{(\ell)}) \\
&\leq \frac{1}{K^2} \max_{\|\alpha\|_2=1} \sum_{j_1=1}^K \sum_{j_2=1}^K \left[\mathbb{E}_{\mathbf{x}} (\alpha_{j_1}^\top \mathbf{h}^{(\ell)})^4 \cdot \mathbb{E}(\phi'(\mathbf{w}_{\ell,j_1}^{*\top} \mathbf{h}^{(\ell)}))^4 \cdot \mathbb{E}_{\mathbf{x}} (\alpha_{j_2}^\top \mathbf{h}^{(\ell)})^4 \cdot \mathbb{E}_{\mathbf{x}} (\phi'(\mathbf{w}_{\ell,j_2}^{*\top} \mathbf{h}^{(\ell)}))^4 \right]^{1/4} \\
&\leq \frac{1}{K^2} \max_{\|\alpha\|_2=1} \sum_{j_1=1}^K \sum_{j_2=1}^K \left[\mathbb{E}_{\mathbf{x}} (\alpha_{j_1}^\top \mathbf{x})^4 \cdot \mathbb{E}_{\mathbf{x}} (\alpha_{j_2}^\top \mathbf{x})^4 \right]^{1/4} \\
&\leq \frac{3}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \|\alpha_{j_1}\|_2 \cdot \|\alpha_{j_2}\|_2 \leq \frac{6}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \frac{1}{2} (\|\alpha_{j_1}\|_2^2 + \|\alpha_{j_2}\|_2^2) \\
&= \frac{6}{K}.
\end{aligned} \tag{77}$$

Therefore, we have

$$\lambda_{\max}(\mathbf{W}^*) = \max_{\|\alpha\|_2=1} \alpha^\top \nabla_\ell^2 f(\mathbf{W}^*; p) \alpha \leq \frac{6}{K}. \tag{78}$$

Then, given (71), we have

$$\|\mathbf{W} - \mathbf{W}^*\|_2 \lesssim \frac{2\rho}{K}. \tag{79}$$

Combining (79) and Lemma 6, we have

$$\|\nabla_\ell^2 f(\mathbf{W}) - \nabla_\ell^2 f(\mathbf{W}^*)\|_2 \lesssim \frac{\rho}{K^2}. \tag{80}$$

Therefore, from (80) and (75), we have

$$\begin{aligned}
\lambda_{\max}(\mathbf{W}) &\leq \lambda_{\max}(\mathbf{W}^*) + \|\nabla_\ell^2 f(\mathbf{W}) - \nabla_\ell^2 f(\mathbf{W}^*)\|_2 \leq \frac{6}{K} + \frac{\rho}{2K^2} \leq \frac{7}{K}, \\
\lambda_{\min}(\mathbf{W}) &\geq \lambda_{\min}(\mathbf{W}^*) - \|\nabla_\ell^2 f(\mathbf{W}) - \nabla_\ell^2 f(\mathbf{W}^*)\|_2 \geq \frac{\rho}{K^2} - \frac{\rho}{2K^2} = \frac{\rho}{2K^2},
\end{aligned} \tag{81}$$

which completes the proof. \square

D Proof of Lemma 3

Before illustrating the whole proof, we first introduce some preliminary lemmas and definitions. Lemma 9 is the concentration theorem for independent random matrices. The definitions of the sub-Gaussian and sub-exponential variables are summarized in Definitions 3 and 4, and it is easy to verify that any bounded variables belong to sub-Gaussian distribution. Lemmas 10 and 11 serve as the technical tools in bounding matrix norms under the framework of the confidence interval.

The error bound between $\|\nabla_\ell f - g_\ell\|_2$ is divided into bounding I_1 , I_2 , and I_3 as shown in (91). I_1 in (92) represent the deviation of the mean of several random variables to their expectation, which can be bounded through concentration inequality, i.e, Chernoff bound. I_2 in (93) come from the inconsistency of "noisy" label in (8) and the "ground truth" label in the population risk function (31). I_3 in (94) come from the data distribution shift defined in Definition 1.

Lemma 9 ([72], Theorem 1.6). Consider a finite sequence $\{\mathbf{Z}_k\}$ of independent, random matrices with dimensions $d_1 \times d_2$. Assume that such a random matrix satisfies

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

Define

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

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

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

Definition 3 (Definition 5.7, [74]). A random variable X is called a sub-Gaussian random variable if it satisfies

$$(\mathbb{E}|X|^p)^{1/p} \leq c_1 \sqrt{p} \quad (82)$$

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

$$\mathbb{E}e^{s(X - \mathbb{E}X)} \leq e^{c_2 \|X\|_{\psi_2}^2 s^2} \quad (83)$$

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 $\boldsymbol{\alpha}^\top \mathbf{X}$ is sub-Gaussian for any $\boldsymbol{\alpha} \in \mathbb{R}^d$, and the sub-Gaussian norm of \mathbf{X} is defined as $\|\mathbf{X}\|_{\psi_2} = \sup_{\|\boldsymbol{\alpha}\|_2=1} \|\boldsymbol{\alpha}^\top \mathbf{X}\|_{\psi_2}$.

Definition 4 (Definition 5.13, [74]). A random variable X is called a sub-exponential random variable if it satisfies

$$(\mathbb{E}|X|^p)^{1/p} \leq c_3 p \quad (84)$$

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

$$\mathbb{E}e^{s(X - \mathbb{E}X)} \leq e^{c_4 \|X\|_{\psi_1}^2 s^2} \quad (85)$$

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 $\|X\|_{\psi_1} = \sup_{p \geq 1} p^{-1} (\mathbb{E}|X|^p)^{1/p}$.

Lemma 10 (Lemma 5.2, [74]). 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}_ξ is called a ξ -net of $\mathcal{B}(0, 1)$ if every point $\mathbf{z} \in \mathcal{B}(0, 1)$ can be approximated to within ξ by some point $\boldsymbol{\alpha} \in \mathcal{S}_\xi$, i.e., $\|\mathbf{z} - \boldsymbol{\alpha}\|_2 \leq \xi$. Then the minimal cardinality of a ξ -net \mathcal{S}_ξ satisfies

$$|\mathcal{S}_\xi| \leq (1 + 2/\xi)^d. \quad (86)$$

Lemma 11 (Lemma 5.3, [74]). Let \mathbf{A} be an $d_1 \times d_2$ matrix, and let $\mathcal{S}_\xi(d)$ be a ξ -net of $\mathcal{B}(0, 1)$ in \mathbb{R}^d for some $\xi \in (0, 1)$. Then

$$\|\mathbf{A}\|_2 \leq (1 - \xi)^{-1} \max_{\boldsymbol{\alpha}_1 \in \mathcal{S}_\xi(d_1), \boldsymbol{\alpha}_2 \in \mathcal{S}_\xi(d_2)} |\boldsymbol{\alpha}_1^\top \mathbf{A} \boldsymbol{\alpha}_2|. \quad (87)$$

Proof of Lemma 3. From (7), we know that

$$\begin{aligned} & g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) \\ &= \frac{1}{N} \sum_{n=1}^N (Q(\mathbf{W}; \mathbf{s}_n, a_n) - y_n^{(t)}) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \\ &= \frac{1}{N} \sum_{n=1}^N \left(Q(\mathbf{W}; \mathbf{s}_n, a_n) - Q(\mathbf{W}^*; \mathbf{s}_n, a_n) + \gamma \cdot \max_a Q(\mathbf{s}_n, a; \mathbf{W}^*) \right. \\ & \quad \left. - \gamma \cdot \max_a Q(\mathbf{s}_n, a; \mathbf{W}^{(t,0)}) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \\ &= \frac{1}{N} \sum_{n=1}^N \left(Q(\mathbf{W}; \mathbf{s}_n, a_n) - Q(\mathbf{W}^*; \mathbf{s}_n, a_n) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \\ & \quad + \frac{1}{N} \sum_{n=1}^N \gamma \cdot \left(\max_a Q(\mathbf{s}_n, a; \mathbf{W}^*) - \max_a Q(\mathbf{s}_n, a; \mathbf{W}^{(t,0)}) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}_n, a_n)}{\partial \mathbf{w}_{\ell,k}}. \end{aligned} \quad (88)$$

From (31), we know that

$$\frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) = \mathbb{E}_{(s,a) \sim \mu^*} \left(Q(\mathbf{W}; s, a) - Q(\mathbf{W}^*; s, a) \right) \cdot \frac{\partial Q(\mathbf{W}; s, a)}{\partial \mathbf{w}_{\ell,k}}. \quad (89)$$

Then, from (88) and (89), we have

$$\begin{aligned} g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) &= g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \mathbb{E}_{(s,a) \sim \mathcal{D}_t} g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) \\ &\quad + \mathbb{E}_{(s,a) \sim \mu_t} g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}), \end{aligned} \quad (90)$$

where \mathcal{D}_t and μ_t are equivalent because of Assumption 2. Then, we have

$$\begin{aligned} &g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) \\ &= \left[\frac{1}{N} \sum_{n=1}^N \left(Q(\mathbf{W}; s_n, a_n) - Q(\mathbf{W}^*; s_n, a_n) \right) \cdot \frac{\partial Q(\mathbf{W}; s_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \right. \\ &\quad \left. - \mathbb{E}_{(s,a) \sim \mu_t} \left(Q(\mathbf{W}; s, a) - Q(\mathbf{W}^*; s, a) \right) \cdot \frac{\partial Q(\mathbf{W}; s, a)}{\partial \mathbf{w}_{\ell,k}} \right] \\ &\quad + \left[\frac{1}{N} \sum_{n=1}^N \gamma \cdot \left(\max_a Q(s_n, a; \mathbf{W}^*) - \max_a Q(s_n, a; \mathbf{W}^{(t,0)}) \right) \cdot \frac{\partial Q(\mathbf{W}; s_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \right] \\ &\quad + \mathbb{E}_{(s,a) \sim \mu_t} g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}). \end{aligned} \quad (91)$$

For convenience, we define \mathbf{I}_1 , \mathbf{I}_2 , and \mathbf{I}_3 in the following ways with $\mathbf{x}_n := (s_n, a_n)$ be the feature mapping of state-action pair (s_n, a_n) .

