

# EMPIRICAL RISK LANDSCAPE ANALYSIS FOR UNDERSTANDING DEEP NEURAL NETWORKS

**Pan Zhou & Jiashi Feng**

Department of Electrical and Computer Engineering  
National University of Singapore  
Singapore, 117583  
{pzhou@u.nus.edu, elefjia@nus.edu.sg}

## ABSTRACT

This work aims to provide comprehensive landscape analysis of empirical risk in deep neural networks (DNNs), including the convergence behavior of its gradient, its stationary points and the empirical risk itself to their corresponding population counterparts, which reveals how various network parameters determine the convergence performance. In particular, for an  $l$ -layer linear neural network consisting of  $d_i$  neurons in the  $i$ -th layer, we prove the gradient of its empirical risk uniformly converges to the one of its population risk, at the rate of  $\mathcal{O}(r^{2l} \sqrt{l \sqrt{\max_i d_i} s \log(d/l)/n})$ . Here  $d$  is the total weight dimension,  $s$  is the number of nonzero entries of all the weights and the magnitude of weights per layer is upper bounded by  $r$ . Moreover, we prove the one-to-one correspondence of the non-degenerate stationary points between the empirical and population risks and provide convergence guarantee for each pair. We also establish the uniform convergence of the empirical risk to its population counterpart and further derive the stability and generalization bounds for the empirical risk. In addition, we analyze these properties for deep *nonlinear* neural networks with sigmoid activation functions. We prove similar results for convergence behavior of their empirical risk gradients, non-degenerate stationary points as well as the empirical risk itself.

To our best knowledge, this work is the first one theoretically characterizing the uniform convergence of the gradient and stationary points of the empirical risk of DNN models, which benefits the theoretical understanding on how the neural network depth  $l$ , the layer width  $d_i$ , the network size  $d$ , the sparsity in weight and the parameter magnitude  $r$  determine the neural network landscape.

## 1 INTRODUCTION

Deep learning has achieved remarkable success in many fields, such as computer vision (Hinton et al., 2006; Szegedy et al., 2015; He et al., 2016), natural language processing (Collobert & Weston, 2008; Bakshi & Stephanopoulos, 1993), and speech recognition (Hinton et al., 2012; Graves et al., 2013). However, theoretical understanding on the properties of deep learning models still lags behind their practical achievements (Shalev-Shwartz et al., 2017; Kawaguchi, 2016) due to their high non-convexity and internal complexity. In practice, parameters of deep learning models are learned by minimizing the *empirical risk* via (stochastic-)gradient descent. Therefore, some recent works (Bartlett & Maass, 2003; Neyshabur et al., 2015) analyzed the convergence of the empirical risk to the population risk, which are however still far from fully understanding the landscape of the empirical risk in deep learning models. Beyond the convergence properties of the empirical risk itself, the convergence and distribution properties of its gradient and stationary points are also essential in landscape analysis. A comprehensive landscape analysis can reveal important information on the optimization behavior and practical performance of deep neural networks, and will be helpful to designing better network architectures. Thus, in this work we aim to provide comprehensive landscape analysis by looking into the gradients and stationary points of the empirical risk.

Formally, we consider a DNN model  $f(\mathbf{w}; \mathbf{x}, \mathbf{y}) : \mathbb{R}^{d_0} \times \mathbb{R}^{d_l} \rightarrow \mathbb{R}$  parameterized by  $\mathbf{w} \in \mathbb{R}^d$  consisting of  $l$  layers ( $l \geq 2$ ) that is trained by minimizing the commonly used squared loss function

over sample pairs  $\{(\mathbf{x}, \mathbf{y})\} \subset \mathbb{R}^{d_0} \times \mathbb{R}^{d_l}$  from an unknown distribution  $\mathcal{D}$ , where  $\mathbf{y}$  is the target output for the sample  $\mathbf{x}$ . Ideally, the model can find its optimal parameter  $\mathbf{w}^*$  by minimizing the *population risk* through (stochastic-)gradient descent by backpropagation:

$$\min_{\mathbf{w}} \mathbf{J}(\mathbf{w}) \triangleq \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} f(\mathbf{w}; \mathbf{x}, \mathbf{y}),$$

where  $f(\mathbf{w}; \mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2$  is the squared loss associated to the sample  $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$  in which  $\mathbf{v}^{(l)}$  is the output of the  $l$ -th layer. In practice, as the sample distribution  $\mathcal{D}$  is usually unknown and only finite training samples  $\{(\mathbf{x}_{(i)}, \mathbf{y}_{(i)})\}_{i=1}^n$  *i.i.d.* drawn from  $\mathcal{D}$  are provided, the network model is usually trained by minimizing the empirical risk:

$$\min_{\mathbf{w}} \hat{\mathbf{J}}_n(\mathbf{w}) \triangleq \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; \mathbf{x}_{(i)}, \mathbf{y}_{(i)}). \quad (1)$$

Understanding the convergence behavior of  $\hat{\mathbf{J}}_n(\mathbf{w})$  to  $\mathbf{J}(\mathbf{w})$  is critical to statistical machine learning algorithms. In this work, we aim to go further and characterize the landscape of the empirical risk  $\hat{\mathbf{J}}_n(\mathbf{w})$  of deep learning models by analyzing the convergence behavior of its gradient and stationary points to their corresponding population counterparts. We provide analysis for both multi-layer linear and nonlinear neural networks. In particular, we obtain following new results.

- We establish the uniform convergence of empirical gradient  $\nabla_{\mathbf{w}} \hat{\mathbf{J}}_n(\mathbf{w})$  to its population counterpart  $\nabla_{\mathbf{w}} \mathbf{J}(\mathbf{w})$ . Specifically, when the sample size  $n$  is not less than  $\mathcal{O}(\max(l^3 r^2 / (\varepsilon^2 s \log(d/l)), s \log(d/l)/l))$ , with probability at least  $1 - \varepsilon$  the convergence rate is  $\mathcal{O}(r^{2l} \sqrt{l \sqrt{\max_i d_i} s \log(d/l)/n})$ , where there are  $s$  nonzero entries in the parameter  $\mathbf{w}$ , the output dimension of the  $i$ -th layer is  $d_i$  and the magnitude of the weight parameter of each layer is upper bounded by  $r$ . This result implies that as long as the training sample size  $n$  is sufficiently large, any stationary point of  $\hat{\mathbf{J}}_n(\mathbf{w})$  is also a stationary point of  $\mathbf{J}(\mathbf{w})$  and vice versa, although both  $\hat{\mathbf{J}}_n(\mathbf{w})$  and  $\mathbf{J}(\mathbf{w})$  are very complex.
- We then prove the exact correspondence of non-degenerate stationary points between  $\hat{\mathbf{J}}_n(\mathbf{w})$  and  $\mathbf{J}(\mathbf{w})$ . Indeed, the corresponding non-degenerate stationary points also uniformly converge to each other at the same convergence rate as the one revealed above with an extra factor  $2/\zeta$ . Here  $\zeta > 0$  accounts for the geometric topology of non-degenerate stationary points (see Definition 1).

Based on the above two new results, we also derive the uniform convergence of the empirical risk  $\hat{\mathbf{J}}_n(\mathbf{w})$  to its population risk  $\mathbf{J}(\mathbf{w})$ , which helps understand the generalization error of deep learning models and stability of their empirical risk. These analyses reveal the role of the depth  $l$  of a neural network model in determining its convergence behavior and performance. Also, the results tell that the width factor  $\sqrt{\max_i d_i}$ , the nonzero entry number  $s$  of weights, and the total network size  $d$  are also critical to the convergence and performance. In addition, controlling magnitudes of the parameters (weights) in DNNs are demonstrated to be important for performance. To our best knowledge, this work is the first one theoretically characterizing the uniform convergence of empirical gradient and stationary points in both deep linear and nonlinear neural networks.

## 2 RELATED WORK

To date, only a few theories have been developed for understanding DNNs which can be roughly divided into following three categories. The first category aims to analyze training error of DNNs. [Baum \(1988\)](#) pointed out that zero training error can be obtained when the last layer of a neural network has more units than training samples. Later, [Soudry & Carmon \(2016\)](#) proved that for DNNs with leaky rectified linear units (ReLU) and a single output, the training error achieves zero at any of their local minima as long as the product of the number of units in the last two layers is larger than the training sample size.

The second category of analysis works ([Dauphin et al., 2014](#); [Choromanska et al., 2015a](#); [Kawaguchi, 2016](#); [Tian, 2017](#)) focus on analyzing loss surfaces of DNNs, *e.g.*, how the stationary points are distributed. Those results are helpful to understanding performance difference of large- and small-size

networks (Choromanska et al., 2015b). Among them, Dauphin et al. (2014) experimentally verified that a large number of saddle points indeed exist for DNNs. With strong assumptions, Choromanska et al. (2015a) connected the loss function of a deep ReLU network with the spherical spin-class model and described locations of the local minima. Later, Kawaguchi (2016) proved the existence of degenerate saddle points for deep linear neural networks with squared loss function. They also showed that any local minimum is also a global minimum. By utilizing techniques from dynamical system analysis, Tian (2017) gave guarantees that for two-layer bias-free networks with ReLUs, the gradient descent algorithm with certain symmetric weight initialization can converge to the ground-truth weights globally, if the inputs follow Gaussian distribution. Recently, Nguyen & Hein (2017) proved that for a fully connected network with squared loss and analytic activation functions, almost all the local minima are globally optimal if one hidden layer has more units than training samples and the network structure after this layer is pyramidal. Besides, some recent works, *e.g.*, (Zhang et al., 2016; 2017), tried to alleviate analysis difficulties by relaxing the involved highly nonconvex functions into ones easier.

In addition, some existing works (Bartlett & Maass, 2003; Neyshabur et al., 2015) analyze the generalization performance of a DNN model. Based on the Vapnik–Chervonenkis (VC) theory, Bartlett & Maass (2003) proved that for a feedforward neural network with one-dimensional output, the best convergence rate of the empirical risk to its population risk on the sample distribution can be bounded by its fat-shattering dimension. Recently, Neyshabur et al. (2015) adopted Rademacher complexity to analyze learning capacity of a fully-connected neural network model with ReLU activation functions and bounded inputs.

However, although gradient descent with backpropagation is the most common optimization technique for DNNs, none of existing works analyzes convergence properties of gradient and stationary points of the DNN empirical risk. For single-layer optimization problems, some previous works analyze their empirical risk but essentially differ from our analysis method. For example, Negahban et al. (2009) proved that for a regularized convex program, the minimum of the empirical risk uniformly converges to the true minimum of the population risk under certain conditions. Gonen & Shalev-Shwartz (2017) proved that for nonconvex problems without degenerated saddle points, the difference between empirical risk and population risk can be bounded. Unfortunately, the loss of DNNs is highly nonconvex and has degenerated saddle points (Fyodorov & Williams, 2007; Dauphin et al., 2014; Kawaguchi, 2016), thus their analysis results are not applicable. Mei et al. (2017) analyzed the convergence behavior of the empirical risk for nonconvex problems, but they only considered the single-layer nonconvex problems and their analysis demands strong sub-Gaussian and sub-exponential assumptions on the gradient and Hessian of the empirical risk respectively. Their analysis also assumes a linearity property on gradient which is difficult to hold or verify. In contrast, our analysis requires much milder assumptions. Besides, we prove that for deep networks which are highly nonconvex, the non-degenerate stationary points of empirical risk can uniformly converge to their corresponding stationary points of population risk at the rate of  $\mathcal{O}(\sqrt{s/n})$  which is faster than the rate  $\mathcal{O}(\sqrt{d/n})$  for single-layer optimization problems in (Mei et al., 2017). Also, Mei et al. (2017) did not analyze the convergence rate of the empirical risk, stability or generalization error of DNNs as this work.

### 3 PRELIMINARIES

Throughout the paper, we denote matrices by boldface capital letters, *e.g.*  $\mathbf{A}$ . Vectors are denoted by boldface lowercase letters, *e.g.*  $\mathbf{a}$ , and scalars are denoted by lowercase letters, *e.g.*  $a$ . We define the  $r$ -radius ball as  $\mathbb{B}^d(r) \triangleq \{\mathbf{z} \in \mathbb{R}^d \mid \|\mathbf{z}\|_2 \leq r\}$ . To explain the results, we also need the vectorization operation  $\text{vec}(\cdot)$ . It is defined as  $\text{vec}(\mathbf{A}) = (\mathbf{A}(:, 1); \cdots; \mathbf{A}(:, t)) \in \mathbb{R}^{st}$  that vectorizes  $\mathbf{A} \in \mathbb{R}^{s \times t}$  along its columns. We use  $d = \sum_{j=1}^l d_j d_{j-1}$  to denote the total dimension of weight parameters, where  $d_j$  denotes the output dimension of the  $j$ -th layer.

In this work, we consider both linear and nonlinear DNNs. Suppose both networks consist of  $l$  layers. We use  $\mathbf{u}^{(j)}$  and  $\mathbf{v}^{(j)}$  to respectively denote the input and output of the  $j$ -th layer,  $\forall j = 1, \dots, l$ .

**Deep linear neural networks:** The function of the  $j$ -th layer is formulated as

$$\mathbf{u}^{(j)} \triangleq \mathbf{W}^{(j)} \mathbf{v}^{(j-1)} \in \mathbb{R}^{d_j}, \quad \mathbf{v}^{(j)} \triangleq \mathbf{u}^{(j)} \in \mathbb{R}^{d_j}, \quad \forall j = 1, \dots, l,$$

where  $\mathbf{v}^{(0)} = \mathbf{x}$  is the input and  $\mathbf{W}^{(j)} \in \mathbb{R}^{d_j \times d_{j-1}}$  is the weight matrix of the  $j$ -th layer.

**Deep nonlinear neural networks:** We adopt the sigmoid function as the non-linear activation function. The function within the  $j$ -th layer can be written as

$$\mathbf{u}^{(j)} \triangleq \mathbf{W}^{(j)} \mathbf{v}^{(j-1)} \in \mathbb{R}^{d_j}, \quad \mathbf{v}^{(j)} \triangleq h_j(\mathbf{u}^{(j)}) = (\sigma(\mathbf{u}_1^{(j)}); \dots; \sigma(\mathbf{u}_{d_j}^{(j)})) \in \mathbb{R}^{d_j}, \quad \forall j = 1, \dots, l,$$

where  $\mathbf{u}_i^{(j)}$  denotes the  $i$ -th entry of  $\mathbf{u}^{(j)}$  and  $\sigma(\cdot)$  is the sigmoid function, i.e.,  $\sigma(a) = 1/(1 + e^{-a})$ .

Following the common practice, both DNN models adopt the squared loss function defined as  $f(\mathbf{w}; \mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2$ , where  $\mathbf{w} = (\mathbf{w}_{(1)}; \dots; \mathbf{w}_{(l)}) \in \mathbb{R}^d$  contains all the weight parameters and  $\mathbf{w}_{(j)} = \text{vec}(\mathbf{W}^{(j)}) \in \mathbb{R}^{d_j d_{j-1}}$ . Then the empirical risk  $\hat{\mathcal{J}}_n(\mathbf{w})$  is  $\hat{\mathcal{J}}_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; \mathbf{x}_{(i)}, \mathbf{y}_{(i)}) = \frac{1}{2n} \sum_{i=1}^n \|\mathbf{v}_{(i)}^{(l)} - \mathbf{y}_{(i)}\|_2^2$ , where  $\mathbf{v}_{(i)}^{(l)}$  is the network's output of  $\mathbf{x}_{(i)}$ .

## 4 RESULTS FOR DEEP LINEAR NEURAL NETWORKS

We first analyze linear neural network models and present following new results: (1) the uniform convergence of the empirical risk gradient to its population counterpart and (2) the convergence properties of non-degenerate stationary points of the empirical risk. As a corollary, we also derive the uniform convergence of the empirical risk to the population one, which further gives stability and generalization bounds. In the next section, we extend the analysis to non-linear neural network models.

We assume the input datum  $\mathbf{x}$  is  $\tau^2$ -sub-Gaussian and has bounded magnitude, as formally stated in Assumption 1.

**Assumption 1.** *The input datum  $\mathbf{x} \in \mathbb{R}^{d_0}$  has zero mean and is  $\tau^2$ -sub-Gaussian, i.e.,*

$$\mathbb{E}[\exp(\langle \boldsymbol{\lambda}, \mathbf{x} \rangle)] \leq \exp\left(\frac{1}{2} \tau^2 \|\boldsymbol{\lambda}\|_2^2\right), \quad \forall \boldsymbol{\lambda} \in \mathbb{R}^{d_0}.$$

*Besides, the magnitude  $\mathbf{x}$  is bounded as  $\|\mathbf{x}\|_2 \leq r_x$ , where  $r_x$  is a positive universal constant.*

Note that any random vector  $\mathbf{z}$  consisting of independent entries with bounded magnitude is sub-Gaussian and satisfies Assumption 1 (Vershynin, 2012). Moreover, for such a random  $\mathbf{z}$ , we have  $\tau = \|\mathbf{z}\|_\infty \leq \|\mathbf{z}\|_2 \leq r_x$ . Such an assumption on bounded magnitude generally holds for natural data, e.g., images and speech signals. Besides, we assume the weight parameters  $\mathbf{w}_{(j)}$  of each layer are bounded as  $\mathbf{w} \in \Omega = \{\mathbf{w} \mid \mathbf{w}_{(j)} \in \mathbb{B}^{d_j d_{j-1}}(r_j), \forall j = 1, \dots, l\}$  where  $r_j$  is a constant. For notational simplicity, we let  $r = \max_j r_j$ . Such an assumption is common (Xu & Mannor, 2012). Here we assume the entry value of  $\mathbf{y}$  falls in  $[0, 1]$ . For any bounded target output  $\mathbf{y}$ , we can always scale it to satisfy such a requirement.

The results presented for linear neural networks here can be generalized to deep ReLU neural networks by applying the results from Choromanska et al. (2015a) and Kawaguchi (2016), which transform deep ReLU neural networks into deep linear neural networks under proper assumptions.

### 4.1 UNIFORM CONVERGENCE OF EMPIRICAL RISK GRADIENT

We first analyze the convergence of gradients for the DNN empirical and population risks. To our best knowledge, these results are the first ones giving guarantees on gradient convergence, which help better understand the landscape of DNNs and their optimization behavior. The results are stated below.

**Theorem 1.** *Suppose Assumption 1 on the input datum  $\mathbf{x}$  holds and the activation functions in a deep neural network are linear. Then the empirical gradient uniformly converges to the population gradient in Euclidean norm. Specifically, there exist two universal constants  $c_{g'}$  and  $c_g$  such that if  $n \geq c_{g'} \max(l^3 r^2 r_x^4 / (c_q s \log(d/l) \varepsilon^2 \tau^4 \log(1/\varepsilon)), s \log(d/l) / (l \tau^2))$  where  $c_q = \sqrt{\max_{0 \leq i \leq l} d_i}$ , then*

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{\mathcal{J}}_n(\mathbf{w}) - \nabla \mathcal{J}(\mathbf{w}) \right\|_2 \leq \epsilon_g \triangleq c_g \tau \omega_g \sqrt{l c_q} \sqrt{\frac{s \log(dn/l) + \log(12/\varepsilon)}{n}}$$

*holds with probability at least  $1 - \varepsilon$ , where  $s$  denotes the number of nonzero entries of all weight parameters and  $\omega_g = \max(\tau r^{2l-1}, r^{2l-1}, r^{l-1})$ .*

From Theorem 1, one can observe that with an increasingly larger sample size  $n$ , the difference between empirical risk and population risk gradients decreases monotonically at the rate of  $\mathcal{O}(1/\sqrt{n})$  (up to a log factor). Theorem 1 also characterizes how the depth  $l$  contributes to obtaining small difference between the empirical and population risk gradients. Specifically, a deeper neural network needs more training samples to mitigate the difference. Also, due to the factor  $d$ , training a network of larger size using gradient descent also requires more training samples. We observe a factor of  $\sqrt{\max_i \bar{d}_i}$  (i.e.  $c_q$ ), which prefers a DNN architecture of balanced layer sizes (without extremely wide layers). This result also matches the trend and empirical performance in deep learning applications advocating deep but thin networks (He et al., 2016; Szegedy et al., 2015).

By observing Theorem 1, imposing certain regularizations on the weight parameters is useful. For example, reducing the number of nonzero entries  $s$  encourages sparsity regularization like  $\|\mathbf{w}\|_1$ . The results also suggest not choosing large-magnitude weights  $\mathbf{w}$  in order for a smaller factor  $r$  by adopting regularization like  $\|\mathbf{w}\|_2^2$ .

Theorem 1 also reveals the point derived from optimizing that the empirical and population risks have similar properties when the sample size  $n$  is sufficiently large. For example, an  $\epsilon/2$ -stationary point  $\tilde{\mathbf{w}}$  of  $\hat{\mathbf{J}}_n(\mathbf{w})$  is also an  $\epsilon$ -stationary point of  $\mathbf{J}(\mathbf{w})$  with probability  $1 - \epsilon$  if  $n \geq c_\epsilon (\tau\omega_g/\epsilon)^2 l c_q s \log(d/l)$  with  $c_\epsilon$  being a constant. Here  $\epsilon$ -stationary point for a function  $\mathbf{F}$  means the point  $\mathbf{w}$  satisfying  $\|\nabla_{\mathbf{w}} \mathbf{F}\|_2 \leq \epsilon$ . Understanding such properties is useful, since in practice one usually computes an  $\epsilon$ -stationary point of  $\hat{\mathbf{J}}_n(\mathbf{w})$ . These results guarantee the computed point is at most a  $2\epsilon$ -stationary point of  $\mathbf{J}(\mathbf{w})$  and is thus close to the optimum.

#### 4.2 UNIFORM CONVERGENCE OF STATIONARY POINTS

We then proceed to analyze the distribution and convergence properties of stationary points of the DNN empirical risk. Here we consider non-degenerate stationary points which are geometrically isolated and thus unique in local regions. Since degenerate stationary points are not unique in a local region, we cannot expect to establish one-to-one corresponding relationship (see below) between them in empirical risk and population risk.

**Definition 1. (Non-degenerate stationary points)** (Gromoll & Meyer, 1969) *If a stationary point  $\mathbf{w}$  is said to be a non-degenerate stationary point of  $\mathbf{J}(\mathbf{w})$ , then it satisfies*

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

where  $\lambda_i(\nabla^2 \mathbf{J}(\mathbf{w}))$  denotes the  $i$ -th eigenvalue of the Hessian  $\nabla^2 \mathbf{J}(\mathbf{w})$  and  $\zeta$  is a positive constant.

Non-degenerate stationary points include local minima/maxima and non-degenerate saddle points, while degenerate stationary points refer to degenerate saddle points. Then we introduce the index of non-degenerate stationary points which can characterize their geometric properties.

**Definition 2. (Index of non-degenerate stationary points)** (Dubrovin et al., 2012) *The index of a symmetric non-degenerate matrix is the number of its negative eigenvalues, and the index of a non-degenerate stationary point  $\mathbf{w}$  of a smooth function  $\mathbf{F}$  is simply the index of its Hessian  $\nabla^2 \mathbf{F}(\mathbf{w})$ .*

Suppose that  $\mathbf{J}(\mathbf{w})$  has  $m$  non-degenerate stationary points that are denoted as  $\{\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(m)}\}$ . We prove following convergence behavior of these stationary points.

**Theorem 2.** *Suppose Assumption 1 on the input datum  $\mathbf{x}$  holds and the activation functions in a deep neural network are linear. Then if  $n \geq c_h \max(l^3 r^2 r_x^4 / (c_q s \log(d/l) \epsilon^2 \tau^4 \log(1/\epsilon)), s \log(d/l) / \zeta^2)$  where  $c_h$  is a constant, for  $k \in \{1, \dots, m\}$ , there exists a non-degenerate stationary point  $\mathbf{w}_n^{(k)}$  of  $\hat{\mathbf{J}}_n(\mathbf{w})$  which corresponds to the non-degenerate stationary point  $\mathbf{w}^{(k)}$  of  $\mathbf{J}(\mathbf{w})$  with probability at least  $1 - \epsilon$ . In addition,  $\mathbf{w}_n^{(k)}$  and  $\mathbf{w}^{(k)}$  have the same non-degenerate index and they satisfy*

$$\|\mathbf{w}_n^{(k)} - \mathbf{w}^{(k)}\|_2 \leq \frac{2c_g \tau \omega_g}{\zeta} \sqrt{l c_q} \sqrt{\frac{s \log(dn/l) + \log(12/\epsilon)}{n}}, \quad (k = 1, \dots, m)$$

with probability at least  $1 - \epsilon$ , where the parameters  $c_q$ ,  $\omega_g$ , and  $c_g$  are given in Theorem 1.

Theorem 2 guarantees the one-to-one correspondence between the non-degenerate stationary points of the empirical risk  $\hat{\mathbf{J}}_n(\mathbf{w})$  and the popular risk  $\mathbf{J}(\mathbf{w})$ . The distances of the corresponding pairs

become smaller as  $n$  increases. In addition, the corresponding pairs have the same non-degenerate index. This implies that the corresponding stationary points have the same geometric properties, such as whether they are saddle points. Accordingly, we can develop more efficient algorithms, *e.g.* escaping saddle points (Ge et al., 2015), since Dauphin et al. (2014) empirically proved that saddle points are usually surrounded by high error plateaus. Also when  $n$  is sufficiently large, the properties of stationary points of  $\hat{\mathbf{J}}_n(\mathbf{w})$  are similar to the points of the population risk  $\mathbf{J}(\mathbf{w})$  in the sense that they have exactly matching local minima/maxima and non-degenerate saddle points. By comparing Theorems 1 and 2, we find that the requirement for sample number in Theorem 2 is more restrict, since establishing exact one-to-one correspondence between the non-degenerate stationary points of  $\hat{\mathbf{J}}_n(\mathbf{w})$  and  $\mathbf{J}(\mathbf{w})$  and bounding their uniform convergence rate to each other are more challenging. From Theorems 1 and 2, we also notice that the uniform convergence rate of non-degenerate stationary points has an extra factor  $1/\zeta$ . This is because bounding stationary points needs to access not only the gradient itself but also the Hessian matrix. See more details in proof.

Kawaguchi (2016) pointed out that degenerate stationary points indeed exist for DNNs. However, since degenerate stationary points are not isolated, such as forming flat regions, it is hard to establish the unique correspondence for them as for non-degenerate ones. Fortunately, by Theorem 1, the gradients at these points of  $\hat{\mathbf{J}}_n(\mathbf{w})$  and  $\mathbf{J}(\mathbf{w})$  are close. This implies that a degenerate stationary point of  $\mathbf{J}(\mathbf{w})$  will also give a near-zero gradient for  $\hat{\mathbf{J}}_n(\mathbf{w})$ , *i.e.*, it is also a stationary point for  $\hat{\mathbf{J}}_n(\mathbf{w})$ .

In the proof, we consider the essential multi-layer architecture of the deep linear network, and do not transform it into a linear regression model and directly apply existing results (see Loh & Wainwright (2015) and Negahban et al. (2011)). This is because we care more about deep ReLU networks which cannot be reduced in this way. Our proof technique is more suitable for analyzing the multi-layer neural networks which paves a way for analyzing deep ReLU networks. Also such an analysis technique can reveal the role of network parameters (dimension, norm, etc.) of each weight matrix in the results which may benefit the design of networks. Besides, the obtained results are more consistent with those for deep nonlinear networks (see Sec. 5).

#### 4.3 UNIFORM CONVERGENCE, STABILITY AND GENERALIZATION OF EMPIRICAL RISK

Based on the above results, we can derive the uniform convergence of empirical risk to population risk easily. In this subsection, we first give the uniform convergence rate of empirical risk for deep linear neural networks in Theorem 3, and then use this result to derive the stability and generalization bounds for DNNs in Corollary 1.

**Theorem 3.** *Suppose Assumption 1 on the input datum  $\mathbf{x}$  holds and the activation functions in a deep neural network are linear. Then there exist two universal constants  $c_{f'}$  and  $c_f$  such that if  $n \geq c_{f'} \max(l^3 r_x^4 / (d_l s \log(d/l) \varepsilon^2 \tau^4 \log(1/\varepsilon)), s \log(d/l) / (\tau^2 d_l))$ , then*

$$\sup_{\mathbf{w} \in \Omega} |\hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w})| \leq \epsilon_f \triangleq c_f \tau \max\left(\sqrt{d_l} \tau r^{2l}, r^l\right) \sqrt{\frac{s \log(dn/l) + \log(8/\varepsilon)}{n}} \quad (2)$$

*holds with probability at least  $1 - \varepsilon$ . Here  $l$  is the number of layers in the neural network,  $n$  is the sample size and  $d_l$  is the dimension of the final layer.*

From Theorem 3, when  $n \rightarrow +\infty$ , we have  $|\hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w})| \rightarrow 0$ . According to the definition of uniform convergence (Vapnik & Vapnik, 1998; Shalev-Shwartz et al., 2010), under the distribution  $\mathcal{D}$ , the empirical risk of a deep linear neural network converges to its population risk *uniformly* at the rate of  $\mathcal{O}(1/\sqrt{n})$ . Theorem 3 also explains the roles of the depth  $l$ , the network size  $d$ , and the number of nonzero weight parameters  $s$  in a DNN model.

Based on VC-dimension techniques, Bartlett & Maass (2003) proved that for a feedforward neural network with polynomial activation functions and one-dimensional output, with probability at least  $1 - \varepsilon$  the convergence bound satisfies  $|\hat{\mathbf{J}}_n(\mathbf{w}) - \inf_f \mathbf{J}(\mathbf{w})| \leq \mathcal{O}(\sqrt{(\gamma \log^2(n) + \log(1/\varepsilon))/n})$ . Here  $\gamma$  is the shattered parameter and can be as large as the VC-dimension of the network model, *i.e.* at the order of  $\mathcal{O}(ld \log(d) + l^2 d)$  (Bartlett & Maass, 2003). Note that Bartlett & Maass (2003) did not reveal the role of the magnitude of weight in their results. In contrast, our uniform convergence bound is  $\sup_{\mathbf{w} \in \Omega} |\hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w})| \leq \mathcal{O}(\sqrt{(s \log(dn/l) + \log(1/\varepsilon))/n})$ . So our convergence rate is tighter.

Neyshabur et al. (2015) proved that the Rademacher complexity of a fully-connected neural network model with ReLU activation functions and one-dimensional output is  $\mathcal{O}(r^l/\sqrt{n})$  (see Corollary 2 in (Neyshabur et al., 2015)). Then by applying Rademacher complexity based argument (Shalev-Shwartz & Ben-David, 2014a), we have  $|\sup_f(\hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}))| \leq \mathcal{O}((r^l + \sqrt{\log(1/\varepsilon)})/\sqrt{n})$  with probability at least  $1 - \varepsilon$  where the loss function is the training error  $g = \mathbf{1}_{(v^{(l)} \neq \mathbf{y})}$  in which  $v^{(l)}$  is the output of the  $l$ -th layer in the network model  $f(\mathbf{w}; \mathbf{x}, \mathbf{y})$ . The convergence rate in our theorem is  $\mathcal{O}(r^{2l} \sqrt{(s \log(d/l) + \log(1/\varepsilon))/n})$  and has the same convergence speed  $\mathcal{O}(1/\sqrt{n})$  w.r.t. sample number  $n$ . Note that our convergence rate involves  $r^{2l}$  since we use squared loss instead of the training error in (Neyshabur et al., 2015). The extra parameters  $s$  and  $d$  are involved since we consider the parameter space rather than the function hypothesis  $f$  in (Neyshabur et al., 2015), which helps people more transparently understand the roles of the network parameters. Besides, the Rademacher complexity cannot be applied to analyzing convergence properties of the empirical risk gradient and stationary points as our techniques.

Based on Theorem 3, we proceed to analyze the stability property of the empirical risk and the convergence rate of the generalization error in expectation. Let  $\mathcal{S} = \{(\mathbf{x}_{(1)}, \mathbf{y}_{(1)}), \dots, (\mathbf{x}_{(n)}, \mathbf{y}_{(n)})\}$  denote the sample set in which the samples are *i.i.d.* drawn from  $\mathcal{D}$ . When the optimal solution  $\mathbf{w}^n$  to problem (1) is computed by deterministic algorithms, the generalization error is defined as  $\epsilon_g = \hat{\mathbf{J}}_n(\mathbf{w}^n) - \mathbf{J}(\mathbf{w}^n)$ . But one usually employs randomized algorithms, *e.g.* stochastic gradient descent (SGD), for computing  $\mathbf{w}^n$ . In this case, stability and generalization error in expectation defined in Definition 3 are more applicable.

**Definition 3. (Stability and generalization in expectation)** (Vapnik & Vapnik, 1998; Shalev-Shwartz et al., 2010; Gonen & Shalev-Shwartz, 2017) Assume a randomized algorithm  $\mathbf{A}$  is employed,  $((\mathbf{x}'_{(1)}, \mathbf{y}'_{(1)}), \dots, (\mathbf{x}'_{(n)}, \mathbf{y}'_{(n)})) \sim \mathcal{D}$  and  $\mathbf{w}^n = \operatorname{argmin}_{\mathbf{w}} \hat{\mathbf{J}}_n(\mathbf{w})$  is the empirical risk minimizer (ERM). For every  $j \in [n]$ , suppose  $\mathbf{w}_*^j = \operatorname{argmin}_{\mathbf{w}} \frac{1}{n-1} \sum_{i \neq j} f_i(\mathbf{w}; \mathbf{x}_{(i)}, \mathbf{y}_{(i)})$ . We say that the ERM is on average stable with stability rate  $\epsilon_k$  under distribution  $\mathcal{D}$  if  $\left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}, (\mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) \sim \mathcal{D}} \frac{1}{n} \sum_{j=1}^n [f_j(\mathbf{w}_*^j; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) - f_j(\mathbf{w}^n; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)})] \right| \leq \epsilon_k$ . The ERM is said to have generalization error with convergence rate  $\epsilon_{k'}$  under distribution  $\mathcal{D}$  if we have  $\left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}} (\mathbf{J}(\mathbf{w}^n) - \hat{\mathbf{J}}_n(\mathbf{w}^n)) \right| \leq \epsilon_{k'}$ .

