

## A. Proof of Theorem 1

In this section, before presenting the proof of Theorem 1, we start with defining some useful notations. Recall that in (3), the empirical risk function for linear regression problem is defined as

$$\min_{\mathbf{W}} : \hat{f}_{\Omega}(\mathbf{W}) = \frac{1}{2|\Omega|} \sum_{n \in \Omega} \left| y_n - g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X}) \right|^2. \quad (22)$$

Population risk function, which is the expectation of the empirical risk function, is defined as

$$\min_{\mathbf{W}} : f_{\Omega}(\mathbf{W}) = \mathbb{E}_{\mathbf{X}} \frac{1}{2|\Omega|} \sum_{n \in \Omega} \left| y_n - g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X}) \right|^2. \quad (23)$$

Then, the road-map of the proof can be summarized in the following three steps.

First, we show the Hessian matrix of the population risk function  $f_{\Omega_t}$  is positive-definite at ground-truth parameters  $\mathbf{W}^*$  and then characterize the local convexity region of  $f_{\Omega_t}$  near  $\mathbf{W}^*$ , which is summarized in Lemma 2.

Second,  $\hat{f}_{\Omega_t}$  is non-smooth because of ReLU activation, but  $f_{\Omega_t}$  is smooth. Hence, we characterize the gradient descent term as  $\nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) = \langle \nabla^2 f_{\Omega_t}(\widehat{\mathbf{W}}^{(t)}), \mathbf{W}^{(t)} - \mathbf{W}^* \rangle + (\hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) - f_{\Omega_t}(\mathbf{W}^{(t)}))$ . During this step, we need to apply concentration theorem to bound  $\nabla \hat{f}_{\Omega_t}$  to its expectation  $\nabla f_{\Omega_t}$ , which is summarized in Lemma 3.

Third, we take the momentum term of  $\beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)})$  into consideration and obtain the following recursive rule:

$$\begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix} = \mathbf{L}(\beta) \begin{bmatrix} \mathbf{W}^{(t)} - \mathbf{W}^* \\ \mathbf{W}^{(t-1)} - \mathbf{W}^* \end{bmatrix}. \quad (24)$$

Then, we know iterates  $\mathbf{W}^{(t)}$  converge to the ground-truth with a linear rate which is the largest singular value of matrix  $\mathbf{L}(\beta)$ . Recall that AGD reduces to GD with  $\beta = 0$ , so our analysis applies to GD method as well. We are able to show the convergence rate of AGD is faster than GD by proving the largest singular value of  $\mathbf{L}(\beta)$  is smaller than  $\mathbf{L}(0)$  for some  $\beta > 0$ . Lemma 4 provides the estimation error of  $\mathbf{W}^{(0)}$  and sample complexity to guarantee  $\|\mathbf{L}(\beta)\|_2$  is less than 1 for  $t = 0$ .

**Lemma 2.** Let  $f_{\Omega_t}$  be the population risk function in (23) for regression problems, then for any  $\mathbf{W}$  that satisfies

$$\|\mathbf{W}^* - \mathbf{W}\|_2 \leq \frac{\varepsilon_0 \sigma_K}{44\kappa^2 \gamma K^2}, \quad (25)$$

the second-order derivative of  $f_{\Omega_t}$  is bounded as

$$\frac{(1 - \varepsilon_0) \sigma_1^2(\mathbf{A})}{11\kappa^2 \gamma K^2} \mathbf{I} \preceq \nabla^2 f_{\Omega_t}(\mathbf{W}) \preceq \frac{4\sigma_1^2(\mathbf{A})}{K} \mathbf{I}. \quad (26)$$

**Lemma 3.** Let  $\hat{f}_{\Omega_t}$  and  $f_{\Omega_t}$  be the empirical and population risk functions in (22) and (23) for regression problems, respectively. Then, for any fixed point  $\mathbf{W}$  satisfies (25), we have<sup>6</sup>

$$\left\| \nabla f_{\Omega_t}(\mathbf{W}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W}) \right\|_2 \lesssim \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2) d \log N}{|\Omega_t|}} \|\mathbf{W} - \mathbf{W}^*\|_2, \quad (27)$$

with probability at least  $1 - K^2 \cdot N^{-10}$ .

**Lemma 4.** Assume the number of samples  $|\Omega_t| \gtrsim \kappa^3 (1 + \delta^2) \sigma_1^4(\mathbf{A}) K d \log^4 N$ , the tensor initialization method via Subroutine 1 outputs  $\mathbf{W}^{(0)}$  such that

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_2 \lesssim \kappa^6 \sigma_1^2(\mathbf{A}) \sqrt{\frac{K^4 (1 + \delta^2) d \log N}{|\Omega_t|}} \|\mathbf{W}^*\|_2 \quad (28)$$

with probability at least  $1 - N^{-10}$ .

<sup>6</sup>We use  $f(d) \gtrsim$  ( or  $\lesssim, \approx$ )  $g(d)$  to denote there exists some positive constant  $C$  such that  $f(d) \geq$  ( or  $\leq, =$ )  $C \cdot g(d)$  when  $d$  is sufficiently large.

The proofs of Lemmas 2 and 3 are included in Appendix A.1 and A.2, respectively, while the proof of Lemma 4 can be found in Appendix D. With these three preliminary lemmas on hand, the proof of Theorem 1 is formally summarized in the following contents.

*Proof of Theorem 1.* The update rule of  $\mathbf{W}^{(t)}$  is

$$\begin{aligned}\mathbf{W}^{(t+1)} &= \mathbf{W}^{(t)} - \eta \nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) \\ &= \mathbf{W}^{(t)} - \eta \nabla f_{\Omega_t}(\mathbf{W}^{(t)}) + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) + \eta(\nabla f_{\Omega_t}(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)})).\end{aligned}\quad (29)$$

Since  $\nabla_{\Omega_t}^2$  is a smooth function, by the intermediate value theorem, we have

$$\begin{aligned}\mathbf{W}^{(t+1)} &= \mathbf{W}^{(t)} - \eta \nabla^2 f_{\Omega_t}(\widehat{\mathbf{W}}^{(t)})(\mathbf{W}^{(t)} - \mathbf{W}^*) \\ &\quad + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) \\ &\quad + \eta(\nabla f_{\Omega_t}(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)})),\end{aligned}\quad (30)$$

where  $\widehat{\mathbf{W}}^{(t)}$  lies in the convex hull of  $\mathbf{W}^{(t)}$  and  $\mathbf{W}^*$ .

Next, we have

$$\begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \eta \nabla^2 f_{\Omega_t}(\widehat{\mathbf{W}}^{(t)}) + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{W}^{(t)} - \mathbf{W}^* \\ \mathbf{W}^{(t-1)} - \mathbf{W}^* \end{bmatrix} + \eta \begin{bmatrix} \nabla f_{\Omega_t}(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) \\ \mathbf{0} \end{bmatrix}.\quad (31)$$

Let  $\mathbf{L}(\beta) = \begin{bmatrix} \mathbf{I} - \eta \nabla^2 f_{\Omega_t}(\widehat{\mathbf{W}}^{(t)}) + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}$ , so we have

$$\left\| \begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix} \right\|_2 = \|\mathbf{L}(\beta)\|_2 \left\| \begin{bmatrix} \mathbf{W}^{(t)} - \mathbf{W}^* \\ \mathbf{W}^{(t-1)} - \mathbf{W}^* \end{bmatrix} \right\|_2 + \eta \left\| \begin{bmatrix} \nabla f_{\Omega_t}(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) \\ \mathbf{0} \end{bmatrix} \right\|_2.$$

From Lemma 3, we know that

$$\eta \left\| \nabla f_{\Omega_t}(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) \right\|_2 \lesssim \eta \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \|\mathbf{W} - \mathbf{W}^*\|_2.\quad (32)$$

Then, we have

$$\begin{aligned}\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 &\lesssim \left( \|\mathbf{L}(\beta)\|_2 + \eta \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \right) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2 \\ &:\approx \nu(\beta) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2.\end{aligned}\quad (33)$$

Let  $\nabla^2 f(\widehat{\mathbf{W}}^{(t)}) = \mathbf{S} \mathbf{\Lambda} \mathbf{S}^T$  be the eigen-decomposition of  $\nabla^2 f(\widehat{\mathbf{W}}^{(t)})$ . Then, we define

$$\tilde{\mathbf{L}}(\beta) := \begin{bmatrix} \mathbf{S}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^T \end{bmatrix} \mathbf{L}(\beta) \begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \eta \mathbf{\Lambda} + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}.\quad (34)$$

Since  $\begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{S}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^T \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$ , we know  $\mathbf{L}(\beta)$  and  $\tilde{\mathbf{L}}(\beta)$  share the same eigenvalues. Let  $\lambda_i$  be the  $i$ -th eigenvalue of  $\nabla^2 f_{\Omega_t}(\widehat{\mathbf{W}}^{(t)})$ , then the corresponding  $i$ -th eigenvalue of  $\mathbf{L}(\beta)$ , denoted by  $\delta_i(\beta)$ , satisfies

$$\delta_i^2 - (1 - \eta \lambda_i + \beta) \delta_i + \beta = 0.\quad (35)$$

Then, we have

$$\delta_i(\beta) = \frac{(1 - \eta \lambda_i + \beta) + \sqrt{(1 - \eta \lambda_i + \beta)^2 - 4\beta}}{2},\quad (36)$$

and

$$|\delta_i(\beta)| = \begin{cases} \sqrt{\beta}, & \text{if } \beta \geq (1 - \sqrt{\eta\lambda_i})^2, \\ \frac{1}{2} \left| (1 - \eta\lambda_i + \beta) + \sqrt{(1 - \eta\lambda_i + \beta)^2 - 4\beta} \right|, & \text{otherwise.} \end{cases} \quad (37)$$

Note that the other root of (35) is abandoned because the root in (36) is always no less than the other root with  $|1 - \eta\lambda_i| < 1$ . By simple calculations, we have

$$\delta_i(0) > \delta_i(\beta), \quad \text{for } \forall \beta \in (0, (1 - \eta\lambda_i)^2). \quad (38)$$

Moreover,  $\delta_i$  achieves the minimum  $\delta_i^* = |1 - \sqrt{\eta\lambda_i}|$  when  $\beta = (1 - \sqrt{\eta\lambda_i})^2$ .

Let us first assume  $\mathbf{W}^{(t)}$  satisfies (25), then from Lemma 2, we know that

$$0 < \frac{(1 - \varepsilon_0)\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2} \leq \lambda_i \leq \frac{4\sigma_1^2(\mathbf{A})}{K}.$$

Let  $\gamma_1 = \frac{(1 - \varepsilon_0)\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2}$  and  $\gamma_2 = \frac{4\sigma_1^2(\mathbf{A})}{K}$ . If we choose  $\beta$  such that

$$\beta^* = \max \{(1 - \sqrt{\eta\gamma_1})^2, (1 - \sqrt{\eta\gamma_2})^2\}, \quad (39)$$

then we have  $\beta \geq (1 - \sqrt{\eta\lambda_i})^2$  and  $\delta_i = \max \{|1 - \sqrt{\eta\gamma_1}|, |1 - \sqrt{\eta\gamma_2}|\}$  for any  $i$ .