Then, \mathbf{I}_1 is defined as

$$\begin{aligned} \mathbf{I}_1 &:= \frac{1}{N} \sum_{n=1}^N \left(Q(\mathbf{W}; s_n, a_n) - Q(\mathbf{W}^*; s_n, a_n) \right) \cdot \frac{\partial Q(\mathbf{W}; s_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \\ &\quad - \mathbb{E}_{(s,a) \sim \mathcal{D}_t} \left(Q(\mathbf{W}; s, a) - Q(\mathbf{W}^*; s, a) \right) \cdot \frac{\partial Q(\mathbf{W}; s, a)}{\partial \mathbf{w}_{\ell,k}}, \end{aligned} \quad (92)$$

\mathbf{I}_2 is defined as

$$\mathbf{I}_2 := \frac{1}{N} \sum_{n=1}^N \gamma \cdot \left(\max_a Q(s'_n, a; \mathbf{W}^*) - \max_a Q(s'_n, a; \mathbf{W}^{(t,0)}) \right) \cdot \frac{\partial Q(\mathbf{W}; s_n, a_n)}{\partial \mathbf{w}_{\ell,k}}, \quad (93)$$

and \mathbf{I}_3 is defined as

$$\mathbf{I}_3 := \mathbb{E}_{(s,a) \sim \mu_t} g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}), \quad (94)$$

where

$$\frac{\partial Q(\mathbf{W}; s_n, a_n)}{\partial \mathbf{w}_{\ell,k}} = \frac{1}{K} \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^\ell) \mathbf{h}^\ell(\mathbf{W}) \quad (95)$$

from (43). Therefore, we have

$$\left\| g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) \right\|_2 \leq \|\mathbf{I}_1\|_2 + \|\mathbf{I}_2\|_2 + \|\mathbf{I}_3\|_2. \quad (96)$$

Next, we will provide the bound for $\|\mathbf{I}_1\|_2$, $\|\mathbf{I}_2\|_2$, and $\|\mathbf{I}_3\|_2$.

Bound of \mathbf{I}_1 . We first divide the data in \mathcal{D}_t into two parts, namely, $\mathcal{D}_{t,1}$ and $\mathcal{D}_{t,2}$. $\mathcal{D}_{t,1}$ includes the state-action pair (s, a) such that a_n is randomly selected from action space \mathcal{A} , and $\mathcal{D}_{t,2}$ includes the state-action pair (s, a) such that a_n is selected based on the greedy policy with respect to $Q(\mathbf{W}^{(t,0)})$.

Then, we define a random variable $Z^{(\ell,1)} = (Q(\mathbf{x}; \mathbf{W}) - Q(\mathbf{x}; \mathbf{W}^*)) \cdot \mathcal{J}_{\ell,k} \cdot \boldsymbol{\alpha}^T \mathbf{h}^{(\ell)}(\mathbf{W})$ with $\mathbf{x} \sim \mathcal{D}_{t,1}$ and $Z_n^{(\ell,1)} = (Q(\mathbf{x}_n; \mathbf{W}) - Q(\mathbf{x}_n; \mathbf{W}^*)) \cdot \mathcal{J}_{\ell,k} \cdot \boldsymbol{\alpha}^T \mathbf{h}_n^{(\ell)}(\mathbf{W})$ as the realization of $Z^{(\ell,1)}$ for $n = 1, 2, \dots, N$, where $\boldsymbol{\alpha} \in \mathbb{R}^d$ is any fixed unit vector with $\|\boldsymbol{\alpha}\|_2 \leq 1$. We know that \mathbf{s} and a are independent for $\mathbf{x} \sim \mathcal{D}_{t,1}$. Let $\boldsymbol{\Sigma}_1$ denote the covariance matrix of $\mathbf{x} \sim \mathcal{D}_{t,1}$. Moreover, $\mathbf{x}(\mathbf{s}, a)$ is bounded by 1, then we have $\|\boldsymbol{\Sigma}_1\|_2 \leq 1$.

Similar to $Z^{(\ell,1)}$, we define a random variable $Z^{(\ell,2)} = (Q(\mathbf{x}; \mathbf{W}) - Q(\mathbf{x}; \mathbf{W}^*)) \cdot \mathcal{J}_{\ell,k} \cdot \boldsymbol{\alpha}^T \mathbf{h}^{(\ell)}(\mathbf{W})$ with $\mathbf{x} \sim \mathcal{D}_{t,2}$ and $Z_n^{(\ell,2)} = (Q(\mathbf{x}_n; \mathbf{W}) - Q(\mathbf{x}_n; \mathbf{W}^*)) \cdot \mathcal{J}_{\ell,k} \cdot \boldsymbol{\alpha}^T \mathbf{h}_n^{(\ell)}(\mathbf{W})$ as the realization of $Z^{(\ell,2)}$ for $n = 1, 2, \dots, N$. Differ from $Z^{(\ell,1)}$, \mathbf{s} and a are dependent for $\mathbf{x} \sim \mathcal{D}_{t,2}$. Let $\boldsymbol{\Sigma}_2$ denote the covariance matrix of $\mathbf{x} \sim \mathcal{D}_{t,1}$. Then, we have $\|\boldsymbol{\Sigma}_2\|_2 \leq 1 + \max_j \rho_{x_j, a} \leq 2$, where $\rho_{x_j, a}$ denotes the correlation between a and x_j .

According to the definition of (92), we can rewrite \mathbf{I}_1 as

$$\begin{aligned}
\mathbf{I}_1 &= \frac{1}{K} \left[\frac{1}{N} \sum_{n=1}^N (Q(\mathbf{W}; \mathbf{x}_n) - Q(\mathbf{W}^*; \mathbf{x}_n)) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}_n^\ell) \mathbf{h}_n^\ell \right. \\
&\quad \left. - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} (Q(\mathbf{W}; \mathbf{x}) - Q(\mathbf{W}^*; \mathbf{x})) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^\ell) \mathbf{h}^\ell \right] \\
&= \frac{1}{K} \left[\frac{1}{N} \left(\sum_{n \in \mathcal{D}_{t,1}} (Q(\mathbf{W}; \mathbf{x}_n) - Q(\mathbf{W}^*; \mathbf{x}_n)) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}_n^\ell) \mathbf{h}_n^\ell \right. \right. \\
&\quad \left. \left. + \sum_{n \in \mathcal{D}_{t,2}} (Q(\mathbf{W}; \mathbf{x}_n) - Q(\mathbf{W}^*; \mathbf{x}_n)) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}_n^\ell) \mathbf{h}_n^\ell \right) \right. \\
&\quad \left. - \left(\varepsilon \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t,1}} (Q(\mathbf{W}; \mathbf{x}) - Q(\mathbf{W}^*; \mathbf{x})) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^\ell) \mathbf{h}^\ell \right. \right. \\
&\quad \left. \left. + (1 - \varepsilon) \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t,2}} (Q(\mathbf{W}; \mathbf{x}) - Q(\mathbf{W}^*; \mathbf{x})) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^\ell) \mathbf{h}^\ell \right) \right] \tag{97} \\
&= \frac{1}{K^2} \left[\varepsilon \cdot \left(\frac{1}{\varepsilon N} \sum_{n \in \mathcal{D}_{t,1}} (Q(\mathbf{W}; \mathbf{x}_n) - Q(\mathbf{W}^*; \mathbf{x}_n)) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}_n^\ell) \mathbf{h}_n^\ell \right. \right. \\
&\quad \left. \left. - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t,1}} (Q(\mathbf{W}; \mathbf{x}) - Q(\mathbf{W}^*; \mathbf{x})) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^\ell) \mathbf{h}^\ell \right) \right. \\
&\quad \left. + (1 - \varepsilon) \left(\frac{1}{(1 - \varepsilon) N} \sum_{n \in \mathcal{D}_{t,2}} (Q(\mathbf{W}; \mathbf{x}_n) - Q(\mathbf{W}^*; \mathbf{x}_n)) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}_n^\ell) \mathbf{h}_n^\ell \right. \right. \\
&\quad \left. \left. - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t,2}} (Q(\mathbf{W}; \mathbf{x}) - Q(\mathbf{W}^*; \mathbf{x})) \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^\ell) \mathbf{h}^\ell \right) \right]
\end{aligned}$$

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

$$\begin{aligned}
(\mathbb{E}|Z^{(1)}|^p)^{1/p} &= \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t,1}} |Q(\mathbf{W}; \mathbf{x}) - Q(\mathbf{W}^*; \mathbf{x})|^p \cdot |\mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{x})| \cdot |\boldsymbol{\alpha}^T \mathbf{h}^\ell|^p \right)^{1/p} \\
&\leq \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t,1}} |Q(\mathbf{W}; \mathbf{x}) - Q(\mathbf{W}^*; \mathbf{x})|^p \cdot |\boldsymbol{\alpha}^T \mathbf{h}^\ell|^p \right)^{1/p} \tag{98} \\
&\leq \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{t,1}} \left[\|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \|\mathbf{x}\|_2 \right]^p \cdot |\boldsymbol{\alpha}^T \mathbf{x}|^p \right)^{1/p} \\
&\leq C_1 \cdot \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot p
\end{aligned}$$

where C_1 is a positive constant.