Stability measures the sensibility of the empirical risk to the input and generalization error measures the effectiveness of ERM on new data. Generalization error in expectation is especially important for applying DNNs considering their internal randomness, *e.g.* from SGD optimization. Now we present the results on stability and generalization performance of deep linear neural networks.

**Corollary 1.** Suppose Assumption 1 on the input datum  $\mathbf{x}$  holds and the activation functions in a deep neural network are linear. Then with probability at least  $1 - \varepsilon$ , both the stability rate and the generalization error rate of ERM of a deep linear neural network are at least  $\epsilon_f$ :

$$\left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}, (\mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) \sim \mathcal{D}} \frac{1}{n} \sum_{j=1}^n (f_j^* - f_j) \right| \leq \epsilon_f \quad \text{and} \quad \left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}} (\mathbf{J}(\mathbf{w}^n) - \hat{\mathbf{J}}_n(\mathbf{w}^n)) \right| \leq \epsilon_f,$$

where  $f_j^*$  and  $f_j$  respectively denote  $f_j(\mathbf{w}_*^j; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)})$  and  $f_j(\mathbf{w}^n; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)})$ , and  $\epsilon_f$  is defined in Eqn. (2).

According to Corollary 1, both the stability rate and the convergence rate of generalization error are  $\mathcal{O}(\epsilon_f)$ . This result indicates that deep learning empirical risk is stable and its output is robust to small perturbation over the training data. When  $n$  is sufficiently large, small generalization error of DNNs is guaranteed.

## 5 RESULTS FOR DEEP NONLINEAR NEURAL NETWORKS

In the above section, we analyze the empirical risk optimization landscape for deep linear neural network models. In this section, we extend our analysis to deep nonlinear neural networks which adopt the sigmoid activation function. Our analysis techniques are also applicable to other third-order

differentiable activation functions, *e.g.*, tanh function with different convergence rate. Here we assume the input data are *i.i.d.* Gaussian variables.

**Assumption 2.** *The input datum  $\mathbf{x}$  is a vector of i.i.d. Gaussian variables from  $\mathcal{N}(0, \tau^2)$ .*

Since for any input, the sigmoid function always maps it to the range  $[0, 1]$ . Thus, we do not require the input  $\mathbf{x}$  to have bounded magnitude. Such an assumption is common. For instance, Tian (2017) and Soudry & Hoffer (2017) both assumed that the entries in the input vector are from Gaussian distribution. We also assume  $\mathbf{w} \in \Omega$  as in (Xu & Mannor, 2012). Here we also assume that the entry value of the target output  $\mathbf{y}$  falls in  $[0, 1]$ . Similar to the analysis of deep linear neural networks, here we also aim to characterize the empirical risk gradient, stationary points and empirical risk for deep nonlinear neural networks.

## 5.1 UNIFORM CONVERGENCE OF GRADIENT AND STATIONARY POINTS

Here we analyze convergence properties of gradients of the empirical risk for deep nonlinear neural networks.

**Theorem 4.** *Assume the input sample  $\mathbf{x}$  obeys Assumption 2 and the activation functions in a deep neural network are sigmoid functions. Then the empirical gradient uniformly converges to the population gradient in Euclidean norm. Specifically, there are two universal constants  $c_y$  and  $c_{y'}$  such that if  $n \geq c_{y'} c_d l^3 r^2 / (s \log(d) \tau^2 \varepsilon^2 \log(1/\varepsilon))$  where  $c_d = \max_{0 \leq i \leq l} \mathbf{d}_i$ , then with probability at least  $1 - \varepsilon$*

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right\|_2 \leq \varepsilon_l \triangleq \tau \sqrt{\frac{512}{729} c_y l (l+2) (l c_r + 1) c_d c_r} \sqrt{\frac{s \log(dn/l) + \log(4/\varepsilon)}{n}},$$

where  $c_r = \max(r^2/16, (r^2/16)^{l-1})$ , and  $s$  denotes the nonzero entry number of all weights.

Similar to deep linear neural networks, the layer number  $l$ , width  $\mathbf{d}_i$ , number of nonzero parameter entries  $s$ , network size  $d$  and magnitude of weights are all critical to the convergence rate. Also, since there is a factor  $\max_i \mathbf{d}_i$  in the convergence rate, it is better to avoid choosing an extremely wide layer. Interestingly, when analyzing the representation ability of deep learning, Eldan & Shamir (2016) also suggested non-extreme-wide layers, though the conclusion was derived from a different perspective. By comparing Theorems 1 and 4, one can observe that there is a factor  $(1/16)^{l-1}$  in the convergence rate in Theorem 4. This is because the convergence rate accesses the Lipschitz constant and when we bound it, sigmoid activation function brings the factor  $1/16$  for each layer.

Now we analyze the non-degenerate stationary points of the empirical risk for deep nonlinear neural networks. Here we also assume that the population risk has  $m$  non-degenerate stationary points denoted by  $\{\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(m)}\}$ .

**Theorem 5.** *Assume the input sample  $\mathbf{x}$  obeys Assumption 2 and the activation functions in a deep neural network are sigmoid functions. Then if  $n \geq c_s \max(c_d l^3 r^2 / (s \log(d) \tau^2 \varepsilon^2 \log(1/\varepsilon)), s \log(d/l) / \zeta^2)$  where  $c_s$  is a constant, for  $k \in \{1, \dots, m\}$ , there exists a non-degenerate stationary point  $\mathbf{w}_n^{(k)}$  of  $\hat{\mathbf{J}}_n(\mathbf{w})$  which corresponds to the non-degenerate stationary point  $\mathbf{w}^{(k)}$  of  $\mathbf{J}(\mathbf{w})$  with probability at least  $1 - \varepsilon$ . Moreover,  $\mathbf{w}_n^{(k)}$  and  $\mathbf{w}^{(k)}$  have the same non-degenerate index and they obey*

$$\left\| \mathbf{w}_n^{(k)} - \mathbf{w}^{(k)} \right\|_2 \leq \frac{2\tau}{\zeta} \sqrt{\frac{512}{729} c_y l (l+2) (l c_r + 1) c_d c_r} \sqrt{\frac{s \log(dn/l) + \log(4/\varepsilon)}{n}}, \quad (k = 1, \dots, m)$$

with probability at least  $1 - \varepsilon$ , where  $c_y$ ,  $c_d$  and  $c_r$  are the same parameters in Theorem 4.

According to Theorem 5, there is one-to-one correspondence between the non-degenerate stationary points of  $\hat{\mathbf{J}}_n(\mathbf{w})$  and  $\mathbf{J}(\mathbf{w})$ . Also the corresponding pair has the same non-degenerate index, implying they have exactly matching local minima/maxima and non-degenerate saddle points. When  $n$  is sufficiently large, the non-degenerate stationary point  $\mathbf{w}_n^{(k)}$  in  $\hat{\mathbf{J}}_n(\mathbf{w})$  is very close to its corresponding non-degenerate stationary point  $\mathbf{w}^{(k)}$  in  $\mathbf{J}(\mathbf{w})$ . As for the degenerate stationary points, Theorem 4 guarantees the gradients at these points of  $\mathbf{J}(\mathbf{w})$  and  $\hat{\mathbf{J}}_n(\mathbf{w})$  are very close to each other.

## 5.2 UNIFORM CONVERGENCE, STABILITY AND GENERALIZATION OF EMPIRICAL RISK

Here we first give the uniform convergence analysis of the empirical risk and then analyze its stability and generalization.

**Theorem 6.** *Assume the input sample  $\mathbf{x}$  obeys Assumption 2 and the activation functions in a deep neural network are the sigmoid functions. If  $n \geq 18l^2\tau^2/(s \log(d)\tau^2\epsilon^2 \log(1/\epsilon))$ , then*

$$\sup_{\mathbf{w} \in \Omega} \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| \leq \epsilon_n \triangleq \tau \sqrt{\frac{9}{8} c_y c_d (1 + c_r(l-1))} \sqrt{\frac{s \log(nd/l) + \log(4/\epsilon)}{n}} \quad (3)$$

holds with probability at least  $1 - \epsilon$ , where  $c_y$ ,  $c_d$  and  $c_r$  are given in Theorem 4.

From Theorem 6, we obtain that under the distribution  $\mathcal{D}$ , the empirical risk of a deep nonlinear neural network converges at the rate of  $\mathcal{O}(1/\sqrt{n})$  (up to a log factor). Theorem 6 also gives similar results as Theorem 3, including the inclination of regularization penalty on weight and suggestion on non-extreme-wide layers. Similar to linear networks, our risk convergence rate is also tighter than the convergence rate on the networks with polynomial activation functions and one-dimensional output in (Bartlett & Maass, 2003) since ours is at the order of  $\mathcal{O}(\sqrt{(l-1)(s \log(dn/l) + \log(1/\epsilon))/n})$ , while the later is  $\mathcal{O}(\sqrt{(\gamma \log^2(n) + \log(1/\epsilon))/n})$  where  $\gamma$  is at the order of  $\mathcal{O}(ld \log(d) + l^2 d)$  (Bartlett & Maass, 2003).

We then establish the stability property and the generalization error of the empirical risk for nonlinear neural networks. By Theorem 6, we can obtain the following results.

**Corollary 2.** *Assume the input sample  $\mathbf{x}$  obeys Assumption 2 and the activation functions in a deep neural network are sigmoid functions. Then with probability at least  $1 - \epsilon$ , we have*

$$\left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}, (\mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) \sim \mathcal{D}} \frac{1}{n} \sum_{j=1}^n (f_j^* - f_j) \right| \leq \epsilon_n \quad \text{and} \quad \left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}} \left( \mathbf{J}(\mathbf{w}^n) - \hat{\mathbf{J}}_n(\mathbf{w}^n) \right) \right| \leq \epsilon_n,$$

where  $\epsilon_n$  is defined in Eqn. (3). The notations  $f_j^*$  and  $f_j$  here are the same in Corollary 1.

By Corollary 2, we know that both the stability convergence rate and the convergence rate of generalization error are  $\mathcal{O}(1/\sqrt{n})$ . This result accords with Theorems 8 and 9 in (Shalev-Shwartz et al., 2010) which implies  $\mathcal{O}(1/\sqrt{n})$  is the bottleneck of the stability and generalization convergence rate for generic learning algorithms. From this result, we have that if  $n$  is sufficiently large, the empirical risk can be expected to be very stable. This also dispels misgivings of the random selection of training samples in practice. Such a result indicates that the deep nonlinear neural network can offer good performance on testing data if it achieves small training error.

## 6 PROOF ROADMAP

Here we briefly introduce our proof roadmap. Due to space limitation, all the proofs of Theorems 1 ~ 6 and Corollaries 1 and 2 as well as technical lemmas are deferred to the supplementary material.

The proofs of Theorems 1 and 4 are similar but essentially differ in some techniques for bounding probability due to their different assumptions. For explanation simplicity, we define four events:  $\mathbf{E} = \{\sup_{\mathbf{w} \in \Omega} \|\nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w})\|_2 > t\}$ ,  $\mathbf{E}_1 = \{\sup_{\mathbf{w} \in \Omega} \|\frac{1}{n} \sum_{i=1}^n (\nabla f(\mathbf{w}, \mathbf{x}_{(i)}) - \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}))\|_2 > t/3\}$ ,  $\mathbf{E}_2 = \{\sup_{\mathbf{w}_{k_w} \in \mathcal{N}_i, i \in [l]} \|\frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E} \nabla f(\mathbf{w}_{k_w}, \mathbf{x})\|_2 > t/3\}$ , and  $\mathbf{E}_3 = \{\sup_{\mathbf{w} \in \Omega} \|\mathbb{E} \nabla f(\mathbf{w}_{k_w}, \mathbf{x}) - \mathbb{E} \nabla f(\mathbf{w}, \mathbf{x})\|_2 > t/3\}$ , where  $\mathbf{w}_{k_w} = [\mathbf{w}_{k_w}^1; \mathbf{w}_{k_w}^2; \dots; \mathbf{w}_{k_w}^l]$  is constructed by selecting  $\mathbf{w}_{k_w}^i \in \mathbb{R}^{d_i d_{i-1}}$  from  $d_i d_{i-1} \epsilon/d$ -net  $\mathcal{N}_i$  such that  $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$ . Note that in Theorems 1 and 4,  $t$  is respectively set to  $\epsilon_g$  and  $\epsilon_l$ . Then we have  $\mathbb{P}(\mathbf{E}) \leq \mathbb{P}(\mathbf{E}_1) + \mathbb{P}(\mathbf{E}_2) + \mathbb{P}(\mathbf{E}_3)$ . So we only need to separately bound  $\mathbb{P}(\mathbf{E}_1)$ ,  $\mathbb{P}(\mathbf{E}_2)$  and  $\mathbb{P}(\mathbf{E}_3)$ . For  $\mathbb{P}(\mathbf{E}_1)$  and  $\mathbb{P}(\mathbf{E}_3)$ , we use the gradient Lipschitz constant and the properties of  $\epsilon$ -net to prove  $\mathbb{P}(\mathbf{E}_1) \leq \epsilon/2$  and  $\mathbb{P}(\mathbf{E}_3) = 0$ , while bounding  $\mathbb{P}(\mathbf{E}_2)$  needs more efforts. Here based on the assumptions, we prove that  $\mathbb{P}(\mathbf{E}_2)$  has sub-exponential tail associated to the sample number  $n$  and the networks parameters, and it satisfies  $\mathbb{P}(\mathbf{E}_2) \leq \epsilon/2$  with proper conditions. Finally, combining the bounds of the three terms, we obtain the desired results.

To prove Theorems 2 and 5, we first prove the uniform convergence of the empirical Hessian to its population Hessian. Then, we define such a set  $D = \{\mathbf{w} \in \Omega : \|\nabla \mathbf{J}(\mathbf{w})\|_2 < \epsilon \text{ and } \inf_i |\lambda_i(\nabla^2 \mathbf{J}(\mathbf{w}))| \geq \zeta\}$ . In this way,  $D$  can be decomposed into countably components, with each component containing either exactly one or zero non-degenerate stationary point. For each component, the uniform convergence of gradient and the results in differential topology guarantee that if  $\mathbf{J}(\mathbf{w})$  has no stationary points, then  $\hat{\mathbf{J}}_n(\mathbf{w})$  also has no stationary points and vice versa. Similarly, for each component, the uniform convergence of Hessian and the results in differential topology guarantee that if  $\mathbf{J}(\mathbf{w})$  has a unique non-degenerate stationary point, then  $\hat{\mathbf{J}}_n(\mathbf{w})$  also has a unique non-degenerate stationary point with the same index. After establishing exact correspondence between the non-degenerate stationary points of empirical risk and population risk, we use the uniform convergence of gradient and Hessian to bound the distance between the corresponding pairs.

We adopt a similar strategy to prove Theorems 3 and 6. Specifically, we divide the event  $\sup_{\mathbf{w} \in \Omega} |\hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w})| > t$  into  $\mathbf{E}_1$ ,  $\mathbf{E}_2$  and  $\mathbf{E}_3$  which have the same forms as their counterparts in the proofs of Theorem 1 with the gradient replaced by the loss function. To prove  $\mathbb{P}(\mathbf{E}_1) \leq \epsilon/2$  and  $\mathbb{P}(\mathbf{E}_3) = 0$ , we can use the Lipschitz constant of the loss function and the  $\epsilon$ -net properties. The remaining is to prove  $\mathbb{P}(\mathbf{E}_2)$ . We also prove that it has sub-exponential tail associated to the sample number  $n$  and the networks parameters and it obeys  $\mathbb{P}(\mathbf{E}_2) \leq \epsilon/2$  with proper conditions. Then we utilize the uniform convergence of  $\hat{\mathbf{J}}_n(\mathbf{w})$  to prove the stability and generalization bounds of  $\hat{\mathbf{J}}_n(\mathbf{w})$  (i.e. Corollaries 1 and 2).

## 7 CONCLUSION

In this work, we provided theoretical analysis on the landscape of empirical risk optimization for deep linear/nonlinear neural networks with (stochastic-)gradient descent, including the properties of the gradient and stationary points of empirical risk as well as the uniform convergence, stability, and generalization of the empirical risk itself. To our best knowledge, most of the results are new to deep learning community. These results also reveal that the depth  $l$ , the nonzero entry number  $s$  of all weights, the network size  $d$  and the width of a network are critical to the convergence rates. We also prove that the weight parameter magnitude is important to the convergence rate. Indeed, small magnitude of the weights is suggested. All the results are consistent with the widely used network architectures in practice.

## ACKNOWLEDGMENT

This work is partially supported by National University of Singapore startup grant R-263-000-C08-133, Ministry of Education of Singapore AcRF Tier One grant R-263-000-C21-112, NUS IDS R-263-000-C67-646 and ECRA R-263-000-C87-133.

## REFERENCES

- R. Alessandro. Lecture notes of advanced statistical theory I, CMU. [http://www.stat.cmu.edu/~arinaldo/36755/F16/Scribed\\_Lectures/LEC0914.pdf](http://www.stat.cmu.edu/~arinaldo/36755/F16/Scribed_Lectures/LEC0914.pdf), 2016.
- B. Bakshi and G. Stephanopoulos. Wave-net: A multiresolution, hierarchical neural network with localized learning. *AIChE Journal*, 39(1):57–81, 1993.
- P. Bartlett and W. Maass. Vapnik-chervonenkis dimension of neural nets. *The handbook of brain theory and neural networks*, pp. 1188–1192, 2003.
- E. Baum. On the capabilities of multilayer perceptrons. *Journal of complexity*, 4(3):193–215, 1988.
- A. Choromanska, M. Henaff, M. Mathieu, G. Arous, and Y. LeCun. The loss surfaces of multilayer networks. In *AISTATS*, 2015a.
- A. Choromanska, Y. LeCun, and G. Arous. Open problem: The landscape of the loss surfaces of multilayer networks. In *COLT*, pp. 1756–1760, 2015b.
- R. Collobert and J. Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *ICML*, pp. 160–167, 2008.

- Y. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *NIPS*, pp. 2933–2941, 2014.
- B. Dubrovin, A. Fomenko, and S. Novikov. *Modern geometry—methods and applications: Part II: The geometry and topology of manifolds*, volume 104. Springer Science & Business Media, 2012.
- R. Eldan and O. Shamir. The power of depth for feedforward neural networks. In *COLT*, pp. 907–940, 2016.
- Y. Fyodorov and I. Williams. Replica symmetry breaking condition exposed by random matrix calculation of landscape complexity. *Journal of Statistical Physics*, 129(5-6):1081–1116, 2007.
- R. Ge, F. Huang, C. Jin, and Y. Yuan. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *COLT*, pp. 797–842, 2015.
- A. Gonen and S. Shalev-Shwartz. Fast rates for empirical risk minimization of strict saddle problems. *COLT*, 2017.
- A. Graves, A. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In *ICASSP*, pp. 6645–6649, 2013.
- D. Gromoll and W. Meyer. On differentiable functions with isolated critical points. *Topology*, 8(4):361–369, 1969.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pp. 770–778, 2016.
- G. Hinton, S. Osindero, and Y. Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.
- G. Hinton, L. Deng, D. Yu, G. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. Sainath, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6):82–97, 2012.
- K. Kawaguchi. Deep learning without poor local minima. In *NIPS*, pp. 1097–1105, 2016.
- P. Loh and M. J. Wainwright. Regularized m-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *JMLR*, 16(Mar):559–616, 2015.
- S. Mei, Y. Bai, and A. Montanari. The landscape of empirical risk for non-convex losses. *Annals of Statistics*, 2017.
- S. Negahban, B. Yu, M. Wainwright, and P. Ravikumar. A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. In *NIPS*, pp. 1348–1356, 2009.
- S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers. In *NIPS*, 2011.
- B. Neyshabur, R. Tomioka, and N. Srebro. Norm-based capacity control in neural networks. In *COLT*, pp. 1376–1401, 2015.
- Q. Nguyen and M. Hein. The loss surface of deep and wide neural networks. In *ICML*, 2017.
- P. Rigollet. Statistic s997 lecture notes, MIT mathematics. *MIT OpenCourseWare*, pp. 23–24, 2015.
- M. Rudelson and R. Vershynin. Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18(82):1–9, 2013.
- S. Shalev-Shwartz and S. Ben-David. Understanding machine learning: From theory to algorithms. *Cambridge Univ. Press, Cambridge*, pp. 375–382, 2014a.
- S. Shalev-Shwartz and S. Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014b.
- S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Learnability, stability and uniform convergence. *JMLR*, 11:2635–2670, 2010.
- S. Shalev-Shwartz, O. Shamir, and S. Shammah. Failures of deep learning. *ICML*, 2017.
- D. Soudry and Y. Carmon. No bad local minima: Data independent training error guarantees for multilayer neural networks. *arXiv preprint arXiv:1605.08361*, 2016.

- D. Soudry and E. Hoffer. Exponentially vanishing sub-optimal local minima in multilayer neural networks. *arXiv preprint arXiv:1702.05777*, 2017.
- C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, pp. 1–9, 2015.
- Y. Tian. An analytical formula of population gradient for two-layered relu network and its applications in convergence and critical point analysis. *ICML*, 2017.
- V. N. Vapnik and V. Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.
- R. Vershynin. Introduction to the non-asymptotic analysis of random matrices, compressed sensing. *Cambridge Univ. Press, Cambridge*, pp. 210–268, 2012.
- H. Xu and S. Mannor. Robustness and generalization. *Machine Learning*, 86(3):391–423, 2012.
- Y. Zhang, J. Lee, and M. Jordan.  $\ell_1$ -regularized neural networks are improperly learnable in polynomial time. In *ICML*, pp. 993–1001, 2016.
- Y. Zhang, P. Liang, and M. Wainwright. Convexified convolutional neural networks. *ICML*, 2017.

## SUPPLEMENTARY MATERIAL OF EMPIRICAL RISK LANDSCAPE ANALYSIS FOR UNDERSTANDING DEEP NEURAL NETWORKS

### A STRUCTURE OF THIS DOCUMENT

This document gives some other necessary notations and preliminaries for our analysis in Sec. B. Then we prove Theorems 1~3 and Corollary 1 for deep linear neural networks in Sec. C. Then we present the proofs of Theorems 4~6 and Corollary 2 for deep nonlinear neural networks in Sec. D.

In both Sec. C and D, we first present the technical lemmas for proving our final results and subsequently present the proofs of these lemmas. Then we utilize these technical lemmas to prove our desired results. Finally, we give the proofs of other auxiliary lemmas.

### 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.  $\mathbf{A} \otimes \mathbf{C}$  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  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ , its operation norm is computed as

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

where  $\mathcal{A}_{ijk}$  denotes the  $(i, j, k)$ -th entry of  $\mathcal{A}$ . Also we denote the vectorization of  $\mathbf{W}^{(j)}$  (the weight matrix of the  $j$ -th layer) as

$$\mathbf{w}_{(j)} = \text{vec}(\mathbf{W}^{(j)}) \in \mathbb{R}^{d_j d_{j-1}}.$$

We denote  $\mathbf{I}_k$  as the identity matrix of size  $k \times k$ .

For notational simplicity, we further define  $\mathbf{e} \triangleq \mathbf{v}^{(l)} - \mathbf{y}$  as the output error vector. Then the squared loss is defined as  $f(\mathbf{w}; \mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{e}\|_2^2$ , where  $\mathbf{w} = (\mathbf{w}_{(1)}; \dots; \mathbf{w}_{(l)}) \in \mathbb{R}^d$  contains all the weight parameters.

#### B.2 TECHNICAL LEMMAS

We first introduce Lemmas 1 and 2 which are respectively used for bounding the  $\ell_2$ -norm of a vector and the operation norm of a matrix. Then we introduce Lemmas 3 and 4 which discuss the topology of functions. In Lemma 5, we give the relationship between the stability and generalization of empirical risk.

**Lemma 1.** (Vershynin, 2012) For any vector  $\mathbf{x} \in \mathbb{R}^d$ , its  $\ell_2$ -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  $\epsilon$ -covering net of  $\mathbb{B}^d(1)$ .

**Lemma 2.** (Vershynin, 2012) For any symmetric matrix  $\mathbf{X} \in \mathbb{R}^{d \times d}$ , its operator norm can be bounded 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  $\boldsymbol{\lambda}_\epsilon = \{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k_w}\}$  be an  $\epsilon$ -covering net of  $\mathbb{B}^d(1)$ .

**Lemma 3.** (Mei et al., 2017) Let  $D \subseteq \mathbb{R}^d$  be a compact set with a  $C^2$  boundary  $\partial D$ , 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 4.** (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

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

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

**Lemma 5.** (Shalev-Shwartz & Ben-David, 2014b; Gonen & Shalev-Shwartz, 2017) Assume that  $\mathcal{D}$  is a sample distribution and randomized algorithm  $\mathbf{A}$  is employed for optimization. Suppose that  $((\mathbf{x}'_{(1)}, \mathbf{y}'_{(1)}), \dots, (\mathbf{x}'_{(n)}, \mathbf{y}'_{(n)})) \sim \mathcal{D}$  and  $\mathbf{w}^n = \operatorname{argmin}_{\mathbf{w}} \hat{\mathbf{J}}_n(\mathbf{w})$ . For every  $j \in \{1, \dots, n\}$ , suppose  $\mathbf{w}_*^j = \operatorname{argmin}_{\mathbf{w}} \frac{1}{n-1} \sum_{i \neq j} f_i(\mathbf{w}; \mathbf{x}_{(i)}, \mathbf{y}_{(i)})$ . For arbitrary distribution  $\mathcal{D}$ , we have

$$\left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}, (\mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) \sim \mathcal{D}} \frac{1}{n} \sum_{j=1}^n (f_j^* - f_j) \right| = \left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}} \left( \mathbf{J}(\mathbf{w}^n) - \hat{\mathbf{J}}_n(\mathbf{w}^n) \right) \right|.$$

where  $f_j^*$  and  $f_j$  respectively denote  $f_j(\mathbf{w}_*^j; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)})$  and  $f_j(\mathbf{w}^n; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)})$ .

## C PROOFS FOR DEEP LINEAR NEURAL NETWORKS

In this section, we first present the technical lemmas in Sec. C.1 and then we give the proofs of these lemmas in Sec. C.2. Next, we utilize these lemmas to prove the results in Theorems 1~3 and Corollary 1 in Sec. C.3. Finally, we give the proofs of other lemmas in Sec. C.4.

### C.1 TECHNICAL LEMMAS

Here we present the technical lemmas for proving our desired results. For brevity, we also define  $\mathbf{B}_{j:s}$  as follows:

$$\mathbf{B}_{s:t} \triangleq \mathbf{W}^{(s)} \mathbf{W}^{(s-1)} \dots \mathbf{W}^{(t)} \in \mathbb{R}^{d_s \times d_{t-1}}, \quad (s \geq t); \quad \mathbf{B}_{s:t} \triangleq \mathbf{I}, \quad (s < t). \quad (4)$$

**Lemma 6.** Assume that the activation functions in the deep neural network  $f(\mathbf{w}, \mathbf{x})$  are linear functions. Then the gradient of  $f(\mathbf{w}, \mathbf{x})$  with respect to  $\mathbf{w}_{(j)}$  can be written as

$$\nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}, \mathbf{x}) = ((\mathbf{B}_{j-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:j+1}^T) \mathbf{e}, \quad (j = 1, \dots, l),$$

where  $\otimes$  denotes the Kronecke product. Then we can compute the Hessian matrix as follows:

$$\nabla^2 f(\mathbf{w}, \mathbf{x}) = \begin{bmatrix} \nabla_{\mathbf{w}_{(1)}} (\nabla_{\mathbf{w}_{(1)}} f(\mathbf{w}, \mathbf{x})) & \cdots & \nabla_{\mathbf{w}_{(1)}} (\nabla_{\mathbf{w}_{(l)}} f(\mathbf{w}, \mathbf{x})) \\ \nabla_{\mathbf{w}_{(2)}} (\nabla_{\mathbf{w}_{(1)}} f(\mathbf{w}, \mathbf{x})) & \cdots & \nabla_{\mathbf{w}_{(2)}} (\nabla_{\mathbf{w}_{(l)}} f(\mathbf{w}, \mathbf{x})) \\ \vdots & \ddots & \vdots \\ \nabla_{\mathbf{w}_{(l)}} (\nabla_{\mathbf{w}_{(1)}} f(\mathbf{w}, \mathbf{x})) & \cdots & \nabla_{\mathbf{w}_{(l)}} (\nabla_{\mathbf{w}_{(l)}} f(\mathbf{w}, \mathbf{x})) \end{bmatrix},$$

where  $\mathbf{Q}_{st} \triangleq \nabla_{\mathbf{w}_{(s)}} (\nabla_{\mathbf{w}_{(t)}} f(\mathbf{w}, \mathbf{x}))$  is defined as

$$\mathbf{Q}_{st} = \begin{cases} (\mathbf{B}_{t-1:s+1}^T) \otimes (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:t+1}^T) + (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1}), & \text{if } s < t, \\ (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{s-1:1}) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:s+1}), & \text{if } s = t, \\ (\mathbf{B}_{l:s+1}^T \mathbf{e} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes \mathbf{B}_{s-1:t+1} + (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1}), & \text{if } s > t. \end{cases}$$

**Lemma 7.** Suppose Assumption 1 on the input data  $\mathbf{x}$  holds and the activation functions in deep neural network are linear functions. Then for any  $t > 0$ , the objective  $f(\mathbf{w}, \mathbf{x})$  obeys

$$\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n(f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{x}_{(i)}))) > t\right) \leq 2 \exp\left(-c_{f'}n \min\left(\frac{t^2}{\omega_f^2 \max(\mathbf{d}_l \omega_f^2 \tau^4, \tau^2)}, \frac{t}{\omega_f^2 \tau^2}\right)\right),$$

where  $c_{f'}$  is a positive constant and  $\omega_f = r^l$ .

**Lemma 8.** Suppose Assumption 1 on the input data  $\mathbf{x}$  holds and the activation functions in deep neural network are linear functions. 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

$$\begin{aligned} \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n(\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E}\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}_{(i)}) \rangle) > t\right) \\ \leq 3 \exp\left(-c_{g'}n \min\left(\frac{t^2}{l \max(\omega_g \tau^2, \omega_g \tau^4, \omega_{g'} \tau^2)}, \frac{t}{\sqrt{l\omega_g} \max(\tau, \tau^2)}\right)\right), \end{aligned}$$

where  $c_{g'}$  is a constant;  $\omega_g = c_q r^{2(2l-1)}$  and  $\omega_{g'} = c_q r^{2(l-1)}$  in which  $c_q = \sqrt{\max_{0 \leq i \leq l} \mathbf{d}_i}$ .

**Lemma 9.** Suppose Assumption 1 on the input data  $\mathbf{x}$  holds and the activation functions in deep neural network are linear functions. Then for any  $t > 0$  and arbitrary unit vector  $\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$ , the Hessian  $\nabla^2 f(\mathbf{w}, \mathbf{x})$  obeys

$$\begin{aligned} \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n(\langle \boldsymbol{\lambda}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E}\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x}_{(i)}))\boldsymbol{\lambda} \rangle) > t\right) \\ \leq 5 \exp\left(-c_{h'}n \min\left(\frac{t^2}{\tau^2 l^2 \max(\omega_g, \omega_g \tau^2, \omega_h)}, \frac{t}{\sqrt{\omega_g} l \max(\tau, \tau^2)}\right)\right), \end{aligned}$$

where  $\omega_g = r^{4(l-1)}$  and  $\omega_h = r^{2(l-2)}$ .

**Lemma 10.** Suppose the activation functions in deep neural network are linear functions. Then for any  $\mathbf{w} \in \mathbb{B}^d(r)$  and  $\mathbf{x} \in \mathbb{B}^{d_0}(r_x)$ , we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})\|_2 \leq \sqrt{\alpha_g}, \quad \text{where } \alpha_g = c_l r_x^4 r^{4l-2}.$$

in which  $c_l$  is a constant. Further, for any  $\mathbf{w} \in \mathbb{B}^d(r)$  and  $\mathbf{x} \in \mathbb{B}^{d_0}(r_x)$ , we also have

$$\|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_{op} \leq \|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_F \leq l\sqrt{\alpha_l}, \quad \text{where } \alpha_l \triangleq c_{l'} r_x^4 r^{4l-2}.$$

in which  $c_{l'}$  is a constant. With the same condition, we can bound the operation norm of  $\nabla^3 f(\mathbf{w}, \mathbf{x})$ . That is, there exists a universal constant  $\alpha_p$  such that  $\|\nabla^3 f(\mathbf{w}, \mathbf{x})\|_{op} \leq \alpha_p$ .

**Lemma 11.** Suppose Assumption 1 on the input data  $\mathbf{x}$  holds and the activation functions in deep neural network are linear functions. Then there exist two universal constant  $c_g$  and  $c_h$  such that the sample Hessian converges uniformly to the population Hessian in operator norm. Specifically, there exist two universal constants  $c_{h_1}$  and  $c_{h_2}$  such that if  $n \geq c_{h_2} \max(\frac{\alpha_p^2 r^2}{\tau^2 l^2 \omega_h^2 \varepsilon^2 s \log(d/l)}, s \log(d/l)/(l\tau^2))$ , then

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w}) \right\|_{op} \leq c_{h_1} \tau l \omega_h \sqrt{\frac{d \log(nl) + \log(20/\varepsilon)}{n}}$$

holds with probability at least  $1 - \varepsilon$ , where  $\omega_h = \max(\tau r^{2(l-1)}, r^{2(l-2)}, r^{l-2})$ .

## C.2 PROOFS OF TECHNICAL LEMMAS

To prove the above lemmas, we first introduce some useful results.

