
Understanding Generalization and Optimization Performance for Deep CNNs (Supplementary File)

Pan Zhou*

Jiashi Feng*

* National University of Singapore, Singapore

pzhou@u.nus.edu

elefjia@nus.edu.sg

A Structure of This Document

This document gives some other necessary notations and preliminaries for our analysis in Sec. B.1 and provides auxiliary lemmas in Sec. B.2. Then in Sec. C, we present the technical lemmas for proving our final results and their proofs. Next, in Sec. D, we utilize these technical lemmas to prove our desired results. Finally, we give the proofs of other auxiliary lemmas in Sec. E.

As for the results in manuscript, the proofs of Lemma 1 and Theorem 1 in Sec. 3.1 in the manuscript are respectively provided in Sec. D.1 and Sec. D.2. As for the results in Sec. 4 in the manuscript, Sec. D.3 and D.4 present the proofs of Theorem 2 and Corollary 1, respectively. Finally, we respectively introduce the proofs of Theorem 3 and Corollary 2 in Sec. D.5 and D.6.

B Notations and Preliminary Tools

Beyond the notations introduced in the manuscript, we need some other notations used in this document. Then we introduce several lemmas that will be used later.

B.1 Notations

Throughout this document, we use $\langle \cdot, \cdot \rangle$ to denote the inner product and use $\tilde{\otimes}$ to denote the convolution operation with stride 1. $\mathbf{A} \otimes \mathbf{C}$ denotes the Kronecker product between \mathbf{A} and \mathbf{C} . Note that \mathbf{A} and \mathbf{C} in $\mathbf{A} \otimes \mathbf{C}$ can be matrices or vectors. For a matrix $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$, we use $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} \mathbf{A}_{ij}^2}$ to denote its Frobenius norm, where \mathbf{A}_{ij} is the (i, j) -th entry of \mathbf{A} . We use $\|\mathbf{A}\|_{\text{op}} = \max_i |\lambda_i(\mathbf{A})|$ to denote the operation norm of a matrix $\mathbf{A} \in \mathbb{R}^{n_1 \times n_1}$, where $\lambda_i(\mathbf{A})$ denotes the i -th eigenvalue of the matrix \mathbf{A} . For a 3-way tensor $\mathbf{A} \in \mathbb{R}^{s \times t \times q}$, its operation norm is computed as

$$\|\mathbf{A}\|_{\text{op}} = \sup_{\|\boldsymbol{\lambda}\|_2 \leq 1} \langle \boldsymbol{\lambda}^{\otimes 3}, \mathbf{A} \rangle = \sum_{i,j,k} \mathbf{A}_{ijk} \lambda_i \lambda_j \lambda_k,$$

where \mathbf{A}_{ijk} denotes the (i, j, k) -th entry of \mathbf{A} .

For brevity, in this document we use $f(\mathbf{w}, \mathbf{D})$ to denote $f(g(\mathbf{w}; \mathbf{D}), \mathbf{y})$ in the manuscript. Let $\mathbf{w}_{(i)} = (\mathbf{w}_{(i)}^1; \dots; \mathbf{w}_{(i)}^{d_i}) \in \mathbb{R}^{k_i^2 d_{i-1} d_i}$ ($i = 1, \dots, l$) be the parameter of the i -th layer where $\mathbf{w}_{(i)}^k = \text{vec}(\mathbf{W}_{(i)}^k) \in \mathbb{R}^{k_i^2 d_{i-1}}$ is the vectorization of $\mathbf{W}_{(i)}^k$. Similarly, let $\mathbf{w}_{(l+1)} = \text{vec}(\mathbf{W}_{(l+1)}) \in \mathbb{R}^{r_l c_l d_l d_{l+1}}$. Then, we further define $\mathbf{w} = (\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(l)}, \mathbf{w}_{(l+1)}) \in \mathbb{R}^{r_1 c_1 d_1 d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i}$ which contains all the parameter in the network. Here we use $\mathbf{W}_{(i)}^k$ to denote the k -th kernel in the i -th convolutional layer. For brevity, let $\mathbf{W}_{(i)}^{k,j}$ denotes the j -th slice of $\mathbf{W}_{(i)}^k$, i.e. $\mathbf{W}_{(i)}^{k,j} = \mathbf{W}_{(i)}^k(:, :, j)$.

For a matrix $\mathbf{M} \in \mathbb{R}^{s \times t}$, $\widehat{\mathbf{M}}$ denotes the matrix which is obtained by rotating the matrix \mathbf{M} by 180 degrees. Then we use δ_i to denote the gradient of $f(\mathbf{w}, \mathbf{D})$ w.r.t. $\mathbf{X}_{(i)}$:

$$\delta_i = \nabla_{\mathbf{X}_{(i)}} f(\mathbf{w}, \mathbf{D}) \in \mathbb{R}^{r_i \times c_i \times d_i}, \quad (i = 1, \dots, l),$$

Based on δ_i , we further define $\tilde{\delta}_i \in \mathbb{R}^{(\tilde{r}_{i-1}-k_i+1) \times (\tilde{c}_{i-1}-k_i+1) \times d_i}$. Each slice $\tilde{\delta}_{i+1}^k$ can be computed as follows. Firstly, let $\tilde{\delta}_{i+1}^k = \delta_{i+1}^k$. Then, we pad zeros of $s_i - 1$ rows between the neighboring rows in $\tilde{\delta}_{i+1}^k$ and similarly we pad zeros of $s_i - 1$ columns between the neighboring columns in $\tilde{\delta}_{i+1}^k$. Accordingly, the size of $\tilde{\delta}_{i+1}^k$ is $(s_i(r_i - 1) + 1) \times (s_i(c_i - 1) + 1)$. Finally, we pad zeros of width $k_i - 1$ around $\tilde{\delta}_{i+1}^k$ to obtain new $\tilde{\delta}_{i+1}^k \in \mathbb{R}^{(s_i(r_i-1)+2k_i-1) \times (s_i(c_i-1)+2k_i-1)}$. Note that since $r_{i+1} = (\tilde{r}_i - k_{i+1})/s_{i+1} + 1$ and $r_{i+1} = (\tilde{r}_i - k_{i+1})/s_{i+1} + 1$, we have $s_i(r_i - 1) + 2k_i - 1 = \tilde{r}_{i-1} - k_i + 1$ and $s_i(c_i - 1) + 2k_i - 1 = \tilde{c}_{i-1} - k_i + 1$.

Then we define four operators $\mathbf{G}(\cdot)$, $\mathbf{Q}(\cdot)$, $\text{up}(\cdot)$ and $\text{vec}(\cdot)$, which are useful for explaining the following analysis.

Operation $\mathbf{G}(\cdot)$: It maps an arbitrary vector $\mathbf{z} \in \mathbb{R}^d$ into a diagonal matrix $\mathbf{G}(\mathbf{z}) \in \mathbb{R}^{d \times d}$ with its i -th diagonal entry equal to $\sigma_1(\mathbf{z}_i)(1 - \sigma_1(\mathbf{z}_i))$ in which \mathbf{z}_i denotes the i -th entry of \mathbf{z} .

Operation $\mathbf{Q}(\cdot)$: it maps a vector $\mathbf{z} \in \mathbb{R}^d$ into a matrix of size $d^2 \times d$ whose $((i-1)d + i, i)$ ($i = 1, \dots, d$) entry equal to $\sigma_1(\mathbf{z}_i)(1 - \sigma_1(\mathbf{z}_i))(1 - 2\sigma_1(\mathbf{z}_i))$ and rest entries are all 0.

Operation $\text{up}(\cdot)$: $\text{up}(\mathbf{M})$ represents conducting upsampling on $\mathbf{M} \in \mathbb{R}^{s \times t \times q}$. Let $\mathbf{N} = \text{up}(\mathbf{M}) \in \mathbb{R}^{ps \times pt \times pq}$. Specifically, for each slice $\mathbf{N}(:, :, i)$ ($i = 1, \dots, q$), we have $\mathbf{N}(:, :, i) = \text{up}(\mathbf{M}(:, :, i))$. It actually upsamples each entry $M(g, h, i)$ into a matrix of p^2 same entries $\frac{1}{p^2}M(g, h, i)$.

Operation $\text{vec}(\cdot)$: For a matrix $\mathbf{A} \in \mathbb{R}^{s \times t}$, $\text{vec}(\mathbf{A})$ is defined as $\text{vec}(\mathbf{A}) = (\mathbf{A}(:, 1); \dots; \mathbf{A}(:, t)) \in \mathbb{R}^{st}$ that vectorizes $\mathbf{A} \in \mathbb{R}^{s \times t}$ along its columns. If $\mathbf{A} \in \mathbb{R}^{t \times s \times q}$ is a 3-way tensor, $\text{vec}(\mathbf{A}) = [\text{vec}(\mathbf{A}(:, :, 1)); \text{vec}(\mathbf{A}(:, :, 2)); \dots; \text{vec}(\mathbf{A}(:, :, q))] \in \mathbb{R}^{stq}$.

B.2 Auxiliary Lemmas

We first introduce Lemmas 2 and 3 which are respectively used for bounding the ℓ_2 -norm of a vector and the operation norm of a matrix. Then we introduce Lemmas 4 and 5 which discuss the topology of functions. In Lemma 6, we introduce the Hoeffding's inequality which provides an upper bound on the probability that the sum of independent random variables deviates from its expected value. In Lemma 7, we provide the covering number of an ϵ -net for a low-rank matrix. Finally, several commonly used inequalities are presented in Lemma 8 for our analysis.

Lemma 2. (Vershynin, 2012) For any vector $\mathbf{x} \in \mathbb{R}^d$, its ℓ_2 -norm can be bounded as

$$\|\mathbf{x}\|_2 \leq \frac{1}{1 - \epsilon} \sup_{\boldsymbol{\lambda} \in \boldsymbol{\lambda}_\epsilon} \langle \boldsymbol{\lambda}, \mathbf{x} \rangle.$$

where $\boldsymbol{\lambda}_\epsilon = \{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k_w}\}$ be an ϵ -covering net of $\mathbb{B}^d(1)$.

Lemma 3. (Vershynin, 2012) For any symmetric matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$, its operator norm can be bounded as

$$\|\mathbf{X}\|_{op} \leq \frac{1}{1 - 2\epsilon} \sup_{\boldsymbol{\lambda} \in \boldsymbol{\lambda}_\epsilon} |\langle \boldsymbol{\lambda}, \mathbf{X}\boldsymbol{\lambda} \rangle|.$$

where $\boldsymbol{\lambda}_\epsilon = \{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k_w}\}$ be an ϵ -covering net of $\mathbb{B}^d(1)$.

Lemma 4. (Mei et al., 2017) Let $D \subseteq \mathbb{R}^d$ be a compact set with a C^2 boundary ∂D , and $f, g : A \rightarrow \mathbb{R}$ be C^2 functions defined on an open set A , with $D \subseteq A$. Assume that for all $\mathbf{w} \in \partial D$ and all $t \in [0, 1]$, $t\nabla f(\mathbf{w}) + (1-t)\nabla g(\mathbf{w}) \neq \mathbf{0}$. Finally, assume that the Hessian $\nabla^2 f(\mathbf{w})$ is non-degenerate and has index equal to r for all $\mathbf{w} \in D$. Then the following properties hold:

- (1) If g has no critical point in D , then f has no critical point in D .
- (2) If g has a unique critical point \mathbf{w} in D that is non-degenerate with an index of r , then f also has a unique critical point \mathbf{w}' in D with the index equal to r .

Lemma 5. (Mei et al., 2017) Suppose that $F(\mathbf{w}) : \Theta \rightarrow \mathbb{R}$ is a C^2 function where $\mathbf{w} \in \Theta$. Assume that $\{\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(m)}\}$ is its non-degenerate critical points and let $D = \{\mathbf{w} \in \Theta : \|\nabla F(\mathbf{w})\|_2 < \epsilon \text{ and } \inf_i |\lambda_i(\nabla^2 F(\mathbf{w}))| \geq \zeta\}$. Then D can be decomposed into (at most) countably components, with each component containing either exactly one critical point, or no critical point. Concretely, there exist disjoint open sets $\{D_k\}_{k \in \mathbb{N}}$, with D_k possibly empty for $k \geq m + 1$, such that

$$D = \bigcup_{k=1}^{\infty} D_k.$$

Furthermore, $\mathbf{w}_{(k)} \in D_k$ for $1 \leq k \leq m$ and each D_i , $k \geq m + 1$ contains no stationary points.

Lemma 6. (Hoeffding, 1963) Let x_1, \dots, x_n be independent random variables where x_i is bounded by the interval $[a_i, b_i]$. Suppose $s_n = \frac{1}{n} \sum_{i=1}^n x_i$, then the following properties hold:

$$\mathbb{P}(s_n - \mathbb{E}s_n \geq t) \leq \exp\left(-\frac{2n^2 t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

where t is an arbitrary positive value.

Lemma 7. (Candes & Plan, 2009) Let $\mathbb{S}_r = \{\mathbf{X} \in \mathbb{R}^{n_1 \times n_2} : \text{rank}(\mathbf{X}) \leq r, \|\mathbf{X}\|_F = b\}$. Then there exists an ϵ -net $\tilde{\mathbb{S}}_r$ for the Frobenius norm obeying

$$|\tilde{\mathbb{S}}_r| \leq \left(\frac{9b}{\epsilon}\right)^{r(n_1+n_2+1)}.$$

Then we give a lemma to summarize the properties of $\mathbf{G}(\cdot)$, $\mathbf{Q}(\cdot)$ and $\text{up}(\cdot)$ defined in Section B.1, the convolutional operation \otimes defined in manuscript.

Lemma 8. For $\mathbf{G}(\cdot)$, $\mathbf{Q}(\cdot)$ and $\text{up}(\cdot)$ defined in Section B.1, the convolutional operation \otimes defined in manuscript, we have the following properties:

(1) For arbitrary vector \mathbf{z} , and arbitrary matrices \mathbf{M} and \mathbf{N} of proper sizes, we have

$$\|\mathbf{G}(\mathbf{z})\mathbf{M}\|_F^2 \leq \frac{1}{16}\|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N}\mathbf{G}(\mathbf{z})\|_F^2 \leq \frac{1}{16}\|\mathbf{N}\|_F^2.$$

(2) For arbitrary vector \mathbf{z} , and arbitrary matrices \mathbf{M} and \mathbf{N} of proper sizes, we have

$$\|\mathbf{Q}(\mathbf{z})\mathbf{M}\|_F^2 \leq \frac{2^6}{38}\|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N}\mathbf{Q}(\mathbf{z})\|_F^2 \leq \frac{2^6}{38}\|\mathbf{N}\|_F^2.$$

(3) For any tensor $\mathbf{M} \in \mathbb{R}^{s \times t \times q}$, we have

$$\|\text{up}(\mathbf{M})\|_F^2 \leq \frac{1}{p^2}\|\mathbf{M}\|_F^2.$$

(4) For any kernel $\mathbf{W} \in \mathbb{R}^{k_i \times k_i \times d_i}$ and $\tilde{\boldsymbol{\delta}}_{i+1} \in \mathbb{R}^{(\tilde{r}_{i-1}-k_i+1) \times (\tilde{c}_{i-1}-k_i+1) \times d_i}$ defined in Sec. B.1, then we have

$$\|\tilde{\boldsymbol{\delta}}_{i+1} \otimes \mathbf{W}\|_F^2 \leq (k_i - s_i + 1)^2 \|\mathbf{W}\|_F^2 \|\tilde{\boldsymbol{\delta}}_{i+1}\|_F^2.$$

(5) For softmax activation function σ_2 , we can bound the norm of difference between output \mathbf{v} and its corresponding one-hot label \mathbf{y}

$$0 \leq \|\mathbf{v} - \mathbf{y}\|_2^2 \leq 2.$$

It should be pointed out that we defer the proof of Lemma 8 to Sec. E.

C Technical Lemmas and Their Proofs

Here we present the key lemmas and theorems for proving our desired results. For brevity, in this document we use $f(\mathbf{w}, \mathbf{D})$ to denote $f(g(\mathbf{w}; \mathbf{D}), \mathbf{y})$ in the manuscript.