From Definition 4, we know that $Z^{(\ell,1)}$ belongs to sub-exponential distribution with $\|Z^{(\ell,1)}\|_{\psi_1} \leq C_1 \|\mathbf{W} - \mathbf{W}^*\|_2$. Therefore, by Chernoff inequality, we have

$$\mathbb{P} \left\{ \left| \frac{1}{N} \sum_{n=1}^N Z_n^{(\ell,1)}(j) - \mathbb{E} Z^{(\ell,1)}(j) \right| < t \right\} \leq 1 - \frac{e^{-C(C_1 \|\mathbf{W} - \mathbf{W}^*\|_2)^2 \cdot N s^2}}{e^{N s t}} \tag{99}$$

for some positive constant C and any $s \in \mathbb{R}$.

Let $t = C_1 \|\mathbf{W} - \mathbf{W}^*\|_2 \sqrt{\frac{d \log q}{N}}$ and $s = \frac{2}{C \|\mathbf{W} - \mathbf{W}^*\|_2} \cdot t$ for some large constant $q > 0$. Then, we have

$$\left| \frac{1}{N} \sum_{n=1}^N Z_n^{(\ell,1)}(j) - \mathbb{E} Z^{(\ell,1)}(j) \right| \lesssim C_1 \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \sqrt{\frac{d \log q}{N}} \quad (100)$$

with probability at least $1 - q^{-d}$.

Similar to (98), we have

$$\left(\mathbb{E} |Z^{(\ell,2)}|^p \right)^{1/p} \leq C_2 \cdot \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot p, \quad (101)$$

where $C_2 = 2 \cdot C_1$. Then, we have

$$\left| \frac{1}{N} \sum_{n=1}^N Z_n^{(\ell,2)}(j) - \mathbb{E} Z^{(\ell,2)}(j) \right| \lesssim 2C_1 \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \sqrt{\frac{d \log q}{N}} \quad (102)$$

with probability at least $1 - q^{-d}$.

From Lemma 11 and (97), we have

$$\begin{aligned} \|\mathbf{I}_1\|_2 &\leq 2 \cdot \frac{1}{K^2} \left[\varepsilon \cdot \left| \frac{1}{\varepsilon N} \sum_{n \in \mathcal{D}_{t,1}} Z_n^{(\ell,1)}(j) - \mathbb{E} Z^{(\ell,1)}(j) \right| \right. \\ &\quad \left. + (1 - \varepsilon) \cdot \left| \frac{1}{(1 - \varepsilon)N} \sum_{n \in \mathcal{D}_{t,2}} Z_n^{(\ell,2)}(j) - \mathbb{E} Z^{(\ell,2)}(j) \right| \right] \\ &\lesssim \frac{2 - \varepsilon}{K^2} \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \sqrt{\frac{d \log q}{N}} \end{aligned} \quad (103)$$

with probability at least $1 - |\mathcal{S}_{\frac{1}{2}}(d)| \cdot q^{-d}$.

From Lemma 10, we know that $|\mathcal{S}_{\frac{1}{2}}(d)| \leq 5^d$. Therefore, the probability for (103) 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.

Bound of \mathbf{I}_2 . Let $\mathbf{a}_n^* = \arg \max_{\mathbf{a} \in \mathcal{A}} Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a})$. While for $Q(\mathbf{W})$, we have

$$\max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \geq Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}^*). \quad (104)$$

Then, we have

$$\begin{aligned} \max_{\mathbf{a}} Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}) - \max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) &= Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}_n^*) - \max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \\ &\leq Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}_n^*) - Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}_n^*). \end{aligned} \quad (105)$$

Similarly to (105), let us define $\tilde{\mathbf{a}}_n^* = \arg \max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}_n, \mathbf{a})$. Then, we have

$$\max_{\mathbf{a}} Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}) - \max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \geq Q(\mathbf{W}^*; \mathbf{s}'_n, \tilde{\mathbf{a}}_n^*) - Q(\mathbf{W}; \mathbf{s}'_n, \tilde{\mathbf{a}}_n^*). \quad (106)$$

Combining (105) and (106), we have

$$\left| \max_{\mathbf{a}} Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}) - \max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \right| \leq \max_{\mathbf{a}} \left| Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}) - Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \right|. \quad (107)$$

Following the definition of $Z^{(\ell,1)}$ in (98), we define

$$Z^{(\ell,3)}(j) = \left(\max_{\mathbf{a}} Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}) - \max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \right) \cdot \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^{(\ell)}) \cdot \boldsymbol{\alpha}^\top \mathbf{h}^{(\ell)}.$$

Therefore, from (105) and (106), we know

$$\begin{aligned} \left(\mathbb{E} |Z^{(3)}|^p \right)^{1/p} &\leq \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} \left| \max_{\mathbf{a}} Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}) - \max_{\mathbf{a}} Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \right|^p \right. \\ &\quad \left. \cdot \left| \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,k}^\top \mathbf{h}^{(\ell)}) \right|^p \cdot \left| \boldsymbol{\alpha}^\top \mathbf{h}_n^{(\ell)} \right|^p \right)^{1/p} \\ &\leq \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} \max_{\mathbf{a}} \left| Q(\mathbf{W}^*; \mathbf{s}'_n, \mathbf{a}) - Q(\mathbf{W}; \mathbf{s}'_n, \mathbf{a}) \right|^p \cdot \left| \boldsymbol{\alpha}^\top \mathbf{h}_n^{(\ell)} \right|^p \right)^{1/p} \\ &\lesssim (2 - \varepsilon) \cdot \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \log |\mathcal{A}| \cdot p. \end{aligned} \quad (108)$$

Following the steps in (98) to (100), we have

$$\begin{aligned}
\|I_2\|_2 &\lesssim \frac{(1-\varepsilon/2)\gamma}{K} \cdot \left(\|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \sqrt{\frac{d \cdot \log q \cdot \log |\mathcal{A}|}{N}} + \mathbb{E}Z^{(\ell,3)} \right) \\
&\lesssim \frac{(1-\varepsilon/2)\gamma}{K} \cdot \left(\|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \left(\sqrt{\frac{d \cdot \log q \cdot \log |\mathcal{A}|}{N}} + C \right) \right) \\
&\lesssim \frac{(1-\varepsilon/2)\gamma}{K} \cdot \|\mathbf{W} - \mathbf{W}^*\|_2
\end{aligned} \tag{109}$$

with probability at least $1 - q^{-d}$, where the last inequality holds when $N \gtrsim d \cdot \log q \cdot \log |\mathcal{A}|$.

Bound of I_3 . We have

$$\begin{aligned}
&I_3 \\
&= \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t} g_t(\mathbf{w}_{\ell, k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell, k}}(\mathbf{W}) \\
&= \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t} \left(Q(\mathbf{W}; \mathbf{s}, a) - Q(\mathbf{W}^*; \mathbf{s}, a) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \\
&\quad - \mathbb{E}_{(\mathbf{s}, a) \sim \mu^*} \left(Q(\mathbf{W}; \mathbf{s}, a) - Q(\mathbf{W}^*; \mathbf{s}, a) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \\
&= \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t} \left(Q(\mathbf{W}; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \mathbb{E}_{\mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \max_{a'} Q(\mathbf{W}^*; \mathbf{s}', a') \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \\
&\quad - \mathbb{E}_{(\mathbf{s}, a) \sim \mu^*} \left(Q(\mathbf{W}; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \mathbb{E}_{\mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \max_{a'} Q(\mathbf{W}^*; \mathbf{s}', a') \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \\
&= \mathbb{E}_{(\mathbf{s}, a) \sim \mu_t, \mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \left(Q(\mathbf{W}; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \max_{a'} Q(\mathbf{W}^*; \mathbf{s}', a') \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \\
&\quad - \mathbb{E}_{(\mathbf{s}, a) \sim \mu^*, \mathbf{s}' \sim p_{\mathbf{s}, \mathbf{s}'}}^a \left(Q(\mathbf{W}; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \max_{a'} Q(\mathbf{W}^*; \mathbf{s}', a') \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}}
\end{aligned} \tag{110}$$

Then, we have

$$\begin{aligned}
|I_3| &= \left| \int_{(\mathbf{s}, a)} \int_{\mathbf{s}'} \left(Q(\mathbf{W}; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \max_{a'} Q(\mathbf{W}^*; \mathbf{s}', a') \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \right. \\
&\quad \left. \cdot (\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| \\
&\leq \left| Q(\mathbf{W}; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \max_{a'} Q(\mathbf{W}^*; \mathbf{s}', a') \right| \cdot \left| \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \right| \\
&\quad \cdot \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| Q(\mathbf{W}; \mathbf{s}, a) - r(\mathbf{s}, a) - \gamma \cdot \max_{a'} Q(\mathbf{W}^*; \mathbf{s}', a') \right| \cdot \left| \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell, k}} \right| \\
&\quad \cdot \left[(1-\varepsilon) \cdot \left| \int_{(\mathbf{s}, a)} \int_{\mathbf{s}'} (\mu^*(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mu_{t,1}(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \right. \\
&\quad \left. + \varepsilon \cdot \left| \int_{(\mathbf{s}, a)} \int_{\mathbf{s}'} (\mu^*(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mu_{t,2}(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \right] \\
&\leq \frac{R_{\max}}{1-\gamma} \cdot \left[(1-\varepsilon) \cdot \left| \int_{(\mathbf{s}, a)} \int_{\mathbf{s}'} (\mu^*(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mu_{t,1}(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \right. \\
&\quad \left. + \varepsilon \cdot \left| \int_{(\mathbf{s}, a)} \int_{\mathbf{s}'} (\mu^*(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mu_{t,2}(d\mathbf{s}, da)\mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \right].
\end{aligned} \tag{111}$$