Let  $\eta = \frac{1}{2\gamma_2}$ , then  $\beta^*$  equals to  $(1 - \sqrt{\frac{\gamma_1}{2\gamma_2}})^2$ . Then, for any  $\varepsilon_0 \in (0, 1/2)$ , we have

$$\|\mathbf{L}(\beta^*)\|_2 = \max_i \delta_i(\beta^*) = 1 - \sqrt{\frac{\gamma_1}{2\gamma_2}} = 1 - \sqrt{\frac{1 - \varepsilon_0}{88\kappa^2\gamma K}} \leq 1 - \frac{1 - (3/4) \cdot \varepsilon_0}{\sqrt{88\kappa^2\gamma K}}. \quad (40)$$

Then, let

$$\eta\sigma_1^2(\mathbf{A})\sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \lesssim \frac{\varepsilon_0}{4\sqrt{88\kappa^2\gamma K}}, \quad (41)$$

we need  $|\Omega_t| \gtrsim \varepsilon_0^{-2}\kappa^2\gamma M(1 + \delta^2)\sigma_1^2(\mathbf{A})K^3d \log N$ . Combining (40) and (41), we have

$$\nu(\beta^*) \leq 1 - \frac{1 - \varepsilon_0}{\sqrt{88\kappa^2\gamma K}}. \quad (42)$$

Let  $\beta = 0$ , we have

$$\begin{aligned} \nu(0) &\geq \|\mathbf{A}(0)\|_2 = 1 - \frac{1 - \varepsilon_0}{88\kappa^2\gamma K}, \\ \nu(0) &\lesssim \|\mathbf{A}(0)\|_2 + \eta\sigma_1^2(\mathbf{A})\sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \leq 1 - \frac{1 - 2\varepsilon_0}{88\kappa^2\gamma K} \end{aligned}$$

if  $|\Omega_t| \gtrsim \varepsilon_0^{-2}\kappa^2\gamma M(1 + \delta^2)\sigma_1^2(\mathbf{A})K^3d \log N$ .

Hence, with  $\eta = \frac{1}{2\gamma_2}$  and  $\beta = (1 - \frac{\gamma_1}{2\gamma_2})^2$ , we have

$$\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 \leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{88\kappa^2\gamma K}}\right) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2, \quad (43)$$

provided  $\mathbf{W}^{(t)}$  satisfies (25), and

$$|\Omega_t| \gtrsim \varepsilon_0^{-2}\kappa^2\gamma(1 + \delta^2)\sigma_1^4(\mathbf{A})K^3d \log N. \quad (44)$$

Then, we can start mathematical induction of (43) over  $t$ .

**Base case:** According to Lemma 4, we know that (25) holds for  $\mathbf{W}^{(0)}$  if

$$|\Omega_1| \gtrsim \varepsilon_0^{-2}\kappa^9\gamma^2(1 + \delta^2)\sigma_1^4(\mathbf{A})K^8d \log N. \quad (45)$$

According to Theorem 1, it is clear that the number of samples  $|\Omega_t|$  satisfies (45), then (25) indeed holds for  $t = 0$ . Since (25) holds for  $t = 0$  and  $|\Omega_t|$  in Theorem 1 satisfies (44) as well, we have (43) holds for  $t = 0$ .

**Induction step:** Assuming (43) holds for  $\mathbf{W}^{(t)}$ , we need to show that (43) holds for  $\mathbf{W}^{(t+1)}$ . That is to say, we need  $|\Omega_t|$  satisfies (44), which holds naturally from Theorem 1.

Therefore, when  $|\Omega_t| \gtrsim \varepsilon_0^{-2} \kappa^9 \gamma^2 (1 + \delta^2) \sigma_1^4(\mathbf{A}) K^8 d \log N$ , we know that (43) holds for all  $0 \leq t \leq T - 1$  with probability at least  $1 - K^2 T \cdot N^{-10}$ . By simple calculations, we can obtain

$$\|\mathbf{W}^{(T)} - \mathbf{W}^*\|_2 \leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{88\kappa^2\gamma K}}\right)^T \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_2 \quad (46)$$

□

### A.1. Proof of Lemma 2

In this section, we provide the proof of Lemma 2 which shows the local convexity of  $f_{\Omega_t}$  in a small neighborhood of  $\mathbf{W}^*$ . The roadmap is to first bound the smallest eigenvalue of  $\nabla^2 f_{\Omega_t}$  in the ground truth as shown in Lemma 5, then show that the difference of  $\nabla^2 f_{\Omega_t}$  between any fixed point  $\mathbf{W}$  in this region and the ground truth  $\mathbf{W}^*$  is bounded in terms of  $\|\mathbf{W} - \mathbf{W}^*\|_2$  by Lemma 6.

**Lemma 5.** *The second-order derivative of  $f_{\Omega_t}$  at the ground truth  $\mathbf{W}^*$  satisfies*

$$\frac{\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2} \preceq \nabla^2 f_{\Omega_t}(\mathbf{W}^*) \preceq \frac{3\sigma_1^2(\mathbf{A})}{K}. \quad (47)$$

**Lemma 6.** *Suppose  $\mathbf{W}$  satisfies (25), we have*

$$\|\nabla^2 f_{\Omega_t}(\mathbf{W}) - \nabla^2 f_{\Omega_t}(\mathbf{W}^*)\|_2 \leq 4\sigma_1^2(\mathbf{A}) \frac{\|\mathbf{W}^* - \mathbf{W}\|_2}{\sigma_K}. \quad (48)$$

The proofs of Lemmas 5 and 6 can be found in Sec. A.3. With these two preliminary lemmas on hand, the proof of Lemma 2 is formally summarized in the following contents.

*Proof of Lemma 2.* By the triangle inequality, we have

$$\left| \|\nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 - \|\nabla^2 f_{\Omega_t}(\mathbf{W}^*)\|_2 \right| \leq \|\nabla^2 f_{\Omega_t}(\mathbf{W}^*) - \nabla^2 f_{\Omega_t}(\mathbf{W})\|_2,$$

and

$$\begin{aligned} \|\nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 &\leq \|\nabla^2 f_{\Omega_t}(\mathbf{W}^*)\|_2 + \|\nabla^2 f_{\Omega_t}(\mathbf{W}^*) - \nabla^2 f_{\Omega_t}(\mathbf{W})\|_2, \\ \|\nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 &\geq \|\nabla^2 f_{\Omega_t}(\mathbf{W}^*)\|_2 - \|\nabla^2 f_{\Omega_t}(\mathbf{W}^*) - \nabla^2 f_{\Omega_t}(\mathbf{W})\|_2. \end{aligned}$$

The error bound of  $\|\nabla^2 f_{\Omega_t}(\mathbf{W}^*) - \nabla^2 f_{\Omega_t}(\mathbf{W})\|_2$  can be derived from Lemma 6, and the error bound of  $\nabla^2 f_{\Omega_t}(\mathbf{W}^*)$  is provided in Lemma 5.

Therefore, for any  $\mathbf{W}$  satisfies (25), we have

$$\frac{(1 - \varepsilon_0)\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2} \leq \|\nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 \leq \frac{4\sigma_1^2(\mathbf{A})}{K}. \quad (49)$$

□

### A.2. Proof of Lemma 3

The proof of Lemma 3 is mainly to bound the concentration error of random variables  $z_n(j, k)$  as shown in (60). We first show that  $z_n(j, k)$  is a sub-exponential random variable, and the definitions of sub-Gaussian and sub-exponential random variables are provided in Definitions 1 and 2. Though Hoeffding's inequality provides the concentration error for sum of independent random variables, random variables  $z_n(j, k)$  with different  $j, k$  are not independent. Hence, we introduce Lemma 7 to provide the upper bound for the moment generation function of the sum of partly dependent random variables and then apply standard Chernoff inequality. Lemmas 8 and 9 are standard tools in analyzing spectral norms of high-dimensional random matrices.

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

$$(\mathbb{E}|X|^p)^{1/p} \leq c_1 \sqrt{p} \quad (50)$$

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

$$\mathbb{E}e^{s(X-\mathbb{E}X)} \leq e^{c_2 \|X\|_{\psi_2}^2 s^2} \quad (51)$$

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

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

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

$$(\mathbb{E}|X|^p)^{1/p} \leq c_3 p \quad (52)$$

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

$$\mathbb{E}e^{s(X-\mathbb{E}X)} \leq e^{c_4 \|X\|_{\psi_1}^2 s^2} \quad (53)$$

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 7.** Given a sampling set  $\mathcal{X} = \{x_n\}_{n=1}^N$  that contains  $N$  partly dependent random variables, for each  $n \in [N]$ , suppose  $x_n$  is dependent with at most  $d_{\mathcal{X}}$  random variables in  $\mathcal{X}$  (including  $x_n$  itself), and the moment generate function of  $x_n$  satisfies  $\mathbb{E}_{x_n} e^{sx_n} \leq e^{Cs^2}$  for some constant  $C$  that may depend on  $x_n$ . Then, the moment generation function of  $\sum_{n=1}^N x_n$  is bounded as

$$\mathbb{E}_{\mathcal{X}} e^{s \sum_{n=1}^N x_n} \leq e^{Cd_{\mathcal{X}}Ns^2}. \quad (54)$$

**Lemma 8** (Lemma 5.2, (Vershynin, 2010)). Let  $\mathcal{B}(0, 1) \in \{\boldsymbol{\alpha} \mid \|\boldsymbol{\alpha}\|_2 = 1, \boldsymbol{\alpha} \in \mathbb{R}^d\}$  denote a unit ball in  $\mathbb{R}^d$ . Then, a subset  $\mathcal{S}_{\xi}$  is called a  $\xi$ -net of  $\mathcal{B}(0, 1)$  if every point  $\mathbf{z} \in \mathcal{B}(0, 1)$  can be approximated to within  $\xi$  by some point  $\boldsymbol{\alpha} \in \mathcal{S}_{\xi}$ , i.e.  $\|\mathbf{z} - \boldsymbol{\alpha}\|_2 \leq \xi$ . Then the minimal cardinality of a  $\xi$ -net  $\mathcal{S}_{\xi}$  satisfies

$$|\mathcal{S}_{\xi}| \leq (1 + 2/\xi)^d. \quad (55)$$

**Lemma 9** (Lemma 5.3, (Vershynin, 2010)). Let  $\mathbf{A}$  be an  $N \times d$  matrix, and let  $\mathcal{S}_{\xi}$  be a  $\xi$ -net of  $\mathcal{B}(0, 1)$  in  $\mathbb{R}^d$  for some  $\xi \in (0, 1)$ . Then

$$\|\mathbf{A}\|_2 \leq (1 - \xi)^{-1} \max_{\boldsymbol{\alpha} \in \mathcal{S}_{\xi}} |\boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha}|. \quad (56)$$

The proof of Lemma 7 can be found in Appendix A.3. With these preliminary Lemmas and definition on hand, the proof of Lemma 3 is formally summarized in the following contents.