Lemma 9. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Then the gradient of $f(\mathbf{w}, \mathbf{D})$ with respect to $\mathbf{w}_{(j)}$ can be written as

$$\begin{aligned} \nabla_{\mathbf{w}_{(l+1)}} f(\mathbf{w}, \mathbf{D}) &= \mathbf{S}(\mathbf{v} - \mathbf{y}) \mathbf{z}_{(l)}^T \in \mathbb{R}^{d_{l+1} \times \tilde{r}_l \tilde{c}_l d_l}, \\ \nabla_{\mathbf{w}_{(i)}^{k,j}} f(\mathbf{w}, \mathbf{D}) &= \mathbf{Z}_{(i-1)}^j \otimes \tilde{\boldsymbol{\delta}}_i^k \in \mathbb{R}^{k_i \times k_i}, \quad (j = 1, \dots, d_{i-1}; k = 1, \dots, d_i; i = 1, \dots, l), \end{aligned}$$

where δ_i^k is the k -slice (i.e. $\delta_i(:, :, k)$) of δ_i which is defined as

$$\delta_i = \nabla_{\mathbf{X}_{(i)}} f(\mathbf{w}, \mathbf{D}) \in \mathbb{R}^{r_i \times c_i \times d_i}, \quad (i = 1, \dots, l),$$

and further satisfies

$$\delta_i^j = \text{up} \left(\sum_{k=1}^{d_{i+1}} \tilde{\delta}_{i+1}^k \otimes \widehat{\mathbf{W}}_{(i+1)}^{k,j} \right) \odot \sigma'_1(\mathbf{X}_{(i)}^j) \in \mathbb{R}^{r_i \times c_i}, \quad (j = 1, \dots, d_i; i = 1, \dots, l-1),$$

where the matrix $\widehat{\mathbf{W}}_{(i+1)}^{k,j} \in \mathbb{R}^{k_{i+1} \times k_{i+1}}$ is obtained by rotating $\mathbf{W}_{(i+1)}^{k,j}$ with 180 degrees. Also, δ_i is computed as follows:

$$\begin{aligned} \nabla_{\mathbf{u}} f(\mathbf{w}, \mathbf{D}) &= \mathbf{S}(\mathbf{v} - \mathbf{y}) \in \mathbb{R}^{d_{l+1}}, & \nabla_{\mathbf{z}_{(l)}} f(\mathbf{w}, \mathbf{D}) &= (\mathbf{W}_{(l+1)})^T \mathbf{S}(\mathbf{v} - \mathbf{y}) \in \mathbb{R}^{\tilde{r}_l \tilde{c}_l d_l}, \\ \nabla_{\mathbf{Z}_{(l)}} f(\mathbf{w}, \mathbf{D}) &= \text{reshape}(\nabla_{\mathbf{z}_{(l)}} f(\mathbf{w}, \mathbf{D})) \in \mathbb{R}^{\tilde{r}_l \times \tilde{c}_l \times d_l}, & \delta_l &= \sigma'_1(\mathbf{X}_{(l)}) \odot \text{up} \left(\frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{Z}_{(l)}} \right) \in \mathbb{R}^{r_l \times c_l \times d_l}. \end{aligned}$$

where $\mathbf{S} = \mathbf{G}(\mathbf{u})$.

Proof. We use chain rule to compute the gradient of $f(\mathbf{w}, \mathbf{D})$ with respect to $\mathbf{Z}_{(i)}$. We first compute several basis gradient. According to the relationship between $\mathbf{X}_{(i)}$, $\mathbf{Y}_{(i)}$, $\mathbf{Z}_{(i)}$ and $f(\mathbf{w}, \mathbf{D})$, we have

$$\begin{aligned} \nabla_{\mathbf{u}} f(\mathbf{w}, \mathbf{D}) &= \mathbf{S}(\mathbf{v} - \mathbf{y}) \in \mathbb{R}^{d_{l+1}}, \\ \nabla_{\mathbf{z}_{(l)}} f(\mathbf{w}, \mathbf{D}) &= (\mathbf{W}_{(l+1)})^T \mathbf{S}(\mathbf{v} - \mathbf{y}) \in \mathbb{R}^{\tilde{r}_l \tilde{c}_l d_l}, \\ \nabla_{\mathbf{Z}_{(l)}} f(\mathbf{w}, \mathbf{D}) &= \text{reshape}(\nabla_{\mathbf{z}_{(l)}} f(\mathbf{w}, \mathbf{D})) \in \mathbb{R}^{\tilde{r}_l \times \tilde{c}_l \times d_l}, \\ \nabla_{\mathbf{X}_{(l)}} f(\mathbf{w}, \mathbf{D}) &= \frac{\partial \mathbf{Y}_{(l)}}{\partial \mathbf{X}_{(l)}} \frac{\partial \mathbf{Z}_{(l)}}{\partial \mathbf{Y}_{(l)}} \frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{Z}_{(l)}} = \sigma'_1(\mathbf{X}_{(l)}) \odot \text{up} \left(\frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{Z}_{(l)}} \right) \triangleq \delta_l \in \mathbb{R}^{r_l \times c_l \times d_l}. \end{aligned}$$

Then we can further obtain

$$\delta_i^j = \text{up} \left(\sum_{k=1}^{d_{i+1}} \tilde{\delta}_{i+1}^k \otimes \widehat{\mathbf{W}}_{(i+1)}^{k,j} \right) \odot \sigma'_1(\mathbf{X}_{(i)}^j) \in \mathbb{R}^{r_i \times c_i}, \quad (j = 1, \dots, d_i; i = 1, \dots, l-1).$$

where $\widehat{\mathbf{W}}_{(i)}^{k,j}$ denotes the j -th slice $\widehat{\mathbf{W}}_{(i)}^k(:, : j)$ of $\widehat{\mathbf{W}}_{(i)}^k$. Note, we clockwise rotate the matrix $\mathbf{W}_{(i+1)}^{k,j}$ by 180 degrees to obtain $\widehat{\mathbf{W}}_{(i+1)}^{k,j}$. Finally, we can compute the gradient w.r.t. $\mathbf{W}_{(l+1)}$ and $\mathbf{W}_{(i)}$ ($i = 1, \dots, l$):

$$\begin{aligned} \nabla_{\mathbf{W}_{(l+1)}} f(\mathbf{w}, \mathbf{D}) &= \mathbf{S}(\mathbf{v} - \mathbf{y}) \mathbf{z}_{(l)}^T \in \mathbb{R}^{d_{l+1} \times \tilde{r}_l \tilde{c}_l d_l}, \\ \nabla_{\mathbf{W}_{(i)}^{k,j}} f(\mathbf{w}, \mathbf{D}) &= \mathbf{X}_{(i-1)}^j \otimes \tilde{\delta}_i^k \in \mathbb{R}^{k_i \times k_i}, \quad (j = 1, \dots, d_{i-1}; k = 1, \dots, d_i; i = 1, \dots, l). \end{aligned}$$

The proof is completed. \square

Lemma 10. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Then the gradient of $f(\mathbf{w}, \mathbf{D})$ with respect to $\mathbf{w}_{(j)}$ can be written as

$$\|\delta_l\|_F^2 \leq \frac{\vartheta}{16p^2} b_{l+1}^2, \quad \|\delta_i\|_F^2 \leq \frac{d_{i+1} b_{i+1}^2 (k_{i+1} - s_{i+1} + 1)^2}{16p^2} \|\delta_{i+1}\|_F^2, \quad \|\delta_i\|_F^2 \leq \frac{\vartheta b_{l+1}^2}{16p^2} \prod_{s=i+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2},$$

where $\vartheta = 1/8$.

Proof. We first bound δ_l . By Lemma 9, we have

$$\begin{aligned} \|\delta_l\|_F^2 &= \left\| \sigma'_1(\mathbf{X}_{(l)}) \odot \text{up} \left(\frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{Z}_{(l)}} \right) \right\|_F^2 \stackrel{\textcircled{1}}{\leq} \frac{1}{16p^2} \left\| \frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{Z}_{(l)}} \right\|_F^2 = \frac{1}{16p^2} \left\| \frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{z}_{(l)}} \right\|_2^2 \\ &\leq \frac{1}{16p^2} \left\| (\mathbf{W}_{(l+1)})^T \mathbf{S}(\mathbf{v} - \mathbf{y}) \right\|_2^2 \stackrel{\textcircled{2}}{\leq} \frac{\vartheta}{16p^2} b_{l+1}^2, \end{aligned}$$

where ① holds since the values of the entries in the tensor $\sigma'_1(\mathbf{X}_{(l)}) \in \mathbb{R}^{r_l \times c_l \times d_l}$ belong to $[0, 1/4]$, and the $\text{up}(\cdot)$ operation does repeat one entry into p^2 entries but the entry value becomes $1/p^2$ of the original entry value. ② holds since we have $\|\mathbf{S}(\mathbf{v} - \mathbf{y})\|_2^2 \leq 1/8$ in Lemma 8.

Also by Lemma 9, we can further bound $\|\delta_i^j\|_F^2$ as follows:

$$\begin{aligned} \|\delta_i\|_F^2 &= \sum_{j=1}^{d_i} \|\delta_i^j\|_F^2 = \sum_{j=1}^{d_i} \left\| \text{up} \left(\sum_{k=1}^{d_{i+1}} \tilde{\delta}_{i+1}^k \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{k,j} \right) \odot \sigma'_1(\mathbf{X}_{(i)}^j) \right\|_F^2 \leq \frac{1}{16p^2} \sum_{j=1}^{d_i} \left\| \sum_{k=1}^{d_{i+1}} \tilde{\delta}_{i+1}^k \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{k,j} \right\|_F^2 \\ &\leq \frac{d_{i+1}}{16p^2} \sum_{j=1}^{d_i} \sum_{k=1}^{d_{i+1}} \left\| \tilde{\delta}_{i+1}^k \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{k,j} \right\|_F^2 \leq \frac{d_{i+1}(k_{i+1} - s_{i+1} + 1)^2}{16p^2} \sum_{j=1}^{d_i} \sum_{k=1}^{d_{i+1}} \left\| \tilde{\delta}_{i+1}^k \right\|_F^2 \left\| \widehat{\mathbf{W}}_{(i+1)}^{k,j} \right\|_F^2 \\ &\stackrel{\textcircled{1}}{=} \frac{d_{i+1}(k_{i+1} - s_{i+1} + 1)^2}{16p^2} \sum_{k=1}^{d_{i+1}} \left\| \tilde{\delta}_{i+1}^k \right\|_F^2 \left\| \mathbf{W}_{(i+1)}^k \right\|_F^2 \leq \frac{d_{i+1}b_{i+1}^2(k_{i+1} - s_{i+1} + 1)^2}{16p^2} \|\delta_{i+1}\|_F^2, \end{aligned}$$

where $\widehat{\mathbf{W}}_{(i+1)}^{k,j}$ denotes the j -th slice $\widehat{\mathbf{W}}_{(i+1)}^k(:, :, j)$ of the tensor $\widehat{\mathbf{W}}_{(i+1)}^k$. ① holds since we rotate the matrix $\mathbf{W}_{(i+1)}^{k,j}$ by 180 degrees to obtain $\widehat{\mathbf{W}}_{(i+1)}^k$, indicating $\|\mathbf{W}_{(i+1)}^{k,j}\|_F^2 = \|\widehat{\mathbf{W}}_{(i+1)}^k\|_F^2$ and $\sum_{j=1}^{d_i} \|\mathbf{W}_{(i+1)}^{k,j}\|_F^2 = \|\mathbf{W}_{(i+1)}^k\|_F^2 \leq b_{i+1}^2$. Accordingly, the above inequality gives

$$\begin{aligned} \|\delta_i\|_F^2 &\leq \frac{d_{i+1}b_{i+1}^2(k_{i+1} - s_{i+1} + 1)^2}{16p^2} \|\delta_{i+1}\|_F^2 \leq \dots \leq \|\delta_l\|_F^2 \prod_{s=i+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \\ &\leq \frac{\vartheta b_{l+1}^2}{16p^2} \prod_{s=i+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}. \end{aligned}$$

The proof is completed. \square

Lemma 11. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Then the gradient of $f(\mathbf{w}, \mathbf{D})$ with respect to $\mathbf{W}^{(l+1)}$ and \mathbf{w} can be respectively bounded as follows:

$$\|\nabla_{\mathbf{W}^{(l+1)}} f(\mathbf{w}, \mathbf{D})\|_F^2 \leq \vartheta \tilde{r}_l \tilde{c}_l d_l \quad \text{and} \quad \|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D})\|_2^2 \leq \beta^2,$$

where $\vartheta = 1/8$ and $\beta \triangleq \left[\vartheta \tilde{r}_l \tilde{c}_l d_l + \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right]^{1/2}$.

Proof. By utilizing Lemma 10, we can bound

$$\begin{aligned} \sum_{i=1}^l \|\nabla_{\mathbf{w}^{(i)}} f(\mathbf{w}, \mathbf{D})\|_F^2 &= \sum_{i=1}^l \sum_{k=1}^{d_i} \sum_{j=1}^{d_{i-1}} \left\| \nabla_{\mathbf{W}_{(i)}^{k,j}} f(\mathbf{w}, \mathbf{D}) \right\|_F^2 = \sum_{i=1}^l \sum_{k=1}^{d_i} \sum_{j=1}^{d_{i-1}} \left\| \mathbf{Z}_{(i-1)}^j \tilde{\otimes} \delta_i^k \right\|_F^2 \\ &\stackrel{\textcircled{1}}{\leq} \sum_{i=1}^l \sum_{k=1}^{d_i} \sum_{j=1}^{d_{i-1}} (k_i - s_i + 1)^2 \left\| \mathbf{Z}_{(i-1)}^j \right\|_F^2 \|\delta_i^k\|_F^2 \\ &\stackrel{\textcircled{2}}{\leq} \sum_{i=1}^l \sum_{k=1}^{d_i} \sum_{j=1}^{d_{i-1}} \tilde{r}_{i-1} \tilde{c}_{i-1} (k_i - s_i + 1)^2 \|\delta_i^k\|_F^2 \\ &\leq \sum_{i=1}^l \tilde{r}_{i-1} \tilde{c}_{i-1} d_{i-1} (k_i - s_i + 1)^2 \|\delta_i\|_F^2 \\ &\leq \sum_{i=1}^l \tilde{r}_{i-1} \tilde{c}_{i-1} d_{i-1} (k_i - s_i + 1)^2 \frac{\vartheta b_{l+1}^2}{16p^2} \prod_{s=i+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \\ &= \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}, \end{aligned}$$

where ① holds since $\left\| \mathbf{Z}_{(i-1)}^j \tilde{\otimes} \delta_i^k \right\|_F^2 \leq (k_i - s_i + 1)^2 \left\| \mathbf{Z}_{(i-1)}^j \right\|_F^2 \left\| \delta_i^k \right\|_F^2$; and ② holds since the values of entries in $\mathbf{Z}_{(i-1)}^j$ belong to $[0, 1]$.

On the other hand, we can bound

$$\left\| \nabla_{\mathbf{W}_{(l+1)}} f(\mathbf{w}, \mathbf{D}) \right\|_F^2 = \left\| \mathbf{S}(\mathbf{v} - \mathbf{y}) \mathbf{z}_{(l)}^T \right\|_F^2 \leq \vartheta \tilde{r}_l \tilde{c}_l d_l.$$

So we can bound the ℓ_2 norm of the gradient as follows:

$$\begin{aligned} \left\| \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}) \right\|_F^2 &= \left\| \nabla_{\mathbf{W}_{(l+1)}} f(\mathbf{w}, \mathbf{x}) \right\|_F^2 + \sum_{i=1}^l \left\| \nabla_{\mathbf{w}_{(i)}} f(\mathbf{w}, \mathbf{D}) \right\|_F^2 \\ &\leq \vartheta \tilde{r}_l \tilde{c}_l d_l + \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16 p^2}. \end{aligned}$$

The proof is completed. \square

Lemma 12. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Then for both cases, the gradient of $\mathbf{x}_{(i)}$ with respect to $\mathbf{w}_{(j)}$ can be bounded as follows:

$$\left\| \frac{\partial \text{vec}(\mathbf{X}_{(i)})}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \left\| \frac{\partial \mathbf{x}_{(i)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq d_i r_i c_i \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^i \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16 p^2}$$

and

$$\max_s \left\| \frac{\partial \text{vec}(\mathbf{X}_{(i)}^s)}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \max_s \left\| \frac{\partial \mathbf{x}_{(i)}^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq r_i c_i \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^i \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16 p^2}.$$

Proof. For brevity, let $\mathbf{X}_{(i)}^k(s, t)$ denotes the (s, t) -th entry in the matrix $\mathbf{X}_{(i)}^k \in \mathbb{R}^{r_i \times c_i}$. We also let $\phi_{(i,m)} = \frac{\partial \mathbf{X}_{(i)}^k(s, t)}{\partial \mathbf{X}_{(m)}} \in \mathbb{R}^{r_m \times c_m \times d_m}$. So similar to Lemma 9, we have

$$\phi_{(i,m)}^q = \text{up} \left(\sum_{k=1}^{d_{m+1}} \tilde{\phi}_{(i,m+1)}^k \tilde{\otimes} \widehat{\mathbf{W}}_{(m+1)}^{k,q} \right) \odot \sigma_1'(\mathbf{X}_{(m)}^q) \in \mathbb{R}^{r_m \times c_m}, \quad (q = 1, \dots, d_m),$$

where the matrix $\widehat{\mathbf{W}}_{(m+1)}^{k,q} \in \mathbb{R}^{k_{m+1} \times k_{m+1}}$ is obtained by rotating $\mathbf{W}_{(m+1)}^{k,q}$ with 180 degrees. Then according to the relationship between $\mathbf{X}_{(j)}$ and $\mathbf{W}_{(j)}$, we can compute

$$\frac{\partial \mathbf{X}_{(i)}^k(s, t)}{\partial \mathbf{W}_{(j)}^{g,h}} = \mathbf{Z}_{(j-1)}^h \tilde{\otimes} \phi_{(i,j)}^g \in \mathbb{R}^{k_j \times k_j}, \quad (h = 1, \dots, d_{j-1}; g = 1, \dots, d_j).$$

Therefore, we can further obtain

$$\begin{aligned} \left\| \frac{\partial \mathbf{X}_{(i)}^k(s, t)}{\partial \mathbf{w}_{(j)}} \right\|_F^2 &= \sum_{g=1}^{d_j} \sum_{h=1}^{d_{j-1}} \left\| \frac{\partial \mathbf{X}_{(i)}^k(s, t)}{\partial \mathbf{W}_{(j)}^{g,h}} \right\|_F^2 = \sum_{g=1}^{d_j} \sum_{h=1}^{d_{j-1}} \left\| \frac{\partial \mathbf{X}_{(i)}^k(s, t)}{\partial \mathbf{W}_{(j)}^{g,h}} \right\|_F^2 = \sum_{g=1}^{d_j} \sum_{h=1}^{d_{j-1}} \left\| \mathbf{Z}_{(j-1)}^h \tilde{\otimes} \phi_{(i,j)}^g \right\|_F^2 \\ &\leq \sum_{g=1}^{d_j} \sum_{h=1}^{d_{j-1}} (k_j - s_j + 1)^2 \left\| \mathbf{Z}_{(j-1)}^h \right\|_F^2 \left\| \phi_{(i,j)}^g \right\|_F^2 \leq (k_j - s_j + 1)^2 \left\| \mathbf{Z}_{(j-1)} \right\|_F^2 \left\| \phi_{(i,j)} \right\|_F^2 \\ &\leq \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \left\| \phi_{(i,j)} \right\|_F^2. \end{aligned}$$