Then, we have

$$\begin{aligned}
& \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} (\mu^*(d\mathbf{s}, da) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mu_{t,1}(d\mathbf{s}, da) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \\
&= \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} (\mathcal{P}^*(d\mathbf{s}) \pi^*(da|\mathbf{s}) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) - \mathcal{P}_{t,1}(d\mathbf{s}) \pi_{t,1}(da|\mathbf{s}) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a)) \right| \\
&\leq \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} (\mathcal{P}^*(d\mathbf{s}) - \mathcal{P}_{t,1}(d\mathbf{s})) \pi^*(da|\mathbf{s}) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) \right| \\
&\quad + \left| \int_{(\mathbf{s},a)} \int_{\mathbf{s}'} \mathcal{P}_{t,1}(d\mathbf{s}) (\pi_{t,1}(da|\mathbf{s}) - \pi^*(da|\mathbf{s})) \mathcal{P}(d\mathbf{s}'|\mathbf{s}, a) \right| \\
&\leq |\mathcal{A}| \cdot C_t.
\end{aligned} \tag{112}$$

Therefore, the bound of I_3 can be found as

$$\begin{aligned}
|I_3| &\lesssim \frac{R_{\max}}{1-\gamma} \cdot |\mathcal{A}| \cdot ((1-\varepsilon)C_t + \varepsilon \cdot C_t) \\
&= C_d \cdot (C_t + (1-C_t)\varepsilon) \cdot \frac{R_{\max}}{1-\gamma},
\end{aligned} \tag{113}$$

where $C_d = |\mathcal{A}|$.

In conclusion, let $\boldsymbol{\alpha} \in \mathbb{R}^{Kd}$ and $\boldsymbol{\alpha}_j \in \mathbb{R}^d$ with $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1^T, \boldsymbol{\alpha}_2^T, \dots, \boldsymbol{\alpha}_K^T]^T$, we have

$$\begin{aligned}
& \|g_t(\mathbf{W}) - \nabla f_t(\mathbf{W})\|_2 \\
&= \left| \boldsymbol{\alpha}^T (g_t(\mathbf{W}) - \nabla f_t(\mathbf{W})) \right| \\
&\leq \sum_{k=1}^K \left| \boldsymbol{\alpha}_k^T (g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W})) \right| \\
&\leq \sum_{k=1}^K \left\| g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) \right\|_2 \cdot \|\boldsymbol{\alpha}_k\|_2 \\
&\leq \sum_{k=1}^K (\|\mathbf{I}_1\|_2 + \|\mathbf{I}_2\|_2 + \|\mathbf{I}_3\|_2) \cdot \|\boldsymbol{\alpha}_k\|_2 \\
&\leq \frac{2-\varepsilon}{K} \sqrt{\frac{d \log q}{N}} \cdot \|\mathbf{W} - \mathbf{W}^*\|_2 + \frac{(1-\varepsilon/2)\gamma}{K} \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_2 \\
&\quad + C_d \cdot (C_t + (1-C_t)\varepsilon) \cdot \frac{R_{\max}}{1-\gamma}
\end{aligned} \tag{114}$$

with probability at least $1 - q^{-d}$. □

E Additional proof of the lemmas in Appendix C

E.1 Proof of Lemma 6

The distance of the second order derivatives of the population risk function $f(\cdot)$ at point \mathbf{W} and \mathbf{W}^* can be converted into bounding $\mathbf{P}_1, \mathbf{P}_2$, which are defined in (116). The major idea in proving \mathbf{P}_1 is to connect the error bound to the angle between \mathbf{W} and \mathbf{W}^* given $\mathbf{h}^{(\ell)}$ belongs to the sub-Gaussian distribution.

Proof of Lemma 6. From the definition of f in (31), we have

$$\begin{aligned}
\frac{\partial^2 f}{\partial \mathbf{w}_{\ell,j_1} \partial \mathbf{w}_{\ell,j_2}}(\mathbf{W}^*) &= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{j_1}^{*\top} \mathbf{h}) \cdot \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{j_2}^{*\top} \mathbf{h}) \cdot \mathbf{h}^* \mathbf{h}^{*\top}, \\
\text{and } \frac{\partial^2 f}{\partial \mathbf{w}_{\ell,j_1} \partial \mathbf{w}_{\ell,j_2}}(\mathbf{W}) &= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \phi' \mathcal{J}_{\ell,k}^*(\mathbf{w}_{j_1}^\top \mathbf{h}) \cdot \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{j_2}^\top \mathbf{h}) \cdot \mathbf{h} \mathbf{h}^\top,
\end{aligned} \tag{115}$$

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

Then, we have

$$\begin{aligned}
& \frac{\partial^2 f}{\partial \mathbf{w}_{\ell,j_1} \partial \mathbf{w}_{\ell,j_2}}(\mathbf{W}^*) - \frac{\partial^2 f}{\partial \mathbf{w}_{\ell,j_1} \partial \mathbf{w}_{\ell,j_2}}(\mathbf{W}) \\
&= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} [\mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_1}^{*T} \mathbf{h}^*) \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) \mathbf{h}^* \mathbf{h}^{*\top} - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_1}^\top \mathbf{h}) \mathcal{J}_{\ell,k} \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) \mathbf{h} \mathbf{h}^\top] \\
&= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} [\mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_1}^{*T} \mathbf{h}^*) (\mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) \mathbf{h}^* \mathbf{h}^{*\top} - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) \mathbf{h} \mathbf{h}^\top) \\
&\quad + \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) (\mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_1}^{*T} \mathbf{h}^*) \mathbf{h}^* \mathbf{h}^{*\top} - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_1}^\top \mathbf{h}) \mathbf{h} \mathbf{h}^\top)] \\
&:= \frac{1}{K^2} (\mathbf{P}_1 + \mathbf{P}_2).
\end{aligned} \tag{116}$$

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

$$\mathbf{a}^\top \mathbf{P}_1 \mathbf{a} = \mathbb{E}_{\mathbf{x}} \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_1}^{*T} \mathbf{h}^*) \left(\mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right). \tag{117}$$

Then, we have

$$\begin{aligned}
|\mathbf{a}^\top \mathbf{P}_1 \mathbf{a}| &= \left| \mathbb{E}_{\mathbf{x}} \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_1}^{*T} \mathbf{h}^*) \left(\mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right) \right| \\
&\leq \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right| \\
&\leq \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h}^*)^2 - \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h})^2 \right| \\
&\quad + \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k}^* \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h}^*) (\mathbf{a}^\top \mathbf{h})^2 - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right| \\
&\quad + \mathbb{E}_{\mathbf{x}} \left| \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 - \mathcal{J}_{\ell,k} \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) (\mathbf{a}^\top \mathbf{h})^2 \right| \\
&\lesssim \|\mathbf{W} - \mathbf{W}^*\|_2 + \|\mathbf{W} - \mathbf{W}^*\|_2 \\
&\quad + \mathbb{E}_{\mathbf{x}} \left| (\phi'(\mathbf{w}_{\ell,j_2}^{*\top} \mathbf{h}) - \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h})) \cdot (\mathbf{a}^\top \mathbf{h})^2 \right| \\
&\lesssim \|\mathbf{W} - \mathbf{W}^*\|_2 + \mathbb{E}_{\mathbf{x}} \left| (\phi'(\mathbf{w}_{\ell,j_2}^{*\top} \mathbf{h}) - \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h})) \cdot (\mathbf{a}^\top \mathbf{h})^2 \right|.
\end{aligned} \tag{118}$$

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

$$\begin{aligned}
\mathbf{w}_{\ell,j_2}^{*\top} \mathbf{h} &= \mathbf{w}_{\ell,j_2}^{*\top} (z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c} + \dots + z_d \mathbf{a}_d^\perp) \\
&= z_1 \mathbf{w}_{\ell,j_2}^{*\top} \mathbf{a} + z_2 \mathbf{w}_{\ell,j_2}^{*\top} \mathbf{b} + z_3 \mathbf{w}_{\ell,j_2}^{*\top} \mathbf{c} + \dots + z_d \mathbf{w}_{\ell,j_2}^{*\top} \mathbf{a}_d^\perp \\
&= z_1 \mathbf{w}_{\ell,j_2}^{*\top} \mathbf{a} + z_2 \mathbf{w}_{\ell,j_2}^{*\top} \mathbf{b} + z_3 \mathbf{w}_{\ell,j_2}^{*\top} \mathbf{c} + 0 \\
&= \mathbf{w}_{\ell,j_2}^{*\top} (z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c}) \\
&:= \mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{h}}.
\end{aligned} \tag{119}$$

where $\tilde{\mathbf{h}} = z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c}$. Similar to (119), we have $\mathbf{w}_{\ell,j_2}^\top \mathbf{h} = \mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{h}}$ and $\mathbf{a}^\top \mathbf{h} = \mathbf{a}^\top \tilde{\mathbf{h}}$.

