

Guaranteed Recovery of One-Hidden-Layer Neural Networks via Cross Entropy

Haoyu Fu, *Student Member*, Yuejie Chi, *Senior Member*, Yingbin Liang, *Senior Member*

Abstract—We study model recovery for data classification, where the training labels are generated from a one-hidden-layer neural network with sigmoid activations, also known as a single-layer feedforward network, and the goal is to recover the weights of the neural network. We consider two network models, the fully-connected network (FCN) and the non-overlapping convolutional neural network (CNN). We prove that with Gaussian inputs, the empirical risk based on cross entropy exhibits strong convexity and smoothness *uniformly* in a local neighborhood of the ground truth, as soon as the sample complexity is sufficiently large. This implies that if initialized in this neighborhood, gradient descent converges linearly to a critical point that is provably close to the ground truth. Furthermore, we show such an initialization can be obtained via the tensor method. This establishes the global convergence guarantee for empirical risk minimization using cross entropy via gradient descent for learning one-hidden-layer neural networks, at the near-optimal sample and computational complexity with respect to the network input dimension without unrealistic assumptions such as requiring a fresh set of samples at each iteration.

I. INTRODUCTION

Neural networks have attracted a significant amount of research interest in recent years due to the success of deep neural networks [1] in practical domains such as computer vision and artificial intelligence [2], [3], [4]. However, the theoretical underpinnings behind such success remains mysterious to a large extent. Efforts have been taken to understand which classes of functions can be represented by deep neural networks [5], [6], [7], [8], when (stochastic) gradient descent is effective for optimizing a nonconvex loss function [9], and why these networks generalize well [10], [11], [12].

One important line of research that has attracted extensive attention is the model-recovery problem, which is important for the network to generalize well [13]. Specifically, it is shown in [13] that in a model-recovery setting, a network cannot generalize well if the underlying parameters cannot be recovered accurately, therefore linking model recovery to generalization. In addition, the problem of model recovery provides a framework to leverage statistical nature of the input data in an intuitive manner, which allows shedding more light to the understanding of optimization of complex neural networks.

H. Fu and Y. Liang are with Dept. of ECE, The Ohio State University, Columbus, OH 43210, USA. Emails: {fu.436, liang.889}@osu.edu.

Y. Chi is with Dept. of ECE, Carnegie Mellon University, Pittsburgh, PA 15213, USA. Email: yuejiechi@cmu.edu.

The work of H. Fu and Y. Liang is supported in part by U.S. National Science Foundation under the grants CCF-1761506, CCF-1801855 and CCF-1900145. The work of Y. Chi is supported in part by AFOSR under the grant FA9550-15-1-0205, by ONR under the grant N00014-18-1-2142, by ARO under the grant W911NF-18-1-0303, and by NSF under the grants CAREER ECCS-1818571, ECCS-1833553, CCF-1806154 and CCF-1901199.

Let the training samples $(\mathbf{x}_i, y_i) \sim (\mathbf{x}, y)$, $i = 1, \dots, n$, be generated independently and identically distributed (i.i.d.) from a distribution \mathcal{D} based on a neural network model with the ground truth parameter \mathbf{W}^* , and the goal is to recover \mathbf{W}^* using the training samples given the network architecture. Consider a network whose output is given as $H(\mathbf{W}^*, \mathbf{x})$. Previous studies along this topic can be mainly divided into two cases of data generations, with the input $\mathbf{x} \in \mathbb{R}^d$ being drawn from the Gaussian distribution.

- *Regression*, where each sample $y \in \mathbb{R}$ is generated as

$$y = H(\mathbf{W}^*, \mathbf{x}).$$

This type of regression problem has been studied in various settings. In particular, [14] studied the single-neuron model under the Rectified Linear Unit (ReLU) activation, [15] studied the one-hidden-layer multi-neuron network model, and [16] studied a two-layer feedforward network with ReLU activations and identity mapping.

- *Classification*, where a label $y \in \{0, 1\}$ is drawn according to the conditional distribution

$$\mathbb{P}(y = 1 | \mathbf{x}) = H(\mathbf{W}^*, \mathbf{x}).$$

Such a problem has been studied in [17] when the network contains only a single neuron.