On the other hand, by Lemma 9, we can further bound $\|\phi_{(i,j)}\|_F^2$ as follows:

$$\begin{aligned}
 \|\phi_{(i,m)}\|_F^2 &= \sum_{q=1}^{d_m} \|\phi_{(i,m)}^q\|_F^2 = \sum_{q=1}^{d_m} \left\| \text{up} \left(\sum_{k=1}^{d_{m+1}} \tilde{\phi}_{(i,m+1)}^k \otimes \widehat{\mathbf{W}}_{(m+1)}^{k,q} \right) \odot \sigma'_1(\mathbf{X}_{(m)}^q) \right\|_F^2 \\
 &\leq \frac{1}{16p^2} \sum_{q=1}^{d_m} \left\| \sum_{k=1}^{d_{m+1}} \tilde{\phi}_{(i,m+1)}^k \otimes \widehat{\mathbf{W}}_{(m+1)}^{k,q} \right\|_F^2 \leq \frac{d_{m+1}}{16p^2} \sum_{q=1}^{d_m} \sum_{k=1}^{d_{m+1}} \|\tilde{\phi}_{(i,m+1)}^k \otimes \widehat{\mathbf{W}}_{(m+1)}^{k,q}\|_F^2 \\
 &\leq \frac{d_{m+1}(k_{m+1} - s_{m+1} + 1)^2}{16p^2} \sum_{q=1}^{d_m} \sum_{k=1}^{d_{m+1}} \|\tilde{\phi}_{(i,m+1)}^k\|_F^2 \|\widehat{\mathbf{W}}_{(m+1)}^{k,q}\|_F^2 \\
 &\stackrel{\textcircled{1}}{=} \frac{d_{m+1}(k_{m+1} - s_{m+1} + 1)^2}{16p^2} \sum_{k=1}^{d_{m+1}} \|\tilde{\phi}_{(i,m+1)}^k\|_F^2 \|\widehat{\mathbf{W}}_{(m+1)}^k\|_F^2 \\
 &\leq \frac{d_{m+1}b_{m+1}^2(k_{m+1} - s_{m+1} + 1)^2}{16p^2} \|\phi_{(i,m+1)}\|_F^2,
 \end{aligned}$$

where $\textcircled{1}$ holds since $\|\mathbf{W}_{(m+1)}^{k,q}\|_F^2 = \|\widehat{\mathbf{W}}_{(m+1)}^{k,q}\|_F^2$. It further yields

$$\|\phi_{(i,m)}\|_F^2 \leq \|\phi_{(i,i)}\|_F^2 \prod_{s=m+1}^i \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \stackrel{\textcircled{1}}{=} \prod_{s=m+1}^i \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}.$$

where $\textcircled{1}$ holds since we have $\|\phi_{(i,i)}\|_F^2 = \left\| \frac{\partial \mathbf{X}_{(i)}^k(s,t)}{\partial \mathbf{X}_{(i)}} \right\|_F^2 = 1$.

Therefore, we have

$$\begin{aligned}
 \left\| \frac{\partial \mathbf{x}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 &= \sum_{s=1}^{r_i} \sum_{t=1}^{c_i} \left\| \frac{\partial \mathbf{X}_{(i)}^k(s,t)}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \sum_{s=1}^{r_i} \sum_{t=1}^{c_i} \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \|\phi_{(i,j)}\|_F^2 \\
 &= r_i c_i \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \|\phi_{(i,j)}\|_F^2 \\
 &\leq r_i c_i \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^i \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}.
 \end{aligned}$$

It further gives

$$\left\| \frac{\partial \mathbf{x}_{(i)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \sum_{s=1}^{d_i} \left\| \frac{\partial \mathbf{x}_{(i)}^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq d_i r_i c_i \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^i \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}.$$

The proof is completed. \square

Lemma 13. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Then the gradient of δ_l^s with respect to $\mathbf{w}_{(j)}$ can be bounded as follows:

$$\left\| \frac{\partial \delta_l^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \frac{\tilde{\vartheta} b_{l+1}^2}{16p^2} \|\mathbf{W}_{(l+1)}^s\|_F^2 d_l r_l c_l \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}$$

and

$$\left\| \frac{\partial \delta_l}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \frac{\tilde{\vartheta} b_{l+1}^4}{16p^2} d_l r_l c_l \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2},$$

where $\tilde{\vartheta} = \frac{3}{64}$.

Proof. Assume that $\mathbf{W}_{(l+1)} = [\mathbf{W}_{(l+1)}^1, \mathbf{W}_{(l+1)}^2, \dots, \mathbf{W}_{(l+1)}^{d_l}]$ where $\mathbf{W}_{(l+1)}^i \in \mathbb{R}^{d_{l+1} \times \tilde{r}_l \tilde{c}_l}$ is a submatrix in $\mathbf{W}_{(l+1)}$. Then we have $\mathbf{v} = \sigma_2(\sum_{k=1}^{d_l} \mathbf{W}_{(l+1)}^k \mathbf{z}_{(l)}^k)$. For brevity, we further define a matrix $\mathbf{G}_{(k)}$ as follows:

$$\mathbf{G}_{(k)} = \underbrace{\left[\sigma_1'(\mathbf{x}_{(k)}), \sigma_1'(\mathbf{x}_{(k)}), \dots, \sigma_1'(\mathbf{x}_{(k)}) \right]}_{r_k c_k \text{ columns}} \in \mathbb{R}^{r_k c_k d_k \times r_k c_k},$$

Then we have

$$\frac{\partial}{\partial \mathbf{z}_{(l)}} \left(\frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{z}_{(l)}^s} \right) = \left[(\mathbf{v} - \mathbf{y})^T \otimes (\mathbf{W}_{(l+1)}^s)^T \right] \mathbf{Q}(\mathbf{u}) \mathbf{W}_{(l+1)} + (\mathbf{W}_{(l+1)}^s)^T \mathbf{G}(\mathbf{u}) \mathbf{G}(\mathbf{u}) \mathbf{W}_{(l+1)},$$

where $\mathbf{Q}(\mathbf{u})$ is a matrix of size $d_{l+1}^2 \times d_{l+1}$ whose $(s, (s-1)d_{l+1} + s)$ entry equal to $\sigma_1(\mathbf{u}_s)(1 - \sigma_1(\mathbf{u}_s))(1 - 2\sigma_1(\mathbf{u}_s))$ and rest entries are all 0. Accordingly, we have

$$\begin{aligned} \left\| \frac{\partial}{\partial \mathbf{z}_{(l)}} \left(\frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{z}_{(l)}^s} \right) \right\|_F^2 &\leq 2 \left(\left\| \left[(\mathbf{v} - \mathbf{y})^T \otimes (\mathbf{W}_{(l+1)}^s)^T \right] \mathbf{Q}(\mathbf{u}) \mathbf{W}_{(l+1)} \right\|_F^2 + \left\| (\mathbf{W}_{(l+1)}^s)^T \mathbf{G}(\mathbf{u}) \mathbf{G}(\mathbf{u}) \mathbf{W}_{(l+1)} \right\|_F^2 \right) \\ &\leq 2 \left(\frac{2^6}{3^8} \|\mathbf{v} - \mathbf{y}\|_F^2 \|\mathbf{W}_{(l+1)}^s\|_F^2 \|\mathbf{W}_{(l+1)}\|_F^2 + \frac{1}{16^2} \|\mathbf{W}_{(l+1)}^s\|_F^2 \|\mathbf{W}_{(l+1)}\|_F^2 \right) \\ &\stackrel{\textcircled{1}}{\leq} 2 \left(\frac{2^7}{3^8} + \frac{1}{16^2} \right) b_{l+1}^2 \|\mathbf{W}_{(l+1)}^s\|_F^2 \\ &\leq \frac{3}{64} b_{l+1}^2 \|\mathbf{W}_{(l+1)}^s\|_F^2, \end{aligned}$$

where we have $\|\mathbf{v} - \mathbf{y}\|_F^2 \leq 2$ by Lemma 8. Then by similar way, we can have

$$\begin{aligned} \left\| \frac{\partial \delta_l^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 &= \left\| \frac{\partial}{\partial \mathbf{z}_{(l)}} \left(\frac{\partial f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{z}_{(l)}^s} \right) \frac{\partial \mathbf{z}_{(l)}}{\partial \mathbf{y}_l} \frac{\partial \mathbf{y}_l}{\partial \mathbf{x}_{(l)}} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \frac{3}{64 * 16p^2} b_{l+1}^2 \|\mathbf{W}_{(l+1)}^s\|_F^2 \left\| \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\ &\leq \frac{3}{64 * 16p^2} b_{l+1}^2 \|\mathbf{W}_{(l+1)}^s\|_F^2 d_l r_l c_l \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}. \end{aligned}$$

Therefore, we can further obtain:

$$\left\| \frac{\partial \delta_l}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \sum_{s=1}^{d_{l+1}} \left\| \frac{\partial \delta_l^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \frac{3}{64 * 16p^2} b_{l+1}^4 d_l r_l c_l \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}.$$

The proof is completed. \square

Lemma 14. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Then the Hessian of $f(\mathbf{w}, \mathbf{x})$ with respect to \mathbf{w} can be bounded as follows:

$$\|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D})\|_F^2 \leq \mathcal{O}(\gamma^2),$$

where $\gamma = \left(\frac{\partial b_{l+1}^2 d_0^2 l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2}p^2} \right]^2 \right)^{1/2}$. With the same condition, we can bound the operation norm of $\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{D})$. That is, there exists a universal constant ν such that $\|\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{D})\|_{op} \leq \nu$.

Proof. From Lemma 9, we can further compute the Hessian matrix $\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D})$. Recall that $\mathbf{w}_{(i)}^k \in \mathbb{R}^{k_i^2 d_{i-1}}$ ($k = 1, \dots, d_i$) is the vectorization of $\mathbf{W}_{(i)}^k \in \mathbb{R}^{k_i \times k_i \times d_{i-1}}$, i.e. $\mathbf{w}_{(i)}^k = \left[\text{vec} \left(\mathbf{W}_{(i)}^k(:, :, 1) \right); \dots; \text{vec} \left(\mathbf{W}_{(i)}^k(:, :, d_{i-1}) \right) \right]$.

Let $\mathbf{w}_{(i)} = [\mathbf{w}_{(i)}^1; \dots; \mathbf{w}_{(i)}^{d_i}] \in \mathbb{R}^{k_i^2 d_{i-1} \times d_i}$ ($i = 1, \dots, l$). Also, $\mathbf{w}_{(l+1)} \in \mathbb{R}^{\tilde{r}_l \tilde{c}_l d_{l+1}}$ is the vectorization of the weight matrix $\mathbf{W}_{(l+1)}$. Then if $1 \leq i, j \leq l$, we can have

$$\begin{aligned} \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}^k} &= \begin{bmatrix} \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}^{k,1}} \\ \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}^{k,2}} \\ \vdots \\ \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}^{k,d_{i-1}}} \end{bmatrix} = \begin{bmatrix} \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^1 \tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \\ \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^2 \tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \\ \vdots \\ \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^{d_{i-1}} \tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{P}_1(\tilde{\delta}_i^k) \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^1))}{\partial \mathbf{w}_{(j)}} + \mathbf{P}_2(\mathbf{Z}_{(i-1)}^1) \frac{\partial(\text{vec}(\tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \\ \mathbf{P}_1(\tilde{\delta}_i^k) \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^2))}{\partial \mathbf{w}_{(j)}} + \mathbf{P}_2(\mathbf{Z}_{(i-1)}^2) \frac{\partial(\text{vec}(\tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \\ \vdots \\ \mathbf{P}_1(\tilde{\delta}_i^k) \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^{d_{i-1}}))}{\partial \mathbf{w}_{(j)}} + \mathbf{P}_2(\mathbf{Z}_{(i-1)}^{d_{i-1}}) \frac{\partial(\text{vec}(\tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \end{bmatrix} \in \mathbb{R}^{k_i^2 d_{i-1} \times k_j^2 d_j d_{j-1}}, \end{aligned} \quad (4)$$

where $\mathbf{P}_1(\tilde{\delta}_i^k) \in \mathbb{R}^{k_i^2 \times \tilde{r}_{i-1} \tilde{c}_{i-1} d_{i-1}}$ and $\mathbf{P}_2(\mathbf{Z}_{(i-1)}^{d_{i-1}}) \in \mathbb{R}^{k_i^2 \times (\tilde{r}_i - k_i + 1)(\tilde{c}_i - k_i + 1)}$ satisfy: each row in $\mathbf{P}_1(\tilde{\delta}_i^k)$ contains the vectorization of $(\tilde{\delta}_i^k)^T$ at the right position and the remaining entries are 0s, and each row in $\mathbf{P}_2(\mathbf{Z}_{(i-1)}^{d_{i-1}})$ is the submatrix in $\mathbf{Z}_{(i-1)}^{d_{i-1}}$ that need to conduct inner product with $\tilde{\delta}_i^k$ in turn. Note that there are $s_i - 1$ rows and columns between each neighboring nonzero entries in \mathbf{N} which is decided by the definition of $\tilde{\delta}_{i+1}^k$ in Sec. B.1. Accordingly, we have

$$\left\| \mathbf{P}_1(\tilde{\delta}_i^k) \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^{d_{i-1}}))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq (k_i - s_i + 1)^2 \|\tilde{\delta}_i^k\|_F^2 \left\| \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^{d_{i-1}}))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = (k_i - s_i + 1)^2 \|\tilde{\delta}_i^k\|_F^2 \left\| \frac{\partial(\text{vec}(\mathbf{Z}_{(i-1)}^{d_{i-1}}))}{\partial \mathbf{w}_{(j)}} \right\|_F^2$$

and

$$\left\| \mathbf{P}_2(\mathbf{Z}_{(i-1)}^{d_{i-1}}) \frac{\partial(\text{vec}(\tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq (k_i - s_i + 1)^2 \|\mathbf{Z}_{(i-1)}^{d_{i-1}}\|_F^2 \left\| \frac{\partial(\text{vec}(\tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = (k_i - s_i + 1)^2 \|\mathbf{Z}_{(i-1)}^{d_{i-1}}\|_F^2 \left\| \frac{\partial(\text{vec}(\tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2.$$

Then in order to bound

$$\|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D})\|_F^2 = \sum_{i=1}^{l+1} \sum_{j=1}^{l+1} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}} \right\|_F^2,$$

we try to bound each term separately. So we consider the following five cases: $l \geq i \geq j$, $i \leq j \leq l$, $l+1 = i > j$, $l+1 = j > i$ and $l+1 = i = j$.

Case 1: $l \geq i \geq j$

In the following, we first consider the first case, *i.e.* $i \geq j$, and bound

$$\begin{aligned}
 & \left\| \frac{\partial (\text{vec}(\tilde{\delta}_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \left\| \frac{\partial (\text{vec}(\delta_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \left\| \frac{\partial}{\partial \mathbf{w}_{(j)}} \text{vec} \left[\text{up} \left(\sum_{s=1}^{d_{i+1}} \tilde{\delta}_{i+1}^s \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{s,k} \right) \odot \sigma'_1(\mathbf{X}_{(i)}^k) \right] \right\|_F^2 \\
 & \stackrel{\textcircled{1}}{\leq} \frac{2}{16} \left\| \frac{\partial}{\partial \mathbf{w}_{(j)}} \text{vec} \left[\text{up} \left(\sum_{s=1}^{d_{i+1}} \tilde{\delta}_{i+1}^s \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{s,k} \right) \right] \right\|_F^2 + 2 \left\| \text{up} \left(\sum_{s=1}^{d_{i+1}} \tilde{\delta}_{i+1}^s \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{s,k} \right) \right\|_F^2 \left\| \frac{\partial \sigma'_1(\mathbf{X}_{(i)}^k)}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & = \frac{2}{16p^2} \left\| \frac{\partial}{\partial \mathbf{w}_{(j)}} \text{vec} \left[\sum_{s=1}^{d_{i+1}} \tilde{\delta}_{i+1}^s \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{s,k} \right] \right\|_F^2 + \frac{2}{p^2} \left\| \sum_{s=1}^{d_{i+1}} \tilde{\delta}_{i+1}^s \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{s,k} \right\|_F^2 \left\| \frac{\partial \text{vec}(\sigma'(\mathbf{x}_{(i)}^k))}{\partial \mathbf{x}_{(i)}^k} \frac{\partial \mathbf{x}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \stackrel{\textcircled{2}}{\leq} \frac{2}{16p^2} \left\| \frac{\partial}{\partial \mathbf{w}_{(j)}} \text{vec} \left[\sum_{s=1}^{d_{i+1}} \tilde{\delta}_{i+1}^s \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{s,k} \right] \right\|_F^2 + \frac{2 \cdot 2^6}{3^8 p^2} \left\| \sum_{s=1}^{d_{i+1}} \tilde{\delta}_{i+1}^s \tilde{\otimes} \widehat{\mathbf{W}}_{(i+1)}^{s,k} \right\|_F^2 \left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \stackrel{\textcircled{3}}{\leq} \frac{2d_{i+1}}{16p^2} (k_{i+1} - s_{i+1} + 1)^2 \sum_{s=1}^{d_{i+1}} \|\widehat{\mathbf{W}}_{(i+1)}^{s,k}\|_F^2 \left\| \frac{\partial \tilde{\delta}_{i+1}^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{2 \cdot 2^6}{3^8 p^2} d_{i+1} (k_{i+1} - s_{i+1} + 1)^2 \sum_{s=1}^{d_{i+1}} \|\widehat{\mathbf{W}}_{(i+1)}^{s,k}\|_F^2 \|\tilde{\delta}_{i+1}^s\|_F^2 \left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \stackrel{\textcircled{4}}{=} \frac{2d_{i+1}}{16p^2} (k_{i+1} - s_{i+1} + 1)^2 \sum_{s=1}^{d_{i+1}} \|\mathbf{W}_{(i+1)}^{s,k}\|_F^2 \left\| \frac{\partial \delta_{i+1}^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{2 \cdot 2^6}{3^8 p^2} d_{i+1} (k_{i+1} - s_{i+1} + 1)^2 \sum_{s=1}^{d_{i+1}} \|\mathbf{W}_{(i+1)}^{s,k}\|_F^2 \|\delta_{i+1}^s\|_F^2 \left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2. \tag{5}
 \end{aligned}$$

① holds since $\mathbf{X}_{(i)}^k$ is independent on $\mathbf{w}_{(j)}$ and the values of entries in $\sigma'_1(\mathbf{X}_{(i)}^k)$ is not larger than $1/4$ since for any constant a , $\sigma'(a) = \sigma(a)(1 - \sigma(a)) \leq 1/4$. ② holds since for arbitrary tensor M , we have $\|\text{up}(M)\|_F^2 \leq \|M\|_F^2/p^2$ in Lemma 8, and we also have

$$\left\| \frac{\partial \text{vec}(\sigma'(\mathbf{x}_{(i)}^k))}{\partial \mathbf{x}_{(i)}^k} \frac{\partial \mathbf{x}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \left\| \mathbf{Q}(\mathbf{x}_{(i)}^k) \frac{\partial \mathbf{x}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \frac{2^6}{3^8} \left\| \frac{\partial \mathbf{x}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \frac{2^6}{3^8} \left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2.$$

③ holds since we can just adopt similar strategy in Eqn. (4) to separate $\widehat{\mathbf{W}}_{(i+1)}^{s,k}$ and the conclusion in Lemma 8; ④ holds since the difference between $\tilde{\delta}_{i+1}^s$ and δ_{i+1}^s is that we pad 0 around δ_{i+1}^s to obtain $\tilde{\delta}_{i+1}^s$, indicating $\|\delta_{i+1}^s\|_F^2 = \|\tilde{\delta}_{i+1}^s\|_F^2$.