Then, we define I_4 as

$$\begin{aligned}
I_4 &:= \mathbb{E}_{\mathbf{h}} \left| (\phi'(\mathbf{w}_{\ell,j_2}^{*\top} \mathbf{h}) - \phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h})) \cdot (\mathbf{a}^\top \mathbf{h}) \right| \\
&= \int_{\mathcal{R}_{\mathbf{h}}} |\phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) - \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h})| \cdot |\mathbf{a}^\top \mathbf{h}|^2 \cdot f_{\mathbf{H}}(\mathbf{h}) d\mathbf{h} \\
&= \int_{\mathcal{R}_{\mathbf{z}}} |\phi'(\mathbf{w}_{\ell,j_2}^\top \mathbf{h}) - \phi'(\mathbf{w}_{\ell,j_2}^{*T} \mathbf{h})| \cdot |\mathbf{a}^\top \mathbf{h}|^2 \cdot f_{\mathbf{Z}}(\mathbf{z}) \cdot |\mathbf{J}_{\mathbf{h}}(\mathbf{z})| dz
\end{aligned} \tag{120}$$

where $|\mathbf{J}_h(\mathbf{z})|$ 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 $|\mathbf{J}_h(\mathbf{z})| = 1$. According to (119), I_4 can be rewritten as

$$\begin{aligned} I_4 &= \int_{\mathcal{R}_z} |\phi'(\mathbf{w}_{\ell, j_2}^\top \tilde{\mathbf{h}}) - \phi'(\mathbf{w}_{\ell, j_2}^{*T} \tilde{\mathbf{h}})| \cdot |\mathbf{a}^\top \tilde{\mathbf{h}}|^2 \cdot f_Z(\mathbf{z}) d\mathbf{z} \\ &= \int_{\mathcal{R}_z} |\phi'(\mathbf{w}_{\ell, j_2}^\top \tilde{\mathbf{h}}) - \phi'(\mathbf{w}_{\ell, j_2}^{*T} \tilde{\mathbf{h}})| \cdot |\mathbf{a}^\top \tilde{\mathbf{h}}|^2 \cdot f_Z(z_1, z_2, z_3) dz_1 dz_2 dz_3 \end{aligned} \quad (121)$$

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}_z .

Next, we show that \mathbf{z} is rotational invariant over \mathcal{R}_z . Let $\mathbf{R} = [\mathbf{a} \ \mathbf{b} \ \mathbf{c} \ \dots \ \mathbf{a}_d^\perp]$, 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 $z_1 = R \cos \phi_1$, $z_2 = R \sin \phi_1 \sin \phi_2$, $z_3 = R \sin \phi_1 \cos \phi_2$. Hence, we have

$$I_4 = \int |\phi'(\mathbf{w}_{\ell, j_2}^\top \tilde{\mathbf{h}}) - \phi'(\mathbf{w}_{\ell, j_2}^{*T} \tilde{\mathbf{h}})| \cdot |R \cos \phi_1|^2 \cdot f_Z(R, \phi_1, \phi_2) \cdot R^2 \sin \phi_1 \cdot dR d\phi_1 d\phi_2. \quad (122)$$

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

$$f_Z(R, \phi_1, \phi_2) = f_Z(R). \quad (123)$$

Then, we have

$$\begin{aligned} I_4 &= \int |\phi'(\mathbf{w}_{\ell, j_2}^\top (\tilde{\mathbf{h}}/R)) - \phi'(\mathbf{w}_{\ell, j_2}^{*T} (\tilde{\mathbf{h}}/R))| \cdot |R \cos \phi_1|^2 \cdot f_Z(R) R^2 \sin \phi_1 dR d\phi_1 d\phi_2 \\ &= \int_0^\infty R^4 f_Z(R) dR \int_0^{\psi_1(R)} \int_0^{\psi_2(R)} |\cos \phi_1|^2 \cdot \sin \phi_1 \\ &\quad \cdot |\phi'(\mathbf{w}_{\ell, j_2}^\top (\tilde{\mathbf{h}}/R)) - \phi'(\mathbf{w}_{\ell, j_2}^{*T} (\tilde{\mathbf{h}}/R))| d\phi_1 d\phi_2 \\ &\leq \int_0^\infty R^4 f_Z(R) dR \int_0^\pi \int_0^{2\pi} \sin \phi_1 \cdot |\phi'(\mathbf{w}_{\ell, j_2}^\top \bar{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell, j_2}^{*T} \bar{\mathbf{x}})| d\phi_1 d\phi_2, \end{aligned} \quad (124)$$

where the first equality holds because $\phi'(\mathbf{w}_{\ell, j_2}^\top \mathbf{h})$ only depends on the direction of \mathbf{h} , and $\bar{\mathbf{x}} := \mathbf{h}/R = (\cos \phi_1, \sin \phi_1 \sin \phi_2, \sin \phi_1 \cos \phi_2)$ 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

$$\begin{aligned} \int_0^\infty R^4 f_Z(R) dR &= \int_0^\infty R^4 d(1 - F_z(R)) \\ &\leq \int_0^\infty 4R^3 (1 - F_z(R)) dR \\ &\leq \int_0^\infty 8R^3 e^{-\frac{R^2}{\sigma^2}} dR \\ &\leq \frac{32}{\sqrt{2\pi}} \sigma \int_0^\infty R^2 e^{-\frac{R^2}{\sigma^2}} dR \\ &= 32\sigma^2 \int_0^\infty R^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{R^2}{\sigma^2}} dR, \end{aligned} \quad (125)$$

where the last inequality comes from the calculation that

$$\begin{aligned} \int_0^\infty 2R^2 e^{-\frac{R^2}{\sigma^2}} dR &= \sqrt{2\pi}\sigma^3, \\ \int_0^\infty 2R^3 e^{-\frac{R^2}{\sigma^2}} dR &= 4\sigma^4. \end{aligned} \quad (126)$$

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

$$\begin{aligned} I_4 &\leq 32\sigma^2 \cdot \int_0^\infty R^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{R^2}{\sigma^2}} dR \int_0^\pi \int_0^{2\pi} \sin \phi_1 \cdot |\phi'(\mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})| d\phi_1 d\phi_2 \\ &= 32\sigma^2 \cdot \mathbb{E}_{z_1, z_2, z_3} |\phi'(\mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})| \\ &\approx \mathbb{E}_{\tilde{\mathbf{x}}} |\phi'(\mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})|, \end{aligned} \quad (127)$$

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 $\mathbb{E}_{\tilde{\mathbf{x}}} |\phi'(\mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})|$.

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

$$|\phi'(\mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})| = \begin{cases} 1, & \text{if } \tilde{\mathbf{x}} \in \mathcal{A}_1 \\ 0, & \text{if } \tilde{\mathbf{x}} \in \mathcal{A}_1^c \end{cases}. \quad (128)$$

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

$$|\mathbf{w}_{\ell,j_2}^{*T} \tilde{\mathbf{x}}| \leq |\mathbf{w}_{\ell,j_2}^{*T} \tilde{\mathbf{x}} - \mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}| \leq \|\mathbf{w}_{\ell,j_2}^* - \mathbf{w}_{\ell,j_2}\|_2 \cdot \|\tilde{\mathbf{x}}\|_2. \quad (129)$$

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

$$\begin{aligned} \mathcal{A}_2 &= \left\{ \tilde{\mathbf{x}} \left| \frac{|\mathbf{w}_{\ell,j_2}^{*T} \tilde{\mathbf{x}}|}{\|\mathbf{w}_{\ell,j_2}^*\|_2 \|\tilde{\mathbf{x}}\|_2} \leq \frac{\|\mathbf{w}_{\ell,j_2}^* - \mathbf{w}_{\ell,j_2}\|_2}{\|\mathbf{w}_{\ell,j_2}^*\|_2} \right. \right\} \\ &= \left\{ \theta_{\tilde{\mathbf{x}}, \mathbf{w}_{\ell,j_2}^*} \left| \left| \cos \theta_{\tilde{\mathbf{x}}, \mathbf{w}_{\ell,j_2}^*} \right| \leq \frac{\|\mathbf{w}_{\ell,j_2}^* - \mathbf{w}_{\ell,j_2}\|_2}{\|\mathbf{w}_{\ell,j_2}^*\|_2} \right. \right\}. \end{aligned} \quad (130)$$

Hence, we have that

$$\begin{aligned} \mathbb{E}_{\tilde{\mathbf{x}}} |\phi'(\mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})|^2 &= \mathbb{E}_{\tilde{\mathbf{x}}} |\phi'(\mathbf{w}_{\ell,j_2}^\top \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{\ell,j_2}^{*\top} \tilde{\mathbf{x}})| \\ &= \text{Prob}(\tilde{\mathbf{x}} \in \mathcal{A}_1) \\ &\leq \text{Prob}(\tilde{\mathbf{x}} \in \mathcal{A}_2). \end{aligned} \quad (131)$$

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

$$\begin{aligned} \text{Prob}(\tilde{\mathbf{x}} \in \mathcal{A}_2) &= \frac{\pi - \arccos \frac{\|\mathbf{w}_{\ell,j_2}^* - \mathbf{w}_{\ell,j_2}\|_2}{\|\mathbf{w}_{\ell,j_2}^*\|_2}}{\pi} \leq \frac{1}{\pi} \tan(\pi - \arccos \frac{\|\mathbf{w}_{\ell,j_2}^* - \mathbf{w}_{\ell,j_2}\|_2}{\|\mathbf{w}_{\ell,j_2}^*\|_2}) \\ &= \frac{1}{\pi} \cot(\arccos \frac{\|\mathbf{w}_{\ell,j_2}^* - \mathbf{w}_{\ell,j_2}\|_2}{\|\mathbf{w}_{\ell,j_2}^*\|_2}) \\ &\leq \frac{2}{\pi} \frac{\|\mathbf{w}_{\ell,j_2}^* - \mathbf{w}_{\ell,j_2}\|_2}{\|\mathbf{w}_{\ell,j_2}^*\|_2} \\ &\leq \|\mathbf{W}_\ell^* - \mathbf{W}_\ell\|_2 \end{aligned} \quad (132)$$

Hence, (124) and (132) suggest that

$$\begin{aligned} I_4 &\lesssim \|\mathbf{W}_i - \mathbf{W}_i^*\|_2 \cdot \|\mathbf{a}\|_2^2, \\ \text{and } \|\mathbf{P}_1\|_2 &\leq \|\mathbf{W} - \mathbf{W}^*\|_2 + I_4 \lesssim \|\mathbf{W} - \mathbf{W}^*\|_2, \end{aligned} \quad (133)$$

The same bound that is shown in (133) holds for \mathbf{P}_2 as well.