*Proof of Lemma 3.* We have

$$\hat{f}_{\Omega_t}(\mathbf{W}) = \frac{1}{2|\Omega_t|} \sum_{n \in \Omega_t} \left| y_n - g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X}) \right|^2 = \frac{1}{2|\Omega_t|} \sum_{n \in \Omega_t} \left| y_n - \sum_{j=1}^K \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j) \right|^2, \quad (57)$$

and

$$f_{\Omega_t}(\mathbf{W}) = \mathbb{E}_{\mathbf{X}} \hat{f}_{\Omega_t}(\mathbf{W}) = \frac{1}{2|\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{x}} \left| y_n - \sum_{j=1}^K \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j) \right|^2. \quad (58)$$

The gradients of  $\hat{f}_{\Omega_t}$  are

$$\begin{aligned}
 \frac{\partial \hat{f}_{\Omega_t}}{\partial \mathbf{w}_k}(\mathbf{W}) &= \frac{1}{K^2|\Omega_t|} \sum_{n \in \Omega_t} \left( y_n - \sum_{j=1}^K \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j) \right) \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_k) \\
 &= \frac{1}{K^2|\Omega_t|} \sum_{n \in \Omega_t} \left( \sum_{j=1}^K \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j^*) - \sum_{j=1}^K \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j) \right) \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_k) \\
 &= \sum_{j=1}^K \frac{1}{K^2|\Omega_t|} \sum_{n \in \Omega_t} (\phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j^*) - \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j)) \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_k).
 \end{aligned} \tag{59}$$

Let us define

$$\mathbf{z}_n(k, j) = \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_k) (\phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j^*) - \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j)), \tag{60}$$

then for any normalized  $\boldsymbol{\alpha} \in \mathbb{R}^d$ , we have

$$\begin{aligned}
 & p^{-1} \left( \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_k) (\phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j^*) - \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j))|^p \right)^{1/p} \\
 & \leq p^{-1} \left( \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^{2p} \cdot \mathbb{E}_{\mathbf{X}} |\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_k) (\phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j^*) - \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j))|^{2p} \right)^{1/2p} \\
 & \leq p^{-1} \left( \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^{2p} \right)^{1/2p} \cdot \left( \mathbb{E}_{\mathbf{X}} |\mathbf{a}_n^T \mathbf{X} (\mathbf{w}_j^* - \mathbf{w}_j)|^{2p} \right)^{1/2p}
 \end{aligned} \tag{61}$$

where the first inequality comes from the Cauchy-Schwarz inequality. Furthermore,  $\mathbf{a}_n^T \mathbf{X}$  belongs to the Gaussian distribution and thus is a sub-Gaussian random vector as well. Then, from Definition 1, we have

$$\begin{aligned}
 & \left( \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^{2p} \right)^{1/2p} \leq (2p)^{1/2} \|\mathbf{X}^T \mathbf{a}_n\|_{\psi_2} \leq (2p)^{1/2} \|\mathbf{a}_n\|_2, \\
 & \text{and } \left( \mathbb{E}_{\mathbf{X}} |\mathbf{a}_n^T \mathbf{X} (\mathbf{w}_j^* - \mathbf{w}_j)|^{2p} \right)^{1/2p} \leq (2p)^{1/2} \|\mathbf{a}_n\|_2 \cdot \|\mathbf{w}_j^* - \mathbf{w}_j\|_2.
 \end{aligned} \tag{62}$$

Then, we have

$$\begin{aligned}
 & p^{-1} \left( \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_k) (\phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j^*) - \phi(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_j))|^p \right)^{1/p} \\
 & \leq p^{-1} \cdot 2p \|\mathbf{a}_n\|_2^2 \cdot \|\mathbf{w}_j^* - \mathbf{w}_j\|_2 \\
 & \leq 2\sigma_1^2(\mathbf{A}) \cdot \|\mathbf{w}_j^* - \mathbf{w}_j\|_2.
 \end{aligned} \tag{63}$$

Therefore, from Definition 2,  $\mathbf{z}_n(k, j)$  belongs to the sub-exponential distribution with

$$\|\mathbf{z}_n\|_{\phi_1} \leq 2\sigma_1^2(\mathbf{A}) \cdot \|\mathbf{w}_j^* - \mathbf{w}_j\|_2. \tag{64}$$

Recall that each node is connected with at most  $\delta$  other nodes. Hence, for any fixed  $\mathbf{z}_n$ , there are at most  $(1 + \delta^2)$  (including  $\mathbf{z}_n$  itself) elements in  $\{\mathbf{z}_l | l \in \Omega_t\}$  are dependant with  $\mathbf{z}_n$ . From Lemma 7, the moment generation function of  $\sum_{n \in \Omega_t} (\mathbf{z}_n - \mathbb{E}_{\mathbf{X}} \mathbf{z}_n)$  satisfies

$$\mathbb{E}_{\mathbf{X}} e^{s \sum_{n \in \Omega_t} (\mathbf{z}_n - \mathbb{E}_{\mathbf{X}} \mathbf{z}_n)} \leq e^{C(1+\delta^2)|\Omega_t|s^2}. \tag{65}$$

By Chernoff inequality, we have

$$\text{Prob} \left\{ \left\| \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} (\mathbf{z}_n(k, j) - \mathbb{E}_{\mathbf{X}} \mathbf{z}_n(k, j)) \right\|_2 > t \right\} \leq \frac{e^{C(1+\delta^2)|\Omega_t|s^2}}{e^{|\Omega_t|ts}} \tag{66}$$

for any  $s > 0$ .

Let  $s = t/(C(1 + \delta^2)\|z_n\|_{\phi_1}^2)$  and  $t = \sqrt{\frac{(1+\delta^2)d \log N}{|\Omega_t|}} \|z_n\|_{\phi_1}$ , we have

$$\begin{aligned} \left\| \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} (z_n(k, j) - \mathbb{E}_{\mathbf{X}} z_n(k, j)) \right\|_2 &\leq C \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \sigma_1^2(\mathbf{A}) \cdot \|w_j^* - w_j\|_2 \\ &\leq C \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \cdot \|\mathbf{W}^* - \mathbf{W}\|_2 \end{aligned} \quad (67)$$

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

In conclusion, by selecting  $\xi = \frac{1}{2}$  in Lemmas 8 and 9, we have

$$\begin{aligned} \left\| \frac{\partial \hat{f}_{\Omega_t}}{\partial w_k}(\mathbf{W}) - \frac{\partial f_{\Omega_t}}{\partial w_k}(\mathbf{W}) \right\|_2 &\leq \sum_{k=1}^K \sum_{j=1}^K \frac{1}{K^2} \left\| \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} z_n(k, j) - \mathbb{E}_{\mathbf{X}} z_n(k, j) \right\|_2 \\ &\leq C \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \cdot \|\mathbf{W}^* - \mathbf{W}\|_2 \end{aligned} \quad (68)$$

with probability at least  $1 - \left(\frac{5}{N}\right)^d$ . □

### A.3. Proof of auxiliary lemmas for regression problems

#### A.3.1. PROOF OF LEMMA 5

*Proof of Lemma 5.* For any normalized  $\alpha \in \mathbb{R}^{Kd}$ , the lower bound of  $\nabla^2 f_{\Omega_t}(\mathbf{W}^*)$  is derived from

$$\begin{aligned} \alpha^T \nabla^2 f(\mathbf{W}^*) \alpha &= \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} \left[ \left( \sum_{j=1}^K \alpha_j^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} w_j^*) \right)^2 \right] \\ &\geq \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \frac{\|\mathbf{a}_n\|_2^2}{11\kappa^2\gamma} \|\alpha\|_2^2 = \frac{\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2}, \end{aligned} \quad (69)$$

where the last inequality can be derived from Lemma D.6 in (Zhong et al., 2017c). In spite that the error bound in (Zhong et al., 2017c) is given in terms of  $x_n$  instead of  $\mathbf{X}^T \mathbf{a}_n$ , both  $x_n$  and  $\mathbf{X}^T \mathbf{a}_n$  belong to Gaussian distribution. Hence, we can follow the similar steps in (Zhong et al., 2017c) to derive the results for Gaussian random variable  $\mathbf{X}^T \mathbf{a}_n$  with 0 mean and  $\|\mathbf{a}_n\|_2^2$  variance.

Next, the upper bound of  $\nabla^2 f_{\Omega_t}(\mathbf{W}^*)$  is derived from

$$\begin{aligned} &\alpha^T \nabla^2 f(\mathbf{W}^*) \alpha \\ &= \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} \left[ \left( \sum_{j=1}^K \alpha_j^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} w_j^*) \right)^2 \right] \\ &= \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \sum_{j_1=1}^K \sum_{j_2=1}^K \mathbb{E}_{\mathbf{X}} \left[ \alpha_{j_1}^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} w_{j_1}^*) \alpha_{j_2}^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{a}_n^T \mathbf{X} w_{j_2}^*) \right] \\ &\leq \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \sum_{j_1=1}^K \sum_{j_2=1}^K \left[ \mathbb{E}_{\mathbf{X}} |\alpha_{j_1}^T \mathbf{X}^T \mathbf{a}_n|^4 \cdot \mathbb{E}_{\mathbf{X}} |\phi'(\mathbf{a}_n^T \mathbf{X} w_{j_1}^*)|^4 \cdot \mathbb{E}_{\mathbf{X}} |\alpha_{j_2}^T \mathbf{X}^T \mathbf{a}_n|^4 \cdot \mathbb{E}_{\mathbf{X}} |\phi'(\mathbf{a}_n^T \mathbf{X} w_{j_2}^*)|^4 \right]^{\frac{1}{4}} \\ &\leq \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \sum_{j_1=1}^K \sum_{j_2=1}^K 3\sigma_1^2(\mathbf{A}) \|\alpha_{j_1}\|_2 \|\alpha_{j_2}\|_2 \\ &\leq 3\sigma_1^2(\mathbf{A}) \frac{\|\alpha\|_2^2}{K}, \end{aligned} \quad (70)$$

which completes the proof. □

## A.3.2. PROOF OF LEMMA 6

*Proof of Lemma 6.* The second-order derivative of  $f_{\Omega_t}$  is written as

$$\begin{aligned}
 & \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}) - \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}^*) \\
 &= \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T \left[ \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}^*) \right] \\
 &= \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T (\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*)) \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}) \\
 & \quad - \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) (\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}^*) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2})).
 \end{aligned} \tag{71}$$

For any normalized  $\boldsymbol{\alpha} \in \mathbb{R}^d$ , we have

$$\begin{aligned}
 & \left| \boldsymbol{\alpha}^T \left[ \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}) - \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}^*) \right] \boldsymbol{\alpha} \right| \\
 & \leq \left| \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2 (\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*)) \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}) \right| \\
 & \quad + \left| \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2 \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) (\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}^*) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2})) \right| \\
 & \leq \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^2 \cdot \left| \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) \right| \\
 & \quad + \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^2 \cdot \left| \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}^*) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_2}) \right|.
 \end{aligned} \tag{72}$$

It is easy to verify there exists a basis such that  $\mathcal{B} = \{\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\alpha}_4^\perp, \dots, \boldsymbol{\alpha}_d^\perp\}$  with  $\{\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}\}$  spanning a subspace that contains  $\boldsymbol{\alpha}, \mathbf{w}_{j_1}$  and  $\mathbf{w}_{j_1}^*$ . Then, for any  $\mathbf{X}^T \mathbf{a}_n \in \mathbb{R}^d$ , we have a unique  $\mathbf{z} = [z_1 \ z_2 \ \dots \ z_d]^T$  such that

$$\mathbf{X}^T \mathbf{a}_n = z_1 \boldsymbol{\alpha} + z_2 \boldsymbol{\beta} + z_3 \boldsymbol{\gamma} + \dots + z_d \boldsymbol{\alpha}_d^\perp.$$