Accordingly, we can further bound

$$\begin{aligned}
 & \left\| \frac{\partial \tilde{\delta}_i}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \sum_{k=1}^{d_i-1} \left\| \frac{\partial (\text{vec}(\delta_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \leq \frac{2d_{i+1}}{16p^2} (k_{i+1} - s_{i+1} + 1)^2 \sum_{k=1}^{d_{i-1}} \sum_{s=1}^{d_{i+1}} \|\mathbf{W}_{(i+1)}^{s,k}\|_F^2 \left\| \frac{\partial \delta_{i+1}^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{2 \cdot 2^6}{3^8 p^2} d_{i+1} (k_{i+1} - s_{i+1} + 1)^2 \sum_{k=1}^{d_{i-1}} \sum_{s=1}^{d_{i+1}} \|\mathbf{W}_{(i+1)}^{s,k}\|_F^2 \|\delta_{i+1}^s\|_F^2 \left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \leq \frac{2d_{i+1}}{16p^2} (k_{i+1} - s_{i+1} + 1)^2 \sum_{s=1}^{d_{i+1}} \|\mathbf{W}_{(i+1)}^s\|_F^2 \left\| \frac{\partial \delta_{i+1}^s}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{2 \cdot 2^6}{3^8 p^2} d_{i+1} (k_{i+1} - s_{i+1} + 1)^2 \|\delta_{i+1}\|_F^2 \max_s \|\mathbf{W}_{(i+1)}^s\|_F^2 \max_k \left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \stackrel{\textcircled{1}}{\leq} \frac{2d_{i+1}}{16p^2} (k_{i+1} - s_{i+1} + 1)^2 b_{i+1}^2 \left\| \frac{\partial \delta_{i+1}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{2 \cdot 2^6}{3^8 p^2} d_{i+1} (k_{i+1} - s_{i+1} + 1)^2 b_{i+1}^2 \|\delta_{i+1}\|_F^2 \max_k \left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \stackrel{\textcircled{2}}{\leq} \frac{2d_{i+1}}{16p^2} (k_{i+1} - s_{i+1} + 1)^2 b_{i+1}^2 \left\| \frac{\partial \delta_{i+1}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{\vartheta b_{i+1}^2 d_{j-1}}{3p^2 b_j^2 d_j} r_i c_i r_{j-1} c_{j-1} \prod_{s=j}^l d_s b_s^2 (k_s - s_s + 1)^2 \frac{1}{16p^2},
 \end{aligned}$$

where ① holds since we have $\|\mathbf{W}_{(i+1)}^s\|_F \leq r_w$; ② holds due to the bounds of $\|\delta_{i+1}\|_F^2$ and $\left\| \frac{\partial \mathbf{X}_{(i)}^k}{\partial \mathbf{w}_{(j)}} \right\|_F^2$ in Lemma 10 and 12.

Then, we can use the above recursion inequality to further obtain

$$\begin{aligned}
 & \left\| \frac{\partial \delta_i}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \leq \left[\prod_{s=i+1}^{i+1} \frac{2d_s}{16p^2} (k_s - s_s + 1)^2 b_s^2 \right] \left(\frac{2d_{i+2}}{16p^2} (k_{i+2} - s_{i+2} + 1)^2 b_{i+2}^2 \left\| \frac{\partial \delta_{i+2}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{\vartheta b_{l+1}^2 d_{j-1}}{3p^2 b_j^2 d_j} r_{j-1} c_{j-1} r_{i+1} c_{i+1} \prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right) \\
 & \quad + \frac{\vartheta b_{l+1}^2 d_{j-1}}{3p^2 b_j^2 d_j} r_i c_i r_{j-1} c_{j-1} \prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \\
 & \leq \left[\prod_{s=i+1}^l \frac{2d_s}{16p^2} (k_s - s_s + 1)^2 b_s^2 \right] \left\| \frac{\partial \delta_l}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \frac{\vartheta b_{l+1}^2 d_{j-1}}{3p^2 b_j^2 d_j} r_{j-1} c_{j-1} \prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \left[r_i c_i + r_{i+1} c_{i+1} \left[\prod_{s=i+1}^{i+1} \frac{2d_s}{16p^2} (k_s - s_s + 1)^2 b_s^2 \right] \right. \\
 & \quad \left. + r_{i+2} c_{i+2} \left[\prod_{s=i+1}^{i+2} \frac{2d_s}{16p^2} (k_s - s_s + 1)^2 b_s^2 \right] + \dots + r_l c_l \left[\prod_{s=i+1}^l \frac{2d_s}{16p^2} (k_s - s_s + 1)^2 b_s^2 \right] \right].
 \end{aligned}$$

By Lemma 13, we have

$$\left\| \frac{\partial \delta_l}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \frac{\tilde{\vartheta} b_{l+1}^4 d_{j-1}}{p^2 b_j^2 d_j} d_l r_l c_l r_{j-1} c_{j-1} \prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2},$$

where $\tilde{\vartheta} = \frac{3}{64}$. Thus, we can establish

$$\left\| \frac{\partial \delta_i}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \leq \frac{\tilde{\vartheta} b_{l+1}^2 d_{j-1}}{p^2 b_j^2 d_j} r_{j-1} c_{j-1} \left[\frac{\tau}{3} + b_{l+1}^2 d_l r_l c_l \prod_{s=i+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right] \prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}.$$

where $\tau = r_i c_i + r_{i+1} c_{i+1} \left[\prod_{s=i+1}^{i+1} \frac{2d_s}{16p^2} (k_s - s_s + 1)^2 b_s^2 \right] + \dots + r_l c_l \left[\prod_{s=i+1}^l \frac{2d_s}{16p^2} (k_s - s_s + 1)^2 b_s^2 \right]$. It further gives the bound of $\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}} \right\|_F^2$ as follows:

$$\begin{aligned}
 & \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}} \right\|_F^2 = \sum_{k=1}^{d_{i-1}} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}^k} \right\|_F^2 = \sum_{k=1}^{d_{i-1}} \sum_{s=1}^{d_i} \left\| \mathbf{P}_1(\delta_i^k) \frac{\partial (\text{vec}(\mathbf{X}_{(i-1)}^s))}{\partial \mathbf{w}_{(j)}} + \mathbf{P}_2(\mathbf{X}_{(i-1)}^s) \frac{\partial (\text{vec}(\delta_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \leq 2 \sum_{k=1}^{d_{i-1}} \sum_{s=1}^{d_i} \left(\left\| \mathbf{P}_1(\delta_i^k) \frac{\partial (\text{vec}(\mathbf{X}_{(i-1)}^s))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \left\| \mathbf{P}_2(\mathbf{X}_{(i-1)}^s) \frac{\partial (\text{vec}(\delta_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \right) \\
 & \leq 2(k_i - s_i + 1)^2 \sum_{k=1}^{d_{i-1}} \sum_{s=1}^{d_i} \left(\left\| \delta_i^k \right\|_F^2 \left\| \frac{\partial (\text{vec}(\mathbf{X}_{(i-1)}^s))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \left\| \mathbf{X}_{(i-1)}^s \right\|_F^2 \left\| \frac{\partial (\text{vec}(\delta_i^k))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \right) \\
 & \leq 2(k_i - s_i + 1)^2 \left(\left\| \delta_i \right\|_F^2 \left\| \frac{\partial (\text{vec}(\mathbf{X}_{(i-1)}))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 + \left\| \mathbf{X}_{(i-1)} \right\|_F^2 \left\| \frac{\partial (\text{vec}(\delta_i))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \right) \\
 & \stackrel{\textcircled{1}}{\leq} \frac{32\vartheta b_{l+1}^2 d_{i-1}}{b_i^2 b_j^2 d_i d_j} r_{i-1} c_{i-1} r_{j-1} c_{j-1} \prod_{s=i+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} + r_{i-1} c_{i-1} d_{i-1} \left\| \frac{\partial (\text{vec}(\delta_i))}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \stackrel{\textcircled{2}}{\leq} \mathcal{O} \left(\frac{\vartheta b_{l+1}^2 d_{i-1} d_{j-1}}{b_i^2 b_j^2 d_i d_j} r_{i-1} c_{i-1} r_{j-1} c_{j-1} \left[\prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right] \left[\prod_{s=i}^l \frac{2d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right] \right).
 \end{aligned}$$

where ① holds because of Lemma 12, while ② holds due to 13.

Case 2: $i \leq j \leq l$

Since $\frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w} \partial \mathbf{w}^T}$ is symmetrical, we have $\frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} = \left(\frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)}^T \partial \mathbf{w}_{(i)}} \right)^T$ ($1 \leq i, j \leq l$). Thus, it yields

$$\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 = \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)}^T \partial \mathbf{w}_{(i)}} \right\|_F^2.$$

Case 3: $l+1 = i > j$

In the following, we first consider the first case, *i.e.* cross entropy and softmax activation, and bound

$$\begin{aligned} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(l+1)}} \right\|_F^2 &= \left\| \frac{\partial(\mathbf{v} - \mathbf{y}) \mathbf{z}_{(l)}^T}{\partial \mathbf{w}_{(j)}} \right\|_F^2 = \left\| [\mathbf{I}_{\tilde{r}_l \tilde{c}_l d_l} \otimes (\mathbf{v} - \mathbf{y})] \frac{\partial \mathbf{z}_{(l)}^T}{\partial \mathbf{x}_{(l)}} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} + [\mathbf{z}_{(l)} \otimes \mathbf{I}_{d_{l+1}}] \frac{\partial \mathbf{v}}{\partial \mathbf{z}_{(l)}} \frac{\partial \mathbf{z}_{(l)}}{\partial \mathbf{x}_{(l)}} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\ &= \left\| \widetilde{\text{up}}([\mathbf{I}_{\tilde{r}_l \tilde{c}_l d_l} \otimes (\mathbf{v} - \mathbf{y})] + [\mathbf{z}_{(l)} \otimes \mathbf{I}_{d_{l+1}}] \text{diag}(\sigma'_2(\mathbf{u})) \mathbf{W}_{(l+1)}) \odot \mathbf{G}_{(l)} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2, \end{aligned}$$

where $\mathbf{G}_{(l)}$ is defined as

$$\mathbf{G}_{(l)} = \underbrace{\left[\sigma'_1(\mathbf{x}_{(l)}), \sigma'_1(\mathbf{x}_{(l)}), \dots, \sigma'_1(\mathbf{x}_{(l)}) \right]}_{r_l c_l \text{ columns}} \in \mathbb{R}^{r_l c_l d_l \times r_l c_l}.$$

Thus, we can further obtain

$$\begin{aligned} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(l+1)}} \right\|_F^2 &\leq \frac{1}{16p^2} \left\| [\mathbf{I}_{\tilde{r}_l \tilde{c}_l d_l} \otimes (\mathbf{v} - \mathbf{y})] + [\mathbf{z}_{(l)} \otimes \mathbf{I}_{d_{l+1}}] \text{diag}(\sigma'_2(\mathbf{u})) \mathbf{W}_{(l+1)} \right\|_F^2 \left\| \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\ &\leq \frac{2}{16p^2} \left(\left\| \mathbf{I}_{\tilde{r}_l \tilde{c}_l d_l} \otimes (\mathbf{v} - \mathbf{y}) \right\|_F^2 + \left\| [\mathbf{z}_{(l)} \otimes \mathbf{I}_{d_{l+1}}] \text{diag}(\sigma'_2(\mathbf{u})) \mathbf{W}_{(l+1)} \right\|_F^2 \right) \left\| \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\ &\stackrel{\textcircled{1}}{\leq} \frac{2}{16p^2} \left(\left\| \mathbf{I}_{\tilde{r}_l \tilde{c}_l d_l} \otimes (\mathbf{v} - \mathbf{y}) \right\|_F^2 + \left\| \mathbf{z}_{(l)} \otimes [\text{diag}(\sigma'_2(\mathbf{u})) \mathbf{W}_{(l+1)}] \right\|_F^2 \right) \left\| \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\ &\stackrel{\textcircled{2}}{\leq} \frac{1}{8p^2} \tilde{r}_l \tilde{c}_l d_l \left(2 + \frac{1}{16} b_{l+1}^2 \right) d_l r_l c_l \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \\ &= \frac{2d_{j-1}}{p^4 b_j^2 d_j} r_l^2 c_l^2 d_l^2 r_{j-1} c_{j-1} \left(2 + \frac{1}{16} b_{l+1}^2 \right) \prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \end{aligned}$$

where ① holds since for an arbitrary vector $\mathbf{u} \in \mathbb{R}^k$ and an arbitrary matrix $\mathbf{M} \in \mathbb{R}^{k \times k}$, we have $(\mathbf{u} \otimes \mathbf{I}_k) \mathbf{M} = \mathbf{u} \otimes \mathbf{M}$; ② holds since we use Lemma 12 and the assumption that $\left\| \mathbf{W}_{(l+1)} \right\|_F^2 \leq b_{l+1}^2$.

Now we consider the least square loss and softmax activation function. In such a case, we can further obtain:

$$\begin{aligned} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(l+1)}} \right\|_F^2 &= \left\| \frac{\partial(\mathbf{v} - \mathbf{y}) \mathbf{G}(\mathbf{u}) \mathbf{z}_{(l)}^T}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\ &= \left\| [\mathbf{I}_{\tilde{r}_l \tilde{c}_l d_l} \otimes (\mathbf{v} - \mathbf{y})] \frac{\partial \mathbf{z}_{(l)}^T}{\partial \mathbf{x}_{(l)}} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} + [\mathbf{z}_{(l)} \otimes (\mathbf{v} - \mathbf{y})] \frac{\partial \text{vec}(\mathbf{G}(\mathbf{u}))}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{z}_{(l)}} \frac{\partial \mathbf{z}_{(l)}}{\partial \mathbf{x}_{(l)}} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} + [\mathbf{z}_{(l)} \otimes \mathbf{I}_{d_{l+1}}] \frac{\partial \mathbf{v}}{\partial \mathbf{z}_{(l)}} \frac{\partial \mathbf{z}_{(l)}}{\partial \mathbf{x}_{(l)}} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\ &= \left\| \widetilde{\text{up}}([\mathbf{I}_{\tilde{r}_l \tilde{c}_l d_l} \otimes (\mathbf{v} - \mathbf{y})] + [\mathbf{z}_{(l)} \otimes (\mathbf{v} - \mathbf{y})] \mathbf{Q}(\mathbf{u}) \mathbf{W}_{(l+1)} + [\mathbf{z}_{(l)} \otimes \mathbf{I}_{d_{l+1}}] \mathbf{Q}(\mathbf{u}) \mathbf{G}(\mathbf{u}) \mathbf{W}_{(l+1)}) \odot \mathbf{G}_{(l)} \frac{\partial \mathbf{x}_{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2. \end{aligned}$$

Thus, we can further obtain

$$\begin{aligned}
 & \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(l+1)}} \right\|_F^2 \\
 & \leq \frac{1}{16p^2} \left\| [\mathbf{I}_{\tilde{r}_i \tilde{c}_i d_i} \otimes (\mathbf{v} - \mathbf{y})] + [\mathbf{z}^{(l)} \otimes (\mathbf{v} - \mathbf{y})] \mathbf{Q}(\mathbf{u}) \mathbf{W}_{(l+1)} + [\mathbf{z}^{(l)} \otimes \mathbf{I}_{d_{l+1}}] \mathbf{Q}(\mathbf{u}) \mathbf{G}(\mathbf{u}) \mathbf{W}_{(l+1)} \right\|_F^2 \left\| \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \leq \frac{3}{16p^2} \left(\left\| \mathbf{I}_{\tilde{r}_i \tilde{c}_i d_i} \otimes (\mathbf{v} - \mathbf{y}) \right\|_F^2 + \left\| [\mathbf{z}^{(l)} \otimes (\mathbf{v} - \mathbf{y})] \mathbf{Q}(\mathbf{u}) \mathbf{W}_{(l+1)} \right\|_F^2 + \left\| [\mathbf{z}^{(l)} \otimes \mathbf{I}_{d_{l+1}}] \mathbf{Q}(\mathbf{u}) \mathbf{G}(\mathbf{u}) \mathbf{W}_{(l+1)} \right\|_F^2 \right) \left\| \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{w}_{(j)}} \right\|_F^2 \\
 & \stackrel{\textcircled{1}}{\leq} \frac{3}{16p^2} \tilde{r}_i \tilde{c}_i d_i \left(2 + \frac{3}{100} b_{l+1}^2 \right) d_l r_l c_l \tilde{r}_{j-1} \tilde{c}_{j-1} d_{j-1} (k_j - s_j + 1)^2 \prod_{s=j+1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \\
 & = \frac{3d_{j-1}}{p^4 b_j^2 d_j} r_{j-1} c_{j-1} d_l^2 r_l^2 c_l^2 \left(2 + \frac{2^7}{3^8} b_{l+1}^2 + \frac{2^6}{16 \cdot 3^8} b_{l+1}^2 \right) \prod_{s=j}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2},
 \end{aligned}$$

where $\textcircled{1}$ holds since we use Lemma 12 and the fact that $\|\mathbf{W}_{(l+1)}\|_F^2 \leq b_{l+1}^2$.