**Lemma 12.** (Rudelson & Vershynin, 2013) Assume that  $\mathbf{x} = (\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_k) \in \mathbb{R}^k$  is a random vector with independent components  $x_i$  which have zero mean and are independent  $\tau_i^2$ -sub-Gaussian variables. Here  $\max_i \tau_i^2 \leq \tau^2$ . Let  $\mathbf{A}$  be an  $k \times k$  matrix. Then we have

$$\mathbb{E} \exp\left(\lambda \left(\sum_{i,j:i \neq j} \mathbf{A}_{ij} x_i x_j - \mathbb{E}\left(\sum_{i,j:i \neq j} \mathbf{A}_{ij} x_i x_j\right)\right)\right) \leq \exp(2\tau^2 \lambda^2 \|\mathbf{A}\|_F^2), \quad |\lambda| \leq 1/(2\tau \|\mathbf{A}\|_2).$$

**Lemma 13.** Assume that  $\mathbf{x} = (\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_k) \in \mathbb{R}^k$  is a random vector with independent components  $x_i$  which have zero mean and are independent  $\tau_i^2$ -sub-Gaussian variables. Here  $\max_i \tau_i^2 \leq \tau^2$ . Let  $\mathbf{a}$  be an  $n$ -dimensional vector. Then we have

$$\mathbb{E} \exp \left( \lambda \left( \sum_{i=1}^k \mathbf{a}_i x_i^2 - \mathbb{E} \left( \sum_{i=1}^k \mathbf{a}_i x_i^2 \right) \right) \right) \leq \mathbb{E} \exp \left( 128 \lambda^2 \tau^4 \left( \sum_{i=1}^k \mathbf{a}_i^2 \right) \right), \quad |\lambda| \leq \frac{1}{\tau^2 \max_i \mathbf{a}_i}.$$

**Lemma 14.** For  $\mathbf{B}_{j:t}$  defined in Eqn. (4), we have the following properties:

$$\|\mathbf{B}_{s:t}\|_{op} \leq \|\mathbf{B}_{s:t}\|_F \leq \omega_r \quad \text{and} \quad \|\mathbf{B}_{l:1}\|_{op} \leq \|\mathbf{B}_{l:1}\|_F \leq \omega_f,$$

where  $\omega_r = r^{s-t+1} \leq \max(r, r^l)$  and  $\omega_f = r^l$ .

Lemma 13 is useful for bounding probability. The two inequalities in Lemma 14 can be obtained by using  $\|\mathbf{w}_{(j)}\|_2 \leq r$  ( $\forall j = 1, \dots, l$ ). We defer the proofs of Lemmas 13 and 14 to Sec. C.4.2.

### C.2.1 PROOF OF LEMMA 6

*Proof.* When the activation functions are linear functions, we can easily compute the gradient of  $f(\mathbf{w}, \mathbf{x})$  with respect to  $\mathbf{w}_{(j)}$ :

$$\nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}, \mathbf{x}) = ((\mathbf{B}_{j-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:j+1}^T) \mathbf{e}, \quad (j = 1, \dots, l),$$

where  $\otimes$  denotes the Kronecker product. Now we consider the computation of the Hessian matrix. For brevity, let  $\mathbf{Q}_s = ((\mathbf{B}_{s-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:s+1}^T)$ . Then we can compute  $\nabla_{\mathbf{w}_{(s)}}^2 f(\mathbf{w}, \mathbf{x})$  as follows:

$$\begin{aligned} \nabla_{\mathbf{w}_{(s)}}^2 f(\mathbf{w}, \mathbf{x}) &= \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(s)}^T \partial \mathbf{w}_{(s)}} = \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(s)}^T \partial \mathbf{w}_{(s)}} = \frac{\partial(\mathbf{Q}_s \mathbf{e})}{\partial \mathbf{w}_{(s)}^T} = \frac{\partial_{\text{vec}}(\mathbf{Q}_s \mathbf{e})}{\partial \mathbf{w}_{(s)}^T} \\ &= \frac{\partial_{\text{vec}}(\mathbf{Q}_s \mathbf{B}_{l:s+1} \mathbf{W}^{(t)} \mathbf{B}_{s-1:1} \mathbf{x})}{\partial \mathbf{w}_{(s)}^T} \\ &= \frac{\partial((\mathbf{B}_{s-1:1} \mathbf{x})^T \otimes (\mathbf{Q}_s \mathbf{B}_{l:s+1})) \text{vec}(\mathbf{W}^{(s)})}{\partial \mathbf{w}_{(s)}^T} \\ &= (\mathbf{B}_{s-1:1} \mathbf{x})^T \otimes ((\mathbf{B}_{s-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:s+1}^T) \mathbf{B}_{l:s+1} \\ &\stackrel{\textcircled{1}}{=} (\mathbf{B}_{s-1:1} \mathbf{x})^T \otimes ((\mathbf{B}_{s-1:1} \mathbf{x}) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:s+1})) \\ &\stackrel{\textcircled{2}}{=} ((\mathbf{B}_{s-1:1} \mathbf{x})^T \otimes (\mathbf{B}_{s-1:1} \mathbf{x})) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:s+1}) \\ &\stackrel{\textcircled{3}}{=} ((\mathbf{B}_{s-1:1} \mathbf{x})(\mathbf{B}_{s-1:1} \mathbf{x})^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:s+1}), \end{aligned}$$

where  $\textcircled{1}$  holds since  $\mathbf{B}_{j-1:1} \mathbf{x}$  is a vector and for any vector  $\mathbf{x}$ , we have  $(\mathbf{x} \otimes \mathbf{A})\mathbf{B} = \mathbf{x} \otimes (\mathbf{A}\mathbf{B})$ .  $\textcircled{2}$  holds because for any four matrices  $\mathbf{Z}_1 \sim \mathbf{Z}_3$  of proper sizes, we have  $(\mathbf{Z}_1 \otimes \mathbf{Z}_2) \otimes \mathbf{Z}_3 = \mathbf{Z}_1 \otimes (\mathbf{Z}_2 \otimes \mathbf{Z}_3)$ .  $\textcircled{3}$  holds because for any two matrices  $\mathbf{z}_1, \mathbf{z}_2$  of proper sizes, we have  $\mathbf{z}_1 \mathbf{z}_2^T = \mathbf{z}_1 \otimes \mathbf{z}_2^T = \mathbf{z}_2^T \otimes \mathbf{z}_1$ .

Then, we consider the case  $s > t$ :

$$\begin{aligned} \nabla_{\mathbf{w}_{(t)}} (\nabla_{\mathbf{w}_{(s)}} f(\mathbf{w}, \mathbf{x})) &= \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(t)}^T \partial \mathbf{w}_{(s)}} = \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(t)}^T \partial \mathbf{w}_{(s)}} = \frac{\partial(\mathbf{Q}_s \mathbf{e})}{\partial \mathbf{w}_{(t)}^T} = \frac{\partial_{\text{vec}}(\mathbf{Q}_s \mathbf{e})}{\partial \mathbf{w}_{(t)}^T} \\ &= \frac{\partial_{\text{vec}}(\mathbf{Q}_s \mathbf{B}_{l:t+1} \mathbf{W}^{(t)} \mathbf{B}_{t-1:1} \mathbf{x})}{\partial \mathbf{w}_{(t)}^T} + \frac{\partial_{\text{vec}}(((\mathbf{B}_{s-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:s+1}^T) \mathbf{e})}{\partial \mathbf{w}_{(t)}^T}. \end{aligned}$$

Notice, here we just think that  $\mathbf{Q}_s$  in the  $\frac{\partial_{\text{vec}}(\mathbf{Q}_s \mathbf{B}_{l:t+1} \mathbf{W}^{(t)} \mathbf{B}_{t-1:1} \mathbf{x})}{\partial \mathbf{w}_{(t)}^T}$  is a constant matrix and is not related to  $\mathbf{W}^{(t)}$ . Similarly, we also take  $\mathbf{e}$  in  $\frac{\partial_{\text{vec}}(((\mathbf{B}_{s-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:s+1}^T) \mathbf{e})}{\partial \mathbf{w}_{(t)}^T}$  as a constant vector. Since we have

$$\frac{\partial_{\text{vec}}(\mathbf{Q}_s \mathbf{B}_{l:t+1} \mathbf{W}^{(t)} \mathbf{B}_{t-1:1} \mathbf{x})}{\partial \mathbf{w}_{(t)}^T} = (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1}),$$

we only need to consider

$$\begin{aligned}
\frac{\partial_{\text{vec}} \left( (\mathbf{B}_{s-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:s+1}^T \mathbf{e} \right)}{\partial \mathbf{w}_{(t)}^T} &= \frac{\partial_{\text{vec}} \left( (\mathbf{B}_{s-1:1} \mathbf{x}) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{e}) \right)}{\partial \mathbf{w}_{(t)}^T} \\
&= \frac{\partial_{\text{vec}} \left( (\mathbf{B}_{s-1:1} \mathbf{x}) (\mathbf{B}_{l:s+1}^T \mathbf{e})^T \right)}{\partial \mathbf{w}_{(t)}^T} \\
&= \frac{\partial_{\text{vec}} \left( \mathbf{B}_{s-1:t+1} \mathbf{W}^{(t)} (\mathbf{B}_{t-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:s+1}) \right)}{\partial \mathbf{w}_t^T} \\
&= \frac{\partial \left( \mathbf{B}_{t-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:s+1} \right)^T \otimes \mathbf{B}_{s-1:t+1} \text{vec} \left( \mathbf{W}^{(t)} \right)}{\partial \mathbf{w}_t^T} \\
&= (\mathbf{B}_{t-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:s+1})^T \otimes \mathbf{B}_{s-1:t+1}.
\end{aligned}$$

Therefore, for  $s > t$ , by combining the above two terms, we can obtain

$$\nabla_{\mathbf{w}_{(t)}} \left( \nabla_{\mathbf{w}_{(s)}} f(\mathbf{w}, \mathbf{x}) \right) = (\mathbf{B}_{l:s+1}^T \mathbf{e} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes \mathbf{B}_{s-1:t+1} + (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1}).$$

Then, by similar method, we can compute the Hessian for the case  $s < t$  as follows:

$$\nabla_{\mathbf{w}_{(t)}} \left( \nabla_{\mathbf{w}_{(s)}} f(\mathbf{w}, \mathbf{x}) \right) = (\mathbf{B}_{t-1:s+1}^T) \otimes (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:t+1}^T) + (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1}).$$

The proof is completed.  $\square$

### C.2.2 PROOF OF LEMMA 7

*Proof.* We first prove that  $\mathbf{v}^{(l)}$ , which is defined in Eqn. (5), is sub-Gaussian.

$$\mathbf{v}^{(l)} = \mathbf{W}^{(l)} \dots \mathbf{W}^{(1)} \mathbf{x} = \mathbf{B}_{l:1} \mathbf{x}. \quad (5)$$

Then by the convexity in  $\lambda$  of  $\exp(\lambda t)$  and Lemma 14, we can obtain

$$\begin{aligned}
\mathbb{E} \left( \exp \left( \left\langle \boldsymbol{\lambda}, \mathbf{v}^{(l)} - \mathbb{E}(\mathbf{v}^{(l)}) \right\rangle \right) \right) &= \mathbb{E} \left( \exp \left( \left\langle \boldsymbol{\lambda}, \mathbf{B}_{l:1} \mathbf{x} - \mathbb{E} \mathbf{B}_{l:1} \mathbf{x} \right\rangle \right) \right) \\
&\leq \mathbb{E} \left( \exp \left( \left\langle \mathbf{B}_{l:1}^T \boldsymbol{\lambda}, \mathbf{x} \right\rangle \right) \right) \\
&\leq \exp \left( \frac{\|\mathbf{B}_{l:1}^T \boldsymbol{\lambda}\|_2^2 \tau^2}{2} \right) \\
&\stackrel{\textcircled{1}}{\leq} \exp \left( \frac{\omega_f^2 \tau^2 \|\boldsymbol{\lambda}\|_2^2}{2} \right),
\end{aligned} \quad (6)$$

where  $\textcircled{1}$  uses the conclusion that  $\|\mathbf{B}_{l:1}\|_{\text{op}} \leq \|\mathbf{B}_{l:1}\|_F \leq \omega_f$  in Lemma 14. This means that  $\mathbf{v}^{(l)}$  is centered and is  $\omega_f^2 \tau^2$ -sub-Gaussian. Accordingly, we can obtain that the  $k$ -th entry of  $\mathbf{v}^{(l)}$  is also  $z_k \tau^2$ -sub-Gaussian, where  $z_k$  is a universal positive constant. Note that  $\max_k z_k \leq \omega_f^2$ . Let  $\mathbf{v}_i^{(l)}$

denotes the output of the  $i$ -th sample  $\mathbf{x}_{(i)}$ . By Lemma 13, we have that for  $s > 0$ ,

$$\begin{aligned} \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n\left(\|\mathbf{v}_i^{(l)}\|_2^2 - \mathbb{E}\|\mathbf{v}_i^{(l)}\|_2^2\right) > \frac{t}{2}\right) &= \mathbb{P}\left(s\sum_{i=1}^n\left(\|\mathbf{v}_i^{(l)}\|_2^2 - \mathbb{E}\|\mathbf{v}_i^{(l)}\|_2^2\right) > \frac{nst}{2}\right) \\ &\stackrel{\textcircled{1}}{\leq} \exp\left(-\frac{snt}{2}\right)\mathbb{E}\left(s\sum_{i=1}^n\left(\|\mathbf{v}_i^{(l)}\|_2^2 - \mathbb{E}\|\mathbf{v}_i^{(l)}\|_2^2\right)\right) \\ &\stackrel{\textcircled{2}}{\leq} \exp\left(-\frac{snt}{2}\right)\prod_{i=1}^n\mathbb{E}\left(s\left(\|\mathbf{v}_i^{(l)}\|_2^2 - \mathbb{E}\|\mathbf{v}_i^{(l)}\|_2^2\right)\right) \\ &\stackrel{\textcircled{3}}{\leq} \exp\left(-\frac{snt}{2}\right)\prod_{i=1}^n\exp\left(128\mathbf{d}_l s^2\omega_f^4\tau^4\right) \quad |s| \leq \frac{1}{\omega_f^2\tau^2} \\ &\stackrel{\textcircled{4}}{\leq} \exp\left(-c'n\min\left(\frac{t^2}{\mathbf{d}_l\omega_f^4\tau^4}, \frac{t}{\omega_f^2\tau^2}\right)\right). \end{aligned}$$

Note that  $\textcircled{1}$  holds because of Chebyshev's inequality.  $\textcircled{2}$  holds since  $\mathbf{x}_{(i)}$  are independent.  $\textcircled{3}$  is established by applying Lemma 13. We have  $\textcircled{4}$  by optimizing  $s$ . Since  $\mathbf{v}^{(l)}$  is sub-Gaussian, we have

$$\begin{aligned} \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n\left(\mathbf{y}^T\mathbf{v}_i^{(l)} - \mathbb{E}\mathbf{y}^T\mathbf{v}_i^{(l)}\right) > \frac{t}{2}\right) &\leq \mathbb{P}\left(s\sum_{i=1}^n\left(\mathbf{y}^T\mathbf{v}_i^{(l)} - \mathbb{E}\mathbf{y}^T\mathbf{v}_i^{(l)}\right) > \frac{nst}{2}\right) \\ &\leq \exp\left(-\frac{nst}{2}\right)\mathbb{E}\exp\left(s\sum_{i=1}^n\left(\mathbf{y}^T\mathbf{v}_i^{(l)} - \mathbb{E}\mathbf{y}^T\mathbf{v}_i^{(l)}\right)\right) \\ &\leq \exp\left(-\frac{nst}{2}\right)\prod_{i=1}^n\mathbb{E}\exp\left(s\left(\mathbf{y}^T\mathbf{v}_i^{(l)} - \mathbb{E}\mathbf{y}^T\mathbf{v}_i^{(l)}\right)\right) \\ &\stackrel{\textcircled{1}}{\leq} \exp\left(-\frac{nst}{2}\right)\prod_{i=1}^n\exp\left(\frac{\omega_f^2\tau^2s^2\|\mathbf{y}\|_2^2}{2}\right) \\ &\stackrel{\textcircled{2}}{\leq} \exp\left(-\frac{nt^2}{8\omega_f^2\tau^2\|\mathbf{y}\|_2^2}\right), \end{aligned}$$

where  $\textcircled{1}$  holds because of Eqn. (6) and we have  $\textcircled{2}$  since we optimize  $s$ .

Since the loss function  $f(\mathbf{w}, \mathbf{x})$  is defined as  $f(\mathbf{w}, \mathbf{x}) = \|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2$ , we have

$$f(\mathbf{w}, \mathbf{x}) - \mathbb{E}(f(\mathbf{w}, \mathbf{x})) = \|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 - \mathbb{E}(\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2) = \left(\|\mathbf{v}^{(l)}\|_2^2 - \mathbb{E}\|\mathbf{v}^{(l)}\|_2^2\right) + \left(\mathbf{y}^T\mathbf{v}^{(l)} - \mathbb{E}\mathbf{y}^T\mathbf{v}^{(l)}\right).$$

Therefore, we have

$$\begin{aligned} &\mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n\left(f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{x}_{(i)}))\right) > t\right) \\ &\leq \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n\left(\|\mathbf{v}_i^{(l)}\|_2^2 - \mathbb{E}\|\mathbf{v}_i^{(l)}\|_2^2\right) > \frac{t}{2}\right) + \mathbb{P}\left(\frac{1}{n}\sum_{i=1}^n\left(\mathbf{y}^T\mathbf{v}_i^{(l)} - \mathbb{E}\mathbf{y}^T\mathbf{v}_i^{(l)}\right) > \frac{t}{2}\right) \\ &\leq 2\exp\left(-c_{f'}n\min\left(\frac{t^2}{\mathbf{d}_l\omega_f^4\tau^4}, \frac{t^2}{\omega_f^2\tau^2}, \frac{t}{\omega_f^2\tau^2}\right)\right). \end{aligned}$$

where  $c_{f'}$  is a constant. Note that  $\|\mathbf{y}\|_2^2$  is the label of  $\mathbf{x}$ , then it can also be bounded. The proof is completed.  $\square$

### C.2.3 PROOF OF LEMMA 8

*Proof.* For brevity, let  $\mathbf{Q}_j$  denote  $\nabla_{\mathbf{w}_{(j)}}f(\mathbf{w}, \mathbf{x})$ . Then, by Lemma 6 we have

$$\nabla_{\mathbf{w}_{(j)}}f(\mathbf{w}) = \left((\mathbf{B}_{j-1:1}\mathbf{x}) \otimes \mathbf{B}_{i,j+1}^T\right) \mathbf{e} \stackrel{\textcircled{1}}{=} \left(\mathbf{B}_{j-1:1}\mathbf{x}\right) \otimes \left(\mathbf{B}_{i,j+1}^T \mathbf{e}\right) \stackrel{\textcircled{2}}{=} \left(\mathbf{B}_{j-1:1} \otimes \mathbf{B}_{i,j+1}^T\right) (\mathbf{x} \otimes \mathbf{e}), \quad (7)$$

where ① holds since  $\mathbf{B}_{j-1:1}\mathbf{x}$  is a vector, and ② holds because for any four matrices  $\mathbf{Z}_1 \sim \mathbf{Z}_4$  of proper sizes, we have  $(\mathbf{Z}_1\mathbf{Z}_3) \otimes (\mathbf{Z}_2\mathbf{Z}_4) = (\mathbf{Z}_1 \otimes \mathbf{Z}_2)(\mathbf{Z}_3 \otimes \mathbf{Z}_4)$ . Note that  $\mathbf{e} = \mathbf{v}^{(l)} - \mathbf{y} = \mathbf{B}_{l:1}\mathbf{x} - \mathbf{y}$ . Then we know that the  $i$ -th entry  $\mathbf{Q}_j^i$  has the form  $\mathbf{Q}_j^i = \sum_{p,q} z_{pq}^{ij} \mathbf{x}_p \mathbf{x}_q + \sum_p y_p^{ij} \mathbf{x}_p + r^{ij}$  (Step 1 below will give the detailed analysis) where  $\mathbf{x}_p$  denotes the  $p$ -th entry in  $\mathbf{x}$ . Note that  $z_{pq}^{ij}, y_p^{ij}$  and  $r^{ij}$  are constants and independent on  $\mathbf{x}$ .

We divide  $\boldsymbol{\lambda} \in \mathbb{R}^{\sum_{j=1}^l d_j d_{j-1}}$  into  $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1; \dots; \boldsymbol{\lambda}_l)$  where  $\boldsymbol{\lambda}_j \in \mathbb{R}^{d_j d_{j-1}}$ . Let  $\boldsymbol{\lambda}_j^i$  denote the  $i$ -th entry in  $\boldsymbol{\lambda}_j$ . Accordingly, we have

$$\mathbf{E} \triangleq \langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \rangle = \sum_{j=1}^l \langle \boldsymbol{\lambda}_j, \mathbf{Q}_j - \mathbb{E} \mathbf{Q}_j \rangle = \mathbf{E}_1 + \mathbf{E}_2 + \mathbf{E}_3,$$

where  $\mathbf{E}_1, \mathbf{E}_2$ , and  $\mathbf{E}_3$  are defined as

$$\begin{aligned} \mathbf{E}_1 &= \sum_{p,q:p \neq q} \left( \sum_{j=1}^l \sum_{i=1}^{d_j d_{j-1}} \boldsymbol{\lambda}_j^i z_{pq}^{ij} \right) (\mathbf{x}_p \mathbf{x}_q - \mathbb{E} \mathbf{x}_p \mathbf{x}_q), \quad \mathbf{E}_2 = \sum_p \left( \sum_{j=1}^l \sum_{i=1}^{d_j d_{j-1}} \boldsymbol{\lambda}_j^i z_{pp}^{ij} \right) (\mathbf{x}_p^2 - \mathbb{E} \mathbf{x}_p^2), \\ \mathbf{E}_3 &= \sum_p \left( \sum_{j=1}^l \sum_{i=1}^{d_j d_{j-1}} \boldsymbol{\lambda}_j^i y_p^{ij} \right) (\mathbf{x}_p - \mathbb{E} \mathbf{x}_p). \end{aligned} \quad (8)$$

Thus, we can further separate the event as:

$$\mathbb{P}(\mathbf{E} > t) \leq \mathbb{P}\left(\frac{1}{n} \sum_{k=1}^n \mathbf{E}_1^k > \frac{t}{3}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{k=1}^n \mathbf{E}_2^k > \frac{t}{3}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{k=1}^n \mathbf{E}_3^k > \frac{t}{3}\right).$$

Thus, to prove our conclusion, we can respectively establish the upper bounds of the three events. To the end, for each input sample  $\mathbf{x}_{(i)}$ , we divide its corresponding  $\mathbf{Q}_j - \mathbb{E} \mathbf{Q}_j$  into  $\mathbf{E}_1, \mathbf{E}_2$  and  $\mathbf{E}_3$ . Then we bound the three events separately. Before that, we first give several equalities. Since  $\mathbf{B}_{j:s} = \mathbf{W}^{(j)} \mathbf{W}^{(j-1)} \dots \mathbf{W}^{(s)}$  ( $j \geq s$ ), by Lemma 14 we have

$$\|\mathbf{B}_{j:s}\|_F^2 \leq r^{2(j-s+1)} \quad \text{and} \quad \|\mathbf{B}_{l:t+1}\|_F^2 \|\mathbf{B}_{t-1:s+1}\|_F^2 \|\mathbf{B}_{s-1:1}\|_F^2 \leq r^{2(l-2)}, \quad (9)$$

These two inequalities can be obtained by using  $\|\mathbf{W}^{(i)}\|_F^2 = \|\mathbf{w}_{(i)}\|_2^2 \leq r^2$ .

**Step 1. Divide  $\mathbf{Q}_j - \mathbb{E} \mathbf{Q}_j$ :** Note that  $\mathbf{e} = \mathbf{v}^{(l)} - \mathbf{y} = \mathbf{B}_{l:1}\mathbf{x} - \mathbf{y}$ . Let  $\mathbf{H}_j = \mathbf{B}_{j-1:1} \otimes \mathbf{B}_{l:j+1}^T$ . Then we can further write Eqn. (7) as

$$\mathbf{Q}_j = \nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}) = \mathbf{H}_j (\mathbf{x} \otimes (\mathbf{B}_{l:1}\mathbf{x}) - \mathbf{x} \otimes \mathbf{y}) = \mathbf{H}_j ((\mathbf{I}_{d_0} \otimes \mathbf{B}_{l:1}) (\mathbf{x} \otimes \mathbf{x}) - \mathbf{x} \otimes \mathbf{y}), \quad (10)$$

where  $\mathbf{I}_{d_0} \in \mathbb{R}^{d_0 \times d_0}$  is the identity matrix. According to Eqn. (10), we can write the  $i$ -th entry of  $\mathbf{Q}_j$  as the form  $\mathbf{Q}_j^i = \sum_{p,q} z_{pq}^{ij} \mathbf{x}_p \mathbf{x}_q + \sum_p y_p^{ij} \mathbf{x}_p + r^{ij}$  where  $\mathbf{x}_p$  denotes the  $p$ -th entry in  $\mathbf{x}$ . Let  $\mathbf{Z}_j = \mathbf{H}_j (\mathbf{I}_{d_0} \otimes \mathbf{B}_{l:1}) \in \mathbb{R}^{d_j d_{j-1} \times d_0^2}$ . Then, we know that the  $i$ -th entry  $\mathbf{Q}_j^i = \mathbf{Z}_j(i, :) \mathbf{x}'$ , where  $\mathbf{x}' = \mathbf{x} \otimes \mathbf{x} = [\mathbf{x}_1 \mathbf{x}; \mathbf{x}_2 \mathbf{x}; \dots, \mathbf{x}_{d_0} \mathbf{x}] \in \mathbb{R}^{d_0^2}$ . In this way, we have  $z_{pq}^{ij} = \mathbf{Z}_j(i, (p-1)d_0 + q)$  which further implies

$$\sum_{p,q} (z_{pq}^{ij})^2 \leq c_q \|\mathbf{Z}_j(i, :)\|_2^2, \quad (11)$$

where  $c_q = \sqrt{\max_{0 \leq i \leq l} d_i}$ .

We divide the  $i$ -th row  $\mathbf{H}_j(i, :)$  into  $\mathbf{H}_j(i, :) = [\mathbf{H}_{ji}^1, \mathbf{H}_{ji}^2, \dots, \mathbf{H}_{ji}^{d_0}]$  where  $\mathbf{H}_{ji}^p \in \mathbb{R}^{1 \times d_i}$ . Then we have  $y_p^{ij} = \mathbf{y}^T \mathbf{H}_{ji}^p$ . This yields

$$\sum_p (y_p^{ij})^2 \leq c_q \sum_p (\mathbf{y}^T \mathbf{H}_{ji}^p)^2 \leq c_q \sum_p \|\mathbf{y}\|_2^2 \|\mathbf{H}_{ji}^p\|_2^2 = c_q \|\mathbf{y}\|_2^2 \|\mathbf{H}_j(i, :)\|_2^2. \quad (12)$$

Let  $\boldsymbol{\lambda}_j^i$  denote the  $i$ -th entry of  $\boldsymbol{\lambda}_j$ . Then, by Eqn. (8), we can obtain

$$\begin{aligned} \sum_j \langle \boldsymbol{\lambda}_j, (\mathbf{Q}_j - \mathbb{E}(\mathbf{Q}_j)) \rangle &= \sum_{p,q:p \neq q} a_{pq} (\mathbf{x}_p \mathbf{x}_q - \mathbb{E} \mathbf{x}_p \mathbf{x}_q) + \sum_p a_{pp} (\mathbf{x}_p^2 - \mathbb{E} \mathbf{x}_p^2) + \sum_p b_p (\mathbf{x}_p - \mathbb{E} \mathbf{x}_p) \\ &= \mathbf{E}_1 + \mathbf{E}_2 + \mathbf{E}_3, \end{aligned}$$

where  $a_{pq}$  and  $b_p$  are defined as

$$a_{pq} = \sum_{j=1}^l \sum_{i=1}^{d_j d_{j-1}} \lambda_j^i z_{pq}^{ij} \quad \text{and} \quad b_p = \sum_{j=1}^l \sum_{i=1}^{d_j d_{j-1}} \lambda_j^i y_p^{ij}.$$

Note that for any four matrices of proper sizes, we have  $(Q_1 \otimes Q_2)(Q_3 \otimes Q_4) = (Q_1 Q_3) \otimes (Q_2 Q_4)$ , indicating  $Z_j = \left( B_{j-1:1} \otimes B_{l:j+1}^T \right) (I_{d_0} \otimes B_{l:1}) = B_{j-1:1} \otimes \left( B_{l:j+1}^T B_{l:1} \right)$ . This gives

$$c_q \|Z_j\|_F^2 \leq c_q \|B_{j-1:1}\|_F^2 \|B_{l:j+1}\|_F^2 \|B_{l:1}\|_F^2 \stackrel{\textcircled{1}}{\leq} c_q r^{2(l-1)} r^{2l} = c_q r^{2(2l-1)} \triangleq \omega. \quad (13)$$

Note that Eqn. (13) uses the conclusion in Eqn. (9). Therefore, we can have the following bound:

$$\sum_{i=1}^{d_j d_{j-1}} (z_{pq}^{ij})^2 \leq \sum_{i=1}^{d_j d_{j-1}} \sum_{p,q} (z_{pq}^{ij})^2 \stackrel{\textcircled{1}}{\leq} c_q \sum_{i=1}^{d_j d_{j-1}} \|Z_j(i, \cdot)\|_2^2 = c_q \|Z_j\|_F^2 \leq \omega, \quad (14)$$

where  $\textcircled{1}$  uses Eqn. (11). Then we can utilize Eqn. (14) and  $\sum_{j=1}^l \left( \sum_{i=1}^{d_j d_{j-1}} (\lambda_j^i)^2 \right) = 1$  to bound  $a_{pq}$  as follows:

$$a_{pq}^2 \leq l \left( \sum_{j=1}^l \left( \sum_{i=1}^{d_j d_{j-1}} \lambda_j^i z_{pq}^{ij} \right)^2 \right) \leq l \sum_{j=1}^l \left( \sum_{i=1}^{d_j d_{j-1}} (\lambda_j^i)^2 \right) \left( \sum_{i=1}^{d_j d_{j-1}} (z_{pq}^{ij})^2 \right) \leq l\omega.$$

which further gives

$$\sum_{p,q} a_{pq}^2 \leq l \sum_{j=1}^l \left( \sum_{i=1}^{d_j d_{j-1}} (\lambda_j^i)^2 \right) \left( \sum_{i=1}^{d_j d_{j-1}} \sum_{p,q} (z_{pq}^{ij})^2 \right) \stackrel{\textcircled{1}}{\leq} l\omega.$$

where  $\textcircled{1}$  uses Eqn. (14).

Similarly, we can obtain

$$\sum_{i=1}^{d_j d_{j-1}} (y_p^{ij})^2 \leq \sum_{i=1}^{d_j d_{j-1}} c_q \|\mathbf{y}\|_2^2 \|H_j(i, \cdot)\|_2^2 = c_q \|\mathbf{y}\|_2^2 \|H_j\|_F^2 \leq c_q \|\mathbf{y}\|_2^2 r^{2(l-1)}. \quad (15)$$

So we can bound  $b_p$  as

$$b_p^2 \leq l \sum_{j=1}^l \left( \sum_{i=1}^{d_j d_{j-1}} \lambda_j^i y_p^{ij} \right)^2 \leq l \sum_{j=1}^l \left( \sum_{i=1}^{d_j d_{j-1}} (\lambda_j^i)^2 \right) \left( \sum_{i=1}^{d_j d_{j-1}} (y_p^{ij})^2 \right) \leq l\omega',$$

where  $\omega' = c_q \|\mathbf{y}\|_2^2 r^{2(l-1)}$ . Accordingly, we can have

$$\sum_p b_p^2 \leq l \sum_{j=1}^l \left( \sum_{i=1}^{d_j d_{j-1}} (\lambda_j^i)^2 \right) \left( \sum_{i=1}^{d_j d_{j-1}} \sum_p (y_p^{ij})^2 \right) \stackrel{\textcircled{1}}{\leq} l\omega',$$

where  $\textcircled{1}$  uses (15).