Therefore, we have

$$\begin{aligned}
\|\nabla_{\ell}^2 f(\mathbf{W}^*) - \nabla_{\ell}^2 f(\mathbf{W})\|_2 &= \max_{\|\boldsymbol{\alpha}\|_2 \leq 1} \left| \boldsymbol{\alpha}^\top \left(\nabla_{\ell}^2 f(\mathbf{W}^*) - \nabla_{\ell}^2 f(\mathbf{W}) \right) \boldsymbol{\alpha} \right| \\
&\leq \frac{1}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \|\mathbf{P}_1 + \mathbf{P}_2\|_2 \cdot \|\boldsymbol{\alpha}_{j_1}\|_2 \cdot \|\boldsymbol{\alpha}_{j_2}\|_2 \\
&\lesssim \frac{1}{K^2} \cdot \sum_{j_1=1}^K \sum_{j_2=1}^K \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \|\boldsymbol{\alpha}_{j_1}\|_2 \|\boldsymbol{\alpha}_{j_2}\|_2 \quad (134) \\
&\lesssim \frac{1}{K^2} \cdot \sum_{j_1=1}^K \sum_{j_2=1}^K \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \left(\frac{\|\boldsymbol{\alpha}_{j_1}\|_2^2 + \|\boldsymbol{\alpha}_{j_2}\|_2^2}{2} \right) \\
&\lesssim \frac{1}{K} \cdot \|\mathbf{W}^* - \mathbf{W}\|_2,
\end{aligned}$$

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

E.2 Proof of Lemma 7

We aim to prove that $\int_{\mathcal{R}} \left(\sum_{j=1}^K \boldsymbol{\alpha}^\top \mathbf{h} \phi'(\mathbf{w}_{\ell,j}^\top \mathbf{h}) \right)^2 p_H(\mathbf{h}) \cdot d\mathbf{h}$ is strictly greater than zero for any $\boldsymbol{\alpha}$. Therefore, the ρ in (2) is strictly greater than zero. The proof is inspired by Theorem 3.1 in [22]. It is obviously that $\left(\sum_{j=1}^K \boldsymbol{\alpha}^\top \mathbf{h} \phi'(\mathbf{w}_{\ell,j}^\top \mathbf{h}) \right)^2$ is greater or equal to zero. Given $\left(\sum_{j=1}^K \boldsymbol{\alpha}^\top \mathbf{h} \phi'(\mathbf{w}_{\ell,j}^\top \mathbf{h}) \right)^2$ is continuous, we only need to show that $\boldsymbol{\alpha}$ such that $\sum_{j=1}^K \boldsymbol{\alpha}^\top \mathbf{h} \phi'(\mathbf{w}_{\ell,j}^\top \mathbf{h}) \neq 0$ for any $\boldsymbol{\alpha}$, namely, $\{\mathbf{h} \phi'(\mathbf{w}_{\ell,j}^\top \mathbf{h})\}_{j=1}^K$ are linear independent. Compared with Theorem 3.1 in [22], we need to address two challenges: (1) the neuron weights \mathbf{w} is the random variable in [22] while the input \mathbf{h} is the random variable in this paper and (2) the random variable belongs to Gaussian distribution in [22] while the random variable belongs to sub-Gaussian distribution in this paper.

Proof of Lemma 7. Let \mathcal{H} be a Hilbert space on $\mathbb{R}^{K\ell}$, and the inner product of \mathcal{H} is defined as

$$\langle f, g \rangle = \int_{\mathcal{R}} f(\mathbf{h})^\top g(\mathbf{h}) f_H(\mathbf{h}) \cdot d\mathbf{h}, \quad \forall f, g \in \mathcal{H}, \quad (135)$$

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

$$\boldsymbol{\alpha}^\top \mathbf{h} \phi'(\mathbf{w}_k^\top \mathbf{h}) \neq 0 \quad \text{almost everywhere.} \quad (136)$$

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

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

We define a group of functions $\{\psi_j(\mathbf{h})\}_{j=1}^K$, where $\psi_j(\mathbf{h}) = \mathbf{h} \phi'(\mathbf{w}_j^\top \mathbf{h})$. From the assumption in Lemma 7, 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 $\mathcal{X}_i = \{\mathbf{h} \mid \mathbf{w}_i^\top \mathbf{h} = 0\}$ for any $i \in [K]$. For any fixed k , we can justify that \mathcal{X}_k cannot be covered by other sets $\{\mathcal{X}_k\}_{j \neq k}$ as long as \mathbf{w}_k does not parallel to any other weights \mathbf{w}_j with $j \neq k$. Namely, $\mathcal{X}_k \not\subseteq \cup_{j \neq k} \mathcal{X}_j$. 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 λ be the Lebesgue measure on \mathcal{X}_k , then $\lambda(\mathcal{X}_k) > 0$. When \mathbf{w}_j does not parallel to \mathbf{w}_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

$$\lambda(\mathcal{X}_k \cap (\cup_{j \neq k} \mathcal{X}_j)) \leq \sum_{j \neq k} \lambda(\mathcal{X}_k \cap \mathcal{X}_j) = 0. \quad (137)$$

Therefore, we have

$$\lambda(\mathcal{X}_k / (\cup_{j \neq k} \mathcal{X}_j)) = \lambda(\mathcal{X}_k) - \lambda(\mathcal{X}_k \cap (\cup_{j \neq k} \mathcal{X}_j)) = \lambda(\mathcal{X}_k) > 0. \quad (138)$$

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

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

$$\mathcal{B}(\mathbf{z}_k, r) \cap \mathcal{D}_j = \emptyset \quad \text{with} \quad \forall r \leq r_0 \text{ and } j \neq k, \quad (139)$$

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

$$\begin{aligned} \mathcal{B}_r^+ &= \mathcal{B}(\mathbf{z}_k, r) \cap \{\mathbf{h} \mid \mathbf{w}_k^\top \mathbf{h} > 0\}, \\ \mathcal{B}_r^- &= \mathcal{B}(\mathbf{z}_k, r) \cap \{\mathbf{h} \mid \mathbf{w}_k^\top \mathbf{h} < 0\}. \end{aligned} \quad (140)$$

Because \mathbf{z}_k is a boundary point of $\{\mathbf{h} \mid \mathbf{w}_k^\top \mathbf{h} = 0\}$, 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 $\phi_j(\mathbf{w}_k^\top \mathbf{h})$ is continuous at point \mathbf{z}_k since $\mathbf{z}_k \notin \mathcal{X}_j$. Hence, it is easy to verify that

$$\lim_{r \rightarrow 0^+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \psi_k(\mathbf{h}) d\mathbf{h} = \lim_{r \rightarrow 0^-} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \psi_k(\mathbf{h}) d\mathbf{h} = \psi_k(\mathbf{z}_k). \quad (141)$$

While for ψ_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

$$\begin{aligned} \lim_{r \rightarrow 0^+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \psi_k(\mathbf{h}) d\mathbf{h} &= \mathbf{z}_k \\ \lim_{r \rightarrow 0^-} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \psi_k(\mathbf{h}) d\mathbf{h} &= 0. \end{aligned} \quad (142)$$

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

$$\sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) \equiv 0, \quad \forall \mathbf{h}. \quad (143)$$

Then, we have

$$\begin{aligned} \lim_{r \rightarrow 0^+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} &= 0 \\ \lim_{r \rightarrow 0^+} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} &= 0 \end{aligned} \quad (144)$$

Then, we have

$$\begin{aligned} 0 &= \lim_{r \rightarrow 0^+} \frac{1}{\lambda(\mathcal{B}_r^+)} \int_{\mathcal{B}_r^+} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} - \lim_{r \rightarrow 0^+} \frac{1}{\lambda(\mathcal{B}_r^-)} \int_{\mathcal{B}_r^-} \sum_{j=1}^K \alpha_j \psi_j(\mathbf{h}) d\mathbf{h} \\ &= \alpha_k \mathbf{z}_k \end{aligned} \quad (145)$$

where the last equality comes from (141) and (142).