Case 4: $i < j = l + 1$

Similar to the Case 2, we also can have

$$\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 = \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)}^T \partial \mathbf{w}_{(i)}} \right\|_F^2.$$

So in this case, we can just directly use the bound in case 3 to bound $\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2$.

Case 5: $i = j = l + 1$

In the following, we first consider the first case, *i.e.* $i = l + 1$, and bound

$$\begin{aligned}
 \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(l+1)} \partial \mathbf{w}_{(l+1)}} \right\|_F^2 &= \left\| \frac{\partial (\mathbf{v} - \mathbf{y}) \mathbf{z}^{(l)T}}{\partial \mathbf{w}_{(l+1)}} \right\|_F^2 = \left\| [\mathbf{z}^{(l)} \otimes \mathbf{I}_{d_{l+1}}] \frac{\partial \mathbf{v}}{\partial \mathbf{w}_{(l+1)}} \right\|_F^2 \\
 &= \left\| [\mathbf{z}^{(l)} \otimes \mathbf{I}_{d_{l+1}}] \mathbf{G}(\mathbf{u}) [\mathbf{z}^{(l)} \otimes \mathbf{I}_{d_{l+1}}]^T \right\|_F^2 \\
 &\stackrel{\textcircled{1}}{=} \left\| [\mathbf{z}^{(l)} (\mathbf{z}^{(l)} \otimes \mathbf{G}(\mathbf{u}))^T]^T \right\|_F^2 \\
 &\leq \|\mathbf{z}^{(l)}\|_F^4 \|\mathbf{G}(\mathbf{u})\|_F^2 \leq \frac{1}{16} \tilde{r}_l^2 \tilde{c}_l^2 d_l^2 d_{l+1},
 \end{aligned}$$

where $\textcircled{1}$ holds since for an arbitrary vector $\mathbf{u} \in \mathbb{R}^k$ and an arbitrary matrix $\mathbf{M} \in \mathbb{R}^{k \times k}$, we have $(\mathbf{u} \otimes \mathbf{I}_k) \mathbf{M} = \mathbf{u} \otimes \mathbf{M}$.

Now we can bound $\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w} \partial \mathbf{w}} \right\|_F^2$ as follows:

$$\begin{aligned}
 \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w} \partial \mathbf{w}} \right\|_F^2 &= \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(l+1)} \partial \mathbf{w}_{(l+1)}} \right\|_F^2 + 2 \sum_{j=1}^l \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(l+1)}} \right\|_F^2 + 2 \sum_{j=1}^l \sum_{i=j}^l \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{D})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}} \right\|_F^2 \\
 &\leq \mathcal{O} \left(\frac{l^2}{k_1^4} \max_{1 \leq i, j \leq l} \tilde{r}_i \tilde{c}_i r_j c_j \frac{b_{l+1}^4}{16p^2} \left(\frac{r_w^2}{8\sqrt{2}p^2} \right)^{2l-2} \prod_{s=1}^l (d_s k_s^2)^2 \right) \\
 &\leq \mathcal{O} \left(\frac{\partial b_{l+1}^2 d_0^2}{b_1^4 d_1^2} l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2}p^2} \right]^2 \right).
 \end{aligned}$$

On the other hand, if the activation functions σ_1 and σ_2 are respectively sigmoid function and softmax function, $f(\mathbf{w}, \mathbf{D})$ is infinitely differentiable. Also $\sigma(a)$, $\sigma'(a)$, $\sigma''(a)$ and $\sigma'''(a)$ are all bounded. This means that $\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{D})$ exists. Also since input \mathbf{D} and the parameter \mathbf{w} are bounded, we can always find a universal constant ν such that

$$\|\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{D})\|_{\text{op}} = \sup_{\|\boldsymbol{\lambda}\|_2 \leq 1} \left\langle \boldsymbol{\lambda}^{\otimes 3}, \nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{D}) \right\rangle = \sum_{i,j,k} [\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{D})]_{ijk} \lambda_i \lambda_j \lambda_k \leq \nu < +\infty.$$

The proof is completed. \square

Lemma 15. *Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Suppose Assumption 1 on the input data \mathbf{D} holds. Then for any $t > 0$, the objective $f(\mathbf{w}, \mathbf{x})$ obeys*

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{D}^{(i)}))\right) > t\right) \leq 2 \exp\left(-\frac{2nt^2}{\alpha^2}\right),$$

where $\alpha = 1$.

Proof. Since the input $\mathbf{D}^{(i)}$ ($i = 1, \dots, n$) are independent from each other, then the output $f(\mathbf{w}, \mathbf{D}^{(i)})$ ($i = 1, \dots, n$) are also independent. Meanwhile, when the loss is the square loss, we can easily bound $0 \leq f(\mathbf{w}, \mathbf{D}^{(i)}) = \frac{1}{2} \|\mathbf{v} - \mathbf{y}\|_2^2 \leq 1$, since the value of entries in \mathbf{v} belongs to $[0, 1]$ and \mathbf{y} is a one-hot vector label of \mathbf{v} .

Besides, for arbitrary random variable x , $|x - \mathbb{E}x| \leq |x|$. So by Hoeffding's inequality in Lemma 6, we have

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{D}^{(i)}))\right) > t\right) \leq \exp\left(-\frac{2nt^2}{\alpha^2}\right),$$

where $\alpha = 1$. This means that $\frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{D}^{(i)})))$ has exponential tails. \square

Lemma 16. *Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Suppose Assumption 1 on the input data \mathbf{D} holds. Then for any $t > 0$ and arbitrary unit vector $\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$, the gradient $\nabla f(\mathbf{w}, \mathbf{x})$ obeys*

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left(\left\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) \right\rangle\right) > t\right) \leq \exp\left(-\frac{nt^2}{2\beta^2}\right).$$

where $\beta \triangleq \left[\vartheta \tilde{r}_l \tilde{c}_l d_l + \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right]^{1/2}$ in which $\vartheta = 1/8$.

Proof. Since the input $\mathbf{D}^{(i)}$ ($i = 1, \dots, n$) are independent from each other, then the output $\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)})$ ($i = 1, \dots, n$) are also independent. Furthermore, for arbitrary vector \mathbf{x} , $\|\mathbf{x} - \mathbb{E}\mathbf{x}\|_2^2 \leq \|\mathbf{x}\|_2^2$. Hence, for an arbitrary unit vector $\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$ where $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$, we have

$$\begin{aligned} \left\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) \right\rangle &\leq \|\boldsymbol{\lambda}\|_2 \|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)})\|_2 \\ &\leq \|\boldsymbol{\lambda}\|_2 \|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)})\|_2 \stackrel{\textcircled{1}}{\leq} \beta, \end{aligned}$$

where $\textcircled{1}$ holds since $\|\boldsymbol{\lambda}\|_2 = 1$ ($\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$) and by Lemma 11, we have $\|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)})\| \leq \beta$ where $\beta \triangleq \left[\vartheta \tilde{r}_l \tilde{c}_l d_l + \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right]^{1/2}$ in which $\vartheta = 1/8$.

Thus, we can use Hoeffding's inequality in Lemma 6 to bound

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left(\left\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) \right\rangle\right) > t\right) \leq \exp\left(-\frac{nt^2}{2\beta^2}\right).$$

The proof is completed. \square

Lemma 17. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Suppose that Assumption 1 on the input data \mathbf{D} and the parameter \mathbf{w} holds. Then for any $t > 0$ and arbitrary unit vector $\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$, the Hessian $\nabla^2 f(\mathbf{w}, \mathbf{D})$ obeys

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \left(\langle \boldsymbol{\lambda}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})) \boldsymbol{\lambda} \rangle \right) > t \right) \leq 2 \exp \left(-\frac{nt^2}{2\gamma^2} \right).$$

$$\text{where } \gamma = \left(\frac{\vartheta b_{l+1}^2 d_0^2}{b_1^4 d_1^2} l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2} p^2} \right]^2 \right)^{1/2}.$$

Proof. Since the input $\mathbf{D}^{(i)}$ ($i = 1, \dots, n$) are independent from each other, then the output $\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)})$ ($i = 1, \dots, n$) are also independent. On the other hand, for arbitrary random matrix \mathbf{X} , $\|\mathbf{X} - \mathbb{E}\mathbf{X}\|_F^2 \leq \|\mathbf{X}\|_F^2$. Thus, for an arbitrary unit vector $\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$ where $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$, we have

$$\begin{aligned} \langle \boldsymbol{\lambda}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})) \boldsymbol{\lambda} \rangle &\leq \|\boldsymbol{\lambda}\|_2 \|(\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})) \boldsymbol{\lambda}\|_2 \\ &\leq \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})\|_{\text{op}} \|\boldsymbol{\lambda}\|_2 \\ &\leq \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})\|_F \\ &\leq \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})\|_F \\ &\stackrel{\textcircled{1}}{\leq} \gamma, \end{aligned}$$

where $\textcircled{1}$ holds since $\|\boldsymbol{\lambda}\|_2 = 1$ ($\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$) and by Lemma 14, we have $\|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})\| \leq \gamma$ where $\gamma = \left(\frac{\vartheta b_{l+1}^2 d_0^2}{b_1^4 d_1^2} l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2} p^2} \right]^2 \right)^{1/2}$.

Thus, we can use Hoeffding's inequality in Lemma 6 to bound

$$\mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n \left(\langle \boldsymbol{\lambda}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{D}^{(i)})) \boldsymbol{\lambda} \rangle \right) > t \right) \leq \exp \left(-\frac{nt^2}{2\gamma^2} \right).$$

The proof is completed. \square

Lemma 18. Suppose that the activation function σ_1 is sigmoid and σ_2 is softmax, and the loss function $f(\mathbf{w}, \mathbf{D})$ is squared loss. Suppose that Assumption 1 on the input data \mathbf{D} and the parameter \mathbf{w} holds. Then the empirical Hessian converges uniformly to the population Hessian in operator norm. Specifically, there exist two universal constants $c_{v'}$ and c_v such that if $n \geq c_{v'} \frac{\nu^2}{d \varrho \varepsilon^2} \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2} p^2} \right]^{-1}$, then with probability at least $1 - \varepsilon$

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right\|_{\text{op}} \leq c_v \gamma \sqrt{\frac{2d + \theta \varrho + \log \left(\frac{4}{\varepsilon} \right)}{2n}},$$

holds with probability at least $1 - \varepsilon$, where $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$, $\theta = a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)$, $\varrho = \sum_{i=1}^l \log \left(\frac{\sqrt{d_i} b_i (k_i - s_i + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128p^2} \right)$, and $\gamma = \left(\frac{\vartheta b_{l+1}^2 d_0^2}{b_1^4 d_1^2} l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2} p^2} \right]^2 \right)^{1/2}$.

Proof. Recall that the weight of each kernel and the feature maps has magnitude bound separately, i.e. $\mathbf{w}_{(i)}^k \in \mathcal{B}^{k_i^2 d_{i-1}}(r_w)$ ($i = 1, \dots, l; k = 1, \dots, d_i$) and $\mathbf{w}_{(l+1)} \in \mathcal{B}^{\tilde{r}_l \tilde{c}_l d_l d_{l+1}}(b_{l+1})$. Since $\tilde{\mathbf{W}}_{(i)} = [\text{vec}(\mathbf{W}_{(i)}^1), \text{vec}(\mathbf{W}_{(i)}^2), \dots, \text{vec}(\mathbf{W}_{(i)}^{d_i-1})] \in \mathbb{R}^{k_i^2 d_i \times d_{i-1}}$, we have $\|\tilde{\mathbf{W}}_{(i)}\|_F \leq d_i b_i$.

So here we assume $\widetilde{\mathbf{W}}_{(i,\epsilon)}$ is the $d_i b_i \epsilon / (b_{l+1} + \sum_{i=1}^l d_i b_i)$ -covering net of the matrix $\widetilde{\mathbf{W}}_{(i)}$ which is the set of all parameters in the i -th layer. Then by Lemma 7, we have the covering number

$$n_\epsilon^i \leq \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)},$$

since the rank of $\widetilde{\mathbf{W}}_{(i)}$ obeys $\text{rank}(\widetilde{\mathbf{W}}_{(i)}) \leq a_i$ for $1 \leq i \leq l$. For the last layer, we also can construct an $b_{l+1} \epsilon / (b_{l+1} + \sum_{i=1}^l d_i b_i)$ -covering net for the weight matrix $\mathbf{W}_{(l+1)}$. Here we have

$$n_\epsilon^{l+1} \leq \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1)},$$

since the rank of $\mathbf{W}_{(l+1)}$ obeys $\text{rank}(\mathbf{W}_{(l+1)}) \leq a_{l+1}$. Finally, we arrange them together to construct a set Θ and claim that there is always an ϵ -covering net \mathbf{w}_ϵ in Θ for any parameter \mathbf{w} . Accordingly, we have

$$|\Theta| \leq \prod_{i=1}^{l+1} n_\epsilon^i = \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)} = \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta,$$

where $\theta = a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)$ which is the total freedom degree of the network. So we can always find a vector $\mathbf{w}_{k_w} \in \Theta$ such that $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$. Now we use the decomposition strategy to bound our goal:

$$\begin{aligned} & \left\| \nabla^2 \widetilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right\|_{\text{op}} \\ &= \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}, \mathbf{D})) \right\|_{\text{op}} \\ &= \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) \right) + \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D})) \right. \\ & \quad \left. + \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}, \mathbf{D})) \right\|_{\text{op}} \\ &\leq \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) \right) \right\|_{\text{op}} + \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D})) \right\|_{\text{op}} \\ & \quad + \left\| \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}, \mathbf{D})) \right\|_{\text{op}}. \end{aligned}$$

Here we also define four events \mathbf{E}_0 , \mathbf{E}_1 , \mathbf{E}_2 and \mathbf{E}_3 as

$$\begin{aligned} \mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \widetilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right\|_{\text{op}} \geq t \right\}, \\ \mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) \right) \right\|_{\text{op}} \geq \frac{t}{3} \right\}, \\ \mathbf{E}_2 &= \left\{ \sup_{\mathbf{w}_{k_w} \in \Theta} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}, \\ \mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{D})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla^2 f(\mathbf{w}, \mathbf{D})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}. \end{aligned}$$

Accordingly, we have

$$\mathbb{P}(\mathbf{E}_0) \leq \mathbb{P}(\mathbf{E}_1) + \mathbb{P}(\mathbf{E}_2) + \mathbb{P}(\mathbf{E}_3).$$

So we can respectively bound $\mathbb{P}(\mathbf{E}_1)$, $\mathbb{P}(\mathbf{E}_2)$ and $\mathbb{P}(\mathbf{E}_3)$ to bound $\mathbb{P}(\mathbf{E}_0)$.

Step 1. Bound $\mathbb{P}(\mathbf{E}_1)$: We first bound $\mathbb{P}(\mathbf{E}_1)$ as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_1) &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right\|_2 \geq \frac{t}{3}\right) \\ &\stackrel{\textcircled{1}}{\leq} \frac{3}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right\|_2 \right) \\ &\leq \frac{3}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 f(\mathbf{w}, \mathbf{D}) - \nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}) \right\|_2 \right) \\ &\leq \frac{3}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \frac{\left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right\|_2}{\|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2 \right) \\ &\stackrel{\textcircled{2}}{\leq} \frac{3\nu\epsilon}{t}, \end{aligned}$$

where $\textcircled{1}$ holds since by Markov inequality and $\textcircled{2}$ holds because of Lemma 14. Therefore, we can set

$$t \geq \frac{6\nu\epsilon}{\epsilon}.$$

Then we can bound $\mathbb{P}(\mathbf{E}_1)$:

$$\mathbb{P}(\mathbf{E}_1) \leq \frac{\epsilon}{2}.$$

Step 2. Bound $\mathbb{P}(\mathbf{E}_2)$: By Lemma 3, we know that for any matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$, its operator norm can be computed as

$$\|\mathbf{X}\|_{\text{op}} \leq \frac{1}{1-2\epsilon} \sup_{\boldsymbol{\lambda} \in \boldsymbol{\lambda}_\epsilon} |\langle \boldsymbol{\lambda}, \mathbf{X} \boldsymbol{\lambda} \rangle|.$$

where $\boldsymbol{\lambda}_\epsilon = \{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k_{\mathbf{w}}}\}$ be an ϵ -covering net of $\mathbb{B}^d(1)$.