For both cases, previous studies attempted to recover \mathbf{W}^* , by minimizing an empirical loss function using the squared loss, i.e. $\min_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n (y_i - H(\mathbf{W}, \mathbf{x}_i))^2$, given the training data. Two types of statistical guarantees were provided for such model recovery problems using the squared loss. More specifically, [15] showed that in the local neighborhood of the ground truth \mathbf{W}^* , the *empirical* loss function is strongly convex for each *given* point under *independent* high probability event, which implies that *fresh samples* are required at *every* iteration for gradient descent to converge linearly with well-designed initializations. On the other hand, studies such as [17] established strong convexity in the entire local neighborhood of the ground truth in a uniform sense, so that resampling per iteration is not needed for gradient descent to have guaranteed linear convergence as long as it enters such a local neighborhood. Here, one weakness of the pointwise strong convexity in [15], compared to the uniform strong convexity in [17], is that independent fresh samples are required at each iteration to guarantee the linear convergence of gradient descent. Consequently, the sample complexity of [15] grows with respect to the recovery accuracy ϵ , typically with an extra factor of $\log(1/\epsilon)$ under linear convergence, which can be large when the desired accuracy is high. Therefore, the latter type of uniform

strong convexity *without requiring per-iteration resampling* is much stronger and more desirable.

In this paper, we focus on the classification setting by minimizing the empirical loss using the cross entropy objective, which is a popular choice in training practical neural networks. The geometry as well as the optimization of the model recovery problem based on the cross-entropy loss function have not yet been understood even for one-hidden-layer networks. The main focus of this paper is to develop technical analysis for guaranteed model recovery under the challenging cross-entropy loss function for the classification problem for two types of one-hidden-layer network structures.

A. Problem Formulation

We consider two popular types of one-hidden-layer nonlinear neural networks illustrated in Fig. 1, i.e., a Fully-Connected Network (FCN) [15] and a non-overlapping Convolutional Neural Network (CNN) [18]. For both cases, we let $\mathbf{x} \in \mathbb{R}^d$ be the input, $K \geq 1$ be the number of neurons, and the activation function be the sigmoid function

$$\phi(x) = \frac{1}{1 + \exp(-x)}.$$

- *FCN*: the network parameter is $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K] \in \mathbb{R}^{d \times K}$, and

$$H_{\text{FCN}}(\mathbf{W}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}_k^\top \mathbf{x}). \quad (1)$$

- *Non-overlapping CNN*: for simplicity we let $d = mK$ for some integers m . Let $\mathbf{w} \in \mathbb{R}^m$ be the network parameter, and the k th stride of \mathbf{x} be given as $\mathbf{x}^{(k)} = [x_{m(k-1)+1}, \dots, x_{m \cdot k}]^\top \in \mathbb{R}^m$. Then,

$$H_{\text{CNN}}(\mathbf{w}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}^\top \mathbf{x}^{(k)}). \quad (2)$$

The non-overlapping CNN model can be viewed as a highly structured instance of the FCN, where the weight matrix can be written as:

$$\mathbf{W}_{\text{CNN}} = \begin{bmatrix} \mathbf{w} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{w} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{w} \end{bmatrix} \in \mathbb{R}^{d \times K}.$$

In a model recovery setting, we are given n training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim (\mathbf{x}, y)$ that are drawn i.i.d. from certain distribution regarding the ground truth network parameter \mathbf{W}^* (or resp. \mathbf{w}^* for CNN). Suppose the network input $\mathbf{x} \in \mathbb{R}^d$ is drawn from a standard Gaussian distribution $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. This assumption has been used a lot in previous literature [14], [19], [18], [20], to name a few. Then, conditioned on $\mathbf{x} \in \mathbb{R}^d$, the output y is mapped to $\{0, 1\}$ via the output of the neural network, i.e.,

$$\mathbb{P}(y = 1 | \mathbf{x}) = H(\mathbf{W}^*, \mathbf{x}). \quad (3)$$

Our goal is to recover the network parameter, i.e., \mathbf{W}^* . One natural choice is to maximize the log-likelihood function, which turns out to be equivalent to minimizing

$$f_n(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{W}; \mathbf{x}_i, y_i), \quad (4)$$

where $\ell(\mathbf{W}; \mathbf{x}, y)$ is the cross-entropy loss function, i.e.,

$$\begin{aligned} \ell(\mathbf{W}; \mathbf{x}, y) &= -y \cdot \log(H(\mathbf{W}, \mathbf{x})) - (1 - y) \cdot \log(1 - H(\mathbf{W}, \mathbf{x})), \end{aligned} \quad (5)$$

where $H(\mathbf{W}, \mathbf{x})$ can subsume either H_{FCN} or H_{CNN} . Although the squared loss has been used in [17] to study the classification problem with a single neuron, the cross-entropy loss is a more natural and popular choice in practice for classification data, due to its natural connection to the principle of maximum likelihood estimation.