Also, since  $\mathbf{X}^T \mathbf{a}_n \sim \mathcal{N}(\mathbf{0}, \|\mathbf{a}_n\|_2^2 \mathbf{I}_d)$ , we have  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \|\mathbf{a}_n\|_2^2 \mathbf{I}_d)$ . Then, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^2 \cdot \left| \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) \right| \\
 &= \mathbb{E}_{z_1, z_2, z_3} |\phi'(\mathbf{w}_{j_1}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_1}^{*T} \tilde{\mathbf{x}})| \cdot |\mathbf{a}^T \tilde{\mathbf{x}}|^2 \\
 &= \int |\phi'(\mathbf{w}_{j_1}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_1}^{*T} \tilde{\mathbf{x}})| \cdot |\mathbf{a}^T \tilde{\mathbf{x}}|^2 \cdot f_Z(z_1, z_2, z_3) dz_1 dz_2 dz_3,
 \end{aligned}$$

where  $\tilde{\mathbf{x}} = z_1 \boldsymbol{\alpha} + z_2 \boldsymbol{\beta} + z_3 \boldsymbol{\gamma}$  and  $f_Z(z_1, z_2, z_3)$  is the probability density function of  $(z_1, z_2, z_3)$ . Next, 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,

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^2 \cdot \left| \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) \right| \\
 &= \int \int \int |\phi'(\mathbf{w}_{j_1}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_1}^{*T} \tilde{\mathbf{x}})| \cdot |r \cos \phi_1|^2 \cdot f_Z(r, \phi_1, \phi_2) r^2 \sin \phi_1 dr d\phi_1 d\phi_2.
 \end{aligned} \tag{73}$$

It is easy to verify that  $\phi'(\mathbf{w}_{j_1}^T \tilde{\mathbf{x}})$  only depends on the direction of  $\tilde{\mathbf{x}}$  and

$$f_Z(r, \phi_1, \phi_2) = \frac{1}{(2\pi \|\mathbf{a}_n\|_2^2)^{\frac{3}{2}}} e^{-\frac{x_1^2 + x_2^2 + x_3^2}{2\|\mathbf{a}_n\|_2^2}} = \frac{1}{(2\pi \|\mathbf{a}_n\|_2^2)^{\frac{3}{2}}} e^{-\frac{r^2}{2\|\mathbf{a}_n\|_2^2}}$$



only depends on  $r$ . Then, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^2 \cdot \left| \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) \right| \\
 &= \int \int \int |\phi'(\mathbf{w}_{j_1}^T(\tilde{\mathbf{x}}/r)) - \phi'(\mathbf{w}_{j_1}^{*T}(\tilde{\mathbf{x}}/r))| \cdot |r \cos \phi_1|^2 \cdot f_Z(r) r^2 \sin \phi_1 dr d\phi_1 d\phi_2 \\
 &= \int_0^\infty r^4 f_Z(r) dr \int_0^\pi \int_0^{2\pi} |\cos \phi_1|^2 \cdot \sin \phi_1 \cdot |\phi'(\mathbf{w}_{j_2}^T(\tilde{\mathbf{x}}/r)) - \phi'(\mathbf{w}_{j_2}^{*T}(\tilde{\mathbf{x}}/r))| d\phi_1 d\phi_2 \\
 &\leq 3 \|\mathbf{a}_n\|_2^2 \cdot \int_0^\infty r^2 f_Z(r) dr \int_0^\pi \int_0^{2\pi} \sin \phi_1 \cdot |\phi'(\mathbf{w}_{j_2}^T(\tilde{\mathbf{x}}/r)) - \phi'(\mathbf{w}_{j_2}^{*T}(\tilde{\mathbf{x}}/r))| d\phi_1 d\phi_2 \\
 &= 3 \|\mathbf{a}_n\|_2^2 \cdot \mathbb{E}_{z_1, z_2, z_3} |\phi'(\mathbf{w}_{j_1}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_1}^{*T} \tilde{\mathbf{x}})| \\
 &= 3 \|\mathbf{a}_n\|_2^2 \cdot \mathbb{E}_{\mathbf{X}} |\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*)|
 \end{aligned} \tag{74}$$

Define a set  $\mathcal{A}_1 = \{\mathbf{x} | (\mathbf{w}_{j_1}^{*T} \mathbf{x})(\mathbf{w}_{j_1}^T \mathbf{x}) < 0\}$ . If  $\mathbf{x} \in \mathcal{A}_1$ , then  $\mathbf{w}_{j_1}^{*T} \mathbf{x}$  and  $\mathbf{w}_{j_1}^T \mathbf{x}$  have different signs, which means the value of  $\phi'(\mathbf{w}_{j_1}^T \mathbf{x})$  and  $\phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x})$  are different. This is equivalent to say that

$$|\phi'(\mathbf{w}_{j_1}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x})| = \begin{cases} 1, & \text{if } \mathbf{x} \in \mathcal{A}_1 \\ 0, & \text{if } \mathbf{x} \in \mathcal{A}_1^c \end{cases}. \tag{75}$$

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

$$\|\mathbf{w}_{j_1}^{*T} \mathbf{x}\| \leq \|\mathbf{w}_{j_1}^{*T} \mathbf{x} - \mathbf{w}_{j_1}^T \mathbf{x}\| \leq \|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\| \cdot \|\mathbf{x}\|. \tag{76}$$

Define a set  $\mathcal{A}_2$  such that

$$\mathcal{A}_2 = \left\{ \mathbf{x} \mid \frac{|\mathbf{w}_{j_1}^{*T} \mathbf{x}|}{\|\mathbf{w}_{j_1}^*\| \|\mathbf{x}\|} \leq \frac{\|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\|}{\|\mathbf{w}_{j_1}^*\|} \right\} = \left\{ \theta_{\mathbf{x}, \mathbf{w}_{j_1}^*} \mid |\cos \theta_{\mathbf{x}, \mathbf{w}_{j_1}^*}| \leq \frac{\|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\|}{\|\mathbf{w}_{j_1}^*\|} \right\}. \tag{77}$$

Hence, we have that

$$\mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_1}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x})| = \text{Prob}(\mathbf{x} \in \mathcal{A}_1) \leq \text{Prob}(\mathbf{x} \in \mathcal{A}_2). \tag{78}$$

Since  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\theta_{\mathbf{x}, \mathbf{w}_{j_1}^*}$  belongs to the uniform distribution on  $[-\pi, \pi]$ , we have

$$\begin{aligned}
 \text{Prob}(\mathbf{x} \in \mathcal{A}_2) &= \frac{\pi - \arccos \frac{\|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\|}{\|\mathbf{w}_{j_1}^*\|}}{\pi} \\
 &\leq \frac{1}{\pi} \tan\left(\pi - \arccos \frac{\|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\|}{\|\mathbf{w}_{j_1}^*\|}\right) \\
 &= \frac{1}{\pi} \cot\left(\arccos \frac{\|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\|}{\|\mathbf{w}_{j_1}^*\|}\right) \\
 &\leq \frac{2}{\pi} \frac{\|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\|}{\|\mathbf{w}_{j_1}^*\|}.
 \end{aligned} \tag{79}$$

Hence, (81) and (79) suggest that

$$\mathbb{E}_{\mathbf{X}} |\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*)| \leq \frac{6}{\pi} \frac{\|\mathbf{w}_{j_1}^* - \mathbf{w}_{j_1}\|}{\|\mathbf{w}_{j_1}^*\|}. \tag{80}$$

Then, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{X}} |\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n|^2 \cdot \left| \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*) \right| \\
 &= 3 \|\mathbf{a}_n\|_2^2 \cdot \mathbb{E}_{\mathbf{X}} |\phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}) - \phi'(\mathbf{a}_n^T \mathbf{X} \mathbf{w}_{j_1}^*)| \\
 &\leq \frac{6 \|\mathbf{a}_n\|_2^2}{\pi} \cdot \frac{\|\mathbf{w}_{j_1} - \mathbf{w}_{j_1}^*\|_2}{\|\mathbf{w}_{j_1}^*\|_2},
 \end{aligned} \tag{81}$$

All in all, we have

$$\begin{aligned}
 \|\nabla^2 f_{\Omega_t}(\mathbf{W}) - \nabla^2 f_{\Omega_t}(\mathbf{W}^*)\|_2 &\leq \sum_{j_1}^K \sum_{j_2}^K \left\| \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}) - \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}^*) \right\|_2 \\
 &\leq K^2 \max_{j_1, j_2} \left\| \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}) - \frac{\partial^2 f_{\Omega_t}}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}}(\mathbf{W}^*) \right\|_2 \\
 &\leq K^2 \cdot \frac{12 \|\mathbf{a}_n\|_2^2}{\pi} \max_j \frac{\|\mathbf{w}_j - \mathbf{w}_j^*\|_2}{\|\mathbf{w}_j^*\|_2} \\
 &\leq 4\sigma_1^2(\mathbf{A}) \frac{\|\mathbf{W}^* - \mathbf{W}\|_2}{\sigma_K}.
 \end{aligned} \tag{82}$$

□

### A.3.3. PROOF OF LEMMA 7

*Proof of Lemma 7.* According to the Definitions in (Janson, 2004), there exists a family of  $\{(\mathcal{X}_j, w_j)\}_j$ , where  $\mathcal{X}_j \subseteq \mathcal{X}$  and  $w_j \in [0, 1]$ , such that  $\sum_j w_j \sum_{x_{n_j} \in \mathcal{X}_j} x_{n_j} = \sum_{n=1}^N x_n$ , and  $\sum_j w_j \leq d_{\mathcal{X}}$  by equations (2.1) and (2.2) in (Janson, 2004). Then, let  $p_j$  be any positive numbers with  $\sum_j p_j = 1$ . By Jensen's inequality, for any  $s \in \mathbb{R}$ , we have

$$e^{s \sum_{n=1}^N x_n} = e^{\sum_j p_j \frac{sw_j}{p_j} X_j} \leq \sum_j p_j e^{\frac{sw_j}{p_j} X_j}, \tag{83}$$

where  $X_j = \sum_{x_{n_j} \in \mathcal{X}_j} x_{n_j}$ .

Then, we have

$$\begin{aligned}
 \mathbb{E}_{\mathcal{X}} e^{s \sum_{n=1}^N x_n} &\leq \mathbb{E}_{\mathcal{X}} \sum_j p_j e^{\frac{sw_j}{p_j} X_j} = \sum_j p_j \prod_{\mathcal{X}_j} \mathbb{E}_{\mathcal{X}} e^{\frac{sw_j}{p_j} x_{n_j}} \\
 &\leq \sum_j p_j \prod_{\mathcal{X}_j} e^{\frac{Cw_j^2}{p_j^2} s^2} \\
 &\leq \sum_j p_j e^{\frac{C|\mathcal{X}_j|w_j^2}{p_j^2} s^2}.
 \end{aligned} \tag{84}$$

Let  $p_j = \frac{w_j |\mathcal{X}_j|^{1/2}}{\sum_j w_j |\mathcal{X}_j|^{1/2}}$ , then we have

$$\mathbb{E}_{\mathcal{X}} e^{s \sum_{n=1}^N x_n} \leq \sum_j p_j e^{C(\sum_j w_j |\mathcal{X}_j|^{1/2})^2 s^2} = e^{C(\sum_j w_j |\mathcal{X}_j|^{1/2})^2 s^2}. \tag{85}$$

By Cauchy-Schwarz inequality, we have

$$\left( \sum_j w_j |\mathcal{X}_j|^{1/2} \right)^2 \leq \sum_j w_j \sum_j w_j |\mathcal{X}_j| \leq d_{\mathcal{X}} N. \tag{86}$$

Hence, we have

$$\mathbb{E}_{\mathcal{X}} e^{s \sum_{n=1}^N x_n} \leq e^{C d_{\mathcal{X}} N s^2}. \tag{87}$$