Let $\boldsymbol{\lambda}_{1/4}$ be the $\frac{1}{4}$ -covering net of $\mathbb{B}^d(1)$, where $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$. Recall that we use Θ to denote the ϵ -net of $\mathbf{w}_{k_{\mathbf{w}}}$ and we have $|\Theta| \leq \prod_{i=1}^{l+1} n_\epsilon^i = \left(\frac{3(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta$. Then we can bound $\mathbb{P}(\mathbf{E}_2)$ as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P}\left(\sup_{\mathbf{w}_{k_{\mathbf{w}}} \in \Theta} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}) \right) \right\|_2 \geq \frac{t}{3}\right) \\ &\leq \mathbb{P}\left(\sup_{\mathbf{w}_{k_{\mathbf{w}}} \in \Theta, \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/4}} 2 \left| \left\langle \boldsymbol{\lambda}, \left(\frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}) \right) \right) \boldsymbol{\lambda} \right\rangle \right| \geq \frac{t}{3}\right) \\ &\leq 12^d \left(\frac{3(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta \sup_{\mathbf{w}_{k_{\mathbf{w}}} \in \Theta, \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/4}} \mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n \left\langle \boldsymbol{\lambda}, \left(\nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}) \right) \right) \boldsymbol{\lambda} \right\rangle \right| \geq \frac{t}{6}\right) \\ &\stackrel{\textcircled{1}}{\leq} 12^d \left(\frac{3(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta 2 \exp\left(-\frac{nt^2}{72\gamma^2}\right), \end{aligned}$$

where $\textcircled{1}$ holds since by Lemma 17, we have

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left(\left\langle \boldsymbol{\lambda}, \left(\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) \right) \boldsymbol{\lambda} \right\rangle \right) > t \right) \leq \exp\left(-\frac{nt^2}{2\gamma^2}\right).$$

where $\gamma = \left(\frac{\vartheta b_{l+1}^2 d_0^2}{b_1^4 d_1^2} l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2}p^2} \right]^2 \right)^{1/2}$.

Thus, if we set

$$t \geq \sqrt{\frac{72\gamma^2 \left(d \log(12) + \theta \log \left(\frac{3(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right) + \log \left(\frac{4}{\epsilon} \right) \right)}{n}},$$

then we have

$$\mathbb{P}(\mathbf{E}_2) \leq \frac{\epsilon}{2}.$$

Step 3. Bound $\mathbb{P}(\mathbf{E}_3)$: We first bound $\mathbb{P}(\mathbf{E}_3)$ as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_3) &= \mathbb{P} \left(\sup_{\mathbf{w} \in \Omega} \left\| \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} (\nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} (\nabla^2 f(\mathbf{w}, \mathbf{D})) \right\|_2 \geq \frac{t}{3} \right) \\ &\leq \mathbb{P} \left(\mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \sup_{\mathbf{w} \in \Omega} \left\| (\nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D})) - \nabla^2 f(\mathbf{w}, \mathbf{D}) \right\|_2 \geq \frac{t}{3} \right) \\ &\leq \mathbb{P} \left(\sup_{\mathbf{w} \in \Omega} \frac{\left| \frac{1}{n} \sum_{i=1}^n (\nabla^2 f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla^2 f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)})) \right|}{\|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2 \geq \frac{t}{3} \right) \\ &\stackrel{\textcircled{1}}{\leq} \mathbb{P} \left(\nu \epsilon \geq \frac{t}{3} \right), \end{aligned}$$

where $\textcircled{1}$ holds because of Lemma 14. We set ϵ enough small such that $\nu \epsilon < t/3$ always holds. Then it yields $\mathbb{P}(\mathbf{E}_3) = 0$.

Step 4. Final result: To ensure $\mathbb{P}(\mathbf{E}_0) \leq \epsilon$, we just set $\epsilon = \frac{36(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\vartheta^2 n b_{l+1}} \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2}p^2} \right]^{-\frac{1}{2}}$. Note that $\frac{6\epsilon\nu}{\epsilon} > 3\epsilon\nu$. Thus we can obtain

$$t \geq \max \left(\frac{6\nu\epsilon}{\epsilon}, \sqrt{\frac{72\gamma^2 \left(d \log(12) + \theta \log \left(\frac{3(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right) + \log \left(\frac{4}{\epsilon} \right) \right)}{n}} \right).$$

Thus, if $n \geq c_{v'} \frac{\nu^2}{d\epsilon^2} \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2}p^2} \right]^{-1}$ where $c_{v'}$ is a constant, there exists a universal constant c_v such that

$$\begin{aligned} \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right\|_{\text{op}} &\leq \hat{c}_v \gamma \sqrt{\frac{d \log(12) + \theta \left(\sum_{i=1}^l \log \left(\frac{\sqrt{d_s} b_s (k_s - s_s + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128p^2} \right) \right) + \log \left(\frac{4}{\epsilon} \right)}{n}} \\ &= c_v \gamma \sqrt{\frac{2d + \theta \varrho + \log \left(\frac{4}{\epsilon} \right)}{2n}} \end{aligned}$$

holds with probability at least $1 - \epsilon$, where $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$, $\theta = a_{l+1} (d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i (k_i^2 d_i + d_{i-1} - 2a_i + 1)$, $\varrho = \sum_{i=1}^l \log \left(\frac{\sqrt{d_i} b_i (k_i - s_i + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128p^2} \right)$, and

$\gamma = \left(\frac{\vartheta b_{l+1}^2 d_0^2}{b_1^4 d_1^2} l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2}p^2} \right]^2 \right)^{1/2}$. The proof is completed. \square

D Proofs of Main Theorems

D.1 Proof of Lemma 1

Proof. Recall that the weight of each kernel and the feature maps has magnitude bound separately, i.e. $\mathbf{w}_{(i)}^k \in \mathbb{B}^{k_i^2 d_{i-1}}(r_w)$ ($i = 1, \dots, l; k = 1, \dots, d_i$) and $\mathbf{w}_{(l+1)} \in \mathbb{B}^{\tilde{r}_l \tilde{c}_l d_l d_{l+1}}(b_{l+1})$. Since $\tilde{\mathbf{W}}_{(i)} = [\text{vec}(\mathbf{W}_{(i)}^1), \text{vec}(\mathbf{W}_{(i)}^2), \dots, \text{vec}(\mathbf{W}_{(i)}^{d_i-1})] \in \mathbb{R}^{k_i^2 d_i \times d_{i-1}}$, we have $\|\tilde{\mathbf{W}}_{(i)}\|_F \leq d_i b_i$.

So here we assume $\widetilde{\mathbf{W}}_{(i,\epsilon)}$ is the $d_i b_i \epsilon / (b_{l+1} + \sum_{i=1}^l d_i b_i)$ -covering net of the matrix $\widetilde{\mathbf{W}}_{(i)}$ which is the set of all parameters in the i -th layer. Then by Lemma 7, we have the covering number

$$n_\epsilon^i \leq \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)},$$

since the rank of $\widetilde{\mathbf{W}}_{(i)}$ obeys $\text{rank}(\widetilde{\mathbf{W}}_{(i)}) \leq a_i$ for $1 \leq i \leq l$. For the last layer, we also can construct an $b_{l+1} \epsilon / (b_{l+1} + \sum_{i=1}^l d_i b_i)$ -covering net for the weight matrix $\mathbf{W}_{(l+1)}$. Here we have

$$n_\epsilon^{l+1} \leq \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1)},$$

since the rank of $\mathbf{W}_{(l+1)}$ obeys $\text{rank}(\mathbf{W}_{(l+1)}) \leq a_{l+1}$. Finally, we arrange them together to construct a set Θ and claim that there is always an ϵ -covering net \mathbf{w}_ϵ in Θ for any parameter \mathbf{w} . Accordingly, we have

$$|\Theta| \leq \prod_{i=1}^{l+1} n_\epsilon^i = \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)} = \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta,$$

where $\theta = a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)$ which is the total freedom degree of the network. So we can always find a vector $\mathbf{w}_{k_w} \in \Theta$ such that $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$. Now we use the decomposition strategy to bound our goal:

$$\begin{aligned} \left| \tilde{Q}_n(\mathbf{w}) - Q(\mathbf{w}) \right| &= \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}, \mathbf{D})) \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{D}^{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)})) + \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E} f(\mathbf{w}_{k_w}, \mathbf{D}) + \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} f(\mathbf{w}_{k_w}, \mathbf{D}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} f(\mathbf{w}, \mathbf{D}) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{D}^{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)})) \right| + \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} f(\mathbf{w}_{k_w}, \mathbf{D}) \right| \\ &\quad + \left| \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} f(\mathbf{w}_{k_w}, \mathbf{D}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} f(\mathbf{w}, \mathbf{D}) \right|. \end{aligned}$$

Then, we define four events \mathbf{E}_0 , \mathbf{E}_1 , \mathbf{E}_2 and \mathbf{E}_3 as

$$\begin{aligned} \mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \tilde{Q}_n(\mathbf{w}) - Q(\mathbf{w}) \right| \geq t \right\}, \\ \mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{D}^{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}^{(i)})) \right| \geq \frac{t}{3} \right\}, \\ \mathbf{E}_2 &= \left\{ \sup_{\mathbf{w}_{k_w} \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}_{k_w}, \mathbf{D})) \right| \geq \frac{t}{3} \right\}, \\ \mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}_{k_w}, \mathbf{D})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}, \mathbf{D})) \right| \geq \frac{t}{3} \right\}. \end{aligned}$$

Accordingly, we have

$$\mathbb{P}(\mathbf{E}_0) \leq \mathbb{P}(\mathbf{E}_1) + \mathbb{P}(\mathbf{E}_2) + \mathbb{P}(\mathbf{E}_3).$$

So we can respectively bound $\mathbb{P}(\mathbf{E}_1)$, $\mathbb{P}(\mathbf{E}_2)$ and $\mathbb{P}(\mathbf{E}_3)$ to bound $\mathbb{P}(\mathbf{E}_0)$.

Step 1. Bound $\mathbb{P}(\mathbf{E}_1)$: We first bound $\mathbb{P}(\mathbf{E}_1)$ as follows:

$$\begin{aligned}
 \mathbb{P}(\mathbf{E}_1) &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{w}, \mathbf{D}^{(i)}) - f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right| \geq \frac{t}{3}\right) \\
 &\stackrel{\textcircled{1}}{\leq} \frac{3}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{w}, \mathbf{D}^{(i)}) - f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right| \right) \\
 &\leq \frac{3}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \frac{\left| \frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{w}, \mathbf{D}^{(i)}) - f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right|}{\|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2 \right) \\
 &\leq \frac{3\epsilon}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}, \mathbf{D}) \right\|_2 \right),
 \end{aligned}$$

where $\textcircled{1}$ holds since by Markov inequality, we have that for an arbitrary nonnegative random variable x , then

$$\mathbb{P}(x \geq t) \leq \frac{\mathbb{E}(x)}{t}.$$

Now we only need to bound $\mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}, \mathbf{D}) \right\|_2 \right)$. Therefore, by Lemma 11, we have

$$\mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}, \mathbf{D}) \right\|_2 \right) = \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}, \mathbf{D}^{(i)}) \right\|_2 \right) \leq \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \|\nabla f(\mathbf{w}, \mathbf{D})\|_2 \right) \leq \beta.$$

where $\beta \triangleq \left[\vartheta \tilde{r}_l \tilde{c}_l d_l + \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s + 1)^2}{16 p^2} \right]^{1/2}$ in which $\vartheta = 1/8$. Therefore, we have

$$\mathbb{P}(\mathbf{E}_1) \leq \frac{3\epsilon\beta}{t}.$$

We further let

$$t \geq \frac{6\epsilon\beta}{\epsilon}.$$

Then we can bound $\mathbb{P}(\mathbf{E}_1)$:

$$\mathbb{P}(\mathbf{E}_1) \leq \frac{\epsilon}{2}.$$

Step 2. Bound $\mathbb{P}(\mathbf{E}_2)$: Recall that we use Θ to denote the index of $\mathbf{w}_{k_{\mathbf{w}}}$ and we have $|\Theta| \leq \prod_{i=1}^{l+1} n_{\epsilon}^i = \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta$. We can bound $\mathbb{P}(\mathbf{E}_2)$ as follows:

$$\begin{aligned}
 \mathbb{P}(\mathbf{E}_2) &= \mathbb{P}\left(\sup_{\mathbf{w}_{k_{\mathbf{w}}} \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D})) \right| \geq \frac{t}{3}\right) \\
 &\leq \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta \sup_{\mathbf{w}_{k_{\mathbf{w}}} \in \Theta} \mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D})) \right| \geq \frac{t}{3}\right) \\
 &\stackrel{\textcircled{1}}{\leq} \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^\theta 2 \exp\left(-\frac{2nt^2}{\alpha^2}\right),
 \end{aligned}$$

where $\textcircled{1}$ holds because in Lemma 15, we have

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left(f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{D}^{(i)})) \right) > t\right) \leq \exp\left(-\frac{2nt^2}{\alpha^2}\right),$$

where $\alpha = 1$. Thus, if we set

$$t \geq \sqrt{\frac{\alpha^2 \left(\theta \log \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right) + \log \left(\frac{4}{\epsilon} \right) \right)}{2n}},$$

then we have

$$\mathbb{P}(\mathbf{E}_2) \leq \frac{\epsilon}{2}.$$

Step 3. Bound $\mathbb{P}(\mathbf{E}_3)$: We first bound $\mathbb{P}(\mathbf{E}_3)$ as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_3) &= \mathbb{P} \left(\sup_{\mathbf{w} \in \Omega} \|\mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}, \mathbf{D}))\|_2 \geq \frac{t}{3} \right) \\ &= \mathbb{P} \left(\sup_{\mathbf{w} \in \Omega} \frac{\|\mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}) - f(\mathbf{w}, \mathbf{D}))\|_2}{\|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2 \geq \frac{t}{3} \right) \\ &\leq \mathbb{P} \left(\epsilon \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \sup_{\mathbf{w} \in \Omega} \|\nabla \mathbf{Q}_{\mathbf{w}}(\mathbf{w}, \mathbf{D})\|_2 \geq \frac{t}{3} \right) \\ &\stackrel{\textcircled{1}}{\leq} \mathbb{P} \left(\beta \epsilon \geq \frac{t}{3} \right), \end{aligned}$$

where $\textcircled{1}$ holds since we utilize Lemma 11. We set ϵ enough small such that $\beta \epsilon < t/3$ always holds. Then it yields $\mathbb{P}(\mathbf{E}_3) = 0$.

Step 4. Final result: To ensure $\mathbb{P}(\mathbf{E}_0) \leq \epsilon$, we just set $\epsilon = \frac{18p^2(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\vartheta^2 n b_{l+1}} \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2} \right]^{-\frac{1}{2}}$. Note that $\frac{6\epsilon\beta}{\epsilon} > 3\epsilon\beta$ due to $\epsilon \leq 1$. Thus we can obtain

$$t \geq \max \left(\frac{6\epsilon\beta}{\epsilon}, \sqrt{\frac{\alpha^2 \left(\theta \log \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right) + \log \left(\frac{4}{\epsilon} \right) \right)}{2n}} \right).$$

By comparing the values of α , we can observe that if $n \geq c_f \frac{l^2(b_{l+1} + \sum_{i=1}^l d_i b_i)^2 \max_i \sqrt{r_i c_i}}{\theta \varrho \epsilon^2}$ where c_f is a constant, there exists such a universal constant c_f such that

$$\sup_{\mathbf{w} \in \Omega} \left| \tilde{\mathbf{Q}}_n(\mathbf{w}) - \mathbf{Q}(\mathbf{w}) \right| \leq \alpha \sqrt{\frac{\theta \left(\sum_{i=1}^l \log \left(\frac{\sqrt{d_i} b_i (k_i - s_i + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128p^2} \right) \right) + \log \left(\frac{4}{\epsilon} \right)}{2n}} = \sqrt{\frac{\theta \varrho + \log \left(\frac{4}{\epsilon} \right)}{2n}}$$

holds with probability at least $1 - \epsilon$, where $\theta = a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i(k_i^2 d_i + d_{i-1} - 2a_i + 1)$, $\varrho = \sum_{i=1}^l \log \left(\frac{\sqrt{d_i} b_i (k_i - s_i + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128p^2} \right)$, and $\alpha = 1$. The proof is completed. \square

D.2 Proof of Theorem 1

Proof. By Lemma 1 in the manuscript, we know that if $n \geq c_f l^2(b_{l+1} + \sum_{i=1}^l d_i b_i)^2 \max_i \sqrt{r_i c_i} / (\theta \varrho \epsilon^2)$ where c_f is a universal constant, then with probability at least $1 - \epsilon$, we have

$$\sup_{\mathbf{w} \in \Omega} \left| \tilde{\mathbf{Q}}_n(\mathbf{w}) - \mathbf{Q}(\mathbf{w}) \right| \leq \sqrt{\frac{\theta \varrho + \log \left(\frac{4}{\epsilon} \right)}{2n}},$$

where the total freedom degree θ of the network is $\theta = a_{l+1}(d_{l+1} + \tilde{r}_l \tilde{c}_l d_l + 1) + \sum_{i=1}^l a_i(k_i^2 d_{i-1} + d_i + 1)$ and $\varrho = \sum_{i=1}^l \log \left(\frac{\sqrt{d_i} b_i (k_i - s_i + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128p^2} \right)$.