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 (145), 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 (143) 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 $\int_{\mathcal{R}} \left(\sum_{j=1}^K \alpha^\top \mathbf{h} \phi'(\mathbf{w}_{\ell,j}^\top \mathbf{h}) \right)^2 p_H(\mathbf{h}) \cdot d\mathbf{h}$ is strictly greater than zero. \square

E.3 Proof of Lemma 8

Proof of Lemma 8. From the definition of (39), we have

$$\begin{aligned}
& \|\mathbf{h}^{(\ell)}(\mathbf{W}) - \mathbf{h}^{(\ell)}(\mathbf{W}^*)\|_2 \\
&= \|\phi(\mathbf{W}_{\ell-1}^\top \mathbf{h}^{(\ell-1)}(\mathbf{W})) - \phi(\mathbf{W}_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\mathbf{W}^*))\|_2 \\
&= \|\phi(\mathbf{W}_{\ell-1}^\top \mathbf{h}^{(\ell-1)}(\mathbf{W})) - \phi(\mathbf{W}_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\mathbf{W})) \\
&\quad + \phi(\mathbf{W}_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\mathbf{W})) - \phi(\mathbf{W}_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\mathbf{W}^*))\|_2 \\
&\leq \|\phi(\mathbf{W}_{\ell-1}^\top \mathbf{h}^{(\ell-1)}(\mathbf{W})) - \phi(\mathbf{W}_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\mathbf{W}))\|_2 \\
&\quad + \|\phi(\mathbf{W}_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\mathbf{W})) - \phi(\mathbf{W}_{\ell-1}^{*\top} \mathbf{h}^{(\ell-1)}(\mathbf{W}^*))\|_2 \\
&\leq \|\mathbf{W}_{\ell-1} - \mathbf{W}_{\ell-1}^*\|_2 \cdot \|\mathbf{h}^{(\ell-1)}(\mathbf{W})\|_2 + \|\mathbf{h}^{(\ell-1)}(\mathbf{W}) - \mathbf{h}^{(\ell-1)}(\mathbf{W}^*)\|_2.
\end{aligned} \tag{146}$$

With the assumption in the Lemma 8 such that \mathbf{W} is close enough to \mathbf{W}^* , we have

$$\|\mathbf{W}_i\|_2 \leq \|\mathbf{W}_i^*\|_2 + \|\mathbf{W}_i - \mathbf{W}_i^*\|_2 \lesssim 1. \tag{147}$$

Therefore, we have

$$\|\mathbf{h}^{(i)}(\mathbf{W})\|_2 \leq \|\mathbf{W}_i\|_2 \cdots \|\mathbf{W}_1\|_2 \cdot \|\mathbf{x}\|_2 \lesssim \|\mathbf{x}\|_2. \tag{148}$$

Then, we have

$$\begin{aligned}
& \|\mathbf{h}^{(\ell)}(\mathbf{W}) - \mathbf{h}^{(\ell)}(\mathbf{W}^*)\|_2 \\
&\leq \|\mathbf{W}_{\ell-1} - \mathbf{W}_{\ell-1}^*\|_2 \cdot \|\mathbf{x}\|_2 + \|\mathbf{h}^{(\ell-1)}(\mathbf{W}) - \mathbf{h}^{(\ell-1)}(\mathbf{W}^*)\|_2 \\
&\leq \sum_{i=1}^{\ell-1} \|\mathbf{W}_i - \mathbf{W}_i^*\|_2 \cdot \|\mathbf{x}\|_2 + \|\mathbf{h}^{(1)}(\mathbf{W}) - \mathbf{h}^{(1)}(\mathbf{W}^*)\|_2 \\
&= \sum_{i=1}^{\ell-1} \|\mathbf{W}_i - \mathbf{W}_i^*\|_2 \cdot \|\mathbf{x}\|_2 + \|\mathbf{x} - \mathbf{x}\|_2 \\
&= \sum_{i=1}^{\ell-1} \|\mathbf{W}_i - \mathbf{W}_i^*\|_2 \cdot \|\mathbf{h}^{(i-1)}(\mathbf{W})\|_2 \\
&\leq \|\mathbf{W} - \mathbf{W}^*\|_2 \cdot \|\mathbf{x}\|_2,
\end{aligned} \tag{149}$$

which completes the proof. \square

F Additional experiments

In this section, we provide numerical justification that our theoretical findings are aligned with DDQN through the Atari Breakout game. The neural network follows the same architecture as the one used in Section 5. The algorithm terminates if the average score over the recent 100 episodes does not improve or the algorithm reaches the maximum episode set as 200, which is around 4×10^5 training steps. The testing score is calculated based on a similar setup as the training process by fixing the maximum memory size N as 2000 and greedy policy, i.e., $\varepsilon = 0$. Each point in the plot is averaged over 10 experiments with an error bar representing the standard deviation.

Estimation errors with respect to the sample complexity N . We follow the setup in Section 5 to use the expected cumulative reward as the estimation error of the learned model to the optimal Q-value function. The ε_t in ε -greedy policy decreases geometrically from 1 to 0.01. We vary the number of samples in the replay buffer from 3000 to 10000. Figure 4 shows that the test error is almost linear in $1/\sqrt{N}$, which is consistent with our characterization in (20). In addition, experiments with a large N have a shorter error bar indicating a more stable learning performance with a large sample complexity as shown in (12).

Convergence with different selections of ε . Figure 5 illustrates the convergence rate when ε_t in the ε -greedy policy changes. For each point, ε_0 is selected as the value in the x-axis, and we decrease ε_t geometrically as the iteration t increases. Each point is averaged over 10 independent trials. We can see that the convergence rate is a linear function of c_ε , matching our findings in (19).

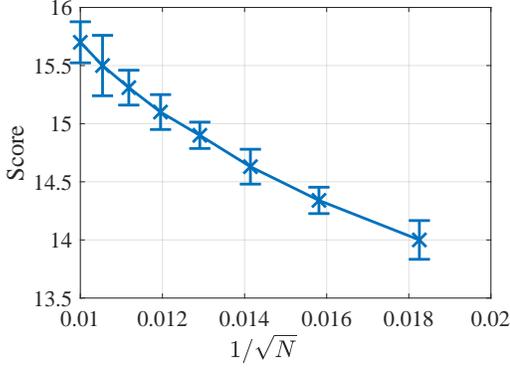


Figure 4: Test error in scores against the number of samples.

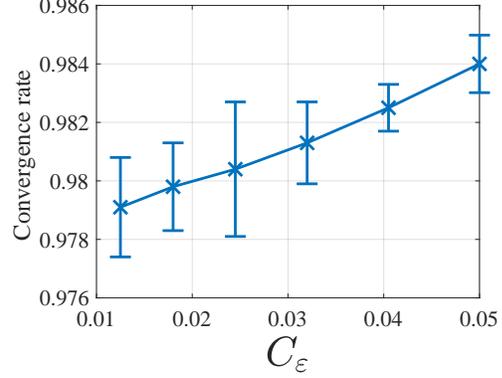


Figure 5: The convergence rate against the value of c_ε .

G Extension to non i.i.d. samples

Assumption 3. At any fixed outer iteration t , the behavior policy π_t and transition kernel \mathcal{P}_t satisfy

$$\sup_{\mathbf{s} \in \mathcal{S}} d_{TV}(\mathbb{P}(\mathbf{s}_\tau \in \cdot) \mid \mathbf{s}_0 = \mathbf{s}, \mathcal{P}_t) \leq \lambda \nu^\tau, \quad \forall \tau \geq 0 \quad (150)$$

for some constant $\lambda > 0$ and $\nu \in (0, 1)$, where d_{TV} denotes the total-variation distance between the probability measures.

Assumption 3 assumes the Markov chain $\{\mathbf{s}_t\}$ induced by the behavior policy, i.e., ε_t -greedy policy at t -th outer loop, is uniformly ergodic with the corresponding invariant measure \mathcal{P}_t . Compared with i.i.d. cases, we need to handle an additional error term when bounding the distance between the g_t and ∇f as shown in (91). Therefore, the upper bound in Lemma 3 changes, which suggests an additional term in the final bound.

We present the major theoretical findings for non-i.i.d. samples in Theorem 2. The major proofs in this context follow similar steps to the proof of Theorem 1, with slight changes in the error bound between the sequences g_t and ∇f . In this section, we omit the details of the proof for Theorem 2 but provide the proof for Lemma 3 under the assumptions outlined in Assumption 2 to simplify the presentation.

Theorem 2 (Convergence for non-i.i.d. case). *Suppose Assumption 1 and (143) hold, the buffer size N satisfies (13). Let us define C_{\max} be a constant that is larger than C_t for $1 \leq t \leq T$ and $C_d = |\mathcal{A}| \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1-\nu})$, when ε_t satisfy*

$$\varepsilon_t = \frac{c_\varepsilon \cdot \Theta(\sqrt{N}) \cdot e_t}{(1 - C_{\max}) \cdot C_d \cdot R_{\max}} - \frac{C_{\max}}{1 - C_{\max}} \quad (151)$$

for a fixed constant $c_\varepsilon \in (0, (1 - \gamma)^2]$, and the initialization satisfies

$$\|\mathbf{W}^{(0,0)} - \mathbf{W}^*\|_F \leq \mathcal{O}\left(1 - \frac{1 - c_\varepsilon}{\Theta(\sqrt{N})}\right) \cdot \frac{\rho \cdot \|\mathbf{W}^*\|_F}{K}. \quad (152)$$

Then, with the high probability of at least $1 - T \cdot q^{-d}$, we have

(C1) The learned weights decay geometrically with

$$\|\mathbf{W}^{(t+1,0)} - \mathbf{W}^*\|_F \leq (\gamma + c_\varepsilon \cdot (1 - \gamma)) \cdot \|\mathbf{W}^{(t,0)} - \mathbf{W}^*\|_F + \frac{(2 + \gamma)R_{\max}\tau^*}{(1 - \gamma)\Theta(N)}, \quad (153)$$

(C2) the returned model $Q(\mathbf{W}^{(T,0)})$ exhibits an estimation error as

$$\sup_{(s,a)} |Q(\mathbf{W}^{(T,0)}) - Q^*| \leq \frac{C_{\max} \cdot C_d \cdot R_{\max}}{(1 - \gamma)^2 \cdot \Theta(\sqrt{N} \cdot T)} + \frac{(2 + \gamma)R_{\max}\tau^*}{(1 - \gamma)\Theta(N \cdot T)}, \quad (154)$$

where $\tau^* = \min\{t \mid \lambda \nu^t \leq 1/(N \cdot T)\}$.