□

## B. Proof of Theorem 2

Recall that the empirical risk function in (4) is defined as

$$\min_{\mathbf{W}} : \hat{f}_{\Omega}(\mathbf{W}) = \frac{1}{|\Omega|} \sum_{n \in \Omega} -y_n \log(g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X})) - (1 - y_n) \log(1 - g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X})). \quad (88)$$

The population risk function is defined as

$$\begin{aligned} f_{\Omega}(\mathbf{W}) &:= \mathbb{E}_{\mathbf{X}, y_n} \hat{f}_{\Omega}(\mathbf{W}) \\ &= \mathbb{E}_{\mathbf{X}} \mathbb{E}_{y_n | \mathbf{X}} \left[ \frac{1}{|\Omega|} \sum_{n \in \Omega} -y_n \log(g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X})) - (1 - y_n) \log(1 - g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X})) \right] \\ &= \mathbb{E}_{\mathbf{X}} \frac{1}{|\Omega|} \sum_{n \in \Omega} -g(\mathbf{W}^*; \mathbf{a}_n^T \mathbf{X}) \log(g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X})) - (1 - g(\mathbf{W}^*; \mathbf{a}_n^T \mathbf{X})) \log(1 - g(\mathbf{W}; \mathbf{a}_n^T \mathbf{X})). \end{aligned} \quad (89)$$

The road-map of proof for Theorem 2 follows the similar three steps as those for Theorem 1. The major differences lie in three aspects: (i) in the second step, the objective function  $\hat{f}_{\Omega_t}$  is smooth since the activation function  $\phi(\cdot)$  is sigmoid. Hence, we can directly apply the mean value theorem as  $\nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) = \langle \nabla^2 \hat{f}_{\Omega_t}(\widehat{\mathbf{W}}^{(t)}), \mathbf{W}^{(t)} - \mathbf{W}^* \rangle$  to characterize the effects of the gradient descent term in each iteration, and the error bound of  $\nabla^2 \hat{f}_{\Omega_t}$  is provided in Lemma 10; (ii) the objective function is the sum of cross-entry loss functions, which have more complex structure of derivatives than those of square loss functions; (iii) as the convergent point may not be the critical point of empirical loss function, we need to provide the distance from the convergent point to the ground-truth parameters additionally, where Lemma 11 is used.

Lemmas 10 and 11 are summarized in the following contents. Also, the notations  $\lesssim$  and  $\gtrsim$  follow the same definitions as in (27). The proofs of Lemmas 10 and 11 can be found in Appendix B.1 and B.2, respectively.

**Lemma 10.** For any  $\mathbf{W}$  that satisfies

$$\|\mathbf{W} - \mathbf{W}^*\| \leq \frac{2\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2} \quad (90)$$

then the second-order derivative of the empirical risk function in (88) for binary classification problems is bounded as

$$\frac{2(1 - \varepsilon_0)}{11\kappa^2\gamma K^2} \sigma_1^2(\mathbf{A}) \preceq \nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}) \preceq \sigma_1^2(\mathbf{A}). \quad (91)$$

provided the number of samples satisfies

$$|\Omega_t| \gtrsim \varepsilon_0^{-2} (1 + \delta^2) \kappa^2 \gamma \sigma_1^4(\mathbf{A}) K^6 d \log N. \quad (92)$$

**Lemma 11.** Let  $\hat{f}_{\Omega_t}$  and  $f_{\Omega_t}$  be the empirical and population risk function in (88) and (89) for binary classification problems, respectively, then the first-order derivative of  $\hat{f}_{\Omega_t}$  is close to its expectation  $f_{\Omega_t}$  with an upper bound as

$$\|\nabla f_{\Omega_t}(\mathbf{W}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W})\|_2 \lesssim K^2 \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2) d \log d}{|\Omega_t|}} \quad (93)$$

with probability at least  $1 - K^2 N^{-10}$ .

With these preliminary lemmas, the proof of Theorem 2 is formally summarized in the following contents.

*Proof of Theorem 2.* The update rule of  $\mathbf{W}^{(t)}$  is

$$\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \nabla \hat{f}_{\Omega_t}(\mathbf{W}^{(t)}) + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) \quad (94)$$

Since  $\widehat{\mathbf{W}}$  is a critical point, then we have  $\nabla \hat{f}_{\Omega_t}(\widehat{\mathbf{W}}) = 0$ . By the intermediate value theorem, we have

$$\begin{aligned} \mathbf{W}^{(t+1)} &= \mathbf{W}^{(t)} - \eta \nabla^2 \hat{f}_{\Omega_t}(\widehat{\mathbf{W}}^{(t)})(\mathbf{W}^{(t)} - \widehat{\mathbf{W}}) \\ &\quad + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) \end{aligned} \quad (95)$$

where  $\widehat{\mathbf{W}}^{(t)}$  lies in the convex hull of  $\mathbf{W}^{(t)}$  and  $\widehat{\mathbf{W}}$ .

Next, we have

$$\begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \eta \nabla^2 \hat{f}_{\Omega_t}(\widehat{\mathbf{W}}^{(t)}) + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{W}^{(t)} - \mathbf{W}^* \\ \mathbf{W}^{(t-1)} - \mathbf{W}^* \end{bmatrix}. \quad (96)$$

Let  $\mathbf{P}(\beta) = \begin{bmatrix} \mathbf{I} - \eta \nabla^2 \hat{f}_{\Omega_t}(\widehat{\mathbf{W}}^{(t)}) + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}$ , so we have

$$\left\| \begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix} \right\|_2 = \|\mathbf{P}(\beta)\|_2 \left\| \begin{bmatrix} \mathbf{W}^{(t)} - \mathbf{W}^* \\ \mathbf{W}^{(t-1)} - \mathbf{W}^* \end{bmatrix} \right\|_2.$$

Then, we have

$$\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 \lesssim \|\mathbf{P}(\beta)\|_2 \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2 \quad (97)$$

Let  $\lambda_i$  be the  $i$ -th eigenvalue of  $\nabla^2 \hat{f}_{\Omega_t}(\widehat{\mathbf{W}}^{(t)})$ , and  $\delta_i$  be the  $i$ -th eigenvalue of matrix  $\mathbf{P}(\beta)$ . Following the similar analysis in proof of Theorem 1, we have

$$\delta_i(0) > \delta_i(\beta), \quad \text{for } \forall \beta \in (0, (1 - \eta \lambda_i)^2). \quad (98)$$

Moreover,  $\delta_i$  achieves the minimum  $\delta_i^* = |1 - \sqrt{\eta \lambda_i}|$  when  $\beta = (1 - \sqrt{\eta \lambda_i})^2$ .

Let us first assume  $\mathbf{W}^{(t)}$  satisfies (90) and the number of samples satisfies (92), then from Lemma 10, we know that

$$0 < \frac{2(1 - \varepsilon_0)\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2} \leq \lambda_i \leq \sigma_1^2(\mathbf{A}).$$

We define  $\gamma_1 = \frac{2(1 - \varepsilon_0)\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2}$  and  $\gamma_2 = \sigma_1^2(\mathbf{A})$ . Also, for any  $\varepsilon_0 \in (0, 1)$ , we have

$$\nu(\beta^*) = \|\mathbf{P}(\beta^*)\|_2 = 1 - \sqrt{\frac{\gamma_1}{2\gamma_2}} = 1 - \sqrt{\frac{1 - \varepsilon_0}{11\kappa^2\gamma K}} \quad (99)$$

Let  $\beta = 0$ , we have

$$\nu(0) = \|\mathbf{A}(0)\|_2 = 1 - \frac{1 - \varepsilon_0}{11\kappa^2\gamma K}.$$

Hence, with probability at least  $1 - K^2 \cdot N^{-10}$ , we have

$$\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 \leq \left(1 - \sqrt{\frac{1 - \varepsilon_0}{11\kappa^2\gamma K}}\right) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2, \quad (100)$$

provided that  $\mathbf{W}^{(t)}$  satisfies (25), and

$$|\Omega_t| \gtrsim \varepsilon_0^{-2} \kappa^2 \gamma (1 + \delta^2) \sigma_1^4(\mathbf{A}) K^6 d \log N. \quad (101)$$

According to Lemma 4, we know that (90) holds for  $\mathbf{W}^{(0)}$  if

$$|\Omega_t| \gtrsim \varepsilon_0^{-2} \kappa^8 \gamma^2 (1 + \delta^2) K^8 d \log N. \quad (102)$$

Combining (101) and (102), we need  $|\Omega_t| \gtrsim \varepsilon_0^{-2} \kappa^8 \gamma^2 (1 + \delta^2) \sigma_1^4(\mathbf{A}) K^8 d \log N$ .

Finally, by the mean value theorem, we have

$$\hat{f}_{\Omega_t}(\widehat{\mathbf{W}}) \leq \hat{f}_{\Omega_t}(\mathbf{W}^*) + \nabla \hat{f}_{\Omega_t}(\mathbf{W}^*)^T (\widehat{\mathbf{W}} - \mathbf{W}^*) + \frac{1}{2} (\widehat{\mathbf{W}} - \mathbf{W}^*)^T \nabla^2 \hat{f}_{\Omega_t}(\widetilde{\mathbf{W}}) (\widehat{\mathbf{W}} - \mathbf{W}^*) \quad (103)$$

for some  $\widetilde{\mathbf{W}}$  between  $\widehat{\mathbf{W}}$  and  $\mathbf{W}^*$ . Since  $\widehat{\mathbf{W}}$  is the local minima, we have  $\hat{f}_{\Omega_t}(\widehat{\mathbf{W}}) \leq \hat{f}_{\Omega_t}(\mathbf{W}^*)$ . That is to say

$$\nabla \hat{f}_{\Omega_t}(\mathbf{W}^*)^T (\widehat{\mathbf{W}} - \mathbf{W}^*) + \frac{1}{2} (\widehat{\mathbf{W}} - \mathbf{W}^*)^T \nabla^2 \hat{f}_{\Omega_t}(\widetilde{\mathbf{W}}) (\widehat{\mathbf{W}} - \mathbf{W}^*) \leq 0 \quad (104)$$

which implies

$$\frac{1}{2} \|\nabla^2 \hat{f}_{\Omega_t}(\widetilde{\mathbf{W}})\|_2 \|\widehat{\mathbf{W}} - \mathbf{W}^*\|_2^2 \leq \|\nabla \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2 \|\widehat{\mathbf{W}} - \mathbf{W}^*\|_2. \quad (105)$$

From Lemma 10, we know that

$$\|\nabla^2 \hat{f}_{\Omega_t}(\widetilde{\mathbf{W}})\|_2 \geq \frac{2(1 - \varepsilon_0)}{11\kappa^2\gamma K^2} \sigma^2(\mathbf{A}). \quad (106)$$

From Lemma 11, we know that

$$\|\nabla \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2 = \|\nabla \hat{f}_{\Omega_t}(\mathbf{W}^*) - \nabla f_{\Omega_t}(\mathbf{W}^*)\|_2 \lesssim K^2 \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}}. \quad (107)$$

Plugging inequalities (106) and (107) back into (105), we have

$$\|\widehat{\mathbf{W}} - \mathbf{W}^*\|_2 \lesssim (1 - \varepsilon_0)^{-1} \kappa^2 \gamma K^4 \sqrt{\frac{(1 + \delta^2)d \log d}{|\Omega_t|}}. \quad (108)$$

□

### B.1. Proof of Lemma 10

The roadmap of proof for Lemma 10 follows the similar steps as those of Lemma 2 for regression problems. Lemmas 12, 13 and 14 are the preliminary lemmas, and their proofs can be found in Appendix B.2. The proof of Lemma 10 is summarized after these preliminary lemmas.