Thus based on such a result, we can derive the following generalization bound:

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \left| \mathbb{E}_{\mathcal{A}} (\mathbf{Q}(\tilde{\mathbf{w}}) - \tilde{\mathbf{Q}}_n(\tilde{\mathbf{w}})) \right| \leq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left| \tilde{\mathbf{Q}}_n(\mathbf{w}) - \mathbf{Q}(\mathbf{w}) \right| \right) \leq \sup_{\mathbf{w} \in \Omega} \left| \tilde{\mathbf{Q}}_n(\mathbf{w}) - \mathbf{Q}(\mathbf{w}) \right| \leq \sqrt{\frac{\theta_{\rho} + \log\left(\frac{4}{\epsilon}\right)}{2n}}.$$

Thus, the conclusion holds. The proof is completed. \square

D.3 Proof of Theorem 2

Proof. Recall that the weight of each kernel and the feature maps has magnitude bound separately, i.e. $\mathbf{w}_{(i)}^k \in \mathcal{B}^{k_i^2 d_{i-1}}(r_w)$ ($i = 1, \dots, l; k = 1, \dots, d_i$) and $\mathbf{w}_{(l+1)} \in \mathcal{B}^{\tilde{r}_l \tilde{c}_l d_l}(b_{l+1})$. Since $\tilde{\mathbf{W}}_{(i)} = [\text{vec}(\mathbf{W}_{(i)}^1), \text{vec}(\mathbf{W}_{(i)}^2), \dots, \text{vec}(\mathbf{W}_{(i)}^{d_i-1})] \in \mathbb{R}^{k_i^2 d_{i-1} \times d_i}$, we have $\|\tilde{\mathbf{W}}_{(i)}\|_F \leq d_i b_i$.

So here we assume $\tilde{\mathbf{W}}_{(i, \epsilon)}$ is the $d_i b_i \epsilon / (b_{l+1} + \sum_{i=1}^l d_i b_i)$ -covering net of the matrix $\tilde{\mathbf{W}}_{(i)}$ which is the set of all parameters in the i -th layer. Then by Lemma 7, we have the covering number

$$n_{\epsilon}^i \leq \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_i (k_i^2 d_{i-1} + d_i - 2a_i + 1)},$$

since the rank of $\tilde{\mathbf{W}}_{(i)}$ obeys $\text{rank}(\tilde{\mathbf{W}}_{(i)}) \leq a_i$ for $1 \leq i \leq l$. For the last layer, we also can construct an $b_{l+1} \epsilon / (b_{l+1} + \sum_{i=1}^l d_i b_i)$ -covering net for the weight matrix $\mathbf{W}_{(l+1)}$. Here we have

$$n_{\epsilon}^{l+1} \leq \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_{l+1} (d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1)},$$

since the rank of $\mathbf{W}_{(l+1)}$ obeys $\text{rank}(\mathbf{W}_{(l+1)}) \leq a_{l+1}$. Finally, we arrange them together to construct a set Θ and claim that there is always an ϵ -covering net \mathbf{w}_{ϵ} in Θ for any parameter \mathbf{w} . Accordingly, we have

$$|\Theta| \leq \prod_{i=1}^{l+1} n_{\epsilon}^i = \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{a_{l+1} (d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i (k_i^2 d_{i-1} + d_i - 2a_i + 1)} = \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon} \right)^{\theta},$$

where $\theta = a_{l+1} (d_{l+1} + \tilde{r}_l \tilde{c}_l d_l - 2a_{l+1} + 1) + \sum_{i=1}^l a_i (k_i^2 d_{i-1} + d_i - 2a_i + 1)$ which is the total freedom degree of the network. So we can always find a vector $\mathbf{w}_{k_w} \in \Theta$ such that $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$. Accordingly, we can decompose $\|\nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla \mathbf{Q}(\mathbf{w})\|_2$ as

$$\begin{aligned} & \left\| \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 \\ &= \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathcal{D} \sim \mathcal{D}} (\nabla f(\mathbf{w}, \mathbf{D})) \right\|_2 \\ &= \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) \right) + \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathcal{D} \sim \mathcal{D}} (\nabla f(\mathbf{w}_{k_w}, \mathbf{D})) \right. \\ & \quad \left. + \mathbb{E}_{\mathcal{D} \sim \mathcal{D}} (\nabla f(\mathbf{w}_{k_w}, \mathbf{D})) - \mathbb{E}_{\mathcal{D} \sim \mathcal{D}} (\nabla f(\mathbf{w}, \mathbf{D})) \right\|_2 \\ &\leq \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) \right) \right\|_2 + \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathcal{D} \sim \mathcal{D}} (\nabla f(\mathbf{w}_{k_w}, \mathbf{D})) \right\|_2 \\ & \quad + \left\| \mathbb{E}_{\mathcal{D} \sim \mathcal{D}} (\nabla f(\mathbf{w}_{k_w}, \mathbf{D})) - \mathbb{E}_{\mathcal{D} \sim \mathcal{D}} (\nabla f(\mathbf{w}, \mathbf{D})) \right\|_2. \end{aligned}$$

Here we also define four events \mathbf{E}_0 , \mathbf{E}_1 , \mathbf{E}_2 and \mathbf{E}_3 as

$$\begin{aligned}\mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 \geq t \right\}, \\ \mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right\|_2 \geq \frac{t}{3} \right\}, \\ \mathbf{E}_2 &= \left\{ \sup_{\mathbf{w}_{k_{\mathbf{w}}} \in \Theta} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D})) \right\|_2 \geq \frac{t}{3} \right\}, \\ \mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}, \mathbf{D})) \right\|_2 \geq \frac{t}{3} \right\}.\end{aligned}$$

Accordingly, we have

$$\mathbb{P}(\mathbf{E}_0) \leq \mathbb{P}(\mathbf{E}_1) + \mathbb{P}(\mathbf{E}_2) + \mathbb{P}(\mathbf{E}_3).$$

So we can respectively bound $\mathbb{P}(\mathbf{E}_1)$, $\mathbb{P}(\mathbf{E}_2)$ and $\mathbb{P}(\mathbf{E}_3)$ to bound $\mathbb{P}(\mathbf{E}_0)$.

Step 1. Bound $\mathbb{P}(\mathbf{E}_1)$: We first bound $\mathbb{P}(\mathbf{E}_1)$ as follows:

$$\begin{aligned}\mathbb{P}(\mathbf{E}_1) &= \mathbb{P} \left(\sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right\|_2 \geq \frac{t}{3} \right) \\ &\stackrel{\textcircled{1}}{\leq} \frac{3}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right\|_2 \right) \\ &\leq \frac{3}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \frac{\left\| \frac{1}{n} \sum_{i=1}^n \left(\nabla f(\mathbf{w}, \mathbf{D}^{(i)}) - \nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{D}^{(i)}) \right) \right\|_2}{\|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2 \right) \\ &\leq \frac{3\epsilon}{t} \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}, \mathbf{D}) \right\|_2 \right),\end{aligned}$$

where $\textcircled{1}$ holds since by Markov inequality, we have that for an arbitrary nonnegative random variable x , then $\mathbb{P}(x \geq t) \leq \frac{\mathbb{E}(x)}{t}$.

Now we only need to bound $\mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}, \mathbf{D}) \right\|_2 \right)$. Here we utilize Lemma 14 to achieve this goal:

$$\mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}, \mathbf{D}) \right\|_2 \right) \leq \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \left(\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 f(\mathbf{w}, \mathbf{D}) - \nabla^2 f(\mathbf{w}^*, \mathbf{D}) \right\|_2 \right) \leq \gamma.$$

where $\gamma = \left(\frac{\partial b_{l+1}^2 d_0^2}{b_1^4 d_1^2} l^2 r_0^2 c_0^2 \left[\prod_{s=1}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{8\sqrt{2} p^2} \right]^2 \right)^{1/2}$. Therefore, we have

$$\mathbb{P}(\mathbf{E}_1) \leq \frac{3\gamma\epsilon}{t}.$$

We further let

$$t \geq \frac{6\gamma\epsilon}{\epsilon}.$$

Then we can bound $\mathbb{P}(\mathbf{E}_1)$:

$$\mathbb{P}(\mathbf{E}_1) \leq \frac{\epsilon}{2}.$$

Step 2. Bound $\mathbb{P}(\mathbf{E}_2)$: By Lemma 2, we know that for any vector $\mathbf{x} \in \mathbb{R}^d$, its ℓ_2 -norm can be computed as

$$\|\mathbf{x}\|_2 \leq \frac{1}{1-\epsilon} \sup_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}_\epsilon} \langle \boldsymbol{\lambda}, \mathbf{x} \rangle.$$

where $\lambda_\epsilon = \{\lambda_1, \dots, \lambda_{k_w}\}$ be an ϵ -covering net of $\mathbb{B}^d(1)$.

Let λ be the $\frac{1}{2}$ -covering net of $\mathbb{B}^d(1)$, where $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$. Recall that we use Θ to denote the index of \mathbf{w}_{k_w} so that $\|\mathbf{w} - \mathbf{w}_{k_w}\| \leq \epsilon$. Besides, $|\Theta| \leq \prod_{i=1}^{l+1} n_\epsilon^i = \left(\frac{3(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon}\right)^\theta$. Then we can bound $\mathbb{P}(\mathbf{E}_2)$ as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P}\left(\sup_{\mathbf{w}_{k_w} \in \Theta} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}_{k_w}, \mathbf{D})) \right\|_2 \geq \frac{t}{3}\right) \\ &= \mathbb{P}\left(\sup_{\mathbf{w}_{k_w} \in \Theta, \lambda \in \lambda_{1/2}} 2 \left\langle \lambda, \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}_{k_w}, \mathbf{D})) \right\rangle \geq \frac{t}{3}\right) \\ &\leq 6^d \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon}\right)^\theta \sup_{\mathbf{w}_{k_w} \in \Theta, \lambda \in \lambda_{1/2}} \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left\langle \lambda, \nabla f(\mathbf{w}_{k_w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}_{k_w}, \mathbf{D})) \right\rangle \geq \frac{t}{6}\right) \\ &\stackrel{\textcircled{1}}{\leq} 6^d \left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon}\right)^\theta 2 \exp\left(-\frac{nt^2}{72\beta^2}\right), \end{aligned}$$

where $\textcircled{1}$ holds since by Lemma 16, we have

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left(\left\langle \lambda, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{D}^{(i)}) \right\rangle\right) > t\right) \leq \exp\left(-\frac{nt^2}{2\beta^2}\right).$$

where $\beta \triangleq \left[\vartheta \tilde{r}_l \tilde{c}_l d_l + \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}\right]^{1/2}$ in which $\vartheta = 1/8$.

Thus, if we set

$$t \geq \sqrt{\frac{72\beta^2 \left(d \log(6) + \theta \log\left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon}\right) + \log\left(\frac{4}{\epsilon}\right)\right)}{n}},$$

then we have

$$\mathbb{P}(\mathbf{E}_2) \leq \frac{\epsilon}{2}.$$

Step 3. Bound $\mathbb{P}(\mathbf{E}_3)$: We first bound $\mathbb{P}(\mathbf{E}_3)$ as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_3) &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \|\mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}, \mathbf{x}))\|_2 \geq \frac{t}{3}\right) \\ &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \frac{\|\mathbb{E}_{\mathbf{D} \sim \mathcal{D}}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x}) - \nabla f(\mathbf{w}, \mathbf{x}))\|_2}{\|\mathbf{w} - \mathbf{w}_{k_w}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \geq \frac{t}{3}\right) \\ &\leq \mathbb{P}\left(\epsilon \mathbb{E}_{\mathbf{D} \sim \mathcal{D}} \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \tilde{\mathcal{Q}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \geq \frac{t}{3}\right) \\ &\leq \mathbb{P}\left(\gamma \epsilon \geq \frac{t}{3}\right). \end{aligned}$$

We set ϵ enough small such that $\gamma \epsilon < t/3$ always holds. Then it yields $\mathbb{P}(\mathbf{E}_3) = 0$.

Step 4. Final result: Note that $\frac{6\beta\epsilon}{\epsilon} \geq 3\beta\epsilon$. Finally, to ensure $\mathbb{P}(\mathbf{E}_0) \leq \epsilon$, we just set $\epsilon = \frac{18p^2(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\vartheta^2 n b_{l+1}} \left[\prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16p^2}\right]^{-\frac{1}{2}}$.

$$t \geq \max\left(\frac{6\gamma\epsilon}{\epsilon}, \sqrt{\frac{72\beta^2 \left(d \log(6) + \theta \log\left(\frac{9(b_{l+1} + \sum_{i=1}^l d_i b_i)}{\epsilon}\right) + \log\left(\frac{4}{\epsilon}\right)\right)}{n}}\right).$$

By comparing the values of β and γ , we have if $n \geq c_{g'} \frac{l^2 b_{l+1}^2 (b_{l+1} + \sum_{i=1}^l d_i b_i)^2 (r_0 c_0 d_0)^4}{d_0^4 b_1^8 (d \log(6) + \theta \varrho) \varepsilon^2 \max_i (r_i c_i)}$ where $c_{g'}$ is a universal constant, then there exists a universal constant c_g such that

$$\begin{aligned} \sup_{\mathbf{w} \in \Omega} \left\| \nabla_{\mathbf{w}} \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla_{\mathbf{w}} \mathbf{Q}(\mathbf{w}) \right\|_2 &\leq c_g \beta \sqrt{\frac{d + \frac{1}{\log(6)} \theta \left[\left(\sum_{i=1}^l \log \left(\frac{\sqrt{d_s} b_s (k_s - s_s + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128 p^2} \right) \right) + \log \left(\frac{4}{\varepsilon} \right)}{n}} \\ &\leq c_g \beta \sqrt{\frac{d + \frac{1}{2} \theta \varrho + \frac{1}{2} \log \left(\frac{4}{\varepsilon} \right)}{n}}, \end{aligned}$$

holds with probability at least $1 - \varepsilon$, where $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$, $\theta = a_{l+1} (d_{l+1} + \tilde{r}_l \tilde{c}_l d_l + 1) + \sum_{i=1}^l a_i (k_i^2 d_{i-1} + d_i + 1)$, $\varrho = \sum_{i=1}^l \log \left(\frac{\sqrt{d_i} b_i (k_i - s_i + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128 p^2} \right)$, and $\beta \triangleq \left[\vartheta \tilde{r}_l \tilde{c}_l d_l + \sum_{i=1}^l \frac{\vartheta b_{l+1}^2 d_{i-1}}{p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{s=i}^l \frac{d_s b_s^2 (k_s - s_s + 1)^2}{16 p^2} \right]^{1/2}$ in which $\vartheta = 1/8$. The proof is completed. \square

D.4 Proof of Corollary 1

Proof. By Theorem 2, we know that there exist universal constants $c_{g'}$ and c_g such that if $n \geq c_{g'} \frac{l^2 b_{l+1}^2 (b_{l+1} + \sum_{i=1}^l d_i b_i)^2 (r_0 c_0 d_0)^4}{d_0^4 b_1^8 (d \log(6) + \theta \varrho) \varepsilon^2 \max_i (r_i c_i)}$, then

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla_{\mathbf{w}} \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla_{\mathbf{w}} \mathbf{Q}(\mathbf{w}) \right\|_2 \leq c_g \beta \sqrt{\frac{2d + \theta \varrho + \log \left(\frac{4}{\varepsilon} \right)}{2n}}$$

holds with probability at least $1 - \varepsilon$, where ϱ is provided in Lemma 1. Here β and d are defined as $\beta = \left[\frac{r_l c_l d_l}{8 p^2} + \sum_{i=1}^l \frac{b_{l+1}^2 d_{i-1}}{8 p^2 b_i^2 d_i} r_{i-1} c_{i-1} \prod_{j=i}^l \frac{d_j b_j^2 (k_j - s_j + 1)^2}{16 p^2} \right]^{1/2}$ and $d = \tilde{r}_l \tilde{c}_l d_l d_{l+1} + \sum_{i=1}^l k_i^2 d_{i-1} d_i$, respectively.

So based on such a result, we can derive that if $n \geq c_g^2 (2d + \theta \varrho + \log(4/\varepsilon)) \beta^2 / (2\varepsilon)$, then we have

$$\left\| \nabla \mathbf{Q}(\tilde{\mathbf{w}}) \right\|_2 \leq \left\| \nabla_{\mathbf{w}} \tilde{\mathbf{Q}}_n(\tilde{\mathbf{w}}) \right\|_2 + \left\| \nabla_{\mathbf{w}} \tilde{\mathbf{Q}}_n(\tilde{\mathbf{w}}) - \nabla_{\mathbf{w}} \mathbf{Q}(\tilde{\mathbf{w}}) \right\|_2 \leq \sqrt{\varepsilon} + c_g \beta \sqrt{\frac{2d + \theta \varrho + \log \left(\frac{4}{\varepsilon} \right)}{2n}} \leq 2\sqrt{\varepsilon}.$$

Thus, we have $\left\| \nabla \mathbf{Q}(\tilde{\mathbf{w}}) \right\|_2^2 \leq 4\varepsilon$, which means that $\tilde{\mathbf{w}}$ is a 4ε -approximate stationary point in population risk with probability at least $1 - \varepsilon$. The proof is completed. \square

D.5 Proof of Theorem 3

Proof. Suppose that $\{\mathbf{w}_{(1)}, \mathbf{w}_{(2)}, \dots, \mathbf{w}_{(m)}\}$ are the non-degenerate critical points of $\mathbf{Q}(\mathbf{w})$. So for any $\mathbf{w}_{(k)}$, it obeys

$$\inf_i |\lambda_i^k (\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)}))| \geq \zeta,$$

where $\lambda_i^k (\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)}))$ denotes the i -th eigenvalue of the Hessian $\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)})$ and ζ is a constant. We further define a set $D = \{\mathbf{w} \in \mathbb{R}^d \mid \left\| \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 \leq \varepsilon \text{ and } \inf_i |\lambda_i (\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)}))| \geq \zeta\}$. According to Lemma 5, $D = \cup_{k=1}^{\infty} D_k$ where each D_k is a disjoint component with $\mathbf{w}_{(k)} \in D_k$ for $k \leq m$ and D_k does not contain any critical point of $\mathbf{Q}(\mathbf{w})$ for $k \geq m + 1$. On the other hand, by the continuity of $\nabla \mathbf{Q}(\mathbf{w})$, it yields $\left\| \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 = \varepsilon$ for $\mathbf{w} \in \partial D_k$. Notice, we set the value of ε blow which is actually a function related to n .