**Step 2. Bound  $\mathbb{P}(\mathbf{E}_1 > t/3)$ ,  $\mathbb{P}(\mathbf{E}_2 > t/3)$  and  $\mathbb{P}(\mathbf{E}_3 > t/3)$ :** Let  $\mathbf{E}_{h1}^k$  denotes the  $\mathbf{E}_{h1}$  which corresponds to the  $k$ -th sample  $\mathbf{x}_{(k)}$ . Therefore, we can bound

$$\begin{aligned}
\mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \mathbf{E}_1^k > \frac{t}{3}\right) &= \mathbb{P}\left(s \sum_{k=1}^n \left(\sum_{p,q:p \neq q} a_{pq}^k (\mathbf{x}_p^k \mathbf{x}_q^k - \mathbb{E} \mathbf{x}_p^k \mathbf{x}_q^k)\right) > \frac{st}{3}\right) \\
&\stackrel{\textcircled{1}}{\leq} \exp\left(-\frac{nst}{3}\right) \mathbb{E} \exp\left(s \sum_{k=1}^n \left(\sum_{p,q:p \neq q} a_{pq}^k (\mathbf{x}_p^k \mathbf{x}_q^k - \mathbb{E} \mathbf{x}_p^k \mathbf{x}_q^k)\right)\right) \\
&\stackrel{\textcircled{2}}{\leq} \exp\left(-\frac{nst}{3}\right) \prod_{k=1}^n \mathbb{E} \exp\left(s \left(\sum_{p,q:p \neq q} a_{pq}^k (\mathbf{x}_p^k \mathbf{x}_q^k - \mathbb{E} \mathbf{x}_p^k \mathbf{x}_q^k)\right)\right) \\
&\stackrel{\textcircled{3}}{\leq} \exp\left(-\frac{nst}{3}\right) \prod_{k=1}^n \exp\left(2\tau^2 s^2 \sum_{p,q:p \neq q} (a_{pq}^k)^2\right) \quad |s| \leq \frac{1}{2\tau\sqrt{l\omega}} \\
&\leq \exp\left(-\frac{nst}{3}\right) \prod_{j=1}^n \exp(2\tau^2 s^2 l\omega) \\
&\stackrel{\textcircled{4}}{\leq} \exp\left(-c'n \min\left(\frac{t^2}{\omega l \tau^2}, \frac{t}{\sqrt{l\omega}\tau}\right)\right),
\end{aligned}$$

where  $\textcircled{1}$  holds because of Chebyshev's inequality.  $\textcircled{2}$  holds since  $\mathbf{x}_{(i)}$  are independent.  $\textcircled{3}$  is established by applying Lemma 12. We have  $\textcircled{4}$  by optimizing  $s$ . Similarly, by Lemma 13 we can bound  $\mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \mathbf{E}_2^k > \frac{t}{3}\right)$  as follows:

$$\begin{aligned}
\mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \mathbf{E}_2^k > \frac{t}{3}\right) &\leq \exp\left(-\frac{nst}{3}\right) \prod_{k=1}^n \mathbb{E} \exp\left(s \left(\sum_p a_{pp}^k ((\mathbf{x}_p^k)^2 - \mathbb{E}(\mathbf{x}_p^k)^2)\right)\right) \\
&\leq \exp\left(-\frac{nst}{3}\right) \prod_{k=1}^n \exp(128\tau^4 s^2 l\omega) \quad |s| \leq \frac{1}{\tau^2\sqrt{l\omega}} \\
&\leq \exp\left(-c''n \min\left(\frac{t^2}{\omega l \tau^4}, \frac{t}{\sqrt{l\omega}\tau^2}\right)\right).
\end{aligned}$$

Finally, since  $\mathbf{x}_{(i)}$  are independent sub-Gaussian, we can use Hoeffding inequality and obtain

$$\mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \mathbf{E}_3^k > \frac{t}{3}\right) \leq \mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \left(\sum_p b_p^k (\mathbf{x}_p^k - \mathbb{E} \mathbf{x}_p^k)\right) > \frac{t}{3}\right) \exp\left(-\frac{c'''nt}{\omega' l \tau^2}\right).$$

**Step 3. Bound  $\mathbb{P}(\mathbf{E} > t)$ :** By comparing the values of  $\omega$  and  $\omega'$ , we can obtain

$$\begin{aligned}
\mathbb{P}(\mathbf{E} > t) &\leq \mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \mathbf{E}_1^k > \frac{t}{3}\right) + \mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \mathbf{E}_2^k > \frac{t}{3}\right) + \mathbb{P}\left(\frac{1}{n}\sum_{k=1}^n \mathbf{E}_3^k > \frac{t}{3}\right) \\
&\leq 3 \exp\left(-c_g n \min\left(\frac{t^2}{l \max(\omega_g \tau^2, \omega_g \tau^4, \omega_{g'} \tau^2)}, \frac{t}{\sqrt{l\omega_g} \max(\tau, \tau^2)}\right)\right),
\end{aligned}$$

where  $\omega_g = c_q r^{2(2l-1)}$  and  $\omega_{g'} = c_q r^{2(l-1)}$  in which  $c_q = \sqrt{\max_{0 \leq i \leq l} \mathbf{d}_i}$ . The proof is completed.  $\square$

#### C.2.4 PROOFS OF LEMMA 9

*Proof.* For brevity, let  $\mathbf{Q}_{ts}$  denote  $\nabla_{\mathbf{w}_{(t)}} (\nabla_{\mathbf{w}_{(s)}} f(\mathbf{w}, \mathbf{x}))$ . Then, by Lemma 6 we have

$$\mathbf{Q}_{ts} = \begin{cases} (\mathbf{B}_{l:s+1}^T \mathbf{e} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes \mathbf{B}_{s-1:t+1} + (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1}), & \text{if } s > t, \\ (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{s-1:1}) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:s+1}), & \text{if } s = t, \\ (\mathbf{B}_{t-1:s+1}^T) \otimes (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:t+1}^T) + (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1}), & \text{if } s < t. \end{cases}$$

Then we know that the  $(i, k)$ -th entry  $Q_{ts}^{ik}$  has the form  $Q_{ts}^{ik} = \sum_{p,q} z_{pq}^{ik} \mathbf{x}_p \mathbf{x}_q + \sum_p y_p^{ik} \mathbf{x}_p + r^{ik}$  (explained in the following Step 1. I) where  $\mathbf{x}_p$  denotes the  $p$ -th entry in  $\mathbf{x}$ . Note that  $z_{pq}^{ik}, y_p^{ik}$  and  $r^{ik}$  are constant and independent on  $\mathbf{x}$ . For convenience, we let  $Q_{ts} = \mathbf{H}_{ts} + \mathbf{G}_{ts}$ , where  $\mathbf{G}_{ts} = (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1})$  and  $\mathbf{H}_{ts}$  is defined as

$$\mathbf{H}_{ts} = \begin{cases} (\mathbf{B}_{l:s+1}^T \mathbf{e} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes \mathbf{B}_{s-1:t+1}, & \text{if } s > t, \\ \mathbf{0}, & \text{if } s = t, \\ (\mathbf{B}_{t-1:s+1}^T) \otimes (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:t+1}^T), & \text{if } s < t. \end{cases}$$

Let

$$\mathbf{E} = \frac{1}{n} \sum_{j=1}^n \langle \boldsymbol{\lambda}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})) \boldsymbol{\lambda} \rangle, \quad \mathbf{E}_h = \frac{1}{n} \sum_{j=1}^n \sum_{t,s} \langle \boldsymbol{\lambda}_t, (\mathbf{H}_{ts} - \mathbb{E}(\mathbf{H}_{ts})) \boldsymbol{\lambda}_s \rangle,$$

$$\mathbf{E}_g = \frac{1}{n} \sum_{j=1}^n \sum_{t,s} \langle \boldsymbol{\lambda}_t, (\mathbf{G}_{ts} - \mathbb{E}(\mathbf{G}_{ts})) \boldsymbol{\lambda}_s \rangle.$$

Then we divide the event as two events:

$$\mathbb{P}(\mathbf{E} > t) = \mathbb{P}(\mathbf{E}_h + \mathbf{E}_g > t) \leq \mathbb{P}(\mathbf{E}_h > t/2) + \mathbb{P}(\mathbf{E}_g > t/2).$$

Now we look each event separately. Similar to  $Q_{ts}$ , the  $(i, k)$ -th entry  $H_{ts}^{ik}$  has the form  $H_{ts}^{ik} = \sum_{p,q} z_{pq}^{ik} \mathbf{x}_p \mathbf{x}_q + \sum_p y_p^{ik} \mathbf{x}_p + r^{ik}$ . We divide the unit vector  $\boldsymbol{\lambda} \in \mathbb{R}^d$  as  $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1; \dots; \boldsymbol{\lambda}_l)$  where  $\boldsymbol{\lambda}_j \in \mathbb{R}^{d_j d_{j-1}}$ . For input vector  $\mathbf{x}$ , let  $\sum_{t,s} \langle \boldsymbol{\lambda}_t, (\mathbf{H}_{ts} - \mathbb{E}(\mathbf{H}_{ts})) \boldsymbol{\lambda}_s \rangle = \mathbf{E}_{h1} + \mathbf{E}_{h2} + \mathbf{E}_{h3}$ , where

$$\mathbf{E}_{h1} = \sum_{p,q:p \neq q} \left( \sum_{t,s} \sum_{i,k} (\boldsymbol{\lambda}_t^i \boldsymbol{\lambda}_s^k) z_{pq}^{ik} \right) (\mathbf{x}_p \mathbf{x}_q - \mathbb{E} \mathbf{x}_p \mathbf{x}_q), \quad \mathbf{E}_{h2} = \sum_p \left( \sum_{t,s} \sum_{i,k} (\boldsymbol{\lambda}_t^i \boldsymbol{\lambda}_s^k) z_{pq}^{ik} \right) (\mathbf{x}_p^2 - \mathbb{E} \mathbf{x}_p^2),$$

$$\mathbf{E}_{h3} = \sum_p \left( \sum_{t,s} \sum_{i,k} (\boldsymbol{\lambda}_t^i \boldsymbol{\lambda}_s^k) y_p^{ik} \right) (\mathbf{x}_p - \mathbb{E} \mathbf{x}_p), \quad (16)$$

where  $\mathbf{x}_p$  denotes the  $p$ -th entry in  $\mathbf{x}$  and  $\boldsymbol{\lambda}_j^i$  denotes the  $i$ -th entry of  $\boldsymbol{\lambda}_j$ . Let  $\mathbf{E}_{h1}^j, \mathbf{E}_{h2}^j$ , and  $\mathbf{E}_{h3}^j$  denote the  $\mathbf{E}_{h1}, \mathbf{E}_{h2}$ , and  $\mathbf{E}_{h3}$  of the  $j$ -th sample. Thus, considering  $n$  samples, we can further separately divide the two events above as:

$$\mathbb{P}\left(\mathbf{E}_h > \frac{t}{2}\right) \leq \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h1}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h2}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h3}^j > \frac{t}{6}\right).$$

Similarly, we can define  $\mathbf{E}_{g1}, \mathbf{E}_{g2}$  and  $\mathbf{E}_{g3}$ .

$$\mathbb{P}\left(\mathbf{E}_g > \frac{t}{2}\right) \leq \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g1}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g2}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g3}^j > \frac{t}{6}\right).$$

Thus, to prove our conclusion, we can respectively establish the upper bounds of  $\mathbb{P}(\mathbf{E}_h > \frac{t}{2})$  and  $\mathbb{P}(\mathbf{E}_g > \frac{t}{2})$ .

**Step 1: Bound  $\mathbb{P}(\mathbf{E}_h > \frac{t}{2})$**

To achieve our goal, for each input sample  $\mathbf{x}_{(i)}$ , we divide its corresponding  $\sum_{t,s} (\mathbf{H}_{ts} - \mathbb{E} \mathbf{H}_{ts})$  as  $\mathbf{E}_{h1}, \mathbf{E}_{h2}$  and  $\mathbf{E}_{h3}$ . Then we bound the three events separately. Before that, we first give two equalities. Since  $\mathbf{B}_{j:s} = \mathbf{W}^{(j)} \mathbf{W}^{(j-1)} \dots \mathbf{W}^{(s)}$  ( $j \geq s$ ), by Lemma 14 we have

$$\|\mathbf{B}_{j:s}\|_F^2 \leq r^{2(j-s+1)} \quad \text{and} \quad \|\mathbf{B}_{l:t+1}\|_F^2 \|\mathbf{B}_{t-1:s+1}\|_F^2 \|\mathbf{B}_{s-1:1}\|_F^2 \leq r^{2(l-2)}, \quad (17)$$

These two inequalities can be obtained by using  $\|\mathbf{W}^{(i)}\|_F^2 = \|\mathbf{w}_{(i)}\|_2^2 \leq r^2$ .

**I. Divide  $\mathbf{H}_{ts} - \mathbb{E}\mathbf{H}_{ts}$ :** For  $t \neq s$ , we can write the  $(i, k)$ -th entry  $\mathbf{H}_{ts}^{ik}$  as the form  $\mathbf{H}_{ts}^{ik} = \sum_{p,q} z_{pq}^{ik} \mathbf{x}_p \mathbf{x}_q + \sum_p y_p^{ik} \mathbf{x}_p + r^{ik}$ . Now we try to bound  $z_{pq}^{ik}$  and  $y_p^{ik}$ . We first consider the case  $s < t$ . Note that  $\mathbf{e} = \mathbf{v}^{(l)} - \mathbf{y} = \mathbf{B}_{l:1} \mathbf{x} - \mathbf{y}$ . Specifically, we have

$$\mathbf{H}_{ts} = (\mathbf{B}_{t-1:s+1}^T) \otimes (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{l:1}^T \mathbf{B}_{l:t+1}^T - \mathbf{B}_{s-1:1} \mathbf{x} \mathbf{y}^T \mathbf{B}_{l:t+1}^T).$$

So the  $(i', k')$ -th entry in the matrix  $\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{l:1}^T \mathbf{B}_{l:t+1}^T$  is  $[\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{l:1}^T \mathbf{B}_{l:t+1}^T]_{i'k'} = (\mathbf{B}_{s-1:1})_{(i', :)} \mathbf{x} (\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k', :)} \mathbf{x} = \mathbf{x}^T ((\mathbf{B}_{s-1:1})_{(i', :)}^T (\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k', :)} \mathbf{x}$ , where  $\mathbf{A}_{(i', :)}$  denotes the  $i'$ -th row of  $\mathbf{A}$ . Let  $i'_i = \text{mod}(i, \mathbf{d}_s)$ ,  $k'_k = \text{mod}(k, \mathbf{d}_{t-1})$ ,  $i''_i = \lfloor i/\mathbf{d}_s \rfloor$  and  $k''_k = \lfloor k/\mathbf{d}_{t-1} \rfloor$ . In this case, the  $(i, k)$ -th entry  $\mathbf{H}_{ts}^{ik} = [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i} \mathbf{x}^T ((\mathbf{B}_{s-1:1})_{(i'_i, :)}^T (\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k'_k, :)} \mathbf{x} + [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i} \mathbf{y}^T (\mathbf{B}_{l:t+1})_{(k'_k, :)}^T (\mathbf{B}_{s-1:1})_{(i'_i, :)} \mathbf{x}$ . Therefore, we have

$$\begin{aligned} \sum_{p,q} (z_{pq}^{ik})^2 &= [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i}^2 \|((\mathbf{B}_{s-1:1})_{(i'_i, :)}^T (\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k'_k, :)} \mathbf{x})\|_F^2 \\ &\leq [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i}^2 \|(\mathbf{B}_{s-1:1})_{(i'_i, :)}\|_2^2 \|(\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k'_k, :)}\|_2^2. \end{aligned}$$

Therefore, we can further establish

$$\begin{aligned} \sum_{i,k} \sum_{p,q} (z_{pq}^{ik})^2 &\leq \sum_{i,k} [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i}^2 \|(\mathbf{B}_{s-1:1})_{(i'_i, :)}\|_2^2 \|(\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k'_k, :)}\|_2^2 \\ &\leq \sum_{i,k} [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i}^2 \|(\mathbf{B}_{s-1:1})_{(i'_i, :)}\|_2^2 \|(\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k'_k, :)}\|_2^2 \\ &= \sum_k \|(\mathbf{B}_{t-1:s+1})_{(k'', :)}\|_2^2 \|\mathbf{B}_{s-1:1}\|_F^2 \|(\mathbf{B}_{l:1} \mathbf{B}_{l:t+1})_{(k'_k, :)}\|_2^2 \\ &= \|\mathbf{B}_{t-1:s+1}\|_F^2 \|\mathbf{B}_{s-1:1}\|_F^2 \|\mathbf{B}_{l:1} \mathbf{B}_{l:t+1}\|_F^2 \\ &\stackrel{\textcircled{1}}{\leq} r^{4(l-1)} \triangleq \omega. \end{aligned} \tag{18}$$

where  $\textcircled{1}$  uses Eqn. (17). Similarly, we can bound

$$\begin{aligned} \sum_p (y_p^{ik})^2 &= [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i}^2 \|\mathbf{y}^T (\mathbf{B}_{l:t+1})_{(k'_k, :)}^T (\mathbf{B}_{s-1:1})_{(i'_i, :)}\|_F^2 \\ &\leq [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i}^2 \|\mathbf{y}\|_2^2 \|(\mathbf{B}_{l:t+1})_{(k'_k, :)}\|_2^2 \|(\mathbf{B}_{s-1:1})_{(i'_i, :)}\|_2^2. \end{aligned}$$

So it further yields

$$\begin{aligned} \sum_{i,k} \sum_p (y_p^{ik})^2 &\leq \sum_{i,k} [\mathbf{B}_{t-1:s+1}]_{k''_k i''_i}^2 \|\mathbf{y}\|_2^2 \|(\mathbf{B}_{l:t+1})_{(k'_k, :)}\|_2^2 \|(\mathbf{B}_{s-1:1})_{(i'_i, :)}\|_2^2 \\ &\leq \|\mathbf{y}\|_2^2 \|\mathbf{B}_{t-1:s+1}\|_F^2 \|\mathbf{B}_{l:t+1}\|_F^2 \|\mathbf{B}_{s-1:1}\|_F^2 \stackrel{\textcircled{1}}{\leq} \|\mathbf{y}\|_2^2 r^{2(l-2)} \triangleq \omega', \end{aligned} \tag{19}$$

where  $\textcircled{1}$  uses Eqn. (17). Note that for the case  $s \geq t$ , Eqn. (18) and (19) also holds. Let  $\lambda_j^i$  denote the  $i$ -th entry of  $\lambda_j$ . Then, by Eqn. (16), we can obtain

$$\begin{aligned} \sum_{t,s} (\langle \lambda_t, (\mathbf{H}_{ts} - \mathbb{E}(\mathbf{H}_{ts})) \lambda_s \rangle) &= \sum_{p,q:p \neq q} a_{pq} (\mathbf{x}_p \mathbf{x}_q - \mathbb{E} \mathbf{x}_p \mathbf{x}_q) + \sum_p a_{pp} (\mathbf{x}_p^2 - \mathbb{E} \mathbf{x}_p^2) + \sum_p b_p (\mathbf{x}_p - \mathbb{E} \mathbf{x}_p) \\ &= \mathbf{E}_{h1} + \mathbf{E}_{h2} + \mathbf{E}_{h3}, \end{aligned}$$

where  $a_{pq}$  and  $b_p$  are defined as

$$a_{pq} = \sum_{t,s} \sum_{i,k} (\lambda_t^i \lambda_s^k) z_{pq}^{ik} \quad \text{and} \quad b_p = \sum_{t,s} \sum_{i,k} (\lambda_t^i \lambda_s^k) y_p^{ik}.$$

Then according to Eqn. (18) and  $\sum_{t,s} (\sum_{i,k} (\lambda_t^i \lambda_s^k)^2) = 1$ , we can bound  $a_{pq}$  as follows:

$$a_{pq}^2 \leq l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k) z_{pq}^{ik} \right)^2 \leq l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k)^2 \right) \left( \sum_{i,k} (z_{pq}^{ik})^2 \right) \leq \omega l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k)^2 \right) \leq \omega l^2.$$

which further yields

$$\sum_{p,q} a_{pq}^2 \leq l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k)^2 \right) \left( \sum_{i,k} \sum_{p,q} (z_{pq}^{ik})^2 \right) \leq \omega l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k)^2 \right) \leq \omega l^2.$$

Similarly, by using Eqn. (19), we have

$$b_p^2 \leq l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k) y_p^{ik} \right)^2 \leq l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k)^2 \right) \left( \sum_{i,k} (y_p^{ik})^2 \right) \stackrel{\textcircled{1}}{\leq} \omega' l^2.$$

Accordingly, we can have

$$\sum_p b_p^2 \leq l^2 \sum_{t,s} \left( \sum_{i,k} (\lambda_t^i \lambda_s^k)^2 \right) \left( \sum_{i,k} \sum_p (y_p^{ik})^2 \right) \leq \omega' l^2.$$

**II. Bound  $\mathbb{P}(\mathbf{E}_{h1} > t/6)$ ,  $\mathbb{P}(\mathbf{E}_{h2} > t/6)$  and  $\mathbb{P}(\mathbf{E}_{h3} > t/6)$ :** Let  $E_{h1}^j$  denotes the  $E_{h1}^j$  which corresponds to the  $j$ -th sample  $\mathbf{x}_{(i)}$ . Therefore, we can bound

$$\begin{aligned} \mathbb{P} \left( \frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h1}^j > \frac{t}{6} \right) &\leq \mathbb{P} \left( s \sum_{j=1}^n \left( \sum_{p,q:p \neq q} a_{pq}^j (\mathbf{x}_p^j \mathbf{x}_q^j - \mathbb{E} \mathbf{x}_p^j \mathbf{x}_q^j) \right) > \frac{snt}{6} \right) \\ &\stackrel{\textcircled{1}}{\leq} \exp \left( -\frac{nst}{6} \right) \mathbb{E} \exp \left( s \sum_{j=1}^n \left( \sum_{p,q:p \neq q} a_{pq}^j (\mathbf{x}_p^j \mathbf{x}_q^j - \mathbb{E} \mathbf{x}_p^j \mathbf{x}_q^j) \right) \right) \\ &\stackrel{\textcircled{2}}{\leq} \exp \left( -\frac{nst}{6} \right) \prod_{j=1}^n \mathbb{E} \exp \left( s \left( \sum_{p,q:p \neq q} a_{pq}^j (\mathbf{x}_p^j \mathbf{x}_q^j - \mathbb{E} \mathbf{x}_p^j \mathbf{x}_q^j) \right) \right) \\ &\stackrel{\textcircled{3}}{\leq} \exp \left( -\frac{nst}{6} \right) \prod_{j=1}^n \exp \left( 2\tau^2 s^2 \sum_{p,q:p \neq q} (a_{pq}^j)^2 \right) \quad |s| \leq \frac{1}{2\tau l \sqrt{\omega}} \\ &\leq \exp \left( -\frac{nst}{6} \right) \prod_{j=1}^n \exp (2\tau^2 s^2 l^2 \omega) \\ &\stackrel{\textcircled{4}}{\leq} \exp \left( -c' n \min \left( \frac{t^2}{\omega l^2 \tau^2}, \frac{t}{\sqrt{\omega} l \tau} \right) \right), \end{aligned}$$

where  $\textcircled{1}$  holds because of Chebyshev's inequality.  $\textcircled{2}$  holds since  $\mathbf{x}_{(i)}$  are independent.  $\textcircled{3}$  is established because of Lemma 12. We have  $\textcircled{4}$  by optimizing  $s$ . Similarly, we can bound  $\mathbb{P} \left( \frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h2}^j > \frac{t}{6} \right)$  as follows:

$$\begin{aligned} \mathbb{P} \left( \frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h2}^j > \frac{t}{6} \right) &\leq \exp \left( -\frac{nst}{6} \right) \prod_{j=1}^n \mathbb{E} \exp \left( s \left( \sum_p a_{pp}^j ((\mathbf{x}_p^j)^2 - \mathbb{E}(\mathbf{x}_p^j)^2) \right) \right) \\ &\leq \exp \left( -\frac{nst}{6} \right) \prod_{j=1}^n \exp (128\tau^4 s^2 l^2 \omega) \quad |s| \leq \frac{1}{\tau^2 l \sqrt{\omega}} \\ &\leq \exp \left( -c'' n \min \left( \frac{t^2}{\omega l^2 \tau^4}, \frac{t}{\sqrt{\omega} l \tau^2} \right) \right). \end{aligned}$$

Finally, since  $\mathbf{x}_{(i)}$  are independent sub-Gaussian, we can use Hoeffding inequality and obtain

$$\mathbb{P} \left( \frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h3}^j > \frac{t}{6} \right) = \mathbb{P} \left( \frac{1}{n} \sum_{j=1}^n \left( \sum_p b_p^j (\mathbf{x}_p^j - \mathbb{E} \mathbf{x}_p^j) \right) > \frac{t}{6} \right) \leq \exp \left( -\frac{c''' n t^2}{\omega' l^2 \tau^2} \right).$$

Since for  $s = t$ ,  $\mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h1}^j > \frac{t}{6}\right) = \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h2}^j > \frac{t}{6}\right) = \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h3}^j > \frac{t}{6}\right) = 0$ , the above upper bounds also hold.

**III: Bound  $\mathbb{P}(\mathbf{E}_h > \frac{t}{2})$**  By comparing the values of  $\omega$  and  $\omega'$ , we can obtain

$$\begin{aligned} \mathbb{P}\left(\mathbf{E}_h > \frac{t}{2}\right) &\leq \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h1}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h2}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{h3}^j > \frac{t}{6}\right) \\ &\leq 3 \exp\left(-c'_2 n \min\left(\frac{t^2}{l^2 \max(\omega\tau^2, \omega\tau^4, \omega_q\tau^2)}, \frac{t}{\sqrt{\omega}l \max(\tau, \tau^2)}\right)\right), \end{aligned}$$

where  $\omega_q = r^{2(l-2)}$ .

**Step 2: Bound  $\mathbb{P}(\mathbf{E}_g > \frac{t}{2})$**  To achieve our goal, for each input sample  $\mathbf{x}_{(i)}$ , we also divide its corresponding  $\sum_{t,s} (\mathbf{G}_{ts} - \mathbb{E}\mathbf{G}_{ts})$  as  $\mathbf{E}_{h1}$ ,  $\mathbf{E}_{h2}$  and  $\mathbf{E}_{h3}$ . Then we bound the three events separately. Before that, we first give several equalities.

**I. Divide  $\mathbf{G}_{ts} - \mathbb{E}\mathbf{G}_{ts}$ :** Dividing  $\mathbf{G}_{ts} - \mathbb{E}\mathbf{G}_{ts}$  is more easy than dividing  $\mathbf{H}_{ts} - \mathbb{E}\mathbf{H}_{ts}$  since the later has more complex form. Since  $\mathbf{G}_{ts} = (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1})$ . we also can write the  $(i, k)$ -th entry  $\mathbf{G}_{ts}^{ik}$  as the form  $\mathbf{G}_{ts}^{ik} = \sum_{p,q} z_{pq}^{ik} \mathbf{x}_p \mathbf{x}_q + \sum_p y_p^{ik} \mathbf{x}_p + r^{ik}$ . But here  $y_p^{ik} = 0$ .

Then similar to the step in dividing  $\mathbf{H}_{ts} - \mathbb{E}\mathbf{H}_{ts}$ , we can bound

$$a_{pq}^2 \leq \omega_g l^2 \quad \text{and} \quad \sum_{p,q} a_{pq}^2 \leq \omega_g l^2 \quad \text{where} \quad \omega_g = r^{4(l-1)}.$$

**II. Bound  $\mathbb{P}(\mathbf{E}_{g1} > t/6)$ ,  $\mathbb{P}(\mathbf{E}_{g2} > t/6)$  and  $\mathbb{P}(\mathbf{E}_{g3} > t/6)$ :** Since  $y_p^{ik} = 0$ ,  $\mathbb{P}(\mathbf{E}_{h3} > t/6) = 0$ . Similar to the above methods, we can bound

$$\mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g1}^j > \frac{t}{6}\right) \leq \exp\left(-c'_1 n \left(\frac{t^2}{\omega_g l^2 \tau^2}, \frac{t}{\sqrt{\omega_g} l \tau}\right)\right),$$

and

$$\mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g2}^j > \frac{t}{6}\right) \leq \exp\left(-c''_1 n \left(\frac{t^2}{\omega_g l^2 \tau^4}, \frac{t}{\sqrt{\omega_g} l \tau^2}\right)\right).$$

**III: Bound  $\mathbb{P}(\mathbf{E}_h > \frac{t}{2})$**  We can obtain  $\mathbb{P}(\mathbf{E}_g > \frac{t}{2})$  as follows:

$$\begin{aligned} \mathbb{P}\left(\mathbf{E}_g > \frac{t}{2}\right) &\leq \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g1}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g2}^j > \frac{t}{6}\right) + \mathbb{P}\left(\frac{1}{n} \sum_{j=1}^n \mathbf{E}_{g3}^j > \frac{t}{6}\right) \\ &\leq 2 \exp\left(-c'_2 n \min\left(\frac{t^2}{\omega_g l^2 \max(\tau^2, \tau^4)}, \frac{t}{\sqrt{\omega_g} l \max(\tau, \tau^2)}\right)\right). \end{aligned}$$

**Step 3: Bound  $\mathbb{P}(\mathbf{E} > t)$**  Finally, we combine the above results and obtain

$$\begin{aligned} \mathbb{P}(\mathbf{E} > t) &\leq \mathbb{P}\left(\mathbf{E}_h > \frac{t}{2}\right) + \mathbb{P}\left(\mathbf{E}_g > \frac{t}{2}\right) \\ &\leq 5 \exp\left(-c_h n \min\left(\frac{t^2}{\tau^2 l^2 \max(\omega_g, \omega_g \tau^2, \omega_h)}, \frac{t}{\sqrt{\omega_g} l \max(\tau, \tau^2)}\right)\right), \end{aligned}$$

where  $\omega_g = r^{4(l-1)}$  and  $\omega_h = r^{2(l-2)}$ . □

### C.2.5 PROOF OF LEMMA 10

*Proof.* Before proving our conclusion, we first give an inequality:

$$\|\mathbf{e}\|_2^2 = \|\mathbf{B}_{l:1} \mathbf{x} - \mathbf{y}\|_2^2 \leq \|\mathbf{B}_{l:1} \mathbf{x}\|_2^2 + 2 |\mathbf{y}^T \mathbf{B}_{l:1} \mathbf{x}| + \|\mathbf{y}\|_2^2 \stackrel{\text{①}}{\leq} r_x^2 \omega_f^2 + 2r_x \omega_f \|\mathbf{y}\|_2 + \|\mathbf{y}\|_2^2,$$

where  $\omega_f = r^l$ . Notice, ① holds since by Lemma 14, we have  $\|\mathbf{B}_{l:1}\|_F^2 \leq r^{2l}$ .

Then we consider  $\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})$ . Firstly, by Lemma 6 we can bound  $\|\nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}, \mathbf{x})\|_2^2$  as follows:

$$\begin{aligned} \|\nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}, \mathbf{x})\|_2^2 &= \|((\mathbf{B}_{j-1:1} \mathbf{x}) \otimes \mathbf{B}_{l:j+1}^T) \mathbf{e}\|_2^2 \leq \|\mathbf{B}_{j-1:1}\|_2^2 \|\mathbf{x}\|_2^2 \|\mathbf{B}_{l:j+1}\|_2^2 \|\mathbf{e}\|_2^2 \\ &\stackrel{\textcircled{1}}{\leq} r_x^2 \omega_{f_1}^2 \left( r_x^2 \omega_f^2 + 2r_x \omega_f \|\mathbf{y}\|_2 + \|\mathbf{y}\|_2^2 \right), \end{aligned}$$

where  $\omega_{f_1} = r^{(l-1)}$ . ① holds since we have  $\|\mathbf{B}_{l:j+1}\|_F^2 \|\mathbf{B}_{j-1:1}\|_F^2 \leq r^{2(l-1)}$  by using  $\|\mathbf{W}^{(i)}\|_F^2 = \|\mathbf{w}_{(i)}\|_2^2 \leq r^2$ . Therefore, we can further obtain

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})\|_2^2 = \sum_{i=1}^l \|\nabla_{\mathbf{w}_{(i)}} f(\mathbf{w}, \mathbf{x})\|_2^2 \leq l r_x^2 \omega_{f_1}^2 \left( r_x^2 \omega_f^2 + 2r_x \omega_f \|\mathbf{y}\|_2 + \|\mathbf{y}\|_2^2 \right).$$

Notice,  $\mathbf{y}$  is the label of sample and the weight magnitude  $r$  is usually larger than 1. Then we have  $\|\mathbf{y}\|_2 \leq r^l$ . Also, the values in input data are usually smaller than  $r^l$ . Thus, we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})\|_2^2 \leq c_t l r_x^4 r^{4l-2} \triangleq \alpha_g,$$

where  $c_t$  is a constant. Then we use the inequality  $\|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_{\text{op}} \leq \|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_F$  to bound  $\|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_{\text{op}}$ . Next we only need to give the upper bound of  $\|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_F$ . Let  $\omega_{f_2} = r^{l-2}$ .