**Lemma 12.** *The second-order derivative of  $f_{\Omega_t}$  at the ground truth  $\mathbf{W}^*$  satisfies*

$$\frac{4\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2} \mathbf{I} \preceq \nabla^2 f_{\Omega_t}(\mathbf{W}^*) \preceq \frac{\sigma_1^2(\mathbf{A})}{4} \mathbf{I}. \quad (109)$$

**Lemma 13.** *Suppose  $f_{\Omega_t}$  is the population loss function with respect to binary classification problems, then we have*

$$\|\nabla^2 f_{\Omega_t}(\mathbf{W}) - \nabla^2 f_{\Omega_t}(\mathbf{W}^*)\|_2 \lesssim \|\mathbf{W} - \mathbf{W}^*\|_2. \quad (110)$$

**Lemma 14.** *Suppose  $\hat{f}_{\Omega_t}$  is the empirical loss function with respect to binary classification problems, then the second-order derivative of  $\hat{f}_{\Omega_t}$  is close to its expectation with an upper bound as*

$$\|\nabla^2 f_{\Omega_t}(\mathbf{W}) - \nabla^2 \hat{f}_{\Omega_t}(\mathbf{W})\|_2 \lesssim K^2 \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log d}{|\Omega_t|}} \quad (111)$$

with probability at least  $1 - K^2 N^{-10}$ .

*Proof of Lemma 10.* For any  $\mathbf{W}$ , we have

$$\left| \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W})\|_2 - \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2 \right| \leq \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}) - \nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2. \quad (112)$$

That is

$$\begin{aligned} \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W})\|_2 &\leq \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2 + \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}) - \nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2 \\ \text{and } \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W})\|_2 &\geq \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2 - \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}) - \nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}^*)\|_2 \end{aligned} \quad (113)$$

Then, for any  $\mathbf{W}$  that satisfies  $\|\mathbf{W} - \mathbf{W}^*\| \leq \frac{2\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2}$ , from Lemmas 12 and 13, we have

$$\frac{2}{11\kappa^2\gamma K^2}\sigma_1^2(\mathbf{A}) \preceq \nabla^2 f_{\Omega_t}(\mathbf{W}) \preceq \frac{1}{2}\sigma_1^2(\mathbf{A}). \quad (114)$$

Next, we have

$$\begin{aligned} \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W})\|_2 &\leq \|\nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 + \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}) - \nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 \\ \text{and } \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W})\|_2 &\geq \|\nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 - \|\nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}) - \nabla^2 f_{\Omega_t}(\mathbf{W})\|_2 \end{aligned} \quad (115)$$

Then, from (114) and Lemma 14, we have

$$\frac{2(1-\varepsilon_0)}{11\kappa^2\gamma K^2}\sigma_1^2(\mathbf{A}) \preceq \nabla^2 \hat{f}_{\Omega_t}(\mathbf{W}) \preceq \sigma_1^2(\mathbf{A}) \quad (116)$$

provided that the sample size  $|\Omega_t| \gtrsim \varepsilon_0^{-2}(1+\delta^2)\kappa^2\gamma\sigma_1^4(\mathbf{A})K^6 d \log N$ .  $\square$

## B.2. Proof of auxiliary lemmas for binary classification problems

### B.2.1. PROOF OF LEMMA 12

*Proof of Lemma 12.* Since  $\mathbb{E}_{\mathbf{X}} y_n = g_n(\mathbf{W}^*; \mathbf{a}_n)$ , then we have

$$\begin{aligned} \frac{\partial^2 f_{\Omega_t}(\mathbf{W}^*)}{\partial \mathbf{w}_j^* \partial \mathbf{w}_k^*} &= \mathbb{E}_{\mathbf{X}} \frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W}^*)}{\partial \mathbf{w}_j^* \partial \mathbf{w}_k^*} \\ &= \mathbb{E}_{\mathbf{X}} \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \frac{1}{g(\mathbf{W}^*; \mathbf{a}_n)(1-g(\mathbf{W}^*; \mathbf{a}_n))} \phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^{*T} \mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T, \end{aligned} \quad (117)$$

for any  $j, k \in [K]$ .

Then, for any  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1^T, \boldsymbol{\alpha}_2^T, \dots, \boldsymbol{\alpha}_K^T]^T \in \mathbb{R}^{dk}$  with  $\boldsymbol{\alpha}_j \in \mathbb{R}^d$ , the lower bound can be obtained from

$$\begin{aligned} \boldsymbol{\alpha}^T \nabla^2 f_{\Omega_t}(\mathbf{W}^*) \boldsymbol{\alpha} &= \mathbb{E}_{\mathbf{X}} \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \frac{\left( \sum_{j=1}^K \boldsymbol{\alpha}_j^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n) \right)^2}{g(\mathbf{W}^*; \mathbf{a}_n)(1-g(\mathbf{W}^*; \mathbf{a}_n))} \\ &\geq \mathbb{E}_{\mathbf{X}} \frac{4}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \left( \sum_{j=1}^K \boldsymbol{\alpha}_j^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n) \right)^2 \\ &\geq \frac{4\sigma_1^2(\mathbf{A})}{11\kappa^2\gamma K^2}. \end{aligned} \quad (118)$$

Also, for the upper bound, we have

$$\begin{aligned} \boldsymbol{\alpha}^T \nabla^2 f_{\Omega_t}(\mathbf{W}^*) \boldsymbol{\alpha} &= \mathbb{E}_{\mathbf{X}} \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \frac{\left( \sum_{j=1}^K \boldsymbol{\alpha}_j^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n) \right)^2}{g(\mathbf{W}^*; \mathbf{a}_n)(1-g(\mathbf{W}^*; \mathbf{a}_n))} \\ &= \mathbb{E}_{\mathbf{X}} \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} \frac{\left( \sum_{j=1}^K \boldsymbol{\alpha}_j^T \mathbf{X}^T \mathbf{a}_n \phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n) \right)^2}{\sum_{j_1=1}^K \phi(\mathbf{w}_{j_1}^{*T} \mathbf{X}^T \mathbf{a}_n) \sum_{j_2=1}^K (1-\phi(\mathbf{w}_{j_2}^{*T} \mathbf{X}^T \mathbf{a}_n))} \\ &\leq \mathbb{E}_{\mathbf{X}} \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} \frac{\sum_{j=1}^K (\boldsymbol{\alpha}_j^T \mathbf{X}^T \mathbf{a}_n)^2 \sum_{j=1}^K (\phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n))^2}{\sum_{j_1=1}^K \phi(\mathbf{w}_{j_1}^{*T} \mathbf{X}^T \mathbf{a}_n) \sum_{j_2=1}^K (1-\phi(\mathbf{w}_{j_2}^{*T} \mathbf{X}^T \mathbf{a}_n))}. \end{aligned} \quad (119)$$

For the denominator item, we have

$$\begin{aligned}
 \sum_{j_1=1}^K \phi(\mathbf{w}_{j_1}^{*T} \mathbf{X}^T \mathbf{a}_n) \sum_{j_2=1}^K (1 - \phi(\mathbf{w}_{j_2}^{*T} \mathbf{X}^T \mathbf{a}_n)) &\geq \sum_{j=1}^K \phi(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n) (1 - \phi(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n)) \\
 &= \sum_{j=1}^K \phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n) \\
 &\geq 4 \sum_{j=1}^K \phi'(\mathbf{w}_j^{*T} \mathbf{X}^T \mathbf{a}_n)^2.
 \end{aligned} \tag{120}$$

Hence, we have

$$\boldsymbol{\alpha}^T \nabla^2 f_{\Omega_t}(\mathbf{W}^*) \boldsymbol{\alpha} \leq \mathbb{E}_{\mathbf{X}} \frac{1}{4|\Omega_t|} \sum_{n \in \Omega_t} \sum_{j=1}^K (\boldsymbol{\alpha}_j^T \mathbf{X}^T \mathbf{a}_n)^2 \leq \frac{1}{4} \sigma_1^2(\mathbf{A}). \tag{121}$$

□

### B.2.2. PROOF OF LEMMA 13

*Proof of Lemma 13.* Recall that

$$\begin{aligned}
 &\frac{\partial^2 f_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} \\
 &= \mathbb{E}_{\mathbf{X}} \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \left( \frac{g(\mathbf{W}^*; \mathbf{a}_n)}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - g(\mathbf{W}^*; \mathbf{a}_n)}{(1 - g(\mathbf{W}; \mathbf{a}_n))^2} \right) \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T,
 \end{aligned} \tag{122}$$

and

$$\begin{aligned}
 \frac{\partial^2 f_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j^2} &= \mathbb{E}_{\mathbf{X}} \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \left( \frac{g(\mathbf{W}^*; \mathbf{a}_n)}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - g(\mathbf{W}^*; \mathbf{a}_n)}{(1 - g(\mathbf{W}; \mathbf{a}_n))^2} \right) \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)^2 (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T \\
 &\quad - \mathbb{E}_{\mathbf{X}} \frac{1}{K |\Omega_t|} \sum_{n \in \Omega_t} \left( -\frac{g(\mathbf{W}^*; \mathbf{a}_n)}{g(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - g(\mathbf{W}^*; \mathbf{a}_n)}{1 - g(\mathbf{W}; \mathbf{a}_n)} \right) \phi''(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T.
 \end{aligned} \tag{123}$$

Let us denote  $A_{j,k}(\mathbf{W}; \mathbf{a}_n)$  as

$$A_{j,k}(\mathbf{W}; \mathbf{a}_n) = \begin{cases} \frac{1}{K^2} \left( \frac{g(\mathbf{W}^*; \mathbf{a}_n)}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - g(\mathbf{W}^*; \mathbf{a}_n)}{(1 - g(\mathbf{W}; \mathbf{a}_n))^2} \right) \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n) \\ \quad - \frac{1}{K} \left( -\frac{g(\mathbf{W}^*; \mathbf{a}_n)}{g(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - g(\mathbf{W}^*; \mathbf{a}_n)}{1 - g(\mathbf{W}; \mathbf{a}_n)} \right) \phi''(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n), & \text{when } j = k; \\ \frac{1}{K^2} \left( \frac{g(\mathbf{W}^*; \mathbf{a}_n)}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - g(\mathbf{W}^*; \mathbf{a}_n)}{(1 - g(\mathbf{W}; \mathbf{a}_n))^2} \right) \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n), & \text{when } j \neq k. \end{cases} \tag{124}$$

Further, let us define  $M(\mathbf{W}; \mathbf{a}_n) = \max \left\{ \frac{2}{K^3} \frac{1}{g^3(\mathbf{W}; \mathbf{a}_n)}, \frac{2}{K^3} \frac{1}{(1 - g(\mathbf{W}; \mathbf{a}_n))^3}, \frac{1}{K^2} \frac{1}{g^2(\mathbf{W}; \mathbf{a}_n)}, \frac{1}{K^2} \frac{1}{(1 - g(\mathbf{W}; \mathbf{a}_n))^2} \right\}$ .