Then by utilizing Theorem 2, we let sample number n sufficient large such that

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 \leq \frac{\varepsilon}{2}$$

holds with probability at least $1 - \varepsilon$, where ε is defined as

$$\frac{\varepsilon}{2} \triangleq c_g \beta \sqrt{\frac{d \log(6) + \theta \left(\sum_{i=1}^l \log \left(\frac{\sqrt{d_s} b_s (k_s - s_s + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128 p^2} \right) \right) + \log \left(\frac{4}{\varepsilon} \right)}{n}}.$$

This further gives that for arbitrary $\mathbf{w} \in D_k$, we have

$$\begin{aligned} \inf_{\mathbf{w} \in D_k} \left\| t \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) + (1-t) \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 &= \inf_{\mathbf{w} \in D_k} \left\| t \left(\nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla \mathbf{Q}(\mathbf{w}) \right) + \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 \\ &\geq \inf_{\mathbf{w} \in D_k} \left\| \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 - \sup_{\mathbf{w} \in D_k} t \left\| \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 \\ &\geq \frac{\epsilon}{2}. \end{aligned} \quad (6)$$

Similarly, by utilizing Lemma 18, let n be sufficient large such that

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right\|_{\text{op}} \leq \frac{\zeta}{2}$$

holds with probability at least $1 - \epsilon$, where ζ satisfies

$$\frac{\zeta}{2} \geq c_v \gamma \sqrt{\frac{d + \theta \varrho + \log\left(\frac{4}{\epsilon}\right)}{n}}.$$

Assume that $\mathbf{b} \in \mathbb{R}^d$ is a vector and satisfies $\mathbf{b}^T \mathbf{b} = 1$. In this case, we can bound $\lambda_i^k \left(\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) \right)$ for arbitrary $\mathbf{w} \in D_k$ as follows:

$$\begin{aligned} \inf_{\mathbf{w} \in D_k} \left| \lambda_i^k \left(\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) \right) \right| &= \inf_{\mathbf{w} \in D_k} \min_{\mathbf{b}^T \mathbf{b} = 1} \left| \mathbf{b}^T \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) \mathbf{b} \right| \\ &= \inf_{\mathbf{w} \in D_k} \min_{\mathbf{b}^T \mathbf{b} = 1} \left| \mathbf{b}^T \left(\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right) \mathbf{b} + \mathbf{b}^T \nabla^2 \mathbf{Q}(\mathbf{w}) \mathbf{b} \right| \\ &\geq \inf_{\mathbf{w} \in D_k} \min_{\mathbf{b}^T \mathbf{b} = 1} \left| \mathbf{b}^T \nabla^2 \mathbf{Q}(\mathbf{w}) \mathbf{b} \right| - \min_{\mathbf{b}^T \mathbf{b} = 1} \left| \mathbf{b}^T \left(\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right) \mathbf{b} \right| \\ &\geq \inf_{\mathbf{w} \in D_k} \min_{\mathbf{b}^T \mathbf{b} = 1} \left| \mathbf{b}^T \nabla^2 \mathbf{Q}(\mathbf{w}) \mathbf{b} \right| - \max_{\mathbf{b}^T \mathbf{b} = 1} \left| \mathbf{b}^T \left(\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right) \mathbf{b} \right| \\ &= \inf_{\mathbf{w} \in D_k} \inf_i \left| \lambda_i^k \left(\nabla^2 f(\mathbf{w}_{(k)}, \mathbf{x}) \right) \right| - \left\| \nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w}) - \nabla^2 \mathbf{Q}(\mathbf{w}) \right\|_{\text{op}} \\ &\geq \frac{\zeta}{2}. \end{aligned} \quad (7)$$

This means that in each set D_k , $\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w})$ has no zero eigenvalues. Then, combine this and Eqn. (6), by Lemma 4 we know that if the population risk $\mathbf{Q}(\mathbf{w})$ has no critical point in D_k , then the empirical risk $\tilde{\mathbf{Q}}_n(\mathbf{w})$ has also no critical point in D_k ; otherwise it also holds.

Now we bound the distance between the corresponding critical points of $\mathbf{Q}(\mathbf{w})$ and $\tilde{\mathbf{Q}}_n(\mathbf{w})$. Assume that in D_k , $\mathbf{Q}(\mathbf{w})$ has a unique critical point $\mathbf{w}_{(k)}$ and $\tilde{\mathbf{Q}}_n(\mathbf{w})$ also has a unique critical point $\mathbf{w}_n^{(k)}$. Then, there exists $t \in [0, 1]$ such that for any $\mathbf{z} \in \partial \mathbf{B}^d(1)$, we have

$$\begin{aligned} \epsilon &\geq \left\| \nabla \mathbf{Q}(\mathbf{w}_n^{(k)}) \right\|_2 \\ &= \max_{\mathbf{z}^T \mathbf{z} = 1} \langle \nabla \mathbf{Q}(\mathbf{w}_n^{(k)}), \mathbf{z} \rangle \\ &= \max_{\mathbf{z}^T \mathbf{z} = 1} \langle \nabla \mathbf{Q}(\mathbf{w}_{(k)}), \mathbf{z} \rangle + \langle \nabla^2 \mathbf{Q}(\mathbf{w}_{(k)} + t(\mathbf{w}_n^{(k)} - \mathbf{w}_{(k)}))(\mathbf{w}_n^{(k)} - \mathbf{w}_{(k)}), \mathbf{z} \rangle \\ &\stackrel{\textcircled{1}}{\geq} \left\langle \left(\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)}) \right)^2 (\mathbf{w}_n^{(k)} - \mathbf{w}_{(k)}), (\mathbf{w}_n^{(k)} - \mathbf{w}_{(k)}) \right\rangle^{1/2} \\ &\stackrel{\textcircled{2}}{\geq} \zeta \left\| \mathbf{w}_n^{(k)} - \mathbf{w}_{(k)} \right\|_2, \end{aligned}$$

where $\textcircled{1}$ holds since $\nabla \mathbf{Q}(\mathbf{w}_{(k)}) = \mathbf{0}$ and $\textcircled{2}$ holds since $\mathbf{w}_{(k)} + t(\mathbf{w}_n^{(k)} - \mathbf{w}_{(k)})$ is in D_k and for any $\mathbf{w} \in D_k$ we have $\inf_i |\lambda_i \left(\nabla^2 \mathbf{Q}(\mathbf{w}) \right)| \geq \zeta$. So if $n \geq c_h \max \left(\frac{l^2 b_{l+1}^2 (b_{l+1} + \sum_{i=1}^l d_i b_i)^2 (r_0 c_0 d_0)^4}{d_0^4 b_1^8 d_{\varrho} \epsilon^2 \max_i (r_i c_i)}, \frac{d + \theta \varrho}{\zeta^2} \right)$ where c_h is a constant, then

$$\left\| \mathbf{w}_n^{(k)} - \mathbf{w}_{(k)} \right\|_2 \leq \frac{2c_g \beta}{\zeta} \sqrt{\frac{d \log(6) + \theta \left(\sum_{i=1}^l \log \left(\frac{\sqrt{d_s} b_s (k_s - s_s + 1)}{4p} \right) + \log(b_{l+1}) + \log \left(\frac{n}{128p^2} \right) \right) + \log \left(\frac{4}{\epsilon} \right)}{n}}$$

holds with probability at least $1 - \varepsilon$. \square

D.6 Proof of Corollary 2

Proof. By Theorem 3, we know that the non-degenerate stationary point $\mathbf{w}_{(k)}$ in the m non-degenerate stationary points in population risk, denoted by $\{\mathbf{w}_{(1)}, \mathbf{w}_{(2)}, \dots, \mathbf{w}_{(m)}\}$ uniquely corresponding to a non-degenerate stationary point $\mathbf{w}_n^{(k)}$ in the empirical risk.

On the other hand, for any $\mathbf{w}_{(k)}$, it obeys

$$\inf_i |\lambda_i^k (\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)}))| \geq \zeta,$$

where $\lambda_i^k (\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)}))$ denotes the i -th eigenvalue of the Hessian $\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)})$ and ζ is a constant. We further define a set $D = \{\mathbf{w} \in \mathbb{R}^d \mid \|\nabla \mathbf{Q}(\mathbf{w})\|_2 \leq \epsilon \text{ and } \inf_i |\lambda_i (\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)}))| \geq \zeta\}$. According to Lemma 5, $D = \cup_{k=1}^{\infty} D_k$ where each D_k is a disjoint component with $\mathbf{w}_{(k)} \in D_k$ for $k \leq m$ and D_k does not contain any critical point of $\mathbf{Q}(\mathbf{w})$ for $k \geq m + 1$. Then $\mathbf{w}_n^{(k)}$ also belong to the component D_k due to the unique corresponding relation between $\mathbf{w}_{(k)}$ and $\mathbf{w}_n^{(k)}$. Then from Eqn. (6) and (7), we know that if the assumptions in Theorem 3 hold, then for arbitrary $\mathbf{w} \in D_k$ and $t \in (0, 1)$,

$$\inf_{\mathbf{w} \in D_k} \left\| t \nabla \tilde{\mathbf{Q}}_n(\mathbf{w}) + (1-t) \nabla \mathbf{Q}(\mathbf{w}) \right\|_2 \geq \frac{\epsilon}{2} \quad \text{and} \quad \inf_{\mathbf{w} \in D_k} \left| \lambda_i^k (\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w})) \right| \geq \frac{\zeta}{2},$$

where ϵ and ζ are constants. This means that in each set D_k , $\nabla^2 \tilde{\mathbf{Q}}_n(\mathbf{w})$ has no zero eigenvalues. Then, combine this and Eqn. (6), we can obtain that in D_k , if $\mathbf{Q}(\mathbf{w})$ has a unique critical point $\mathbf{w}_{(k)}$ with non-degenerate index s_k , then $\tilde{\mathbf{Q}}_n(\mathbf{w})$ also has a unique critical point $\mathbf{w}_n^{(k)}$ in D_k with the same non-degenerate index s_k . Namely, the number of negative eigenvalues of the Hessian matrices $\nabla^2 \mathbf{Q}(\mathbf{w}_{(k)})$ and $\nabla^2 \mathbf{Q}(\mathbf{w}_n^{(k)})$ are the same. This further gives that if one of the pair $(\mathbf{w}_{(k)}, \mathbf{w}_n^{(k)})$ is a local minimum or saddle point, then another one is also a local minimum or a saddle point. The proof is completed. \square

E Proof of Auxiliary Lemmas

E.1 Proof of Lemma 8

Proof. (1) Since $\mathbf{G}(\mathbf{z})$ is a diagonal matrix and its diagonal values are upper bounded by $\sigma_1(\mathbf{z}_i)(1 - \sigma_1(\mathbf{z})) \leq 1/4$ where \mathbf{z}_i denotes the i -th entry of \mathbf{z}_i , we can conclude

$$\|\mathbf{G}(\mathbf{z})\mathbf{M}\|_F^2 \leq \frac{1}{16} \|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N}\mathbf{G}(\mathbf{z})\|_F^2 \leq \frac{1}{16} \|\mathbf{N}\|_F^2.$$

(2) The operator $\mathbf{Q}(\cdot)$ maps a vector $\mathbf{z} \in \mathbb{R}^d$ into a matrix of size $d^2 \times d$ whose $((i-1)d + i, i)$ ($i = 1, \dots, d$) entry equal to $\sigma_1(\mathbf{z}_i)(1 - \sigma_1(\mathbf{z}_i))(1 - 2\sigma_1(\mathbf{z}_i))$ and rest entries are all 0. This gives

$$\begin{aligned} \sigma_1(\mathbf{z}_i)(1 - \sigma_1(\mathbf{z}_i))(1 - 2\sigma_1(\mathbf{z}_i)) &= \frac{1}{3}(3\sigma_1(\mathbf{z}_i))(1 - \sigma_1(\mathbf{z}_i))(1 - 2\sigma_1(\mathbf{z}_i)) \\ &\leq \frac{1}{3} \left(\frac{3\sigma_1(\mathbf{z}_i) + 1 - \sigma_1(\mathbf{z}_i) + 1 - 2\sigma_1(\mathbf{z}_i)}{3} \right)^3 \\ &\leq \frac{2^3}{3^4}. \end{aligned}$$

This means the maximal value in $\mathbf{Q}(\mathbf{z})$ is at most $\frac{2^3}{3^4}$. Consider the structure in $\mathbf{Q}(\mathbf{z})$, we can obtain

$$\|\mathbf{Q}(\mathbf{z})\mathbf{M}\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N}\mathbf{Q}(\mathbf{z})\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{N}\|_F^2.$$

(3) $\text{up}(\mathbf{M})$ represents conducting upsampling on $\mathbf{M} \in \mathbb{R}^{s \times t \times q}$. Let $\mathbf{N} = \text{up}(\mathbf{M}) \in \mathbb{R}^{ps \times pt \times q}$. Specifically, for each slice $\mathbf{N}(:, :, i)$ ($i = 1, \dots, q$), we have $\mathbf{N}(:, :, i) = \text{up}(\mathbf{M}(:, :, i))$. It actually upsamples each entry $\mathbf{M}(g, h, i)$ into a matrix of p^2 same entries $\frac{1}{p^2} \mathbf{M}(g, h, i)$. So it is easy to obtain

$$\|\text{up}(\mathbf{M})\|_F^2 \leq \frac{1}{p^2} \|\mathbf{M}\|_F^2.$$

(4) Let $\mathbf{M} = \mathbf{W}(:, :, i)$ and $\mathbf{N} = \tilde{\delta}_{i+1}(:, : i)$. Assume that $\mathbf{H} = \mathbf{M} \circledast \mathbf{N} \in \mathbb{R}^{m_1 \times m_2}$, where $m_1 = \tilde{r}_{i-1} - 2k_i + 2$ and $m_2 = \tilde{c}_{i-1} - 2k_i + 2$. Then we have

$$\|\mathbf{H}\|_F^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} |\mathbf{H}(i, j)|^2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \langle \mathbf{M}_{\Omega_{i,j}}, \mathbf{N} \rangle^2 \leq \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \|\mathbf{M}_{\Omega_{i,j}}\|_F^2 \|\mathbf{N}\|_F^2,$$

where $\Omega_{i,j}$ denotes the entry index of \mathbf{M} for the (i, j) -th convolution operation (*i.e.* computing the $\mathbf{H}(i, j)$).

Since for each convolution computing, each element in \mathbf{M} is involved at most one time, we can claim that any element in \mathbf{M} in $\sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \|\mathbf{M}_{\Omega_{i,j}}\|_F^2$ occurs at most $(k_i - s_i + 1)^2$ since there are $s_i - 1$ rows and columns between each neighboring nonzero entries in \mathbf{N} which is decided by the definition of $\tilde{\delta}_{i+1}$ in Sec. B.1. Therefore, we have

$$\sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \|\mathbf{M}_{\Omega_{i,j}}\|_F^2 \leq (k_i - s_i + 1)^2 \|\mathbf{M}\|_F^2,$$

which further gives

$$\|\mathbf{M} \circledast \mathbf{N}\|_F^2 \leq (k_i - s_i + 1)^2 \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2.$$

Consider all the slices in $\tilde{\delta}_{i+1}$, we can obtain

$$\|\tilde{\delta}_{i+1} \circledast \mathbf{W}\|_F^2 \leq (k_i - s_i + 1)^2 \|\mathbf{W}\|_F^2 \|\tilde{\delta}_{i+1}\|_F^2.$$

(5) Since for softmax activation function σ_2 , we have $\sum_{i=1}^{d_i+1} \mathbf{v}_i = 1$ ($\mathbf{v}_i \geq 0$) and there is only one nonzero entry (*i.e.* 1) in \mathbf{y} , we can obtain

$$0 \leq \|\mathbf{v} - \mathbf{y}\|_2^2 = \|\mathbf{v}\|_2^2 + \|\mathbf{y}\|_2^2 - 2\langle \mathbf{v}, \mathbf{y} \rangle = 2 - 2\langle \mathbf{v}, \mathbf{y} \rangle \leq 2.$$

The proof is completed. □

References

- Candes, E. and Plan, Y. Tight oracle bounds for low-rank matrix recovery from a minimal number of random measurements. *IEEE TIT*, 57(4):2342–2359, 2009.
- Hoeffding, W. Probability inequalities for sums of bounded random variables. *Journal of the American statistical association*, 58(301): 13–30, 1963.
- Mei, S., Bai, Y., and Montanari, A. The landscape of empirical risk for non-convex losses. *Annals of Statistics*, 2017.
- Vershynin, R. Introduction to the non-asymptotic analysis of random matrices, compressed sensing. *Cambridge Univ. Press, Cambridge*, pp. 210–268, 2012.