We first consider  $\mathbf{Q}_{st} \triangleq \nabla_{\mathbf{w}_{(s)}} (\nabla_{\mathbf{w}_{(t)}} f(\mathbf{w}, \mathbf{x}))$ . By Lemma 6, if  $s < t$ , we have

$$\begin{aligned} \|\mathbf{Q}_{st}\|_F^2 &= \|(\mathbf{B}_{t-1:s+1}^T) \otimes (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:t+1}^T) + (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1})\|_F^2 \\ &\leq 2 \left( \|(\mathbf{B}_{t-1:s+1}^T) \otimes (\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{e}^T \mathbf{B}_{l:t+1}^T)\|_F^2 + \|(\mathbf{B}_{s-1:1} \mathbf{x} \mathbf{x}^T \mathbf{B}_{t-1:1}^T) \otimes (\mathbf{B}_{l:s+1}^T \mathbf{B}_{l:t+1})\|_F^2 \right) \\ &\leq 2 \|\mathbf{B}_{t-1:s+1}\|_F^2 \|\mathbf{B}_{s-1:1}\|_F^2 \|\mathbf{x}\|_2^2 \|\mathbf{e}\|_2^2 \|\mathbf{B}_{l:t+1}\|_F^2 \\ &\quad + 2 \|\mathbf{B}_{s-1:1}\|_F^2 \|\mathbf{x}\|_2^2 \|\mathbf{x}\|_2^2 \|\mathbf{B}_{t-1:1}\|_F^2 \|\mathbf{B}_{l:s+1}\|_F^2 \|\mathbf{B}_{l:t+1}\|_F^2 \\ &\stackrel{\textcircled{1}}{\leq} 2\omega_{f_2}^2 r_x^2 \left( r_x^2 \omega_f^2 + r_x \omega_f \|\mathbf{y}\|_2 + \|\mathbf{y}\|_2^2 \right) + 2\omega_{f_1}^4 r_x^4, \end{aligned}$$

where ① holds since we use  $\|\mathbf{B}_{l:t+1}\|_F^2 \|\mathbf{B}_{t-1:s+1}\|_F^2 \|\mathbf{B}_{s-1:1}\|_F^2 \leq \omega_{f_2}^2$  and  $\|\mathbf{B}_{s-1:1}\|_F^2 \|\mathbf{B}_{l:s+1}\|_F^2 \leq \omega_{f_1}^2$ . Note that when  $s \geq t$ , the above inequality also holds. Similarly, consider the values in input data and the values in label, we have

$$\|\mathbf{Q}_{st}\|_F^2 \leq c_{t'} r_x^4 r^{4l-2} \triangleq \alpha_l,$$

where  $c_{t'}$  is a constant. Therefore, we can bound

$$\|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_{\text{op}} \leq \|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_F \leq \sqrt{\sum_{s=1}^l \sum_{t=1}^l \|\mathbf{Q}_{st}\|_F^2} \leq l \sqrt{\alpha_l}.$$

On the other hand, if the activation functions are linear functions,  $f(\mathbf{w}, \mathbf{x})$  is fourth order differentiable when  $l \geq 2$ . This means that  $\nabla_{\mathbf{x}} \nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{x})$  exists. Also since for any input  $\mathbf{x} \in \mathbb{B}^{d_0}(r_x)$  and  $\mathbf{w} \in \Omega$ , we can always find a universal constant  $\alpha_p$  such that

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

We complete the proofs.  $\square$

### C.2.6 PROOF OF LEMMA 11

*Proof.* Recall that the weight of each layer has magnitude bound separately, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$ . Assume that  $\mathbf{w}_{(j)}$  has  $s_j$  non-zero entries. Then we have  $\sum_{j=1}^l s_j = s$ . So here we separately assume  $\mathbf{w}_{\epsilon}^j = \{\mathbf{w}_{1,\epsilon}^j, \dots, \mathbf{w}_{n_{\epsilon}^j,\epsilon}^j\}$  is the  $d_j d_{j-1} \epsilon / d$ -covering net of the ball  $\mathbb{B}^{d_j d_{j-1}}(r)$  which corresponds

to the weight  $\mathbf{w}_{(j)}$  of the  $j$ -th layer. Let  $n_\epsilon^j$  be the  $\epsilon/l$ -covering number. By  $\epsilon$ -covering theory in (Vershynin, 2012), we can have

$$n_\epsilon^j \leq \binom{\mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{s}_j} \left( \frac{3r}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d} \right)^{\mathbf{s}_j} \leq \exp \left( \mathbf{s}_j \log \left( \frac{3r \mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d} \right) \right) = \exp \left( \mathbf{s}_j \log \left( \frac{3rd}{\epsilon} \right) \right).$$

Let  $\mathbf{w} \in \Omega$  be an arbitrary vector. Since  $\mathbf{w} = [\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(l)}]$  where  $\mathbf{w}_{(j)}$  is the weight of the  $j$ -th layer, we can always find a vector  $\mathbf{w}_{k_j}^j$  in  $\mathbf{w}_\epsilon^j$  such that  $\|\mathbf{w}_{(j)} - \mathbf{w}_{k_j}^j\|_2 \leq \mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ . For brevity, let  $j_w \in [n_\epsilon^j]$  denote the index of  $\mathbf{w}_{k_j}^j$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^j$ . Then let  $\mathbf{w}_{k_w} = [\mathbf{w}_{k_1}^1; \dots; \mathbf{w}_{k_j}^j; \dots; \mathbf{w}_{k_l}^l]$ . This means that we can always find a vector  $\mathbf{w}_{k_w}$  such that  $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$ . Now we use the decomposition strategy to bound our goal:

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

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

$$\begin{aligned} \mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w}) \right\|_{\text{op}} \geq t \right\}, \\ \mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla^2 f(\mathbf{w}, \mathbf{x}_{(i)}) - \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}, \\ \mathbf{E}_2 &= \left\{ \sup_{j_w \in [n_\epsilon^j], j=[l]} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}, \\ \mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(\nabla^2 f(\mathbf{w}, \mathbf{x})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}. \end{aligned}$$

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:

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

where  $\textcircled{1}$  holds since by Markov inequality and  $\textcircled{2}$  holds because of Lemma 10.

Therefore, we can set

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

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

Let  $\boldsymbol{\lambda}_{1/4}$  be the  $\frac{1}{4}$ -covering net of  $\mathbb{B}^d(1)$  but it has only  $s$  nonzero entries. So the size of its  $\epsilon$ -net is

$$\binom{d}{s} \left(\frac{3}{1/4}\right)^s \leq \exp(s \log(12d)).$$

Recall that we use  $j_w$  to denote the index of  $\mathbf{w}_{k_w}^j$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^j$  and we have  $j_w \in [n_\epsilon^j]$ , ( $n_\epsilon^j \leq \exp(s_j \log(\frac{3rd}{\epsilon}))$ ). Then we can bound  $\mathbb{P}(\mathbf{E}_2)$  as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P}\left(\sup_{j_w \in [n_\epsilon^j], j \in [l]} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3}\right) \\ &\leq \mathbb{P}\left(\sup_{j_w \in [n_\epsilon^j], j \in [l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/4}} 2 \left| \left\langle \boldsymbol{\lambda}, \left( \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) \right) \boldsymbol{\lambda} \right\rangle \right| \geq \frac{t}{3}\right) \\ &\leq \exp(s \log(12d)) \exp\left(\sum_{j=1}^l s_j \log\left(\frac{3rd}{\epsilon}\right)\right) \sup_{j_w \in [n_\epsilon^j], j \in [l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/4}} \mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n \left\langle \boldsymbol{\lambda}, \right. \right. \right. \\ &\quad \left. \left. \left. (\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}))) \right\rangle \right| \geq \frac{t}{6}\right) \\ &\stackrel{\textcircled{1}}{\leq} \exp\left(s \log\left(\frac{36rd^2}{\epsilon}\right)\right) 10 \exp\left(-c_{h'} n \min\left(\frac{t^2}{36\tau^2 l^2 \max(\omega_g, \omega_g \tau^2, \omega_h)}, \frac{t}{6\sqrt{\omega_g} l \max(\tau, \tau^2)}\right)\right), \end{aligned}$$

where  $\textcircled{1}$  holds since by Lemma 9, we have

$$\begin{aligned} \mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n \left\langle \boldsymbol{\lambda}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E}_{\mathbf{w}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})) \boldsymbol{\lambda} \right\rangle \right| > t\right) \\ \leq 10 \exp\left(-c_{h'} n \min\left(\frac{t^2}{\tau^2 l^2 \max(\omega_g, \omega_g \tau^2, \omega_h)}, \frac{t}{\sqrt{\omega_g} l \max(\tau, \tau^2)}\right)\right), \end{aligned}$$

where  $\omega_g = r^{4(l-1)}$  and  $\omega_h = r^{2(l-2)}$ .

Let  $d_\epsilon = s \log(36d^2r/\epsilon) + \log(20/\epsilon)$ . Thus, if we set

$$t \geq \max \left( \sqrt{\frac{36\tau^2 l^2 \max(\omega_g, \omega_g \tau^2, \omega_h) d_\epsilon}{c_{h'} n}}, \frac{6\sqrt{\omega_g} l \max(\tau, \tau^2) d_\epsilon}{c_{h'} n} \right),$$

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:

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

where  $\textcircled{1}$  holds because of Lemma 10. We set  $\epsilon$  enough small such that  $\alpha_p \epsilon < t/3$  always holds. Then it yields  $\mathbb{P}(\mathbf{E}_3) = 0$ .

**Step 4. Final result:** For brevity, let  $\omega_2 = 36\tau^2 l^2 \max(\omega_g, \omega_g \tau^2, \omega_h)$  and  $\omega_3 = 6\sqrt{\omega_g} l \max(\tau, \tau^2)$ . To ensure  $\mathbb{P}(\mathbf{E}_0) \leq \epsilon$ , we just set  $\epsilon = 36rl/n$  and

$$\begin{aligned} t &\geq \max \left( \frac{6\alpha_p \epsilon}{\epsilon}, 3\alpha_p \epsilon, \sqrt{\frac{\omega_2 (s \log(36d^2r/\epsilon) + \log(20/\epsilon))}{c_{h'} n}}, \frac{\omega_3 (s \log(36d^2r/\epsilon) + \log(20/\epsilon))}{c_{h'} n} \right) \\ &= \max \left( \frac{216\alpha_p r}{n\epsilon}, \sqrt{\frac{\omega_2 (s \log(d^2nl) + \log(20/\epsilon))}{c_{h'} n}}, \frac{\omega_3 (s \log(36d^2n/l) + \log(20/\epsilon))}{c_{h'} n} \right). \end{aligned}$$

Thus, there exist two universal constants  $c_{h_1}$  and  $c_{h_2}$  such that if  $n \geq c_{h_2} \max(\frac{\alpha_p^2 r^2}{\tau^2 l^2 \omega_h^2 \epsilon^2 s \log(d/l)}, s \log(d/l)/(l\tau^2))$ , then

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w}) \right\|_{\text{op}} \leq c_{h_1} \tau l \omega_h \sqrt{\frac{d \log(nl) + \log(20/\epsilon)}{n}}$$

holds with probability at least  $1 - \epsilon$ , where  $\omega_h = \max(\tau r^{2(l-1)}, r^{2(l-2)}, r^{l-2})$ . The proof is completed.  $\square$

### C.3 PROOFS OF MAIN THEOREMS

#### C.3.1 PROOF OF THEOREM 1

*Proof.* Recall that the weight of each layer has magnitude bound separately, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$ . Assume that  $\mathbf{w}_{(j)}$  has  $s_j$  non-zero entries. Then we have  $\sum_{j=1}^l s_j = s$ . So here we separately assume  $\mathbf{w}_\epsilon^j = \{\mathbf{w}_1^j, \dots, \mathbf{w}_{n_\epsilon^j}^j\}$  is the  $\mathbf{d}_j \mathbf{d}_{j-1} \epsilon/d$ -covering net of the ball  $\mathbf{B}^{\mathbf{d}_j \mathbf{d}_{j-1}}(r)$  which corresponds to the weight  $\mathbf{w}_{(j)}$  of the  $j$ -th layer. Let  $n_\epsilon^j$  be the  $\epsilon/l$ -covering number. By  $\epsilon$ -covering theory in (Vershynin, 2012), we can have

$$n_\epsilon^j \leq \binom{\mathbf{d}_j \mathbf{d}_{j-1}}{s_j} \left( \frac{3r}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon/d} \right)^{s_j} \leq \exp \left( s_j \log \left( \frac{3r \mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon/d} \right) \right) = \exp \left( s_j \log \left( \frac{3rd}{\epsilon} \right) \right).$$

Let  $\mathbf{w} \in \Omega$  be an arbitrary vector. Since  $\mathbf{w} = [\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(l)}]$  where  $\mathbf{w}_{(j)}$  is the weight of the  $j$ -th layer, we can always find a vector  $\mathbf{w}_{k_j}^j$  in  $\mathcal{W}_\epsilon^j$  such that  $\|\mathbf{w}_{(j)} - \mathbf{w}_{k_j}^j\|_2 \leq \mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ . For brevity, let  $j_w \in [n_\epsilon^j]$  denote the index of  $\mathbf{w}_{k_j}^j$  in  $\epsilon$ -net  $\mathcal{W}_\epsilon^j$ . Then let  $\mathbf{w}_{k_w} = [\mathbf{w}_{k_1}^{j_1}; \dots; \mathbf{w}_{k_j}^{j_j}; \dots; \mathbf{w}_{k_l}^{j_l}]$ . This means that we can always find a vector  $\mathbf{w}_{k_w}$  such that  $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$ . Accordingly, we can decompose  $\|\nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w})\|_2$  as

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

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

$$\begin{aligned} \mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right\|_2 \geq t \right\}, \\ \mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f(\mathbf{w}, \mathbf{x}_{(i)}) - \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right\|_2 \geq \frac{t}{3} \right\}, \\ \mathbf{E}_2 &= \left\{ \sup_{j_w \in [n_\epsilon^j], j=[l]} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3} \right\}, \\ \mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(\nabla f(\mathbf{w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3} \right\}. \end{aligned}$$

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:

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

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  $\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right)$ . Now we utilize Lemma 10 to achieve this goal:

$$\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right) \leq \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 f(\mathbf{w}, \mathbf{x}) - \nabla^2 f(\mathbf{w}^*, \mathbf{x}) \right\|_2 \right) \leq l\sqrt{\alpha_l}.$$

where  $\alpha_l = c_l r_x^4 r^{4l-2}$ . Therefore, we have

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

We further let

$$t \geq \frac{6l\sqrt{\alpha_l}\epsilon}{\epsilon}.$$

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 1, we know that for any vector  $\mathbf{x} \in \mathbb{R}^d$ , its  $\ell_2$ -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  $\boldsymbol{\lambda}_\epsilon = \{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k_w}\}$  be an  $\epsilon$ -covering net of  $\mathbb{B}^d(1)$ .

Let  $\boldsymbol{\lambda}_{1/2}$  be the  $\frac{1}{2}$ -covering net of  $\mathbb{B}^d(1)$  but it has only  $s$  nonzero entries. So the size of its  $\epsilon$ -net is

$$\binom{d}{s} \left( \frac{3}{1/2} \right)^s \leq \exp(s \log(6d)).$$

Recall that we use  $j_w$  to denote the index of  $\mathbf{w}_{k_j}^j$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^j$  and we have  $j_w \in [n_\epsilon^j]$ , ( $n_\epsilon^j \leq \exp(s_j \log(\frac{3rd}{\epsilon}))$ ). Then we can bound  $\mathbb{P}(\mathbf{E}_2)$  as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P} \left( \sup_{j_w \in [n_\epsilon^j], j=[l]} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3} \right) \\ &= \mathbb{P} \left( \sup_{j_w \in [n_\epsilon^j], j=[l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/2}} 2 \left\langle \boldsymbol{\lambda}, \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\rangle \geq \frac{t}{3} \right) \\ &\leq \exp(s \log(6d)) \exp \left( \sum_{j=1}^l s_j \log \left( \frac{3rd}{\epsilon} \right) \right) \sup_{j_w \in [n_\epsilon^j], j=[l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/2}} \mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n \left\langle \boldsymbol{\lambda}, \right. \right. \\ &\quad \left. \left. \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\rangle \geq \frac{t}{6} \right) \\ &\stackrel{\textcircled{1}}{\leq} \exp \left( s \log \left( \frac{18rd}{\epsilon} \right) \right) 6 \exp \left( -c_{g'} n \min \left( \frac{t^2}{36l \max(\omega_g \tau^2, \omega_g \tau^4, \omega_{g'} \tau^2)}, \frac{t}{6\sqrt{l\omega_g} \max(\tau, \tau^2)} \right) \right), \end{aligned}$$

where  $\textcircled{1}$  holds since by Lemma 8, we have

$$\begin{aligned} &\mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}_{(i)}) \rangle) > t \right) \\ &\leq 3 \exp \left( -c_{g'} n \min \left( \frac{t^2}{l \max(\omega_g \tau^2, \omega_g \tau^4, \omega_{g'} \tau^2)}, \frac{t}{\sqrt{l\omega_g} \max(\tau, \tau^2)} \right) \right), \end{aligned}$$

where  $c_{g'}$  is a constant;  $\omega_g = c_q r^{2(2l-1)}$  and  $\omega_{g'} = c_q r^{2(l-1)}$  in which  $c_q = \sqrt{\max_{0 \leq i \leq l} d_i}$ .

Let  $\omega_2 = 36l \max(\omega_g \tau^2, \omega_g \tau^4, \omega_{g'} \tau^2)$  and  $\omega_3 = 6\sqrt{l\omega_g} \max(\tau, \tau^2)$ . Thus, if we set

$$t \geq \max \left( \sqrt{\frac{\omega_2(s \log(18dr/\epsilon) + \log(12/\epsilon))}{c_{g'} n}}, \frac{\omega_3(s \log(18dr/\epsilon) + \log(12/\epsilon))}{c_{g'} n} \right),$$

then we have

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

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

$$\begin{aligned} \mathbb{P}(\mathbf{E}_3) &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \|\mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(f(\mathbf{w}, \mathbf{x}))\|_2 \geq \frac{t}{3}\right) \\ &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \frac{\|\mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x}) - f(\mathbf{w}, \mathbf{x}))\|_2}{\|\mathbf{w} - \mathbf{w}_{k_w}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \geq \frac{t}{3}\right) \\ &\leq \mathbb{P}\left(\epsilon \mathbb{E} \sup_{\mathbf{w} \in \Omega} \|\nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x})\|_2 \geq \frac{t}{3}\right) \\ &\leq \mathbb{P}\left(l\sqrt{\alpha_l} \epsilon \geq \frac{t}{3}\right). \end{aligned}$$

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

**Step 4. Final result:** Finally, to ensure  $\mathbb{P}(\mathbf{E}_0) \leq \varepsilon$ , we just set  $\epsilon = 18lr/n$  and

$$\begin{aligned} t &\geq \max\left(\frac{6l\sqrt{\alpha_l}\epsilon}{\varepsilon}, 3l\sqrt{\alpha_l}\epsilon, \sqrt{\frac{\omega_2(s \log(18dr/\epsilon) + \log(12/\varepsilon))}{c_{g'}n}}, \frac{\omega_3(s \log(18dr/\epsilon) + \log(12/\varepsilon))}{c_{g'}n}\right) \\ &= \max\left(\frac{108l^2\sqrt{\alpha_l}r}{n\varepsilon}, \sqrt{\frac{\omega_2(s \log(dn/l) + \log(12/\varepsilon))}{c_{g'}n}}, \frac{\omega_3(s \log(dn/l) + \log(12/\varepsilon))}{c_{g'}n}\right). \end{aligned}$$

Notice, we have  $\alpha_l = c_{l'} r_x^4 r^{4l-2}$  where  $c_{l'}$  is a constant. Therefore, there exists two universal constants  $c_g$  and  $c_{g'}$  such that  $n \geq c_{g'} \max\left(\frac{l^3 r_x^2 r^4}{c_q s \log(d/l) \varepsilon^2 \tau^4 \log(1/\varepsilon)}, s \log(d/l) / (l\tau^2)\right)$ , then

$$\sup_{\mathbf{w} \in \Omega} \|\nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w})\|_2 \leq c_g \tau \omega_g \sqrt{l c_q} \sqrt{\frac{s \log(dn/l) + \log(12/\varepsilon)}{n}}$$

holds with probability at least  $1 - \varepsilon$ , where  $\omega_g = \max(\tau r^{2l-1}, r^{2l-1}, r^{l-1})$ .  $\square$

### C.3.2 PROOF OF THEOREM 2

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

$$\inf_i \left| \lambda_i^k \left( \nabla^2 \mathbf{J}(\mathbf{w}^{(k)}) \right) \right| \geq \zeta,$$

where  $\lambda_i^k \left( \nabla^2 \mathbf{J}(\mathbf{w}^{(k)}) \right)$  denotes the  $i$ -th eigenvalue of the Hessian  $\nabla^2 \mathbf{J}(\mathbf{w}^{(k)})$  and  $\zeta$  is a constant. We further define a set  $D = \{\mathbf{w} \in \mathbb{R}^d \mid \|\nabla \mathbf{J}(\mathbf{w})\|_2 \leq \epsilon \text{ and } \inf_i |\lambda_i \left( \nabla^2 \mathbf{J}(\mathbf{w}^{(k)}) \right)| \geq \zeta\}$ . According to Lemma 4,  $D = \cup_{k=1}^{\infty} D_k$  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  $\mathbf{J}(\mathbf{w})$  for  $k \geq m + 1$ . On the other hand, by the continuity of  $\nabla \mathbf{J}(\mathbf{w})$ , it yields  $\|\nabla \mathbf{J}(\mathbf{w})\|_2 = \epsilon$  for  $\mathbf{w} \in \partial D_k$ . Notice, we set the value of  $\epsilon$  blow which is actually a function related to  $n$ .

Then by utilizing Theorem 1, we let sample number  $n$  sufficient large such that

$$\sup_{\mathbf{w} \in \Omega} \|\nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w})\|_2 \leq z_g \triangleq \frac{\epsilon}{2}$$

holds with probability at least  $1 - \varepsilon$ , where if  $n \geq c_{g'} \max\left(\frac{l^3 r_x^2 r^4}{c_q s \log(d/l) \varepsilon^2 \tau^4 \log(1/\varepsilon)}, \frac{s \log(d/l)}{l\tau^2}\right)$ ,  $z_g = c_g \tau \omega_g \sqrt{l c_q} \sqrt{\frac{s \log(dn/l) + \log(12/\varepsilon)}{n}}$ .

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

$$\begin{aligned} \inf_{\mathbf{w} \in D_k} \left\| t \nabla \hat{\mathbf{J}}_n(\mathbf{w}) + (1-t) \nabla \mathbf{J}(\mathbf{w}) \right\|_2 &= \inf_{\mathbf{w} \in D_k} \left\| t \left( \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right) + \nabla \mathbf{J}(\mathbf{w}) \right\|_2 \\ &\geq \inf_{\mathbf{w} \in D_k} \left\| \nabla \mathbf{J}(\mathbf{w}) \right\|_2 - \sup_{\mathbf{w} \in D_k} t \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right\|_2 \\ &\geq \frac{\epsilon}{2}. \end{aligned} \quad (20)$$

Similarly, by utilizing Lemma 11, let  $n$  be sufficient large such that

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

holds with probability at least  $1 - \epsilon$ , where if  $n \geq c_{h_2} \max\left(\frac{\alpha_p^2 r^2}{\tau^2 l^2 \omega_n^2 \epsilon^2 s \log(d/l)}, s \log(d/l)/(l\tau^2)\right)$ ,

$$z_s = c_{h_1} \tau l \omega_h \sqrt{\frac{s \log(nl) + \log(20/\epsilon)}{n}}.$$

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  $\lambda_i^k \left( \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) \right)$  for arbitrary  $\mathbf{w} \in D_k$  as follows:

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

This means that in each set  $D_k$ ,  $\nabla^2 \hat{\mathbf{J}}_n(\mathbf{w})$  has no zero eigenvalues. Then, combine this and Eqn. (20), by Lemma 3 we know that if the population risk  $\mathbf{J}(\mathbf{w})$  has no critical point in  $D_k$ , then the empirical risk  $\hat{\mathbf{J}}_n(\mathbf{w})$  has also no critical point in  $D_k$ ; otherwise it also holds. By Lemma 3, we can also obtain that in  $D_k$ , if  $\mathbf{J}(\mathbf{w})$  has a unique critical point  $\mathbf{w}^{(k)}$  with non-degenerate index  $s_k$ , then  $\hat{\mathbf{J}}_n(\mathbf{w})$  also has a unique critical point  $\mathbf{w}_n^{(k)}$  in  $D_k$  with the same non-degenerate index  $s_k$ . The first conclusion is proved.

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

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

where  $\textcircled{1}$  holds since  $\nabla \mathbf{J}(\mathbf{w}^{(k)}) = \mathbf{0}$  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  $\inf_i \left| \lambda_i \left( \nabla^2 \mathbf{J}(\mathbf{w}) \right) \right| \geq \zeta$ . Consider the conditions in Lemma 11 and Theorem 1, we can obtain that if  $n \geq c_h \max\left(\frac{l^3 r^2 r^4}{c_q s \log(d/l) \epsilon^2 \tau^4 \log(1/\epsilon)}, s \log(d/l)/\zeta^2\right)$  where  $c_h$  is a constant, then

$$\left\| \mathbf{w}_n^{(k)} - \mathbf{w}^{(k)} \right\|_2 \leq \frac{2c_g \tau \omega_g}{\zeta} \sqrt{l c_q} \sqrt{\frac{s \log(dn/l) + \log(12/\epsilon)}{n}}$$

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

### C.3.3 PROOF OF THEOREM 3

*Proof.* Recall that the weight of each layer has magnitude bound separately, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$ . Assume that  $\mathbf{w}_{(j)}$  has  $s_j$  non-zero entries. Then we have  $\sum_{j=1}^l s_j = s$ . So here we separately assume  $\mathbf{w}_\epsilon^j = \{\mathbf{w}_1^j, \dots, \mathbf{w}_{n_\epsilon^j}^j\}$  is the  $\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ -covering net of the ball  $\mathbb{B}^{\mathbf{d}_j \mathbf{d}_{j-1}}(r)$  which corresponds to the weight  $\mathbf{w}_{(j)}$  of the  $j$ -th layer. Let  $n_\epsilon^j$  be the  $\epsilon/l$ -covering number. By  $\epsilon$ -covering theory in (Vershynin, 2012), we can have

$$n_\epsilon^j \leq \binom{\mathbf{d}_j \mathbf{d}_{j-1}}{s_j} \left( \frac{3r}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d} \right)^{s_j} \leq \exp \left( s_j \log \left( \frac{3r \mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d} \right) \right) = \exp \left( s_j \log \left( \frac{3rd}{\epsilon} \right) \right).$$

Let  $\mathbf{w} \in \Omega$  be an arbitrary vector. Since  $\mathbf{w} = [\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(l)}]$  where  $\mathbf{w}_{(j)}$  is the weight of the  $j$ -th layer, we can always find a vector  $\mathbf{w}_{k_w}^j$  in  $\mathbf{w}_\epsilon^j$  such that  $\|\mathbf{w}_{(j)} - \mathbf{w}_{k_w}^j\|_2 \leq \mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ . For brevity, let  $j_w \in [n_\epsilon^j]$  denote the index of  $\mathbf{w}_{k_w}^j$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^j$ . Then let  $\mathbf{w}_{k_w} = [\mathbf{w}_{k_1}^{j_1}; \dots; \mathbf{w}_{k_j}^{j_j}; \dots; \mathbf{w}_{k_l}^{j_l}]$ . This means that we can always find a vector  $\mathbf{w}_{k_w}$  such that  $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$ . Now we use the decomposition strategy to bound our goal:

$$\begin{aligned} & \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| = \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{x})) \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) + \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x})) + \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x}) - \mathbb{E}f(\mathbf{w}, \mathbf{x}) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right| + \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x})) \right| + \left| \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x}) - \mathbb{E}f(\mathbf{w}, \mathbf{x}) \right|. \end{aligned}$$

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

$$\begin{aligned} \mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| \geq t \right\}, \\ \mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right| \geq \frac{t}{3} \right\}, \\ \mathbf{E}_2 &= \left\{ \sup_{j_w \in [n_\epsilon^j], j=[l]} \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x}))) \right| \geq \frac{t}{3} \right\}, \\ \mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(f(\mathbf{w}, \mathbf{x})) \right| \geq \frac{t}{3} \right\}. \end{aligned}$$

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:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_1) &= \mathbb{P} \left( \sup_{\mathbf{w} \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right| \geq \frac{t}{3} \right) \\ &\leq \frac{3}{t} \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right| \right) \\ &\leq \frac{3}{t} \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \frac{\left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right|}{\|\mathbf{w} - \mathbf{w}_{k_w}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \right) \\ &\leq \frac{3\epsilon}{t} \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right), \end{aligned}$$

where 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  $\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{J}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right)$ . Therefore, by Lemma 10, we have

$$\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{J}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right) = \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}, \mathbf{x}_{(i)}) \right\|_2 \right) = \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \|\nabla f(\mathbf{w}, \mathbf{x})\|_2 \right) \leq \sqrt{\alpha_g}.$$

where  $\alpha_g = c_t l r_x^4 r^{4l-2}$ . Therefore, we have

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

We further let

$$t \geq \frac{6\epsilon\sqrt{\alpha_g}}{\epsilon}.$$

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  $j_w$  to denote the index of  $\mathbf{w}_{k_j}^j$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^j$  and we have  $j_w \in [n_\epsilon^j]$ , ( $n_\epsilon^j \leq \exp(s_j \log(\frac{3rd}{\epsilon}))$ ). We can bound  $\mathbb{P}(\mathbf{E}_2)$  as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P} \left( \sup_{j_w \in [n_\epsilon^j], j \in [l]} \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) \right| \geq \frac{t}{3} \right) \\ &\leq \exp \left( \sum_{j=1}^l s_j \log \left( \frac{3rd}{\epsilon} \right) \right) \sup_{j_w \in [n_\epsilon^j], j \in [l]} \mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) \right| \geq \frac{t}{3} \right) \\ &\stackrel{\textcircled{1}}{\leq} 4 \left( \frac{3dr}{\epsilon} \right)^s \exp \left( -c_{f'} n \min \left( \frac{t^2}{9\omega_f^2 \max(\mathbf{d}_l \omega_f^2 \tau^4, \tau^2)}, \frac{t}{3\omega_f^2 \tau^2} \right) \right), \end{aligned}$$

where ① holds because in Lemma 7, we have

$$\mathbb{P} \left( \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{x}_{(i)}))) > t \right) \leq 2 \exp \left( -c_{f'} n \min \left( \frac{t^2}{\omega_f^2 \max(\mathbf{d}_l \omega_f^2 \tau^4, \tau^2)}, \frac{t}{\omega_f^2 \tau^2} \right) \right),$$

where  $c_{f'}$  is a positive constant and  $\omega_f = r^l$ . Thus, if we set

$$t \geq \max \left( \sqrt{\frac{9\omega_f^2 (s \log(3rd/\epsilon) + \log(8/\epsilon)) \max(\mathbf{d}_l \omega_f^2 \tau^4, \tau^2)}{c_{f'} n}}, \frac{3\omega_f^2 \tau^2 (s \log(3rd/\epsilon) + \log(8/\epsilon))}{c_{f'} n} \right),$$

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:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_3) &= \mathbb{P} \left( \sup_{\mathbf{w} \in \Omega} \|\mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(f(\mathbf{w}, \mathbf{x}))\|_2 \geq \frac{t}{3} \right) \\ &= \mathbb{P} \left( \sup_{\mathbf{w} \in \Omega} \frac{\|\mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) - f(\mathbf{w}, \mathbf{x})\|_2}{\|\mathbf{w} - \mathbf{w}_{k_w}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \geq \frac{t}{3} \right) \\ &\leq \mathbb{P} \left( \epsilon \mathbb{E} \sup_{\mathbf{w} \in \Omega} \|\nabla J_{\mathbf{w}}(\mathbf{w}, \mathbf{x})\|_2 \geq \frac{t}{3} \right) \\ &\stackrel{\textcircled{1}}{\leq} \mathbb{P} \left( \sqrt{\alpha_g} \epsilon \geq \frac{t}{3} \right), \end{aligned}$$

where ① holds since we utilize Lemma 10. We set  $\epsilon$  enough small such that  $\sqrt{\alpha_g}\epsilon < t/3$  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  $\epsilon = 3rl/n$ . Note that  $\frac{6\sqrt{\alpha_g}\epsilon}{\epsilon} > 3\sqrt{\alpha_g}\epsilon$ . Thus we can obtain

$$\begin{aligned} t &\geq \max \left( \frac{6\sqrt{\alpha_g}\epsilon}{\epsilon}, \sqrt{\frac{9\omega_f^2(s \log(3rd/\epsilon) + \log(8/\epsilon)) \max(\mathbf{d}_l \omega_f^2 \tau^4, \tau^2)}{c_f n}}, \frac{3\omega_f^2 \tau^2 (s \log(3rd/\epsilon) + \log(8/\epsilon))}{c_f n} \right) \\ &= \max \left( \frac{18l\sqrt{\alpha_g}r}{n\epsilon}, \sqrt{\frac{9\omega_f^2(s \log(dn/l) + \log(8/\epsilon)) \max(\mathbf{d}_l \omega_f^2 \tau^4, \tau^2)}{c_f n}}, \frac{3\omega_f^2 \tau^2 (s \log(dn/l) + \log(8/\epsilon))}{c_f n} \right). \end{aligned}$$

Note that we have  $\alpha_g = c_t l r_x^4 r^{4l-2}$  where  $c_t$  is a constant. Then there exist four universal constants  $c_f$  and  $c_{f'}$  such that if  $n \geq c_{f'} \max \left( \frac{l^3 r_x^4}{\mathbf{d}_l s \log(d) \epsilon^2 \tau^4 \log(1/\epsilon)}, s \log(d) / (\tau^2 \mathbf{d}_l) \right)$ , then

$$\sup_{\mathbf{w} \in \Omega} \left\| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right\|_2 \leq c_f \omega_f \tau \max \left( \sqrt{\mathbf{d}_l} \omega_f \tau, 1 \right) \sqrt{\frac{s \log(dn/l) + \log(8/\epsilon)}{n}}$$

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

#### C.3.4 PROOF OF COROLLARY 1

*Proof.* By Lemma 5, we know  $\epsilon_s = \epsilon_g$ . Thus, the remaining work is to bound  $\epsilon_s$ . Actually, we can have

$$\begin{aligned} \left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}, (\mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) \sim \mathcal{D}} \frac{1}{n} \sum_{j=1}^n \left( f_j(\mathbf{w}^j; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) - f_j(\mathbf{w}^n; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) \right) \right| &\leq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \left( \sup_{\mathbf{w} \in \Omega} \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| \right) \\ &\leq \sup_{\mathbf{w} \in \Omega} \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| \\ &\leq \epsilon_f. \end{aligned}$$

Thus, we have  $\epsilon_g = \epsilon_s \leq \epsilon_f$ . The proof is completed.  $\square$

### C.4 PROOF OF OTHER LEMMAS

#### C.4.1 PROOF OF LEMMA 13

**Lemma 15.** (Rigollet, 2015) Suppose a random variable  $x$  is  $\tau^2$ -sub-Gaussian, then the random variable  $x^2 - \mathbb{E}x^2$  is sub-exponential and obeys:

$$\mathbb{E} \left( \exp \lambda (x^2 - \mathbb{E}x^2) \right) \leq \exp \left( \frac{256\lambda^2\tau^4}{2} \right), \quad |\lambda| \leq \frac{1}{16\tau^2}. \quad (21)$$

*Proof.* Here we utilize Lemma 15 to prove our conclusion. We have

$$\begin{aligned} \mathbb{E} \exp \left( \lambda \left( \sum_{i=1}^k \mathbf{a}_i \mathbf{x}_i^2 - \mathbb{E} \left( \sum_{i=1}^k \mathbf{a}_i \mathbf{x}_i^2 \right) \right) \right) &\stackrel{\textcircled{1}}{=} \prod_{i=1}^k \mathbb{E} \exp \left( \lambda \mathbf{a}_i (\mathbf{x}_i^2 - \mathbb{E} \mathbf{x}_i^2) \right) \\ &\stackrel{\textcircled{2}}{\leq} \prod_{i=1}^k \mathbb{E} \exp \left( 128\lambda^2 \mathbf{a}_i^2 \tau_i^4 \right), \quad |\lambda| \leq \frac{1}{\max_i \mathbf{a}_i \tau^2} \\ &\leq \mathbb{E} \exp \left( 128\lambda^2 \tau^4 \left( \sum_{i=1}^k \mathbf{a}_i^2 \right) \right), \end{aligned}$$

where ① holds since  $\mathbf{x}_i$  are independent and ② holds because of Lemma 15.  $\square$

### C.4.2 PROOF OF LEMMA 14

*Proof.* Since the  $\ell_2$ -norm of each  $\mathbf{w}_{(j)}$  is bounded, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$  ( $1 \leq j \leq l$ ), we can obtain

$$\|\mathbf{B}_{s:t}\|_F^2 \leq \left\| \mathbf{W}^{(s)} \right\|_F^2 \left\| \mathbf{W}^{(s-1)} \right\|_F^2 \cdots \left\| \mathbf{W}^{(t)} \right\|_F^2 \leq r^{2(t-s+1)} \triangleq \omega_r^2 \stackrel{\textcircled{1}}{\leq} \max(r^2, r^{2l}),$$

where  $\textcircled{1}$  holds since the function  $r^{2x}$  obtains its maximum at two endpoints  $x = 1$  and  $x = l$  for case  $r < 1$  and  $r \geq 1$ , respectively. On the other hand, we have  $\|\mathbf{B}_{s:t}\|_{\text{op}} \leq \|\mathbf{B}_{s:t}\|_F \leq \omega_r$ . Specifically, we have  $\|\mathbf{B}_{l:1}\|_F^2 \leq r^{2l} \triangleq \omega_f^2$ .  $\square$

## D PROOFS FOR DEEP NONLINEAR NEURAL NETWORKS

In this section, we first present the technical lemmas in Sec. D.1. Then in Sec. D.2 we give the proofs of these lemmas. Next, we utilize these technical lemmas to prove the results in Theorems 4 ~ 6 and Corollary 2 in Sec. D.3. Finally, we give the proofs of other lemmas in Sec. D.4.