Then, by the mean value theorem, we have

$$A_{j,k}(\mathbf{W}; \mathbf{a}_n) - A_{j,k}(\mathbf{W}^*; \mathbf{a}_n) = \sum_{l=1}^K \left\langle \frac{\partial A_{j,k}}{\partial \mathbf{w}_l}(\tilde{\mathbf{W}}; \mathbf{a}_n), \mathbf{w}_l - \mathbf{w}_l^* \right\rangle. \tag{125}$$

For  $\frac{\partial A_{j,k}}{\partial \mathbf{w}_l}$ , we have

$$\frac{\partial A_{j,k}}{\partial \mathbf{w}_l}(\tilde{\mathbf{W}}; \mathbf{a}_n) = B_{j,k,l}(\tilde{\mathbf{W}}; \mathbf{a}_n) \mathbf{X}^T \mathbf{a}_n \tag{126}$$

with

$$|B_{j,k,l}(\widetilde{\mathbf{W}}; \mathbf{a}_n)| \leq \frac{2}{K^3} \frac{1}{g^3(\widetilde{\mathbf{W}}; \mathbf{a}_n)} + \frac{2}{K^3} \frac{1}{(1-g(\widetilde{\mathbf{W}}; \mathbf{a}_n))^3} + \frac{1}{K^2} \frac{1}{g(\widetilde{\mathbf{W}}; \mathbf{a}_n)} + \frac{1}{K^2} \frac{1}{(1-g(\widetilde{\mathbf{W}}; \mathbf{a}_n))^2} \leq 4M(\widetilde{\mathbf{W}}; \mathbf{a}_n). \quad (127)$$

for all  $j \in [K], k \in [K], l \in [K]$ .

Therefore, for any  $\boldsymbol{\alpha} \in \mathbb{R}^{Kd}$ , we have

$$\begin{aligned} & \boldsymbol{\alpha}^T \nabla^2 f_{\Omega_t}(\mathbf{W}) \boldsymbol{\alpha} \\ & \leq \frac{1}{|\Omega_t|} \sum_{n \in |\Omega_t|} \sum_{j=1}^K \sum_{k=1}^K \mathbb{E}_{\mathbf{X}} \left| \boldsymbol{\alpha}_j^T \frac{\partial f_{\Omega_t}}{\partial \mathbf{w}_j \partial \mathbf{w}_k}(\mathbf{W}) \boldsymbol{\alpha}_k \right| \\ & = \frac{1}{|\Omega_t|} \sum_{n \in |\Omega_t|} \sum_{j=1}^K \sum_{k=1}^K \mathbb{E}_{\mathbf{X}} \left| \sum_{l=1}^K |B_{j,k,l}(\widetilde{\mathbf{W}}; \mathbf{a}_n)| \langle \mathbf{w}_l - \mathbf{w}_l^*, \mathbf{X}^T \mathbf{a}_n \rangle \langle \boldsymbol{\alpha}_j, \mathbf{X}^T \mathbf{a}_n \rangle \langle \boldsymbol{\alpha}_k, \mathbf{X}^T \mathbf{a}_n \rangle \right| \\ & = \frac{1}{|\Omega_t|} \sum_{n \in |\Omega_t|} \sum_{j=1}^K \sum_{k=1}^K \left( \sum_{l=1}^K \mathbb{E}_{\mathbf{X}} |B_{j,k,l}(\widetilde{\mathbf{W}}; \mathbf{a}_n)|^2 \right)^{\frac{1}{2}} \left( \sum_{l=1}^K \mathbb{E}_{\mathbf{X}} |\langle \mathbf{w}_l - \mathbf{w}_l^*, \mathbf{X}^T \mathbf{a}_n \rangle \langle \boldsymbol{\alpha}_j, \mathbf{X}^T \mathbf{a}_n \rangle \langle \boldsymbol{\alpha}_k, \mathbf{X}^T \mathbf{a}_n \rangle|^2 \right)^{\frac{1}{2}} \\ & \leq \frac{1}{|\Omega_t|} \sum_{n \in |\Omega_t|} \sum_{j=1}^K \sum_{k=1}^K 36K^{\frac{1}{2}} \left( \mathbb{E}_{\mathbf{X}} M^2(\widetilde{\mathbf{W}}; \mathbf{a}_n) \right)^{\frac{1}{2}} \cdot \left( \sum_{l=1}^K \|\mathbf{w}_l - \mathbf{w}_l^*\|_2^2 \right)^{\frac{1}{2}} \|\boldsymbol{\alpha}_j\|_2 \|\boldsymbol{\alpha}_k\|_2 \\ & \leq \frac{1}{|\Omega_t|} \sum_{n \in |\Omega_t|} 36K^3 \left( \mathbb{E}_{\mathbf{X}} M^2(\widetilde{\mathbf{W}}; \mathbf{a}_n) \right)^{\frac{1}{2}} \|\mathbf{W} - \mathbf{W}^*\|_2 \\ & \stackrel{(a)}{\lesssim} e^{\sigma_1^2(\mathcal{A})} \|\mathbf{W} - \mathbf{W}^*\|_2 \\ & \lesssim \|\mathbf{W} - \mathbf{W}^*\|_2, \end{aligned} \quad (128)$$

where (a) comes from Lemma 5 in (Fu et al., 2018).  $\square$

### B.2.3. PROOF OF LEMMA 14

*Proof of Lemma 14.* Recall that

$$\begin{aligned} & \frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} \\ & = \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \left( \frac{y_n}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1-y_n}{(1-g(\mathbf{W}; \mathbf{a}_n))^2} \right) \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T, \end{aligned} \quad (129)$$

and

$$\begin{aligned} \frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j^2} & = \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \left( \frac{y_n}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1-y_n}{(1-g(\mathbf{W}; \mathbf{a}_n))^2} \right) \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)^2 (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T \\ & \quad - \frac{1}{K |\Omega_t|} \sum_{n \in \Omega_t} \left( -\frac{y_n}{g(\mathbf{W}; \mathbf{a}_n)} + \frac{1-y_n}{1-g(\mathbf{W}; \mathbf{a}_n)} \right) \phi''(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T. \end{aligned} \quad (130)$$

When  $y_n = 1$  and  $j \neq k$ , we have

$$\frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} = \frac{1}{K^2 |\Omega_t|} \sum_{n \in \Omega_t} \frac{\phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)}{g^2(\mathbf{W}; \mathbf{a}_n)} (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T, \quad (131)$$



and

$$\begin{aligned} \frac{\phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)}{g^2(\mathbf{W}; \mathbf{a}_n)} &= \frac{\phi(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)(1 - \phi(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)) \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)(1 - \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n))}{\left(\frac{1}{K} \sum_{l=1}^K \phi(\mathbf{w}_l^T \mathbf{X}^T \mathbf{a}_n)\right)^2} \\ &\leq K^2 \frac{\phi(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)(1 - \phi(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)) \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)(1 - \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n))}{\phi(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)} \\ &= K^2 (1 - \phi(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n))(1 - \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)) \leq K^2. \end{aligned} \quad (132)$$

When  $y_n = 1$  and  $j = k$ , we have

$$\frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} = \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} \left[ \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1}{K} \frac{\phi''(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)}{g(\mathbf{W}; \mathbf{a}_n)} \right] (\mathbf{X}^T \mathbf{a}_n) (\mathbf{X}^T \mathbf{a}_n)^T, \quad (133)$$

and

$$\left| \frac{\phi''(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)}{g(\mathbf{W}; \mathbf{a}_n)} \right| = \frac{\phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)(1 - \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)) \cdot |1 - 2\phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)|}{\frac{1}{K} \sum_{l=1}^K \phi(\mathbf{w}_l^T \mathbf{X}^T \mathbf{a}_n)} \leq K. \quad (134)$$

Similar to (132) and (134), we can obtain the following inequality for  $y_n = 0$ .

$$\frac{\phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)}{(1 - g(\mathbf{W}; \mathbf{a}_n))^2} \leq K^2, \quad \text{and} \quad \left| \frac{\phi''(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)}{1 - g(\mathbf{W}; \mathbf{a}_n)} \right| \leq K. \quad (135)$$

Then, for any  $\boldsymbol{\alpha} \in \mathbb{R}^d$ , we have

$$\begin{aligned} \boldsymbol{\alpha}^T \frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} \boldsymbol{\alpha} &= \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} \left[ \frac{1}{K^2} \left( \frac{y_n}{g^2(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - y_n}{(1 - g(\mathbf{W}; \mathbf{a}_n))^2} \right) \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \phi'(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n) \right. \\ &\quad \left. - \frac{\mathbb{1}_{\{j=k\}}}{K} \left( -\frac{y_n}{g(\mathbf{W}; \mathbf{a}_n)} + \frac{1 - y_n}{1 - g(\mathbf{W}; \mathbf{a}_n)} \right) \phi''(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \right] (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2 \\ &:= \frac{1}{|\Omega_t|} \sum_{n \in \Omega_t} H_{j,k}(\mathbf{a}_n) \cdot (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2. \end{aligned} \quad (136)$$

Next, we show that  $H_{j,k}(\mathbf{a}_n) \cdot (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2$  belongs to the sub-exponential distribution. For any  $p \in \mathbb{N}^+$ , we have

$$\begin{aligned} \left( \mathbb{E}_{\mathbf{X}, \mathbf{y}_n} \left[ |H_{j,k}(\mathbf{a}_n) \cdot (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2|^p \right] \right)^{1/p} &\leq \left( \mathbb{E}_{\mathbf{X}} \left[ |4(\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2|^p \right] \right)^{1/p} \\ &\leq 8 \|\mathbf{a}_n\|_{2p}^2 \leq 8 \sigma_1^2(\mathbf{A}) p \end{aligned} \quad (137)$$

Hence,  $H_{j,k}(\mathbf{a}_n) \cdot (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2$  belongs to the sub-exponential distribution with  $\|H_{j,k}(\mathbf{a}_n) (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2\|_{\psi_1} = 8 \sigma_1^2(\mathbf{A})$ . Then, the moment generation function of  $H_{j,k}(\mathbf{a}_n) \cdot (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2$  can be bounded as

$$\mathbb{E} e^{s H_{j,k}(\mathbf{a}_n) \cdot (\boldsymbol{\alpha}^T \mathbf{X}^T \mathbf{a}_n)^2} \leq e^{C \sigma_1^2(\mathbf{A}) s^2} \quad (138)$$

for some positive constant  $C$  and any  $s \in \mathbb{R}$ . From Lemma 7 and Chernoff bound, we have

$$\boldsymbol{\alpha}^T \left( \frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} - \frac{\partial^2 f_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} \right) \boldsymbol{\alpha} \leq C \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2) d \log N}{|\Omega_t|}} \quad (139)$$

with probability at least  $1 - N^{-d}$ . By selecting  $\xi = \frac{1}{2}$  in Lemmas 8 and 9, we have

$$\left\| \frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} - \frac{\partial^2 f_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} \right\|_2 \leq C \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2) d \log N}{|\Omega_t|}} \quad (140)$$

with probability at least  $1 - (\frac{5}{N})^d$ .