Proof of Lemma 3 under Assumption 2. Recall that in (91), we have

$$\begin{aligned}
& g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) \\
&= \left[\frac{1}{N} \sum_{n=1}^N \left(Q(\mathbf{W}; \mathbf{s}_n, a_n) - Q(\mathbf{W}^*; \mathbf{s}_n, a_n) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \right. \\
&\quad \left. - \mathbb{E}_{(s,a) \sim \mathcal{D}_t} \left(Q(\mathbf{W}; \mathbf{s}, a) - Q(\mathbf{W}^*; \mathbf{s}, a) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}, a)}{\partial \mathbf{w}_{\ell,k}} \right] \\
&\quad + \left[\frac{1}{N} \sum_{n=1}^N \gamma \cdot \left(\max_a Q(\mathbf{s}_n, a; \mathbf{W}^*) - \max_a Q(\mathbf{s}_n, a; \mathbf{W}^{(t,0)}) \right) \cdot \frac{\partial Q(\mathbf{W}; \mathbf{s}_n, a_n)}{\partial \mathbf{w}_{\ell,k}} \right] \\
&\quad + \mathbb{E}_{(s,a) \sim \mu_t} g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \frac{\partial f}{\partial \mathbf{w}_{\ell,k}}(\mathbf{W}) \\
&\quad + \mathbb{E}_{(s,a) \sim \mathcal{D}_t, \mathcal{P}} [g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \mathbb{E}_{(s,a) \sim \mu_t, \mathcal{P}} g_t(\mathbf{w}_{\ell,k}; \mathbf{W})] \\
&:= I_1 + I_2 + I_3 + I_4.
\end{aligned} \tag{155}$$

Bound of I_1 and I_2 . Compared with (91), the upper bound for I_1 and I_2 is the same as those shown in (103) and (109), respectively.

Bound of I_3 . Following (111), the upper bound of I_3 can be characterized as

$$\begin{aligned}
\|I_3\|_2 &\leq \frac{R_{\max}}{1-\gamma} \cdot \left[(1-\varepsilon) \cdot \left| \int_{(s,a)} \int_{s'} (\mu^*(ds, da) \mathcal{P}(ds'|s, a) - \mu_{t,1}(ds, da) \mathcal{P}(ds'|s, a)) \right| \right. \\
&\quad \left. + \varepsilon \cdot \left| \int_{(s,a)} \int_{s'} (\mu^*(ds, da) \mathcal{P}(ds'|s, a) - \mu_{t,2}(ds, da) \mathcal{P}(ds'|s, a)) \right| \right].
\end{aligned} \tag{156}$$

and

$$\begin{aligned}
& \left| \int_{(s,a)} \int_{s'} (\mu^*(ds, da) \mathcal{P}(ds'|s, a) - \mu_{t,1}(ds, da) \mathcal{P}(ds'|s, a)) \right| \\
&= \left| \int_{(s,a)} \int_{s'} (\mathcal{P}^*(ds) \pi^*(da|s) \mathcal{P}(ds'|s, a) - \mathcal{P}_{t,1}(ds) \pi_{t,1}(da|ds) \mathcal{P}(ds'|s, a)) \right| \\
&\leq \left| \int_{(s,a)} \int_{s'} (\mathcal{P}^*(ds) - \mathcal{P}_{t,1}(ds)) \pi^*(da|s) \mathcal{P}(ds'|s, a) \right| \\
&\quad + \left| \int_{(s,a)} \int_{s'} \mathcal{P}_{t,1}(ds) (\pi_{t,1}(da|ds) - \pi^*(da|ds)) \mathcal{P}(ds'|s, a) \right|.
\end{aligned} \tag{157}$$

From Theorem 3.1 in [49], we know that

$$\begin{aligned}
\left| \int_{(s,a)} (\mathcal{P}^*(ds) - \mathcal{P}_{t,1}(ds)) \right| &\leq |\mathcal{A}| (\log_\nu \lambda^{-1} + \frac{1}{1-\nu}) C_t \\
\text{and} \quad \|\pi_{t,1}(da|ds) - \pi^*(da|ds)\| &\leq C_t.
\end{aligned} \tag{158}$$

Therefore, the bound of I_3 can be found as

$$\begin{aligned}
\|I_3\|_2 &\leq \frac{R_{\max}}{1-\gamma} \cdot |\mathcal{A}| \cdot ((1-\varepsilon)C_t + \varepsilon \cdot C_t) \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1-\nu}) \\
&= C_d \cdot (C_t + (1-C_t)\varepsilon) \cdot \frac{R_{\max}}{1-\gamma},
\end{aligned} \tag{159}$$

where $C_d = |\mathcal{A}| \cdot (1 + \log_\nu \lambda^{-1} + \frac{1}{1-\nu})$.

Bound of I_4 . I_4 is the bias of the data because the data (s, a) at iteration t depends on the neural network parameters \mathbf{W} . Let us define \bar{g}_t as

$$\bar{g}_t(\mathbf{w}_{\ell,k}; \mathbf{W}) = \mathbb{E}_{\mu_t, \mathcal{P}} g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) \tag{160}$$

and

$$\Delta_t = g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \bar{g}_t(\mathbf{w}_{\ell,k}; \mathbf{W}). \tag{161}$$

It is easy to verify that

$$\begin{aligned} \|g_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - g_t(\tilde{\mathbf{w}}_{\ell,k}; \tilde{\mathbf{W}})\| &\leq (1 + \gamma) \cdot \|\mathbf{W} - \tilde{\mathbf{W}}\|, \\ \|\bar{g}_t(\mathbf{w}_{\ell,k}; \mathbf{W}) - \bar{g}_t(\tilde{\mathbf{w}}_{\ell,k}; \tilde{\mathbf{W}})\| &\leq (1 + \gamma) \cdot \|\mathbf{W} - \tilde{\mathbf{W}}\|, \\ \text{and} \quad \|g_t\| &\lesssim \frac{R_{\max}}{1 - \gamma}. \end{aligned} \quad (162)$$

Then, we have

$$\Delta_t(\mathbf{W}) - \Delta_t(\tilde{\mathbf{W}}) \lesssim (1 + \gamma) \cdot \|\mathbf{W} - \tilde{\mathbf{W}}\|_2. \quad (163)$$

Therefore, we have

$$\Delta_t(\mathbf{W}^{(t,0)}) \leq \Delta_t(\mathbf{W}^{(t-\tau,0)}) + \frac{1 + \gamma}{1 - \gamma} \cdot R_{\max} \cdot \sum_{i=t-\tau}^{t-1} \eta_i. \quad (164)$$

Then, we need to bound $\delta_t(\mathbf{W}^{(t-\tau,0)})$.

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

$$\mathbf{W}^{(t-\tau,0)} \longrightarrow \mathbf{s}_{t-\tau} \longrightarrow \mathbf{s}_t \longrightarrow O_t \quad (165)$$

forms a Markov chain introduced by the policy $\pi_{t-\tau}$.

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

With Lemma 9 in [4], we have

$$\mathbb{E} \Delta_t(\mathbf{W}^{(t-\tau,0)}, O_t) - \mathbb{E} \Delta_t(\tilde{\mathbf{W}}^{(t-\tau,0)}, \tilde{O}_t) \lesssim 2 \sup_{\mathbf{w}, O} |\Delta_t(\mathbf{W}, O)| \cdot \lambda \cdot \nu^\tau. \quad (166)$$

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

$$|\Delta_t(\mathbf{W}, O)| \leq \frac{2 R_{\max}}{1 - \gamma}. \quad (167)$$

Therefore, we have

$$\begin{aligned} \mathbb{E} \Delta_t(\mathbf{W}^{(t,0)}) &\leq \mathbb{E} \Delta_t(\mathbf{W}^{(t-\tau,0)}) + \frac{1 + \gamma}{1 - \gamma} \cdot R_{\max} \cdot \sum_{i=t-\tau}^{t-1} \eta_i \\ &\leq \frac{R_{\max}}{1 - \gamma} \left(\lambda \cdot \nu^\tau + (1 + \gamma) \cdot \tau \cdot \eta_{t-\tau} \right), \end{aligned} \quad (168)$$

where the last inequality comes from the fact that the step size η_t 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

$$\mathbb{E} \Delta_t(\mathbf{W}^{(t,0)}) \leq \frac{R_{\max}}{1 - \gamma} \cdot \tau^* \cdot \eta_0. \quad (169)$$

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

$$\mathbb{E} \Delta_t(\mathbf{W}^{(t,0)}) \leq \frac{R_{\max}}{1 - \gamma} \cdot (1 + \gamma) \tau^* \cdot \eta_{t-\tau^*}. \quad (170)$$

Combining (169) and (170), we have

$$|I_4| \leq \frac{R_{\max}}{1 - \gamma} \cdot (1 + \gamma) \tau^* \cdot \eta_{\max\{0, t-\tau^*\}}, \quad (171)$$

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