### D.1 TECHNICAL LEMMAS

Here we present the key lemmas and theorems for proving our desired results. For brevity, we define an operation  $\mathbb{G}$  which maps an arbitrary vector  $\mathbf{z} \in \mathbb{R}^k$  into a diagonal matrix  $\mathbb{G}(\mathbf{z}) \in \mathbb{R}^{k \times k}$  with its  $i$ -th diagonal entry equal to  $\sigma(z_i)(1 - \sigma(z_i))$  in which  $z_i$  denotes the  $i$ -th entry of  $\mathbf{z}$ . We further define  $\mathbf{A}_i \in \mathbb{R}^{d_{i-1} \times d_i}$  as follows:

$$\mathbf{A}_i = (\mathbf{W}^{(i)})^T \mathbb{G}(\mathbf{u}^{(i)}) \in \mathbb{R}^{d_{i-1} \times d_i} \quad (i = 1, \dots, l), \quad (22)$$

where  $\mathbf{W}^{(i)}$  is the weight matrix in the  $i$ -th layer and  $\mathbf{u}^{(i)}$  is the linear output of the  $i$ -th layer. In this section, we define

$$\mathbf{B}_{s:t} = \mathbf{A}_s \mathbf{A}_{s+1} \cdots \mathbf{A}_t \in \mathbb{R}^{d_{s-1} \times d_t}, \quad (s \leq t) \quad \text{and} \quad \mathbf{B}_{s:t} = \mathbf{I}, \quad (s > t). \quad (23)$$

**Lemma 16.** *Suppose that the activation function in deep neural network are sigmoid functions. Then the gradient of  $f(\mathbf{w}, \mathbf{x})$  with respect to  $\mathbf{w}_{(j)}$  can be formulated as*

$$\nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}, \mathbf{x}) = \text{vec} \left( \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \right) (\mathbf{v}^{(j-1)})^T \right), \quad (j = 1, \dots, l-1),$$

and

$$\nabla_{\mathbf{w}_{(l)}} f(\mathbf{w}, \mathbf{x}) = \text{vec} \left( \left( \mathbb{G}(\mathbf{u}^{(l)}) (\mathbf{v}^{(l)} - \mathbf{y}) \right) (\mathbf{v}^{(l-1)})^T \right).$$

Besides, the loss  $f(\mathbf{w}, \mathbf{x})$  is  $\alpha$ -Lipschitz,

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})\|_2 \leq \alpha,$$

where  $\alpha = \sqrt{\frac{1}{16} c_y c_d (1 + c_r (l-1))}$  in which  $c_y$ ,  $c_d$  and  $c_r$  are defined as

$$\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 \leq c_y < +\infty, \quad c_d = \max(\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_l) \quad \text{and} \quad c_r = \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^{l-1}\right).$$

**Lemma 17.** *Suppose that the activation functions in deep neural network are sigmoid functions. Then there exists two universal constants  $c_{s_1}$  and  $c_{s_2}$  such that*

$$\|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_{\text{op}} \leq \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_F \leq \varsigma,$$

where  $\varsigma = \sqrt{c_{s_1} c_r c_d^2 l^4}$  in which  $c_d = \max_i \mathbf{d}_i$  and  $c_r = \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^{l-1}\right)$ . Moreover, the gradient  $\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})$  is  $\varsigma$ -Lipschitz, i.e.

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}_1, \mathbf{x}) - \nabla_{\mathbf{w}} f(\mathbf{w}_2, \mathbf{x})\|_2 \leq \varsigma \|\mathbf{w}_1 - \mathbf{w}_2\|_2.$$

Similarly, there also exist a universal constant  $\xi$  such that

$$\|\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{x})\|_{\text{op}} \leq \|\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{x})\|_F \leq \xi.$$

**Lemma 18.** Suppose that the activation function in deep neural network are sigmoid functions. Then we have

$$\|\nabla_{\mathbf{w}} \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x})\|_{op} \leq \|\nabla_{\mathbf{w}} \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x})\|_F \leq \beta,$$

where  $\beta = \sqrt{\frac{2^6}{3^8} l(l+2)c_y c_r c_d (lc_r + 1)}$  in which  $c_y$ ,  $c_d$  and  $c_r$  are defined in Lemma 16.

**Lemma 19.** Suppose that the input sample  $\mathbf{x}$  obeys Assumption 2 and the activation functions in deep neural network are sigmoid functions. The gradient of the loss is  $8\beta^2\tau^2$ -sub-Gaussian. Specifically, for any  $\boldsymbol{\lambda} \in \mathbb{R}^d$ , we have

$$\mathbb{E}(\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \rangle) \leq \exp\left(\frac{8\beta^2\tau^2 \|\boldsymbol{\lambda}\|_2^2}{2}\right),$$

where  $\beta = \sqrt{\frac{2^6}{3^8} l(l+2)c_y c_r c_d (lc_r + 1)}$  in which  $c_y$ ,  $c_d$  and  $c_r$  are defined in Lemma 16.

**Lemma 20.** Suppose that the input sample  $\mathbf{x}$  obeys Assumption 2 and the activation functions in deep neural network are sigmoid functions. The Hessian of the loss, evaluated on a unit vector, is sub-Gaussian. Specifically, for any unit  $\boldsymbol{\lambda} \in \mathbb{S}^{d-1}$  (i.e.  $\|\boldsymbol{\lambda}\|_2 = 1$ ), there exist universal constant  $\gamma$  such that

$$\mathbb{E}(t \langle \boldsymbol{\lambda}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})) \boldsymbol{\lambda} \rangle) \leq \exp\left(\frac{8t^2\gamma^2\tau^2}{2}\right).$$

Notice,  $\gamma$  obeys  $\gamma \geq \|\nabla_{\mathbf{x}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_{op}$ .

**Lemma 21.** Assume that the input sample  $\mathbf{x}$  obeys Assumption 2 and the activation functions in deep neural network are sigmoid functions. Then the sample Hessian uniformly converges to the population Hessian in operator norm. That is, there exists such two universal constants  $c_{m'}$  and  $c_m$  such that if  $n \geq \frac{c_{m'} \xi^2 l^2 r^2}{\gamma^2 \tau^2 \varepsilon^2 s \log(d) \log(1/\varepsilon)}$ , then

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w}) \right\|_{op} \leq c_m \gamma \tau \sqrt{\frac{s \log(dn/l) + \log(4/\varepsilon)}{n}}$$

holds with probability at least  $1 - \varepsilon$ . Here  $\gamma$  is the same parameter in Lemma 20.

## D.2 PROOFS OF TECHNICAL LEMMAS

For brevity, we also define

$$\mathbf{D}_{s:t} = \|\mathbf{W}^{(s)}\|_F^2 \cdots \|\mathbf{W}^{(t)}\|_F^2 \quad (s \leq t) \quad \text{and} \quad \mathbf{D}_{s:t} = 1, \quad (s > t).$$

We define a matrix  $\mathbf{P}_k \in \mathbb{R}^{d_k^2 \times d_k}$  whose  $((s-1)d_k + s, s)$  ( $s = 1, \dots, d_k$ ) entry equal to  $\sigma(\mathbf{u}_s^{(k)})(1 - \sigma(\mathbf{u}_s^{(k)}))(1 - 2\sigma(\mathbf{u}_s^{(k)}))$  and rest entries are all 0. On the other hand, since the values in  $\mathbf{v}^{(l)}$  belong to the range  $[0, 1]$  and  $\mathbf{y}$  is the label,  $\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2$  can be bounded:

$$\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 \leq c_y < +\infty,$$

where  $c_y$  is a universal constant. We further define  $c_d = \max(d_0, d_1, \dots, d_l)$ .

Then we give a lemma to summarize the properties of  $\mathbb{G}(\mathbf{u}^{(i)})$  defined in Eqn. (22),  $\mathbf{B}_{s:t}$  defined in Eqn. (23),  $\mathbf{D}_{s:t}$  and  $\mathbf{P}_k$ .

**Lemma 22.** For  $\mathbb{G}(\mathbf{u}^{(i)})$  defined in Eqn. (22),  $\mathbf{B}_{s:t}$  defined in Eqn. (23),  $\mathbf{D}_{s:t}$  and  $\mathbf{P}_k$ , we have the following properties:

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

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

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

$$\|\mathbf{P}_k \mathbf{M}\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N} \mathbf{P}_k\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{N}\|_F^2.$$

(3) For arbitrary matrices  $\mathbf{M}$  and  $\mathbf{N}$  of proper sizes, we have

$$\|\mathbf{B}_{s:t}\|_F^2 \leq \frac{1}{16^{t-s+1}} \mathbf{D}_{s:t} \quad \text{and} \quad \frac{1}{16^{t-s+1}} \mathbf{D}_{s:t} \leq c_{st} \leq c_r,$$

where  $c_{st} = \left(\frac{r}{4}\right)^{2(t-s+1)}$  and  $c_r = \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^{l-1}\right)$ .

(4) For arbitrary matrices  $\mathbf{M}$ ,  $\mathbf{N}$  and  $\mathbf{I}$  of proper sizes, let  $\mathbf{m} = \text{vec}(\mathbf{M})$ . Then we have

$$\|(\mathbf{N} \otimes \mathbf{I})\mathbf{m}\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2 \quad \text{and} \quad \|(\mathbf{I} \otimes \mathbf{N})\mathbf{m}\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2.$$

It should be pointed out that we defer the proof of Lemma 22 to Sec. D.4.

### D.2.1 PROOF OF LEMMA 16

*Proof.* We use chain rule to compute the gradient of  $f(\mathbf{w}, \mathbf{x})$  with respect to  $\mathbf{w}_{(j)}$ . We first compute several basis gradient. According to the relationship between  $\mathbf{u}^{(j)}$ ,  $\mathbf{v}^{(j)}$ ,  $\mathbf{W}^{(j)}$  and  $f(\mathbf{w}, \mathbf{x})$ , we have

$$\begin{aligned} \nabla_{\mathbf{v}^{(l)}} f(\mathbf{w}, \mathbf{x}) &= \mathbf{v}^{(l)} - \mathbf{y}, \\ \nabla_{\mathbf{v}^{(i)}} f(\mathbf{w}, \mathbf{x}) &= \frac{\partial \mathbf{u}^{(i+1)}}{\partial \mathbf{v}^{(i)}} \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i+1)}} = (\mathbf{W}^{(i+1)})^T \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i+1)}}, \quad (i = 1, \dots, l-1), \\ \nabla_{\mathbf{u}^{(i)}} f(\mathbf{w}, \mathbf{x}) &= \frac{\partial \mathbf{v}^{(i)}}{\partial \mathbf{u}^{(i)}} \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{v}^{(i)}} = \mathbb{G}(\mathbf{u}^{(i)}) \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{v}^{(i)}}, \quad (i = 1, \dots, l), \\ \nabla_{\mathbf{W}^{(i)}} f(\mathbf{w}, \mathbf{x}) &= \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{w}^{(i)}} \left( \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}} \right)^T = \mathbf{v}^{(i-1)} \left( \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}} \right)^T, \quad (i = 1, \dots, l). \end{aligned} \quad (24)$$

Then by chain rule, we can easily compute the gradient of  $f(\mathbf{w}, \mathbf{x})$  with respect to  $\mathbf{w}_{(j)}$  which is formulated as

$$\nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}, \mathbf{x}) = \text{vec} \left( \mathbf{v}^{(j-1)} \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \right)^T \right), \quad (j = 1, \dots, l-1),$$

and

$$\nabla_{\mathbf{w}_{(l)}} f(\mathbf{w}, \mathbf{x}) = \text{vec} \left( \mathbf{v}^{(l-1)} \left( \mathbb{G}(\mathbf{u}^{(l)}) (\mathbf{v}^{(l)} - \mathbf{y}) \right)^T \right).$$

Besides, since the values in  $\mathbf{v}^{(l)}$  belong to the range  $[0, 1]$ . Combine with Lemma 22, we can bound  $\|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})\|_2$  as follows:

$$\begin{aligned} \|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})\|_2^2 &= \sum_{j=1}^l \|\nabla_{\mathbf{w}_{(j)}} f(\mathbf{w}, \mathbf{x})\|_2^2 \\ &= \left\| \mathbf{v}^{(l-1)} \left( \mathbb{G}(\mathbf{u}^{(l)}) (\mathbf{v}^{(l)} - \mathbf{y}) \right)^T \right\|_F^2 + \sum_{j=1}^{l-1} \left\| \mathbf{v}^{(j-1)} \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \right)^T \right\|_F^2 \\ &\leq \frac{1}{16} \mathbf{d}_{l-1} \|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 + \frac{1}{16} \|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 \sum_{j=1}^{l-1} \mathbf{d}_{j-1} \|\mathbf{B}_{j+1:l}\|_F^2 \\ &\stackrel{\textcircled{1}}{\leq} \frac{1}{16} c_y c_d + \frac{1}{16} c_y c_d c_r (l-1), \end{aligned}$$

where  $c_y, c_d, c_r$  are defined as

$$\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 \leq c_y, \quad c_d = \max(\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_l) \quad \text{and} \quad c_r = \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^{l-1}\right).$$

Notice,  $\textcircled{1}$  holds since in Lemma 22, we have

$$\|\mathbf{B}_{s:t}\|_F^2 \leq \left(\frac{r}{4}\right)^{2(t-s+1)} \leq \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^{l-1}\right).$$

Thus, we can obtain

$$\|\nabla_{\mathbf{w}} f(w, x)\|_2 \leq \sqrt{\frac{1}{16} c_y c_d (1 + c_r (l - 1))} \triangleq \alpha.$$

The proof is completed.  $\square$

## D.2.2 PROOF OF LEMMA 17

For convenience, we first give the computation of some gradients.

**Lemma 23.** *Assume the activation functions in deep neural network are sigmoid functions. Then the following properties hold:*

(1) We can compute the gradients  $\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}}$  and  $\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{v}^{(i)}}$  as

$$\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}} = \mathcal{G}(\mathbf{u}^{(i)}) \mathbf{B}_{i+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \quad \text{and} \quad \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{v}^{(i)}} = \mathbf{B}_{i+1:l} (\mathbf{v}^{(l)} - \mathbf{y}).$$

(2) We can compute the gradient  $\frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{w}^{(j)}}$  as

$$\frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{w}^{(j)}} = (\mathbf{v}^{(j-1)})^T \otimes \left( \mathcal{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:i-1} (\mathbf{W}^{(i)})^T \right)^T \in \mathbb{R}^{d_i \times d_j d_{j-1}}, \quad (i > j).$$

$$\frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{w}^{(i)}} = (\mathbf{v}^{(i-1)})^T \otimes \mathbf{I}_{d_i} \in \mathbb{R}^{d_i \times d_i d_{i-1}}, \quad (i = j).$$

(3) We can compute the gradient  $\frac{\partial \mathbf{v}^{(i)}}{\partial \mathbf{w}^{(j)}}$  as

$$\frac{\partial \mathbf{v}^{(i)}}{\partial \mathbf{w}^{(j)}} = (\mathbf{v}^{(j-1)})^T \otimes \left( \mathcal{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:i} \right)^T \in \mathbb{R}^{d_i \times d_j d_{j-1}}, \quad (i \geq j).$$

It should be pointed out that the proof of Lemma 23 can be founded Sec. D.4.

*Proof.* To prove our conclusion, we have two steps: computing the Hessian and bounding its operation norm.

**Step 1. Compute the Hessian:** We first consider the computation of  $\frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}}$ :

$$\frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} = \frac{\partial \left( \text{vec} \left( \left( \mathcal{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \right) (\mathbf{v}^{(j-1)})^T \right) \right)}{\partial \mathbf{w}_{(i)}^T}.$$

Recall that we define

$$\mathbf{B}_{s:t} = \mathbf{A}_s \mathbf{A}_{s+1} \cdots \mathbf{A}_t \in \mathbb{R}^{d_{s-1} \times d_t}, \quad (s \leq t) \quad \text{and} \quad \mathbf{B}_{s:t} = \mathbf{I}, \quad (s > t).$$

Then we have

$$\begin{aligned} \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} &= \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{j+1:l}^T \right) \otimes \left( \mathbf{I}_{d_j} \frac{\partial \text{vec} \left( \mathcal{G}(\mathbf{u}^{(j)}) \right)}{\partial \mathbf{w}_{(i)}^T} \right) (\triangleq \mathbf{Q}_1^{ij}) \\ &+ \sum_{k=j+1}^l \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{k+1:l}^T \right) \otimes \left( \mathcal{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:k-1} \mathbf{W}_k^T \right) \frac{\partial \text{vec} \left( \mathcal{G}(\mathbf{u}^{(k)}) \right)}{\partial \mathbf{w}_{(i)}^T} (\triangleq \mathbf{Q}_2^{ij}) \\ &+ \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{i+1:l}^T \mathcal{G}(\mathbf{u}^{(i)}) \right) \otimes \left( \mathcal{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:i-1} \right) \frac{\partial \text{vec} \left( \mathbf{W}_i^T \right)}{\partial \mathbf{w}_{(i)}^T} (\triangleq \mathbf{Q}_3^{ij}) \\ &+ \mathbf{v}^{(j-1)} \otimes \left( \mathcal{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \frac{\partial (\mathbf{v}^{(l)} - \mathbf{y})}{\partial \mathbf{w}_{(i)}^T} (\triangleq \mathbf{Q}_4^{ij}) \\ &+ \mathbf{I}_{d_{j-1}} \otimes \left( \mathcal{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \right) \frac{\partial \mathbf{v}^{(j-1)}}{\partial \mathbf{w}_{(i)}^T} (\triangleq \mathbf{Q}_5^{ij}) \end{aligned}$$

**Case I:**  $i > j$ . We first consider the case that  $i > j$ . In this is case,  $\mathbf{Q}_1^{ij} = \mathbf{0}$  since  $\frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(j)}))}{\partial \mathbf{w}_{(i)}^T} = \mathbf{0}$ .

Computing  $\mathbf{Q}_2^{ij}$  needs more efforts. By utilizing the computation of  $\frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{w}_{(i)}}$  in Lemma 23, we have

$$\frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(k)}))}{\partial \mathbf{w}_{(i)}} = \frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(k)}))}{\partial \mathbf{u}^{(k)}} \frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{w}_{(i)}} = \mathbf{P}_k \left( \mathbf{v}^{(i-1)T} \otimes \left( \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{B}_{i+1:k-1} (\mathbf{W}^{(k)T})^T \right)^T \right), (k > i)$$

where  $\mathbf{P}_k$  is a matrix of size  $\mathbf{d}_k^2 \times \mathbf{d}_k$  whose  $((s-1)\mathbf{d}_k + s, s)$  ( $s = 1, \dots, \mathbf{d}_k$ ) entry equal to  $\sigma(\mathbf{u}_s^{(k)})(1 - \sigma(\mathbf{u}_s^{(k)}))(1 - 2\sigma(\mathbf{u}_s^{(k)}))$  and rest entries are all 0. When  $k = i$ ,

$$\frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(k)}))}{\partial \mathbf{w}_{(k)}} = \frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(k)}))}{\partial \mathbf{u}^{(k)}} \frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{w}_{(k)}} = \mathbf{P}_k \left( (\mathbf{v}^{(k-1)})^T \otimes \mathbf{I}_{\mathbf{d}_k} \right) \in \mathbb{R}^{\mathbf{d}_k^2 \times \mathbf{d}_k \mathbf{d}_{k-1}}.$$

Note that for  $k < i$ , we have  $\frac{\partial \mathbb{G}(\mathbf{u}^{(k)})}{\partial \mathbf{w}_{(i)}} = \mathbf{0}$ . For brevity, let

$$\mathbf{D}_k \triangleq \left( (\mathbf{v}^{(j-1)}(\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{k+1:l}^T) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:k-1} \mathbf{W}_k^T \right) \right) (k = i, \dots, l). \quad (25)$$

Therefore, we have

$$\mathbf{Q}_2^{ij} = \mathbf{D}_i \mathbf{P}_i \left( (\mathbf{v}^{(i-1)})^T \otimes \mathbf{I}_{\mathbf{d}_i} \right) + \sum_{k=i+1}^l \mathbf{D}_k \mathbf{P}_k \left( (\mathbf{v}^{(i-1)})^T \otimes \left( \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{B}_{i+1:k-1} (\mathbf{W}^{(k)T})^T \right)^T \right).$$

Then we consider  $\mathbf{Q}_3^{ij}$ .

$$\mathbf{Q}_3^{ij} = \left( \mathbf{v}^{(j-1)}(\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{i+1:l}^T \mathbb{G}(\mathbf{u}^{(i)}) \right) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:i-1} \right).$$

Also we can use the computation of  $\frac{\partial \mathbf{v}^{(l)}}{\partial \mathbf{w}_{(i)}}$  in Lemma 23 and compute  $\mathbf{Q}_4^{ij}$  as follows:

$$\begin{aligned} \mathbf{Q}_4^{ij} &= \mathbf{v}^{(j-1)} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \frac{\partial (\mathbf{v}^{(l)} - \mathbf{y})}{\partial \mathbf{w}_{(i)}^T} \\ &= \left( \mathbf{v}^{(j-1)} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \right) \left( (\mathbf{v}^{(i-1)})^T \otimes \left( \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{B}_{i+1:l} \right)^T \right). \end{aligned}$$

Finally, since  $i > j$ , we can compute  $\mathbf{Q}_5^{ij} = \mathbf{0}$ .

**Case II:**  $i = j$ . We first consider  $\frac{\partial \mathbb{G}(\mathbf{u}^{(k)})}{\partial \mathbf{w}_{(k)}}$ :

$$\frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(k)}))}{\partial \mathbf{w}_{(k)}^T} = \frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(k)}))}{\partial \mathbf{u}^{(k)}} \frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{w}_{(k)}^T} = \mathbf{P}_k \left( (\mathbf{v}^{(k-1)})^T \otimes \mathbf{I}_{\mathbf{d}_k} \right) \in \mathbb{R}^{\mathbf{d}_k^2 \times \mathbf{d}_k \mathbf{d}_{k-1}},$$

where  $\mathbf{P}_k$  is a matrix of size  $\mathbf{d}_k^2 \times \mathbf{d}_k$  whose  $(s, (s-1)\mathbf{d}_k + s)$  entry equal to  $\sigma(\mathbf{u}_s^{(k)})(1 - \sigma(\mathbf{u}_s^{(k)}))(1 - 2\sigma(\mathbf{u}_s^{(k)}))$  and rest entries are all 0.  $\mathbf{Q}_1^{jj}$  can be computed as

$$\begin{aligned} \mathbf{Q}_1^{jj} &= \left( \mathbf{v}^{(j-1)}(\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{j+1:l}^T \right) \otimes \left( \mathbf{I}_{\mathbf{d}_j} \right) \frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(j)}))}{\partial \mathbf{w}_{(j)}^T} \\ &= \left( \left( \mathbf{v}^{(j-1)}(\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{j+1:l}^T \right) \otimes \left( \mathbf{I}_{\mathbf{d}_j} \right) \right) \left( \mathbf{P}_j \left( (\mathbf{v}^{(j-1)})^T \otimes \mathbf{I}_{\mathbf{d}_j} \right) \right). \end{aligned}$$

As for  $\mathbf{Q}_2^{jj}$ , by Eqn. (25) we have

$$\mathbf{Q}_2^{jj} = \sum_{k=j+1}^l \mathbf{D}_k \mathbf{P}_k \left( (\mathbf{v}^{(j-1)})^T \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:k-1} (\mathbf{W}^{(k)T})^T \right)^T \right).$$

Since  $i = j$ ,  $\mathbf{Q}_3^{jj}$  does not exist. For convenience, we just set  $\mathbf{Q}_3^{jj} = \mathbf{0}$ .

Now we consider  $\mathbf{Q}_4^{jj}$  which can be computed as follows:

$$\begin{aligned}\mathbf{Q}_4^{jj} &= \mathbf{v}^{(j-1)} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \frac{\partial(\mathbf{v}^{(l)} - \mathbf{y})}{\partial \mathbf{w}_{(j)}^T} \\ &= \left( \mathbf{v}^{(j-1)} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \right) \left( (\mathbf{v}^{(j-1)})^T \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right)^T \right).\end{aligned}$$

Finally, since  $i = j$ , we can compute  $\mathbf{Q}_5^{jj} = \mathbf{0}$ .

**Case III:**  $i < j$ . Since  $\frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w} \partial \mathbf{w}^T}$  is symmetrical, we have  $\mathbf{Q}_k^{ij} = \mathbf{Q}_k^{ji}$  ( $k = 1, \dots, 5$ ).

**Step 2. Bound the operation norm of Hessian:** We mainly use Lemma 22 to achieve this goal. From Lemma 22, we have

(1) For arbitrary matrices  $\mathbf{M}$  and  $\mathbf{N}$  of proper size, we have

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

(2) For arbitrary matrices  $\mathbf{M}$  and  $\mathbf{N}$  of proper size, we have

$$\|\mathbf{P}_k \mathbf{M}\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N} \mathbf{P}_k\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{N}\|_F^2.$$

(3) For  $\mathbf{B}_{s:t}$  and  $\mathbf{D}_{s:t}$ , we have

$$\|\mathbf{B}_{s:t}\|_F^2 \leq \frac{1}{16^{t-s+1}} \|\mathbf{D}_{s:t}\|_F^2 \quad \text{and} \quad \frac{1}{16^{t-s+1}} \|\mathbf{D}_{s:t}\|_F^2 \leq c_r,$$

where  $c_r = \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^l\right)$ .

(4) For arbitrary matrices  $\mathbf{M}$ ,  $\mathbf{N}$  and  $\mathbf{I}$  of proper sizes, let  $\mathbf{m} = \text{vec}(\mathbf{M})$ . Then we have

$$\|(\mathbf{N} \otimes \mathbf{I}) \mathbf{m}\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2 \quad \text{and} \quad \|(\mathbf{I} \otimes \mathbf{N}) \mathbf{m}\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2.$$

The values of entries in  $\mathbf{v}^{(h)}$  are bounded by  $0 \leq \sigma(\mathbf{u}_h^{(i)}) \leq 1$  which leads to  $\|\mathbf{v}^{(h)}\|_F^2 \leq \mathbf{d}_h \leq c_d$ , where  $c_d = \max_i \mathbf{d}_i$ . On the other hand, since the values in  $\mathbf{v}^{(l)}$  belong to the range  $[0, 1]$  and  $\mathbf{y}$  is the label,  $\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2$  can be bounded:

$$\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 \leq c_y < +\infty,$$

where  $c_y$  is a universal constant.

We first define

$$\mathbf{C}_k^{ij} = \mathbf{D}_k \mathbf{P}_k \left( \mathbf{v}^{(i-1)} \right)^T \otimes \left( \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{B}_{i+1:k-1} (\mathbf{W}^{(k)})^T \right)^T$$

and

$$\begin{aligned}\mathbf{C}^{ij} &= \mathbf{D}_i \mathbf{P}_i \left( (\mathbf{v}^{(i-1)})^T \otimes \mathbf{I}_{\mathbf{d}_i} \right) = \left( (\mathbf{v}^{(i-1)} \otimes \mathbf{I}_{\mathbf{d}_i}) (\mathbf{D}_i \mathbf{P}_i)^T \right)^T \stackrel{\textcircled{1}}{=} \left( \mathbf{v}^{(i-1)} \otimes (\mathbf{D}_i \mathbf{P}_i)^T \right)^T \\ &= (\mathbf{v}^{(i-1)})^T \otimes (\mathbf{D}_i \mathbf{P}_i),\end{aligned}$$

where  $\mathbf{D}_k$  is defined in Eqn. (25).  $\textcircled{1}$  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  $(\mathbf{u} \otimes \mathbf{I}_k) \mathbf{M} = \mathbf{u} \otimes \mathbf{M}$ .

**Case I:**  $i > j$ . According to the definition of  $\mathbf{C}^{ij}$  and  $\mathbf{C}_k^{ij}$ , we have  $\mathbf{Q}_2^{ij} = \mathbf{C}^{ij} + \sum_{k=i+1}^l \mathbf{C}_k^{ij}$ . So we have

$$\begin{aligned}\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}^T} \right\|_F^2 &= \left\| \mathbf{Q}_1^{ij} + \mathbf{Q}_2^{ij} + \mathbf{Q}_3^{ij} + \mathbf{Q}_4^{ij} + \mathbf{Q}_5^{ij} \right\|_F^2 \\ &= \left\| \mathbf{C}^{ij} + \sum_{k=i+1}^l \mathbf{C}_k^{ij} + \mathbf{Q}_3^{ij} + \mathbf{Q}_4^{ij} \right\|_F^2 \\ &= (l - i + 3) \left( \|\mathbf{C}^{ij}\|_F^2 + \sum_{k=i+1}^l \left( \|\mathbf{C}_k^{ij}\|_F^2 + \|\mathbf{Q}_3^{ij}\|_F^2 + \|\mathbf{Q}_4^{ij}\|_F^2 \right) \right).\end{aligned}$$

Here we bound each term separately:

$$\begin{aligned}
\|C^{ij}\|_F^2 &\leq \|\mathbf{v}^{(j-1)}\|_F^2 \|\mathbf{v}^{(l)} - \mathbf{y}\|_F^2 \|\mathbf{B}_{i+1:l}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:i-1} \mathbf{W}_i^T\|_F^2 \frac{2^6}{3^8} \|\mathbf{v}^{(i-1)}\|_F^2 \\
&\leq \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} \frac{1}{16^{l-i}} \mathbf{D}_{i+1:l} \frac{1}{16^{i-j}} \mathbf{D}_{j+1:i} \\
&\leq \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} \frac{1}{16^{l-j}} \mathbf{D}_{j+1:l} \\
&\leq \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} c_r.
\end{aligned}$$

Similarly, we can bound  $\|C_k^{ij}\|_F^2$  as follows:

$$\begin{aligned}
&\|C_k^{ij}\|_F^2 \\
&\leq \|\mathbf{v}^{(j-1)}\|_F^2 \|\mathbf{v}^{(l)} - \mathbf{y}\|_F^2 \|\mathbf{B}_{k+1:l}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:k-1} \mathbf{W}_k^T\|_F^2 \frac{2^6}{3^8} \|\mathbf{v}^{(i-1)}\|_F^2 \frac{1}{16} \|\mathbf{B}_{i+1:k-1} (\mathbf{W}^{(k)})^T\|_F^2 \\
&\leq \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} \frac{1}{16^{l-k}} \mathbf{D}_{k+1:l} \frac{1}{16^{k-j-1}} \mathbf{D}_{j+1:k} \frac{1}{16^{k-i-1}} \mathbf{D}_{i+1:k} \\
&= \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} \frac{1}{16^{l-j-1}} \mathbf{D}_{j+1:l} \frac{1}{16^{k-i-1}} \mathbf{D}_{i+1:k} \\
&\leq \frac{2^{14}}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} c_r^2.
\end{aligned}$$

We also bound  $\|Q_3^{ij}\|_F^2$  as

$$\|Q_3^{ij}\|_F^2 \leq \|\mathbf{v}^{(j-1)}\|_F^2 \|\mathbf{v}^{(l)} - \mathbf{y}\|_F^2 \frac{1}{16} \|\mathbf{B}_{i+1:l}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:i-1}\|_F^2 \leq \frac{1}{2^8} c_y \mathbf{d}_{j-1} c_r.$$

Finally, we bound  $\|Q_4^{ij}\|_F^2$  as follows:

$$\|Q_4^{ij}\|_F^2 \leq \|\mathbf{v}^{(j-1)}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:l}\|_F^2 \|\mathbf{v}^{(i-1)}\|_F^2 \frac{1}{16} \|\mathbf{B}_{i+1:l}\|_F^2 \leq \frac{1}{2^8} \mathbf{d}_{j-1} \mathbf{d}_{i-1} c_r^2.$$

Note that  $\mathbf{d}_i \leq c_d$ . Thus, we can bound  $\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)}^T \partial \mathbf{w}_{(i)}} \right\|_F^2$  as

$$\begin{aligned}
&\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 \\
&\leq (l-i+3) \left( \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} c_r + \sum_{k=i+1}^l \frac{2^{14}}{3^8} c_y \mathbf{d}_{j-1} \mathbf{d}_{i-1} c_r^2 + \frac{1}{2^8} c_y \mathbf{d}_{j-1} c_r + \frac{1}{2^8} \mathbf{d}_{j-1} \mathbf{d}_{i-1} c_r^2 \right) \\
&\leq (l+1) \left( \frac{64}{6561} c_y c_d^2 c_r + \frac{4096}{6561} c_y (l-2) c_d^2 c_r^2 + \frac{1}{256} c_y c_d c_r + \frac{1}{256} c_d c_r^2 \right).
\end{aligned}$$

**Case II:**  $i = j$ . According to the definition of  $C^{ij}$  and  $C_k^{ij}$ , we have  $Q_2^{jj} = \sum_{k=j+1}^l C_k^{jj}$ .