In conclusion, we have

$$\begin{aligned} \|\nabla^2 f_{\Omega_t}(\mathbf{W}) - \nabla^2 \hat{f}_{\Omega_t}(\mathbf{W})\|_2 &\leq \sum_{j=1}^K \sum_{k=1}^K \left\| \frac{\partial^2 \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} - \frac{\partial^2 f_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_k} \right\|_2 \\ &\leq CK^2 \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log d}{|\Omega_t|}} \end{aligned} \quad (141)$$

with probability at least  $1 - (\frac{5}{d})^d$ .  $\square$

#### B.2.4. PROOF OF LEMMA 11

*Proof of Lemma 11.* Recall that the first-order derivative of  $\hat{f}_{\Omega_t}(\mathbf{W})$  is calculated from

$$\frac{\partial \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j} = -\frac{1}{K|\Omega_t|} \sum_{n \in \Omega} \frac{y_n - g(\mathbf{W}; \mathbf{a}_n)}{g(\mathbf{W}; \mathbf{a}_n)(1 - g(\mathbf{W}; \mathbf{a}_n))} \phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n) \mathbf{X}^T \mathbf{a}_n. \quad (142)$$

Similar to (134), we have

$$\left| \frac{\phi'(\mathbf{w}_j^T \mathbf{X}^T \mathbf{a}_n)}{g(\mathbf{W}; \mathbf{a}_n)} \right| = \frac{\phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n)(1 - \phi(\mathbf{w}_k^T \mathbf{X}^T \mathbf{a}_n))}{\frac{1}{K} \sum_{l=1}^K \phi(\mathbf{w}_l^T \mathbf{X}^T \mathbf{a}_n)} \leq K. \quad (143)$$

Similar to (137), for any fixed  $\boldsymbol{\alpha} \in \mathbb{R}^{dK}$ , we can show that random variable  $\boldsymbol{\alpha}^T \frac{\partial \hat{f}_{\Omega_t}(\mathbf{W})}{\partial \mathbf{w}_j}$  belongs to sub-exponential distribution with the same bounded norm up to a constant. Hence, by applying Lemma 7 and the Chernoff bound, we have

$$\left\| \nabla f_{\Omega_t}(\mathbf{W}) - \nabla \hat{f}_{\Omega_t}(\mathbf{W}) \right\|_2 \lesssim K^2 \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega_t|}} \quad (144)$$

with probability at least  $1 - (\frac{5}{N})^d$ .  $\square$

### C. Proof of Lemma 1

*Proof of Lemma 1.* Let  $\tilde{\mathbf{A}}$  denote the adjacency matrix, then we have

$$\sigma_1(\tilde{\mathbf{A}}) = \max_{\mathbf{z}} \frac{\mathbf{z}^T \tilde{\mathbf{A}} \mathbf{z}}{\mathbf{z}^T \mathbf{z}} \geq \frac{\mathbf{1}^T \tilde{\mathbf{A}} \mathbf{1}}{\mathbf{1}^T \mathbf{1}} = 1 + \frac{\sum_{n=1}^N \delta_n}{N}, \quad (145)$$

where  $\delta_n$  denotes the degree of node  $v_n$ . Let  $\mathbf{z}$  be the eigenvector of the maximum eigenvalue  $\sigma_1(\mathbf{A})$ . Since  $\sigma_1(\mathbf{A}) = \mathbf{D}^{-1/2} \tilde{\mathbf{A}} \mathbf{D}^{-1/2}$  and  $\mathbf{D}$  is diagonal matrix, then  $\mathbf{z}$  is the eigenvector to  $\sigma_1(\tilde{\mathbf{A}})$  as well. Then, let  $n \in [N]$  be the index of the largest value of vector  $\mathbf{z}_n$  as  $z_n = \|\mathbf{z}\|_\infty$ , we have

$$\sigma_1(\tilde{\mathbf{A}}) = \frac{(\tilde{\mathbf{A}} \mathbf{z})_n}{z_n} = \frac{\tilde{\mathbf{a}}_n^T \mathbf{z}}{z_n} \leq \frac{\|\mathbf{a}_n\|_1 \|\mathbf{z}\|_\infty}{z_n} = 1 + \delta. \quad (146)$$

where  $\tilde{\mathbf{a}}_n$  is the  $n$ -th row of  $\tilde{\mathbf{A}}$ .

Since  $\mathbf{D}$  is a diagonal matrix with  $\|\mathbf{D}\|_2 \leq 1 + \delta$ , then we can conclude the inequality in this lemma.  $\square$

### D. Proof of Lemma 4

The proof of Lemma 4 is divided into three major parts to bound  $I_1$ ,  $I_2$  and  $I_3$  in (153). Lemmas 15, 16 and 17 provide the error bounds for  $I_1$ ,  $I_2$  and  $I_3$ , respectively. The proofs of these preliminary lemmas are similar to those of Theorem 5.6 in (Zhong et al., 2017b), the difference is to apply Lemma 7 plus Chernoff inequality instead of standard Hoeffding inequality, and we skip the details of the proofs of Lemmas 15, 16 and 17 here.

**Lemma 15.** Suppose  $M_2$  is defined as in (7) and  $\widehat{M}_2$  is the estimation of  $M_2$  by samples. Then, with probability  $1 - N^{-10}$ , we have

$$\|\widehat{M}_2 - M_2\| \lesssim \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega|}}, \quad (147)$$

provided that  $|\Omega| \gtrsim (1 + \delta^2)d \log^4 N$ .

**Lemma 16.** Let  $\widehat{\mathbf{V}}$  be generated by step 4 in Subroutine 1. Suppose  $M_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  is defined as in (9) and  $\widehat{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  is the estimation of  $M_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  by samples. Further, we assume  $\mathbf{V} \in \mathbb{R}^{d \times K}$  is an orthogonal basis of  $\mathbf{W}^*$  and satisfies  $\|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\| \leq 1/4$ . Then, provided that  $N \gtrsim K^5 \log^6 d$ , with probability at least  $1 - N^{-10}$ , we have

$$\|\widehat{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - M_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})\| \lesssim \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)K^3 \log N}{|\Omega|}}. \quad (148)$$

**Lemma 17.** Suppose  $M_1$  is defined as in (6) and  $\widehat{M}_1$  is the estimation of  $M_1$  by samples. Then, with probability  $1 - N^{-10}$ , we have

$$\|\widehat{M}_1 - M_1\| \lesssim \sigma_1^2(\mathbf{A}) \sqrt{\frac{(1 + \delta^2)d \log N}{|\Omega|}} \quad (149)$$

provided that  $|\Omega| \gtrsim (1 + \delta^2)d \log^4 N$ .

**Lemma 18** ((Zhong et al., 2017b), Lemma E.6). Let  $\mathbf{V} \in \mathbb{R}^{d \times K}$  be an orthogonal basis of  $\mathbf{W}^*$  and  $\widehat{\mathbf{V}}$  be generated by step 4 in Subroutine 1. Assume  $\|\widehat{M}_2 - M_2\|_2 \leq \sigma_K(M_2)/10$ . Then, for some small  $\varepsilon_0$ , we have

$$\|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\|_2 \leq \frac{\|M_2 - \widehat{M}_2\|}{\sigma_K(M_2)}. \quad (150)$$

**Lemma 19** ((Zhong et al., 2017b), Lemma E.13). Let  $\mathbf{V} \in \mathbb{R}^{d \times K}$  be an orthogonal basis of  $\mathbf{W}^*$  and  $\widehat{\mathbf{V}}$  be generated by step 4 in Subroutine 1. Assume  $M_1$  can be written in the form of (6) with some homogeneous function  $\phi_1$ , and let  $\widehat{M}_1$  be the estimation of  $M_1$  by samples. Let  $\widehat{\alpha}$  be the optimal solution of (11) with  $\widehat{\mathbf{w}}_j = \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j$ . Then, for each  $j \in [K]$ , if

$$\begin{aligned} T_1 &:= \|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\|_2 \leq \frac{1}{\kappa^2 \sqrt{K}}, \\ T_2 &:= \|\widehat{\mathbf{u}}_j - \widehat{\mathbf{V}}^T \widehat{\mathbf{w}}_j\|_2 \leq \frac{1}{\kappa^2 \sqrt{K}}, \\ T_3 &:= \|\widehat{M}_1 - M_1\|_2 \leq \frac{1}{4} \|M_1\|_2, \end{aligned} \quad (151)$$

then we have

$$\left| \|\mathbf{w}_j\|_2 - \widehat{\alpha}_j \right| \leq \left( \kappa^4 K^{\frac{3}{2}} (T_1 + T_2) + \kappa^2 K^{\frac{1}{2}} T_3 \right) \|\mathbf{W}^*\|_2. \quad (152)$$

*Proof of Lemma 4.* we have

$$\begin{aligned} \|\mathbf{w}_j^* - \widehat{\alpha}_j \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j\|_2 &\leq \left\| \mathbf{w}_j^* - \|\mathbf{w}_j\|_2 \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j + \|\mathbf{w}_j\|_2 \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j - \widehat{\alpha}_j \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j \right\|_2 \\ &\leq \left\| \mathbf{w}_j^* - \|\mathbf{w}_j\|_2 \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j \right\|_2 + \left\| \|\mathbf{w}_j\|_2 \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j - \widehat{\alpha}_j \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j \right\|_2 \\ &\leq \|\mathbf{w}_j^*\|_2 \|\overline{\mathbf{w}}_j^* - \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j\|_2 + \left| \|\mathbf{w}_j\|_2 - \widehat{\alpha}_j \right| \|\widehat{\mathbf{V}}\widehat{\mathbf{u}}_j\|_2 \\ &\leq \sigma_1 \left( \|\overline{\mathbf{w}}_j^* - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T \overline{\mathbf{w}}_j^*\|_2 + \|\widehat{\mathbf{V}}^T \overline{\mathbf{w}}_j^* - \widehat{\mathbf{u}}_j\|_2 \right) + \left| \|\mathbf{w}_j\|_2 - \widehat{\alpha}_j \right| \\ &:= \sigma_1 (I_1 + I_2) + I_3. \end{aligned} \quad (153)$$

From Lemma 18, we have

$$I_1 = \|\overline{\mathbf{w}}_j^* - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T \overline{\mathbf{w}}_j^*\|_2 \leq \|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\|_2 \leq \frac{\|\widehat{M}_2 - M_2\|_2}{\sigma_K(M_2)}, \quad (154)$$

where the last inequality comes from Lemma 15. Then, from (7), we know that

$$\sigma_K(\mathbf{M}_2) \lesssim \min_{1 \leq j \leq K} \|\mathbf{w}_j^*\|_2 \lesssim \sigma_K(\mathbf{W}^*). \quad (155)$$

From Theorem 3 in (Kuleshov et al., 2015), we have

$$I_2 = \|\widehat{\mathbf{V}}^T \overline{\mathbf{w}}_j^* - \widehat{\mathbf{u}}_j\|_2 \lesssim \frac{\kappa}{\sigma_K(\mathbf{W}^*)} \|\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})\|_2. \quad (156)$$

To guarantee the condition (151) in Lemma 19 hold, according to Lemmas 15 and 16, we need  $|\Omega| \gtrsim \kappa^3(1 + \delta^2)Kd \log N$ . Then, from Lemma 19, we have

$$I_3 = \left( \kappa^4 K^{3/2} (I_1 + I_2) + \kappa^2 K^{1/2} \|\widehat{\mathbf{M}}_1 - \mathbf{M}_1\| \right) \|\mathbf{W}^*\|_2. \quad (157)$$

Since  $d \gg K$ , according to Lemmas 15, 16 and 17, we have

$$\|\mathbf{w}_j^* - \widehat{\alpha}_j \widehat{\mathbf{V}} \widehat{\mathbf{u}}_j\|_2 \lesssim \kappa^6 \sigma_1^2(\mathbf{A}) \sqrt{\frac{K^3(1 + \delta^2)d \log N}{|\Omega|}} \|\mathbf{W}^*\|_2 \quad (158)$$

provided  $|\Omega| \gtrsim (1 + \delta^2)d \log^4 N$ . □