Similarly, we have

$$\begin{aligned}
\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 &= \|Q_1^{jj} + Q_2^{jj} + Q_3^{jj} + Q_4^{jj} + Q_5^{jj}\|_F^2 = \left\| Q_1^{jj} + \sum_{k=j+1}^l C_k^{jj} + Q_4^{jj} \right\|_F^2 \\
&\leq (l-j+2) \left( \|Q_1^{jj}\|_F^2 + \sum_{k=j+1}^l \|C_k^{jj}\|_F^2 + \|Q_4^{jj}\|_F^2 \right).
\end{aligned}$$

Thus, we can bound  $\|Q_1^{jj}\|_F^2$  first:

$$\|Q_1^{jj}\|_F^2 \leq \|\mathbf{v}^{(j-1)}\|_F^2 \|\mathbf{v}^{(l)} - \mathbf{y}\|_F^2 \|\mathbf{B}_{j+1:l}\|_F^2 \frac{2^6}{3^8} \|\mathbf{v}^{(j-1)}\|_F^2 \leq \frac{2^6}{3^8} c_y \mathbf{d}_{j-1}^2 c_r.$$

As for  $Q_2^{jj}$ , we have

$$\begin{aligned} & \|C_k^{ij}\|_F^2 \\ & \leq \|\mathbf{v}^{(j-1)}\|_F^2 \|\mathbf{v}^{(l)} - \mathbf{y}\|_F^2 \|\mathbf{B}_{k+1:l}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:k-1} \mathbf{W}_k^T\|_F^2 \frac{2^6}{3^8} \|\mathbf{v}^{(j-1)}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:k-1} (\mathbf{W}^{(k)})^T\|_F^2 \\ & = \frac{2^6}{3^8} c_y \mathbf{d}_{j-1}^2 \frac{1}{16^{l-k}} \mathbf{D}_{k+1:l} \frac{1}{16^{k-j-1}} \mathbf{D}_{j+1:k} \frac{1}{16^{k-j-1}} \mathbf{D}_{j+1:k} \\ & \leq \frac{2^{14}}{3^8} c_y \mathbf{d}_{j-1}^2 c_r^2. \end{aligned}$$

Then we bound  $\|Q_4^{jj}\|_F^2$ :

$$\|Q_4^{jj}\|_F^2 \leq \|\mathbf{v}^{(j-1)}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:l}\|_F^2 \|\mathbf{v}^{(j-1)}\|_F^2 \frac{1}{16} \|\mathbf{B}_{j+1:l}\|_F^2 \leq \frac{1}{2^8} \mathbf{d}_{j-1}^2 c_r^2.$$

Note that for any input, we have  $c_v = \max_j \|\mathbf{v}^{(j-1)}(\mathbf{v}^{(l)} - \mathbf{y})^T\|_F^2 \leq \max_j \|\mathbf{v}^{(j-1)}\|_F^2 \|\mathbf{v}^{(l)} - \mathbf{y}\|_F^2 \leq c_y c_d$ , where  $\|\mathbf{v}^{(l)} - \mathbf{y}\|_F^2$  can be bounded by a constant  $c_y$ . Thus, we can bound  $\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2$  as

$$\begin{aligned} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 & \leq (l-i+3) \left( \frac{2^6}{3^8} c_y \mathbf{d}_{j-1}^2 c_r + \sum_{k=i+1}^l \frac{2^{14}}{3^8} c_y \mathbf{d}_{j-1}^2 c_r^2 + \frac{1}{2^8} \mathbf{d}_{j-1}^2 c_r^2 \right) \\ & \leq (l+2) \left( \frac{64}{6561} c_y c_d^2 c_r + \frac{4096}{6561} c_y (l-1) c_d^2 c_r^2 + \frac{1}{256} c_d^2 c_r^2 \right). \end{aligned}$$

**Case III:**  $i < j$ . Since  $\frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w} \partial \mathbf{w}^T}$  is symmetrical, we have  $Q_k^{ij} = Q_k^{ji}$  ( $k = 1, \dots, 5$ ). Thus, it yields

$$\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(i)}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 \leq (l+1) \left( \frac{64}{6561} c_y c_d^2 c_r + \frac{4096}{6561} c_y (l-2) c_d^2 c_r^2 + \frac{1}{256} c_y c_d c_r + \frac{1}{256} c_d c_r^2 \right).$$

**Final result:** Thus we can bound

$$\begin{aligned} \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_{\text{op}} & \leq \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_F \\ & \leq \sqrt{(l-1)l \max_{i,j:i \neq j} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(i)}^T} \right\|_F^2 + \sum_{j=1}^l \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)} \partial \mathbf{w}_{(j)}^T} \right\|_F^2} \\ & \leq \left( (l-1)l(l+1) \left( \frac{64}{6561} c_y c_d^2 c_r + \frac{4096}{6561} c_y (l-2) c_d^2 c_r^2 + \frac{1}{256} c_y c_d c_r + \frac{1}{256} c_d c_r^2 \right) \right. \\ & \quad \left. + (l+2) \left( \frac{64}{6561} c_y c_d^2 c_r + \frac{4096}{6561} c_y (l-1) l c_d^2 c_r^2 + \frac{1}{256} l c_d^2 c_r^2 \right) \right)^{\frac{1}{2}} \\ & \leq \sqrt{c_{s_1} c_r c_d^2 l^4}, \end{aligned}$$

where  $c_{s_1}$  and  $c_{s_2}$  are two constants.

Since  $\|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_{\text{op}} \leq \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_F$ , we know that the gradient  $\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})$  is  $\varsigma$ -Lipschitz, where  $\varsigma = \sqrt{c_{s_1} c_r c_d^2 l^4}$ .

On the other hand, since for any input  $\mathbf{x}$ ,  $\sigma(\mathbf{x})$  belongs to  $[0, 1]$ , the values of the entries of  $\nabla_{\mathbf{w}}^3 f(\mathbf{w}, \mathbf{x})$  can be bounded. Thus, we can bound

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

We complete the proof.  $\square$

### D.2.3 PROOF OF LEMMA 18

For convenience, we first give the computation of some gradients.

**Lemma 24.** *Assume the activation functions in deep neural network are sigmoid functions. Then we can compute the gradients  $\frac{\partial \mathbf{u}^{(j)}}{\partial \mathbf{u}^{(1)}}$  and  $\frac{\partial \mathbf{v}^{(j)}}{\partial \mathbf{u}^{(1)}}$  as*

$$\begin{aligned} \frac{\partial \mathbf{u}^{(j)}}{\partial \mathbf{u}^{(1)}} &= \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{A}_2 \cdots \mathbf{A}_{j-1} (\mathbf{W}^j)^T \right)^T \in \mathbb{R}^{d_j \times d_1}, \quad (j > 1). \\ \frac{\partial \mathbf{v}^{(j)}}{\partial \mathbf{u}^{(1)}} &= \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{A}_2 \cdots \mathbf{A}_j \right)^T \in \mathbb{R}^{d_j \times d_1}, \quad (j > 1). \end{aligned}$$

It should be pointed out that the proof of Lemma 24 can be founded Sec. D.4.

*Proof.* To prove our conclusion, we have two steps: computing  $\nabla_{\mathbf{x}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})$  and bounding its operation norm.

**Step 1. Compute  $\nabla_{\mathbf{x}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})$ :**

We first consider the computation of  $\frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{x}^T \partial \mathbf{w}^{(j)}}$ :

$$\frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{x}^T \partial \mathbf{w}^{(j)}} = \frac{\partial \left( \text{vec} \left( \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \right) (\mathbf{v}^{(j-1)})^T \right) \right)}{\partial \mathbf{x}^T}.$$

Recall that we define

$$\begin{aligned} \mathbf{A}_i &= (\mathbf{W}^{(i)})^T \mathbb{G}(\mathbf{u}^{(i)}) \in \mathbb{R}^{d_{i-1} \times d_i}. \\ \mathbf{B}_{s:t} &= \mathbf{A}_s \mathbf{A}_{s+1} \cdots \mathbf{A}_t \in \mathbb{R}^{d_{s-1} \times d_t}, \quad (s \leq t) \quad \text{and} \quad \mathbf{B}_{s:t} = \mathbf{I}, \quad (s > t). \end{aligned}$$

Then we have

$$\begin{aligned} \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{x}^T \partial \mathbf{w}^{(j)}} &= \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{j+1:l}^T \right) \otimes (\mathbf{I}_{d_j}) \frac{\partial \text{vec} \left( \mathbb{G}(\mathbf{u}^{(j)}) \right)}{\partial \mathbf{x}^T} (\triangleq \mathbf{Q}_1^j) \\ &\quad + \sum_{k=j+1}^l \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{k+1:l}^T \right) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:k-1} \mathbf{W}_k^T \right) \frac{\partial \text{vec} \left( \mathbb{G}(\mathbf{u}^{(k)}) \right)}{\partial \mathbf{x}^T} (\triangleq \mathbf{Q}_2^j) \\ &\quad + \mathbf{v}^{(j-1)} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \frac{\partial (\mathbf{v}^{(l)} - \mathbf{y})}{\partial \mathbf{x}^T} (\triangleq \mathbf{Q}_3^j) \\ &\quad + \mathbf{I}_{d_{j-1}} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \right) \frac{\partial \mathbf{v}^{(j-1)}}{\partial \mathbf{x}^T} (\triangleq \mathbf{Q}_4^j) \end{aligned}$$

By using Lemma 24, we can compute  $\mathbf{Q}_1^{ij}$  as

$$\frac{\partial \text{vec} \left( \mathbb{G}(\mathbf{u}^{(k)}) \right)}{\partial \mathbf{x}^T} = \frac{\partial \text{vec} \left( \mathbb{G}(\mathbf{u}^{(k)}) \right)}{\partial \mathbf{u}^{(k)}} \frac{\partial \mathbf{u}^{(k)}}{\partial \mathbf{x}^T} = \mathbf{P}_k \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:k-1} (\mathbf{W}^k)^T \right)^T.$$

Thus, we have

$$\begin{aligned} \mathbf{Q}_1^j &= \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{j+1:l}^T \right) \otimes \mathbf{I}_{d_j} \frac{\partial \text{vec} \left( \mathbb{G}(\mathbf{u}^{(j)}) \right)}{\partial \mathbf{x}^T} \\ &= \left( \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{j+1:l}^T \right) \otimes \mathbf{I}_{d_j} \right) \mathbf{P}_k \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:k-1} (\mathbf{W}^k)^T \right)^T. \end{aligned}$$

As for  $\mathbf{Q}_2^j$ , we also can utilize Lemma 24 to compute it:

$$\begin{aligned}\mathbf{Q}_2^j &= \sum_{k=j+1}^l \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{k+1:l}^T \right) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:k-1} \mathbf{W}_k^T \right) \frac{\partial \text{vec}(\mathbb{G}(\mathbf{u}^{(k)}))}{\partial \mathbf{x}^T} \\ &= \sum_{k=j+1}^l \left( \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{k+1:l}^T \right) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:k-1} \mathbf{W}_k^T \right) \right) \mathbf{P}_k \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:k-1} (\mathbf{W}^k)^T \right)^T.\end{aligned}$$

Then we consider  $\mathbf{Q}_3^{jj}$ .

$$\mathbf{Q}_3^j = \mathbf{v}^{(j-1)} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \frac{\partial (\mathbf{v}^{(l)} - \mathbf{y})}{\partial \mathbf{x}^T} = \left( \mathbf{v}^{(j-1)} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} \right) \right) \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:l} \right)^T.$$

$\mathbf{Q}_4^j$  can be computed as follows:

$$\mathbf{Q}_4^j = \mathbf{I}_{d_{j-1}} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \right) \frac{\partial \mathbf{v}^{(j-1)}}{\partial \mathbf{x}^T} = \left( \mathbf{I}_{d_{j-1}} \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \right) \right) \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:j} \right)^T.$$

**Step 2. Bound the operation norm of Hessian:** We mainly use Lemma 22 to achieve this goal. From Lemma 22, we have

(1) For arbitrary matrices  $\mathbf{M}$  and  $\mathbf{N}$  of proper size, we have

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

(2) For arbitrary matrices  $\mathbf{M}$  and  $\mathbf{N}$  of proper size, we have

$$\|\mathbf{P}_k \mathbf{M}\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N} \mathbf{P}_k\|_F^2 \leq \frac{2^6}{3^8} \|\mathbf{N}\|_F^2.$$

(3) For  $\mathbf{B}_{s:t}$  and  $\mathbf{D}_{s:t}$ , we have

$$\|\mathbf{B}_{s:t}\|_F^2 \leq \frac{1}{16^{t-s+1}} \mathbf{D}_{s:t} \quad \text{and} \quad \frac{1}{16^{t-s+1}} \mathbf{D}_{s:t} \leq c_r,$$

$$\text{where } c_r = \max\left(\frac{r^2}{4}, \left(\frac{r^2}{16}\right)^{l-1}\right).$$

(4) For arbitrary matrices  $\mathbf{M}$ ,  $\mathbf{N}$  and  $\mathbf{I}$  of proper sizes, let  $\mathbf{m} = \text{vec}(\mathbf{M})$ . Then we have

$$\|(\mathbf{N} \otimes \mathbf{I}) \mathbf{m}\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2 \quad \text{and} \quad \|(\mathbf{I} \otimes \mathbf{N}) \mathbf{m}\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2.$$

The values of entries in  $\mathbf{v}^{(h)}$  are bounded by  $0 \leq \sigma(\mathbf{u}_h^{(i)}) \leq 1$  which leads to  $\|\mathbf{v}^{(h)}\|_F^2 \leq \mathbf{d}_h \leq c_d$ , where  $c_d = \max_i \mathbf{d}_i$ . On the other hand, since the values in  $\mathbf{v}^{(l)}$  belong to the range  $[0, 1]$  and  $\mathbf{y}$  is the label,  $\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2$  can be bounded:

$$\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 \leq c_y < +\infty,$$

where  $c_y$  is a universal constant.

We first define

$$\mathbf{C}_k^j = \left( \left( \mathbf{v}^{(j-1)} (\mathbf{v}^{(l)} - \mathbf{y})^T \mathbf{B}_{k+1:l}^T \right) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:k-1} \mathbf{W}_k^T \right) \right) \mathbf{P}_k \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:k-1} (\mathbf{W}^k)^T \right)^T.$$

Then we have  $\mathbf{Q}_2^j = \sum_{k=j+1}^l \mathbf{C}_k^j$ . So we have

$$\begin{aligned}\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{x}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 &= \left\| \mathbf{Q}_1^j + \mathbf{Q}_2^j + \mathbf{Q}_3^j + \mathbf{Q}_4^j \right\|_F^2 = \left\| \mathbf{Q}_1^j + \sum_{k=j+1}^l \mathbf{C}_k^j + \mathbf{Q}_3^j + \mathbf{Q}_4^j \right\|_F^2 \\ &= (l-j+3) \left( \left\| \mathbf{Q}_1^j \right\|_F^2 + \sum_{k=j+1}^l \left\| \mathbf{C}_k^j \right\|_F^2 + \left\| \mathbf{Q}_3^j \right\|_F^2 + \left\| \mathbf{Q}_4^j \right\|_F^2 \right).\end{aligned}$$

Then we bound each term separately:

$$\left\| \mathbf{Q}_1^j \right\|_F^2 \leq \left\| \mathbf{v}^{(j-1)} \right\|_F^2 \left\| \mathbf{v}^{(l)} - \mathbf{y} \right\|_F^2 \left\| \mathbf{B}_{j+1:l} \right\|_F^2 \frac{2^6}{3^8} \frac{1}{16} \left\| \mathbf{B}_{2:k-1} (\mathbf{W}^k)^T \right\|_F^2 \leq \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2.$$

Similarly, we bound  $\left\| \mathbf{C}_k^j \right\|_F^2$ :

$$\begin{aligned} \left\| \mathbf{C}_k^j \right\|_F^2 &= \left\| \mathbf{v}^{(j-1)} \right\|_F^2 \left\| \mathbf{v}^{(l)} - \mathbf{y} \right\|_F^2 \left\| \mathbf{B}_{k+1:l} \right\|_F^2 \frac{1}{16} \left\| \mathbf{B}_{j+1:k-1} \mathbf{W}_k^T \right\|_F^2 \frac{2^6}{3^8} \frac{1}{16} \left\| \mathbf{B}_{2:k-1} (\mathbf{W}^{(k)})^T \right\|_F^2 \\ &= \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} \frac{1}{16^{l-k}} \mathbf{D}_{k+1:l} \frac{1}{16^{k-j-1}} \mathbf{D}_{j+1:k} \frac{1}{16^{k-1}} \mathbf{D}_{2:k} \\ &\leq \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2. \end{aligned}$$

We also bound  $\left\| \mathbf{Q}_3^{ij} \right\|_F^2$  as

$$\left\| \mathbf{Q}_3^{ij} \right\|_F^2 \leq \left\| \mathbf{v}^{(j-1)} \right\|_2^2 \frac{1}{16} \left\| \mathbf{B}_{j+1:l} \right\|_F^2 \frac{1}{16} \left\| \mathbf{B}_{2:l} \right\|_F^2 \leq \frac{1}{2^8} \mathbf{d}_{j-1} c_r^2.$$

Finally, we bound  $\left\| \mathbf{Q}_4^j \right\|_F^2$  as follows:

$$\left\| \mathbf{Q}_4^j \right\|_F^2 = \frac{1}{16} \left\| \mathbf{B}_{j+1:l} \right\|_F^2 \left\| \mathbf{v}^{(l)} - \mathbf{y} \right\|_F^2 \frac{1}{16} \left\| \mathbf{B}_{2:j} \right\|_F^2 \leq \frac{1}{2^8} c_y c_r.$$

Since  $c_d = \max_i \mathbf{d}_i$ , we can bound  $\left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)} \partial \mathbf{x}^T} \right\|_F^2$  as

$$\begin{aligned} \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{x}^T \partial \mathbf{w}_{(j)}} \right\|_F^2 &\leq (l-j+3) \left( \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2 + \sum_{k=j+1}^l \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2 + \frac{1}{2^8} c_y \mathbf{d}_{j-1} c_r + \frac{1}{2^8} c_y c_r \right) \\ &\leq (l+2) \left( \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2 + \sum_{k=j+1}^l \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2 + \frac{1}{2^8} c_y \mathbf{d}_{j-1} c_r + \frac{1}{2^8} c_y c_r \right). \end{aligned}$$

**Final result:** Thus we can bound

$$\begin{aligned} \left\| \nabla_{\mathbf{w}} \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x}) \right\|_{\text{op}} &\leq \left\| \nabla_{\mathbf{w}} \nabla_{\mathbf{x}} f(\mathbf{w}, \mathbf{x}) \right\|_F \leq \sqrt{\sum_{j=1}^l \left\| \frac{\partial^2 f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)} \partial \mathbf{x}^T} \right\|_F^2} \\ &\leq \sqrt{\sum_{j=1}^l (l+2) \left( \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2 + \sum_{k=j+1}^l \frac{2^6}{3^8} c_y \mathbf{d}_{j-1} c_r^2 + \frac{1}{2^8} c_y \mathbf{d}_{j-1} c_r + \frac{1}{2^8} c_y c_r \right)} \\ &\leq \sqrt{\frac{2^6}{3^8} l(l+2) c_y c_r c_d (l c_r + 1)}, \end{aligned}$$

where  $c_d = \max_j \mathbf{d}_j$ . The proof is completed.  $\square$

#### D.2.4 PROOF OF LEMMAS 19 AND 20

**Lemma 25.** (Alessandro, 2016; Rigollet, 2015) *Let  $(\mathbf{x}_1, \dots, \mathbf{x}_k)$  be a vector of i.i.d. Gaussian variables from  $\mathcal{N}(0, \tau^2)$  and let  $f : \mathbb{R}^{d_0} \rightarrow \mathbb{R}$  be  $L$ -Lipschitz. Then the variable  $f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})$  is sub-Gaussian. That is, we have*

$$\mathbb{P}(f(\mathbf{x}) - \mathbb{E}f(\mathbf{x}) > t) \leq \exp\left(-\frac{t^2}{2L^2\tau^2}\right), \quad (\forall t \geq 0),$$

or

$$\mathbb{E}(\lambda(f(\mathbf{x}) - \mathbb{E}f(\mathbf{x}))) \leq \exp(4\lambda^2 L^2 \tau^2), \quad (\forall \lambda \geq 0).$$

Remarkably, this is a dimension free inequality.

*Proof of Lemma 19.* We first define a function  $g(\mathbf{x}) = \mathbf{z}^T \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})$  where  $\mathbf{z} \in \mathbb{R}^d$  is a constant vector. Then we have  $\nabla_{\mathbf{x}} g(\mathbf{x}) = \nabla_{\mathbf{x}} (\mathbf{z}^T \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})) = \nabla_{\mathbf{x}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \mathbf{z}$ . Then by Lemma 18, we can obtain  $\|\nabla_{\mathbf{x}} g(\mathbf{x})\|_2 = \|\nabla_{\mathbf{x}} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \mathbf{z}\|_2 \leq \beta \|\mathbf{z}\|_2$ , where  $\beta = \sqrt{\frac{2^6}{3^8} l(l+2) c_y c_r c_d (lc_r + 1)}$  in which  $c_y$ ,  $c_d$  and  $c_r$  are defined in Lemma 18. This means  $g(\mathbf{x})$  is  $\beta \|\mathbf{z}\|_2$ -Lipschitz. Thus, by Lemma 25, we have

$$\mathbb{E} (t \langle \mathbf{z}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \rangle) = \mathbb{E} (t (g(\mathbf{x}) - \mathbb{E} g(\mathbf{x}))) \leq \exp (4t^2 \beta^2 \|\mathbf{z}\|_2^2 \tau^2).$$

Let  $\boldsymbol{\lambda} = t\mathbf{z}$ . This further gives

$$\mathbb{E} (\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \rangle) \leq \exp (4\beta^2 \tau^2 \|\boldsymbol{\lambda}\|_2^2),$$

which means  $\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \rangle$  is  $8\beta^2 \tau^2$ -sub-Gaussian.  $\square$

*Proof of Lemma 20.* We first define a function  $h(\mathbf{x}) = \mathbf{z}^T \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x}) \mathbf{z}$  where  $\mathbf{z} \in \mathbb{S}^{d-1}$ , i.e.  $\|\mathbf{z}\|_2 = 1$ . Then  $h(\mathbf{x})$  is a  $\gamma$ -Lipschitz function, where  $\gamma = \|\nabla_{\mathbf{x}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_{op}$ . Note that since the sigmoid function is infinitely differentiable function,  $\nabla_{\mathbf{x}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})$  exists. Also since for any input  $x$ ,  $\sigma(x)$  belongs to  $[0, 1]$ . Thus, the values of the entries in  $\nabla_{\mathbf{x}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})$  can be bounded. So according to the definition of the operation norm of a 3-way tensor, the operation norm of  $\nabla_{\mathbf{x}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})$  can be bounded by a constant. Without loss of generality, let  $\|\nabla_{\mathbf{x}} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_{op} \leq \gamma < +\infty$ . Thus, by Lemma 25, we have

$$\mathbb{E} (t \langle \mathbf{z}, (\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})) \mathbf{z} \rangle) = \mathbb{E} (t (h(\mathbf{x}) - \mathbb{E} h(\mathbf{x}))) \leq \exp \left( \frac{8t^2 \gamma^2 \tau^2}{2} \right).$$

This means that the hessian of the loss evaluated on a unit vector is  $8\gamma^2 \tau^2$ -sub-Gaussian.  $\square$

## D.2.5 PROOF OF LEMMA 21

*Proof.* Recall that the weight of each layer has magnitude bound separately, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$ . Assume that  $\mathbf{w}_{(j)}$  has  $s_j$  non-zero entries. Then we have  $\sum_{j=1}^l s_j = s$ . So here we separately assume  $\mathbf{w}_{\epsilon}^j = \{\mathbf{w}_1^j, \dots, \mathbf{w}_{n_{\epsilon}^j}^j\}$  is the  $\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ -covering net of the ball  $\mathbb{B}^{\mathbf{d}_j \mathbf{d}_{j-1}}(r)$  which corresponds to the weight  $\mathbf{w}_{(j)}$  of the  $j$ -th layer. Let  $n_{\epsilon}^j$  be the  $\epsilon/l$ -covering number. By  $\epsilon$ -covering theory in (Vershynin, 2012), we can have

$$n_{\epsilon}^j \leq \binom{\mathbf{d}_j \mathbf{d}_{j-1}}{s_j} \left( \frac{3r}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / (ld)} \right)^{s_j} \leq \exp \left( s_j \log \left( \frac{3r \mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d} \right) \right) = \exp \left( s_j \log \left( \frac{3rd}{\epsilon} \right) \right).$$

Let  $\mathbf{w} \in \Omega$  be an arbitrary vector. Since  $\mathbf{w} = [\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(l)}]$  where  $\mathbf{w}_{(j)}$  is the weight of the  $j$ -th layer, we can always find a vector  $\mathbf{w}_{k_j}^j$  in  $\mathbf{w}_{\epsilon}^j$  such that  $\|\mathbf{w}_{(j)} - \mathbf{w}_{k_j}^j\|_2 \leq \mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ . For brevity, let  $j_w \in [n_{\epsilon}^{j_w}]$  denote the index of  $\mathbf{w}_{k_{j_w}}^{j_w}$  in  $\epsilon$ -net  $\mathbf{w}_{\epsilon}^{j_w}$ . Then let  $\mathbf{w}_{k_w} = [\mathbf{w}_{k_1}^{j_1}; \dots; \mathbf{w}_{k_{j_w}}^{j_{j_w}}; \dots; \mathbf{w}_{k_l}^{j_l}]$ . This means that we can always find a vector  $\mathbf{w}_{k_w}$  such that  $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$ . Accordingly, we can

decompose  $\left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w}) \right\|_{\text{op}}$  as follows:

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

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

$$\begin{aligned}
\mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w}) \right\|_{\text{op}} \geq t \right\}, \\
\mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla^2 f(\mathbf{w}, \mathbf{x}_{(i)}) - \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}, \\
\mathbf{E}_2 &= \left\{ \sup_{j_w \in [n\epsilon^2], j=|l|} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}, \\
\mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(\nabla^2 f(\mathbf{w}, \mathbf{x})) \right\|_{\text{op}} \geq \frac{t}{3} \right\}.
\end{aligned}$$

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:

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

where  $\textcircled{1}$  holds since by Markov inequality and  $\textcircled{2}$  holds because of Lemma 17.

Therefore, we can set

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

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

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

**Step 2. Bound  $\mathbb{P}(\mathbf{E}_2)$ :** By Lemma 2, 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  $\boldsymbol{\lambda}_\epsilon = \{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k_w}\}$  be an  $\epsilon$ -covering net of  $\mathbb{B}^d(1)$ .

Let  $\boldsymbol{\lambda}_{1/4}$  be the  $\frac{1}{4}$ -covering net of  $\mathbb{B}^d(1)$  but it has only  $s$  nonzero entries. So the size of its  $\epsilon$ -net is

$$\binom{d}{s} \left(\frac{3}{1/4}\right)^s \leq \exp(s \log(12d)).$$

Recall that we use  $j_w$  to denote the index of  $\mathbf{w}_{k_j}^j$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^j$  and we have  $j_w \in [n_\epsilon^j]$ ,  $(n_\epsilon^j \leq \exp(s_j \log(\frac{3rd}{\epsilon})))$ . Then we can bound  $\mathbb{P}(\mathbf{E}_2)$  as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P}\left(\sup_{j_w \in [n_\epsilon^j], j=[l]} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3}\right) \\ &= \mathbb{P}\left(\sup_{j_w \in [n_\epsilon^j], j=[l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/4}} 2 \left| \left\langle \boldsymbol{\lambda}, \left( \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) \right) \boldsymbol{\lambda} \right\rangle \right| \geq \frac{t}{3}\right) \\ &\leq \exp(s \log(12d)) \exp\left(\sum_{j=1}^l s_j \log\left(\frac{3rd}{\epsilon}\right)\right) \sup_{j_w \in [n_\epsilon^j], j=[l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/4}} \mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n \left\langle \boldsymbol{\lambda}, \right. \right. \right. \\ &\quad \left. \left. \left. \left( \nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla^2 f(\mathbf{w}_{k_w}, \mathbf{x})) \right) \boldsymbol{\lambda} \right\rangle \right| \geq \frac{t}{6}\right). \end{aligned}$$

Since by Lemma 20,  $\langle \boldsymbol{\lambda}, (\nabla_w^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_w^2 f(\mathbf{w}, \mathbf{x})) \boldsymbol{\lambda} \rangle$  where  $\boldsymbol{\lambda} \in \mathbb{B}^d(1)$  is  $8\gamma^2\tau^2$ -sub-Gaussian, *i.e.*

$$\mathbb{E}\left(t \langle \boldsymbol{\lambda}, (\nabla_w^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_w^2 f(\mathbf{w}, \mathbf{x})) \boldsymbol{\lambda} \rangle\right) \leq \exp\left(\frac{8t^2\gamma^2\tau^2}{2}\right).$$

Thus,  $\frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\lambda}, (\nabla_w^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_w^2 f(\mathbf{w}, \mathbf{x})) \boldsymbol{\lambda} \rangle$  is  $8\gamma^2\tau^2/n$ -sub-Gaussian random variable. So we can obtain

$$\mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{y}, (\nabla_w^2 f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_w^2 f(\mathbf{w}, \mathbf{x})) \mathbf{y} \rangle \right| \geq \frac{t}{6}\right) \leq 2 \exp\left(-\frac{nt^2}{72\gamma^2\tau^2}\right).$$

Note  $d = \sum_j d_j d_{j-1}$ . Then the probability of  $\mathbf{E}_2$  is upper bounded as

$$\mathbb{P}(\mathbf{E}_2) \leq 2 \exp\left(-\frac{nt^2}{72\gamma^2\tau^2} + s \log\left(\frac{36d^2r}{\epsilon}\right)\right).$$

Thus, if we set

$$t \geq \gamma\tau \sqrt{\frac{72(s \log(36d^2r/\epsilon) + \log(4/\epsilon))}{n}},$$

then we have

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

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

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

We set  $\epsilon$  enough small such that  $\xi \epsilon < t/3$  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  $\epsilon = 36rl^2/n$  and

$$t \geq \max\left(\frac{6\xi\epsilon}{\epsilon}, \gamma\tau\sqrt{\frac{72(s \log(36rd^2/\epsilon) + \log(4/\epsilon))}{n}}\right) = \max\left(\frac{108\xi r}{n\epsilon}, c'_4\gamma\tau\sqrt{\frac{d \log(nl) + \log(4/\epsilon)}{n}}\right).$$

Therefore, there exists such two universal constants  $c_{m'}$  and  $c_m$  such that if  $n \geq \frac{c_{m'}\xi^2 l^2 r^2}{\gamma^2 \tau^2 \epsilon^2 s \log(d) \log(1/\epsilon)}$ , then

$$\sup_{\mathbf{w} \in \Omega} \|\nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w})\|_{\text{op}} \leq c_m \gamma \tau \sqrt{\frac{s \log(dn/l) + \log(4/\epsilon)}{n}}$$

holds with probability at least  $1 - \epsilon$ . □

## D.3 PROOFS OF MAIN THEORIES

### D.3.1 PROOF OF THEOREM 4

*Proof.* Recall that the weight of each layer has magnitude bound separately, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$ . Assume that  $\mathbf{w}_{(j)}$  has  $s_j$  non-zero entries. Then we have  $\sum_{j=1}^l s_j = s$ . So here we separately assume  $\mathbf{w}_\epsilon^j = \{\mathbf{w}_1^j, \dots, \mathbf{w}_{n_\epsilon^j}^j\}$  is the  $\mathbf{d}_j \mathbf{d}_{j-1} \epsilon/d$ -covering net of the ball  $\mathbf{B}^{\mathbf{d}_j \mathbf{d}_{j-1}}(r)$  which corresponds to the weight  $\mathbf{w}_{(j)}$  of the  $j$ -th layer. Let  $n_\epsilon^j$  be the  $\epsilon/l$ -covering number. By  $\epsilon$ -covering theory in (Vershynin, 2012), we can have

$$n_\epsilon^j \leq \binom{\mathbf{d}_j \mathbf{d}_{j-1}}{s_j} \left(\frac{3r}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / (ld)}\right)^{s_j} \leq \exp\left(s_j \log\left(\frac{3r \mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d}\right)\right) = \exp\left(s_j \log\left(\frac{3rd}{\epsilon}\right)\right).$$

Let  $\mathbf{w} \in \Omega$  be an arbitrary vector. Since  $\mathbf{w} = [\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(l)}]$  where  $\mathbf{w}_{(j)}$  is the weight of the  $j$ -th layer, we can always find a vector  $\mathbf{w}_{k_j}^j$  in  $\mathbf{w}_\epsilon^j$  such that  $\|\mathbf{w}_{(j)} - \mathbf{w}_{k_j}^j\|_2 \leq \mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ . For brevity, let  $j_w \in [n_\epsilon^{j_w}]$  denote the index of  $\mathbf{w}_{k_{j_w}}^{j_w}$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^{j_w}$ . Then let  $\mathbf{w}_{k_w} = [\mathbf{w}_{k_1}^{j_1}; \dots; \mathbf{w}_{k_{j_w}}^{j_{j_w}}; \dots; \mathbf{w}_{k_l}^{j_l}]$ . This means that we can always find a vector  $\mathbf{w}_{k_w}$  such that  $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$ . Accordingly, we can

decompose  $\left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right\|_2$  as follows:

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

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

$$\begin{aligned}
\mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right\|_2 \geq t \right\}, \\
\mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f(\mathbf{w}, \mathbf{x}_{(i)}) - \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right\|_2 \geq \frac{t}{3} \right\}, \\
\mathbf{E}_2 &= \left\{ \sup_{j_w \in [n^{\epsilon^j}], j=[l]} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3} \right\}, \\
\mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left\| \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(\nabla f(\mathbf{w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3} \right\}.
\end{aligned}$$

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:

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

where  $\textcircled{1}$  holds because of Markov inequality. Then, we bound  $\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right)$  as follows:

$$\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right) \leq \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 f(\mathbf{w}, \mathbf{x}) \right\|_2 \right) = \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \|\nabla^2 f(\mathbf{w}, \mathbf{x})\|_2 \right) \stackrel{\textcircled{1}}{\leq} \varsigma,$$

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

$$\|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_{\text{op}} \leq \|\nabla_{\mathbf{w}}^2 f(\mathbf{w}, \mathbf{x})\|_F \leq \varsigma,$$

where  $\varsigma = \sqrt{c_{s_1} c_r c_d^2 l^4}$  in which  $c_d = \max_i \mathbf{d}_i$  and  $c_r = \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^{l-1}\right)$ . Therefore, we have

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

We further let

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

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 1, we know that for any vector  $\mathbf{x} \in \mathbb{R}^d$ , its  $\ell_2$ -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  $\boldsymbol{\lambda}_\epsilon = \{\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{k_w}\}$  be an  $\epsilon$ -covering net of  $\mathbb{B}^d(1)$ .

Let  $\boldsymbol{\lambda}_{1/2}$  be the  $\frac{1}{2}$ -covering net of  $\mathbb{B}^d(1)$  but it has only  $s$  nonzero entries. So the size of its  $\epsilon$ -net is

$$\binom{d}{s} \left(\frac{3}{1/2}\right)^s \leq \exp(s \log(6d)).$$

Recall that we use  $j_w$  to denote the index of  $\mathbf{w}_{k_j}^j$  in  $\epsilon$ -net  $\mathbf{w}_\epsilon^j$  and we have  $j_w \in [n_\epsilon^j]$ , ( $n_\epsilon^j \leq \exp(s_j \log(\frac{3rd}{\epsilon}))$ ). Then we can bound  $\mathbb{P}(\mathbf{E}_2)$  as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P}\left(\sup_{j_w \in [n_\epsilon^j], j=[l]} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\|_2 \geq \frac{t}{3}\right) \\ &= \mathbb{P}\left(\sup_{j_w \in [n_\epsilon^j], j=[l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/2}} 2 \left\langle \boldsymbol{\lambda}, \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\rangle \geq \frac{t}{3}\right) \\ &\leq \exp(s \log(6d)) \exp\left(\sum_{j=1}^l s_j \log\left(\frac{3rd}{\epsilon}\right)\right) \sup_{j_w \in [n_\epsilon^j], j=[l], \boldsymbol{\lambda} \in \boldsymbol{\lambda}_{1/2}} \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \left\langle \boldsymbol{\lambda}, \right. \right. \\ &\quad \left. \left. \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \right\rangle \geq \frac{t}{6}\right). \end{aligned}$$

Since by Lemma 19,  $\langle \mathbf{y}, \nabla f(\mathbf{w}, \mathbf{x}) \rangle$  is  $8\beta^2\tau^2$ -sub-Gaussian, *i.e.*

$$\mathbb{E}(\langle \boldsymbol{\lambda}, \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) - \mathbb{E} \nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x}) \rangle) \leq \exp\left(\frac{8\beta^2\tau^2 \|\boldsymbol{\lambda}\|_2^2}{2}\right),$$

where  $\beta = \sqrt{\frac{2^6}{3^8} l(l+2) c_y c_r c_d (l c_r + 1)}$  in which  $c_y$ ,  $c_d$  and  $c_r$  are defined in Lemma 16. Thus,  $\frac{1}{n} \sum_{i=1}^n \langle \mathbf{y}, \nabla f(\mathbf{w}, \mathbf{x}) \rangle$  is  $8\beta^2\tau^2/n$ -sub-Gaussian random variable. Thus, we can obtain

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n \langle \mathbf{y}, \nabla f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(\nabla f(\mathbf{w}_{k_w}, \mathbf{x})) \rangle \geq \frac{t}{6}\right) \leq \exp\left(-\frac{nt^2}{72\beta^2\tau^2}\right).$$

Notice,  $\sum_j \mathbf{d}_j \mathbf{d}_{j-1} = d$ . In this case, the probability of  $\mathbf{E}_2$  is upper bounded as

$$\mathbb{P}(\mathbf{E}_2) \leq \exp\left(-\frac{nt^2}{72\beta^2\tau^2} + d \log\left(\frac{18r}{\epsilon}\right)\right).$$

Thus, if we set

$$t \geq \beta\tau \sqrt{\frac{72(s \log(18d^2r/\epsilon) + \log(4/\epsilon))}{n}},$$

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:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_3) &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \|\mathbb{E}(\nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{x})) - \mathbb{E}(\nabla f(\mathbf{w}, \mathbf{x}))\|_2 \geq \frac{t}{3}\right) \\ &= \mathbb{P}\left(\sup_{\mathbf{w} \in \Omega} \frac{\|\mathbb{E}(\nabla f(\mathbf{w}_{k_{\mathbf{w}}}, \mathbf{x})) - \nabla f(\mathbf{w}, \mathbf{x})\|_2}{\|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_{\mathbf{w}}}\|_2 \geq \frac{t}{3}\right) \\ &\leq \mathbb{P}\left(\epsilon \mathbb{E} \sup_{\mathbf{w} \in \Omega} \|\nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x})\|_2 \geq \frac{t}{3}\right) \\ &\stackrel{\textcircled{1}}{\leq} \mathbb{P}\left(\varsigma \epsilon \geq \frac{t}{3}\right). \end{aligned}$$

where  $\textcircled{1}$  holds since by Lemma 17. We set  $\epsilon$  enough small such that  $\varsigma \epsilon < t/3$  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  $\epsilon = 18r/n$  and

$$\begin{aligned} t &\geq \max\left(\frac{6\varsigma\epsilon}{\epsilon}, \beta\tau\sqrt{\frac{72(s\log(18d^2r/\epsilon) + \log(4/\epsilon))}{n}}\right) \\ &= \max\left(\frac{108\varsigma r}{n\epsilon}, \beta\tau\sqrt{\frac{72(s\log(nl) + \log(4/\epsilon))}{n}}\right). \end{aligned}$$

Note that  $\varsigma = \mathcal{O}(\sqrt{lc_d\beta})$ . Therefore, there exists a universal constant  $c_{y'}$  such that if  $n \geq c_{y'}c_d l^3 r^2 / (s \log(d)\tau^2 \epsilon^2 \log(1/\epsilon))$ , then

$$\sup_{\mathbf{w} \in \Omega} \|\nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w})\|_2 \leq \tau \sqrt{\frac{512}{729} c_y l(l+2)(lc_r + 1)c_r c_d \sqrt{\frac{s \log(dn/l) + \log(4/\epsilon)}{n}}}$$

holds with probability at least  $1 - \epsilon$ , where  $c_y, c_d$  and  $c_r$  are defined in Lemma 16.  $\square$

### D.3.2 PROOF OF THEOREM 5

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

$$\inf_i \left| \lambda_i^k \left( \nabla^2 \mathbf{J}(\mathbf{w}^{(k)}) \right) \right| \geq \zeta,$$

where  $\lambda_i^k \left( \nabla^2 \mathbf{J}(\mathbf{w}^{(k)}) \right)$  denotes the  $i$ -th eigenvalue of the Hessian  $\nabla^2 \mathbf{J}(\mathbf{w}^{(k)})$  and  $\zeta$  is a constant. We further define a set  $D = \{\mathbf{w} \in \mathbb{R}^d \mid \|\nabla \mathbf{J}(\mathbf{w})\|_2 \leq \epsilon \text{ and } \inf_i |\lambda_i \left( \nabla^2 \mathbf{J}(\mathbf{w}^{(k)}) \right)| \geq \zeta\}$ . According to Lemma 4,  $D = \cup_{k=1}^{\infty} D_k$  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  $\mathbf{J}(\mathbf{w})$  for  $k \geq m+1$ . On the other hand, by the continuity of  $\nabla \mathbf{J}(\mathbf{w})$ , it yields  $\|\nabla \mathbf{J}(\mathbf{w})\|_2 = \epsilon$  for  $\mathbf{w} \in \partial D_k$ . Notice, we set the value of  $\epsilon$  blow which is actually a function related  $n$ .

Then by utilizing Theorem 4, we let sample number  $n$  sufficient large such that

$$\sup_{\mathbf{w} \in \Omega} \|\nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w})\|_2 \leq \beta \triangleq \frac{\epsilon}{2}$$

holds with probability at least  $1 - \epsilon$ , where  $\beta = \tau \sqrt{\frac{512}{729} c_y l(l+2)(lc_r + 1)c_r c_d \sqrt{\frac{s \log(dn/l) + \log(4/\epsilon)}{n}}}$ .

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

$$\begin{aligned} \inf_{\mathbf{w} \in D_k} \left\| t \nabla \hat{\mathbf{J}}_n(\mathbf{w}) + (1-t) \nabla \mathbf{J}(\mathbf{w}) \right\|_2 &= \inf_{\mathbf{w} \in D_k} \left\| t \left( \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right) + \nabla \mathbf{J}(\mathbf{w}) \right\|_2 \\ &\geq \inf_{\mathbf{w} \in D_k} \|\nabla \mathbf{J}(\mathbf{w})\|_2 - \sup_{\mathbf{w} \in D_k} t \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla \mathbf{J}(\mathbf{w}) \right\|_2 \\ &\geq \frac{\epsilon}{2}. \end{aligned} \tag{26}$$

Similarly, by utilizing Lemma 21, let  $n$  be sufficient large such that

$$\sup_{\mathbf{w} \in \Omega} \left\| \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) - \nabla^2 \mathbf{J}(\mathbf{w}) \right\|_{\text{op}} \leq c_m \gamma \tau \sqrt{\frac{s \log(dn/l) + \log(4/\varepsilon)}{n}} \leq \frac{\zeta}{2}$$

holds with probability at least  $1 - \varepsilon$ . 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  $\lambda_i^k \left( \nabla^2 \hat{\mathbf{J}}_n(\mathbf{w}) \right)$  for arbitrary  $\mathbf{w} \in D_k$  as follows:

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

This means that in each set  $D_k$ ,  $\nabla^2 \hat{\mathbf{J}}_n(\mathbf{w})$  has no zero eigenvalues. Then, combining this and Eqn. (26), by Lemma 3 we know that if the population risk  $\mathbf{J}(\mathbf{w})$  has no critical point in  $D_k$ , then the empirical risk  $\hat{\mathbf{J}}_n(\mathbf{w})$  has also no critical point in  $D_k$ ; otherwise it also holds. By Lemma 3, we can also obtain that in  $D_k$ , if  $\mathbf{J}(\mathbf{w})$  has a unique critical point  $\mathbf{w}^{(k)}$  with non-degenerate index  $s_k$ , then  $\hat{\mathbf{J}}_n(\mathbf{w})$  also has a unique critical point  $\mathbf{w}_n^{(k)}$  in  $D_k$  with the same non-degenerate index  $s_k$ . The first conclusion is proved.

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

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

where  $\textcircled{1}$  holds since  $\nabla \mathbf{J}(\mathbf{w}^{(k)}) = \mathbf{0}$  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  $\inf_i |\lambda_i(\nabla^2 \mathbf{J}(\mathbf{w}))| \geq \zeta$ . Then if  $n \geq c_s \max(c_d l^3 r^2 / (s \log(d) \tau^2 \varepsilon^2 \log(1/\varepsilon)), s \log(d/l) / \zeta^2)$  where  $c_s$  is a constant, then

$$\left\| \mathbf{w}_n^{(k)} - \mathbf{w}^{(k)} \right\|_2 \leq \frac{2\tau}{\zeta} \sqrt{\frac{512}{729} c_y l (l+2) (lc_r + 1) c_r c_d} \sqrt{\frac{s \log(dn/l) + \log(4/\varepsilon)}{n}}$$

holds with probability at least  $1 - \varepsilon$ . The proof is completed.  $\square$

### D.3.3 PROOF OF THEOREM 6

*Proof.* Recall that the weight of each layer has magnitude bound separately, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$ . Assume that  $\mathbf{w}_{(j)}$  has  $s_j$  non-zero entries. Then we have  $\sum_{j=1}^l s_j = s$ . So here we separately assume  $\mathbf{w}_\epsilon^j = \{\mathbf{w}_{1,\epsilon}^j, \dots, \mathbf{w}_{n_\epsilon^j}^j\}$  is the  $d_j d_{j-1} \epsilon / d$ -covering net of the ball  $\mathbb{B}^{d_j d_{j-1}}(r)$  which corresponds

to the weight  $\mathbf{w}_{(j)}$  of the  $j$ -th layer. Let  $n_{\epsilon^j}$  be the  $\epsilon/l$ -covering number. By  $\epsilon$ -covering theory in (Vershynin, 2012), we can have

$$n_{\epsilon^j} \leq \binom{\mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{s}_j} \left( \frac{3r}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / (ld)} \right)^{\mathbf{s}_j} \leq \exp \left( \mathbf{s}_j \log \left( \frac{3r \mathbf{d}_j \mathbf{d}_{j-1}}{\mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d} \right) \right) = \exp \left( \mathbf{s}_j \log \left( \frac{3rd}{\epsilon} \right) \right).$$

Let  $\mathbf{w} \in \Omega$  be an arbitrary vector. Since  $\mathbf{w} = [\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(l)}]$  where  $\mathbf{w}_{(j)}$  is the weight of the  $j$ -th layer, we can always find a vector  $\mathbf{w}_{k_w}^j$  in  $\mathbf{w}_{\epsilon^j}$  such that  $\|\mathbf{w}_{(j)} - \mathbf{w}_{k_w}^j\|_2 \leq \mathbf{d}_j \mathbf{d}_{j-1} \epsilon / d$ . For brevity, let  $j_w \in [n_{\epsilon^j}]$  denote the index of  $\mathbf{w}_{k_w}^j$  in  $\epsilon$ -net  $\mathbf{w}_{\epsilon^j}$ . Then let  $\mathbf{w}_{k_w} = [\mathbf{w}_{k_1}^{j_1}; \dots; \mathbf{w}_{k_j}^{j_j}; \dots; \mathbf{w}_{k_l}^{j_l}]$ . This means that we can always find a vector  $\mathbf{w}_{k_w}$  such that  $\|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \leq \epsilon$ . Accordingly, we can decompose  $|\hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w})|$  as

$$\begin{aligned} |\hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w})| &= \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}, \mathbf{x})) \right| \\ &= \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) + \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x})) + \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x}) - \mathbb{E}f(\mathbf{w}, \mathbf{x}) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right| + \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x})) \right| + \left| \mathbb{E}f(\mathbf{w}_{k_w}, \mathbf{x}) - \mathbb{E}f(\mathbf{w}, \mathbf{x}) \right|. \end{aligned}$$

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

$$\begin{aligned} \mathbf{E}_0 &= \left\{ \sup_{\mathbf{w} \in \Omega} |\hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w})| \geq t \right\}, \\ \mathbf{E}_1 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}, \mathbf{x}_{(i)}) - f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)})) \right| \geq \frac{t}{3} \right\}, \\ \mathbf{E}_2 &= \left\{ \sup_{j_w \in [n_{\epsilon^j}], j=[l]} \left| \frac{1}{n} \sum_{i=1}^n (f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x}))) \right| \geq \frac{t}{3} \right\}, \\ \mathbf{E}_3 &= \left\{ \sup_{\mathbf{w} \in \Omega} \left| \mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(f(\mathbf{w}, \mathbf{x})) \right| \geq \frac{t}{3} \right\}. \end{aligned}$$

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:

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

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

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

Now we only need to bound  $\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right)$ . Then by Lemma 16, we can bound it as follows:

$$\mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \nabla \hat{\mathbf{J}}_n(\mathbf{w}, \mathbf{x}) \right\|_2 \right) \leq \mathbb{E} \left( \sup_{\mathbf{w} \in \Omega} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(\mathbf{w}, \mathbf{x}_{(i)}) \right\|_2 \right) \leq \alpha,$$

where  $\alpha = \sqrt{\frac{1}{16} c_y c_d (1 + c_r (l - 1))}$  in which  $c_y$ ,  $c_d$  and  $c_r$  are defined in Lemma 16.

Therefore, we have

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

We further let

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

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

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

**Step 2. Bound  $\mathbb{P}(\mathbf{E}_2)$ :** Recall that we use  $j_w$  to denote the index of  $\mathbf{w}_{k_j}^j$  in  $\epsilon$ -net  $\mathbf{w}_{\epsilon}^j$  and we have  $j_w \in [n_{\epsilon}^j]$ , ( $n_{\epsilon}^j \leq \exp(s_j \log(\frac{3rd}{\epsilon}))$ ). We can bound  $\mathbb{P}(\mathbf{E}_2)$  as follows:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_2) &= \mathbb{P} \left( \sup_{j_w \in [n_{\epsilon}^j], j=[l]} \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_{k_w}, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) \right| \geq \frac{t}{3} \right) \\ &\leq \exp \left( \sum_{j=1}^l s_j \log \left( \frac{3rd}{\epsilon} \right) \right) \sup_{j_w \in [n_{\epsilon}^j], j=[l]} \mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_j, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_j, \mathbf{x})) \right| \geq \frac{t}{3} \right). \end{aligned}$$

Since when the activation functions are sigmoid functions, the loss  $f(\mathbf{w}, \mathbf{x})$  is  $\alpha$ -Lipschitz. Besides, we assume  $\mathbf{x}$  to be a vector of *i.i.d.* Gaussian variables from  $\mathcal{N}(0, \tau^2)$ . Then by Lemma 25, we know that the variable  $f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})$  is  $8\alpha^2\tau^2$ -sub-Gaussian. Thus, we have

$$\mathbb{P}(|f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})| > t) \leq 2 \exp \left( -\frac{t^2}{2\alpha^2\tau^2} \right), \quad (\forall t \geq 0),$$

where  $\alpha = \sqrt{\frac{1}{16} c_y c_d (1 + c_r (l - 1))}$  in which  $c_y$ ,  $c_d$  and  $c_r$  are defined in Lemma 16. Thus,  $\frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_j, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_j, \mathbf{x}))$  is  $8\alpha^2\tau^2/n$ -sub-Gaussian random variable. Thus, we can obtain

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}_j, \mathbf{x}_{(i)}) - \mathbb{E}(f(\mathbf{w}_j, \mathbf{x})) \right| \geq \frac{t}{3} \right) \leq 2 \exp \left( -\frac{nt^2}{18\alpha^2\tau^2} \right).$$

Notice  $\sum_{j=1}^l s_j = s$ . In this case, the probability of  $\mathbf{E}_2$  is upper bounded as

$$\mathbb{P}(\mathbf{E}_2) \leq 2 \exp \left( -\frac{nt^2}{18\alpha^2\tau^2} + s \log \left( \frac{3dr}{\epsilon} \right) \right).$$

Thus, if we set

$$t \geq \alpha\tau \sqrt{\frac{18(s \log(3dr/\epsilon) + \log(4/\epsilon))}{n}},$$

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:

$$\begin{aligned} \mathbb{P}(\mathbf{E}_3) &= \mathbb{P} \left( \sup_{\mathbf{w} \in \Omega} |\mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) - \mathbb{E}(f(\mathbf{w}, \mathbf{x}))| \geq \frac{t}{3} \right) \\ &= \mathbb{P} \left( \sup_{\mathbf{w} \in \Omega} \frac{|\mathbb{E}(f(\mathbf{w}_{k_w}, \mathbf{x})) - f(\mathbf{w}, \mathbf{x})|}{\|\mathbf{w} - \mathbf{w}_{k_w}\|_2} \sup_{\mathbf{w} \in \Omega} \|\mathbf{w} - \mathbf{w}_{k_w}\|_2 \geq \frac{t}{3} \right) \\ &\leq \mathbb{P} \left( \epsilon \mathbb{E} \sup_{\mathbf{w} \in \Omega} \|\nabla \mathbf{J}_w(\mathbf{w}, \mathbf{x})\|_2 \geq \frac{t}{3} \right) \\ &\stackrel{\textcircled{1}}{\leq} \mathbb{P} \left( \alpha\epsilon \geq \frac{t}{3} \right), \end{aligned}$$

where ① holds since by Lemma 16, for arbitrary  $\mathbf{x}$  and  $\mathbf{w} \in \Omega$ , we have  $\|\nabla_{\mathbf{w}} f(\mathbf{w}, \mathbf{x})\|_2 \leq \alpha$ . We set  $\epsilon$  enough small such that  $\alpha\epsilon < t/3$  always holds. Then it yields  $\mathbb{P}(\mathbf{E}_3) = 0$ .

**Step 4. Final result:** Notice, we have  $\frac{6\alpha\epsilon}{\epsilon} \geq 3\alpha\epsilon$ . To ensure  $\mathbb{P}(\mathbf{E}_0) \leq \epsilon$ , we just set  $\epsilon = 3r/n$  and

$$t \geq \max\left(\frac{6\alpha\epsilon}{\epsilon}, \alpha\tau\sqrt{\frac{18(s\log(3dr/\epsilon) + \log(4/\epsilon))}{n}}\right) = \max\left(\frac{18\alpha r}{n\epsilon}, \alpha\tau\sqrt{\frac{18(s\log(nd/l) + \log(4/\epsilon))}{n}}\right).$$

Therefore, if  $n \geq 18l^2r^2/(s\log(d)\tau^2\epsilon^2\log(1/\epsilon))$ , then

$$\sup_{\mathbf{w} \in \Omega} \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| \leq \tau\sqrt{\frac{9}{8}c_y c_d (1 + c_r(l-1))} \sqrt{\frac{s\log(nd/l) + \log(4/\epsilon)}{n}}$$

holds with probability at least  $1 - \epsilon$ , where  $c_y$ ,  $c_d$ , and  $c_r$  are defined as

$$\|\mathbf{v}^{(l)} - \mathbf{y}\|_2^2 \leq c_y < +\infty, \quad c_d = \max(\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_l) \quad \text{and} \quad c_r = \max\left(\frac{r^2}{16}, \left(\frac{r^2}{16}\right)^{l-1}\right).$$

The proof is completed.  $\square$

#### D.3.4 PROOF OF COROLLARY 2

*Proof.* By Lemma 5, we know  $\epsilon_s = \epsilon_g$ . Thus, the remaining work is to bound  $\epsilon_s$ . Actually, we can have

$$\begin{aligned} \left| \mathbb{E}_{\mathcal{S} \sim \mathcal{D}, \mathbf{A}, (\mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) \sim \mathcal{D}} \frac{1}{n} \sum_{j=1}^n (f_j(\mathbf{w}_*^j; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)}) - f_j(\mathbf{w}^n; \mathbf{x}'_{(j)}, \mathbf{y}'_{(j)})) \right| &\leq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}} \left( \sup_{\mathbf{w} \in \Omega} \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| \right) \\ &\leq \sup_{\mathbf{w} \in \Omega} \left| \hat{\mathbf{J}}_n(\mathbf{w}) - \mathbf{J}(\mathbf{w}) \right| \\ &\leq \epsilon_n. \end{aligned}$$

Thus, we have  $\epsilon_g = \epsilon_s \leq \epsilon_n$ . The proof is completed.  $\square$

#### D.4 PROOF OF OTHER LEMMAS

##### D.4.1 PROOF OF LEMMA 22

*Proof.* Since  $\mathbb{G}(\mathbf{u}^{(i)})$  is a diagonal matrix and its diagonal values are upper bounded by  $\sigma(\mathbf{u}_h^{(i)})(1 - \sigma(\mathbf{u}_h^{(i)})) \leq 1/4$  where  $\mathbf{u}_h^{(i)}$  denotes the  $h$ -th entry of  $\mathbf{u}^{(i)}$ , we can conclude

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

Note that  $\mathbf{P}_k$  is a matrix of size  $\mathbf{d}_k^2 \times \mathbf{d}_k$  whose  $((s-1)\mathbf{d}_k + s, s)$  ( $s = 1, \dots, \mathbf{d}_k$ ) entry equal to  $\sigma(\mathbf{u}_s^{(k)})(1 - \sigma(\mathbf{u}_s^{(k)}))(1 - 2\sigma(\mathbf{u}_s^{(k)}))$  and rest entries are all 0. This gives

$$\begin{aligned} \sigma(\mathbf{u}_s^{(k)})(1 - \sigma(\mathbf{u}_s^{(k)}))(1 - 2\sigma(\mathbf{u}_s^{(k)})) &= \frac{1}{3}(3\sigma(\mathbf{u}_s^{(k)}))(1 - \sigma(\mathbf{u}_s^{(k)}))(1 - 2\sigma(\mathbf{u}_s^{(k)})) \\ &\leq \frac{1}{3} \left( \frac{3\sigma(\mathbf{u}_s^{(k)}) + 1 - \sigma(\mathbf{u}_s^{(k)}) + 1 - 2\sigma(\mathbf{u}_s^{(k)})}{3} \right)^3 \\ &\leq \frac{2^3}{3^4}. \end{aligned}$$

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

$$\|\mathbf{P}_k\mathbf{M}\|_F^2 \leq \frac{2^6}{3^8}\|\mathbf{M}\|_F^2 \quad \text{and} \quad \|\mathbf{N}\mathbf{P}_k\|_F^2 \leq \frac{2^6}{3^8}\|\mathbf{N}\|_F^2.$$

As for  $\mathbf{B}_{s:t}$ , we have

$$\begin{aligned}\|\mathbf{B}_{s:t}\|_F^2 &\leq \|\mathbf{A}_s\|_F^2 \|\mathbf{A}_{s+1}\|_F^2 \cdots \|\mathbf{A}_t\|_F^2 \\ &= \left\| (\mathbf{W}^s)^T \mathbb{G}(\mathbf{u}^{(s)}) \right\|_F^2 \left\| (\mathbf{W}^{(s+1)})^T \mathbb{G}(\mathbf{u}^{(s+1)}) \right\|_F^2 \cdots \left\| (\mathbf{W}^{(t)})^T \mathbb{G}(\mathbf{u}^{(t)}) \right\|_F^2 \\ &\leq \frac{1}{16^{t-s+1}} \left\| \mathbf{W}^{(s)} \right\|_F^2 \left\| \mathbf{W}^{(s+1)} \right\|_F^2 \cdots \left\| \mathbf{W}^{(t)} \right\|_F^2 \\ &= \frac{1}{16^{t-s+1}} \mathbf{D}_{s:t}.\end{aligned}$$

Since the  $\ell_2$ -norm of each  $\mathbf{w}_{(j)}$  is bounded, i.e.  $\|\mathbf{w}_{(j)}\|_2 \leq r$ , we can obtain

$$\frac{1}{16^{t-s+1}} \mathbf{D}_{s:t} \leq \frac{1}{16^{t-s+1}} r^{2(t-s+1)} = \left(\frac{r}{4}\right)^{2(t-s+1)} \triangleq c_{st}.$$

Now we prove the final result. According to the property of Kronecker product that for any matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{X}$  of proper sizes,  $\text{vec}(\mathbf{A}\mathbf{X}\mathbf{B}) = (\mathbf{B}^T \otimes \mathbf{A})\text{vec}(\mathbf{X})$ , we have

$$\text{vec}(\mathbf{M}\mathbf{N}^T) = (\mathbf{N} \otimes \mathbf{I})\text{vec}(\mathbf{M}) = (\mathbf{N} \otimes \mathbf{I})\mathbf{m}.$$

This further yields

$$\|(\mathbf{N} \otimes \mathbf{I})\mathbf{m}\|_F^2 = \|\text{vec}(\mathbf{M}\mathbf{N}^T)\|_F^2 = \|\mathbf{M}\mathbf{N}^T\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2.$$

By similar way, we can obtain

$$\|(\mathbf{I} \otimes \mathbf{N})\mathbf{m}\|_F^2 \leq \|\mathbf{M}\|_F^2 \|\mathbf{N}\|_F^2.$$

The proof is completed.  $\square$

#### D.4.2 PROOF OF LEMMA 23

*Proof.* By utilizing the chain rule in Eqn. (24) in Sec. D.2.1, we can easily compute  $\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}}$  and  $\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{v}^{(i)}}$  as follows:

$$\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}} = \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{A}_{i+1} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) = \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{B}_{i+1:l} (\mathbf{v}^{(l)} - \mathbf{y})$$

and

$$\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{v}^{(i)}} = \mathbf{A}_{i+1} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) = \mathbf{B}_{i+1:l} (\mathbf{v}^{(l)} - \mathbf{y}).$$

Therefore, we can further obtain

$$\begin{aligned}&\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)}} \\ &= \text{vec} \left( \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \right) (\mathbf{v}^{(j-1)})^T \right) \\ &= \text{vec} \left( \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_{i-1} (\mathbf{W}^{(i)})^T \right) \left( \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{A}^{i+1} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \right) (\mathbf{v}^{(j-1)})^T \right) \\ &= (\mathbf{v}^{(j-1)}) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_{i-1} (\mathbf{W}^{(i)})^T \right) \text{vec} \left( \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{A}_{i+1} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \right) \\ &= (\mathbf{v}^{(j-1)}) \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_{i-1} (\mathbf{W}^{(i)})^T \right) \left( \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}} \right).\end{aligned}$$

Note that we have  $\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}_{(j)}} = \frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{w}_{(j)}} \left( \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}} \right)$ . This gives

$$\frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{w}_{(j)}} = (\mathbf{v}^{(j-1)})^T \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:i-1} (\mathbf{W}^{(i)})^T \right)^T \in \mathbb{R}^{\mathbf{d}_i \times \mathbf{d}_j \mathbf{d}_{j-1}} \quad (i > j).$$

When  $i = j$ , we have

$$\frac{\partial \mathbf{u}^{(i)}}{\partial \mathbf{w}_{(i)}} = (\mathbf{v}^{(i-1)})^T \otimes \mathbf{I}_{\mathbf{d}_i} \in \mathbb{R}^{\mathbf{d}_i \times \mathbf{d}_i \mathbf{d}_{i-1}}.$$

Similarly, we can obtain

$$\frac{\partial \mathbf{v}^{(i)}}{\partial \mathbf{w}_{(j)}} = (\mathbf{v}^{(j-1)})^T \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \mathbf{A}_{j+2} \cdots \mathbf{A}_i \right)^T = (\mathbf{v}^{(j-1)})^T \otimes \left( \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{B}_{j+1:i} \right)^T \in \mathbb{R}^{\mathbf{d}_i \times \mathbf{d}_j \mathbf{d}_{j-1}} \quad (i \geq j).$$

The proof is completed.  $\square$

## D.4.3 PROOF OF LEMMA 24

*Proof.* By Lemma 23, we have

$$\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(i)}} = \mathbb{G}(\mathbf{u}^{(i)}) \mathbf{B}_{i+1:l} (\mathbf{v}^{(l)} - \mathbf{y}) \quad \text{and} \quad \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{v}^{(i)}} = \mathbf{B}_{i+1:l} (\mathbf{v}^{(l)} - \mathbf{y}).$$

Therefore, we can further obtain

$$\begin{aligned} \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(1)}} &= \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{A}_2 \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \\ &= \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{A}_2 \cdots \mathbf{A}_{j-1} (\mathbf{W}^j)^T \mathbb{G}(\mathbf{u}^{(j)}) \mathbf{A}_{j+1} \cdots \mathbf{A}_l (\mathbf{v}^{(l)} - \mathbf{y}) \\ &= \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{A}_2 \cdots \mathbf{A}_{j-1} (\mathbf{W}^j)^T \right) \left( \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(j)}} \right). \end{aligned}$$

Note that we have  $\frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(1)}} = \left( \frac{\partial \mathbf{u}^{(j)}}{\partial \mathbf{u}^{(1)}} \right)^T \left( \frac{\partial f(\mathbf{w}, \mathbf{x})}{\partial \mathbf{u}^{(j)}} \right)$ . This gives

$$\frac{\partial \mathbf{u}^{(j)}}{\partial \mathbf{u}^{(1)}} = \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{A}_2 \cdots \mathbf{A}_{j-1} (\mathbf{W}^j)^T \right)^T = \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:j-1} (\mathbf{W}^j)^T \right)^T \in \mathbb{R}^{d_j \times d_1} \quad (j > 1).$$

Similarly, we can obtain

$$\frac{\partial \mathbf{v}^{(j)}}{\partial \mathbf{u}^{(1)}} = \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{A}_2 \cdots \mathbf{A}_j \right)^T = \left( \mathbb{G}(\mathbf{u}^{(1)}) \mathbf{B}_{2:j} \right)^T \in \mathbb{R}^{d_j \times d_1} \quad (j > 1).$$

The proof is completed.  $\square$