B. Our Contributions

Considering the multi-neuron classification problem with either FCN or CNN, the main contributions of this work are summarized as follows. Throughout the discussions below, we assume the number K of neurons is a constant, and state the scaling only in terms of the input dimension d and the number n of samples.

- *Uniform local strong convexity*: If the input is Gaussian, the empirical risk function $f_n(\mathbf{W})$ is *uniformly* strongly convex in a local neighborhood of the ground truth \mathbf{W}^* as soon as the sample size $n = O(d \log^2 d)$.
- *Statistical and computational rate of gradient descent*: consequently, if initialized in this neighborhood, gradient descent converges linearly to a critical point (which we show to exist). Due to the nature of quantized labels here, the recovery of the ground truth is only up to certain statistical accuracy. In particular, gradient descent finds the critical point $\widehat{\mathbf{W}}_n$ with a computation cost of $O(nd \log(1/\epsilon))$, where ϵ denotes the numerical accuracy and $\widehat{\mathbf{W}}_n$ converges to \mathbf{W}^* at a rate of $O(\sqrt{d \log n/n})$ in the Frobenius norm.
- *Tensor initialization*: We adopt the tensor method proposed in [15], and show that it provably provides an initialization in the neighborhood of the ground truth both for FCN and CNN. In particular, we strengthened the guarantee of the tensor method by replacing the homogeneous assumption on activation functions in [15] by a mild condition on the curvature of activation functions around \mathbf{W}^* , which holds for a larger class of activation functions including *sigmoid* and *tanh*.

The cross-entropy loss is much more challenging to analyze than the squared loss, e.g., its gradient and Hessian take much more complicated forms compared with the squared loss; moreover, it is hard to control the values of gradient and Hessian due to the saturation phenomenon, i.e., when $H(\mathbf{W}, \mathbf{x})$ approaches 0 or 1. In order to establish the uniform local strong convexity property for the cross-entropy loss, we first show the *population* loss is smooth regarding to \mathbf{W}^* . Such a property

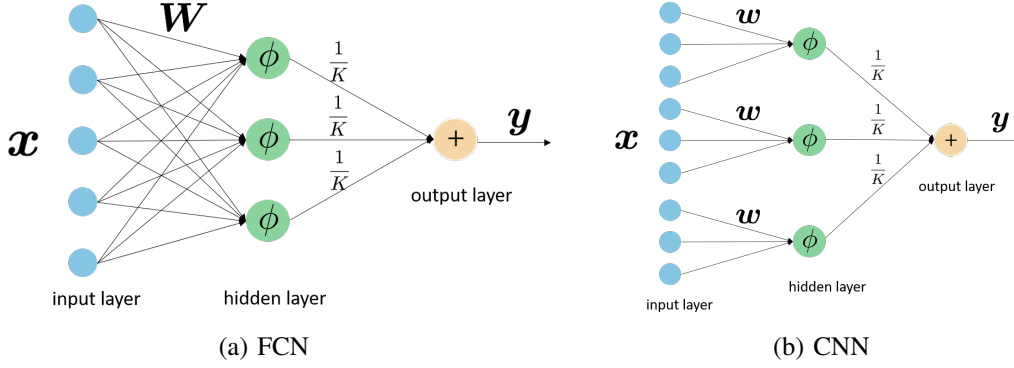


Fig. 1. Illustration of two types of one-hidden-layer neural networks considered in this paper: (a) a fully-connected network (FCN); (b) a non-overlapping convolutional neural network (CNN).

was also established in [15] for the squared loss. However, considering the special form of Hessian under the cross-entropy loss, we need to apply Taylor's approximation together with certain probabilistic upper bounds to control the value of Hessian, and obtain the smooth property. Network-specific quantities to capture the local geometry of the population loss at W^* for FCN and CNN are derived, which imply that the geometry of CNN is more benign than FCN, corroborated by the numerical experiments.

Beyond these two steps, the additional uniform concentration property of the Hessian (Lemma 3) is of key importance for us to obtain the uniform local strong convexity of the *empirical* loss. To show the uniform concentration of the Hessian, we successfully apply a type of covering argument. Different from the arguments in [17], which deal with the squared loss and are facilitated by certain nice assumptions on the activation functions, the cross-entropy loss is more difficult to apply the covering argument, e.g., both the gradient and Hessian no longer have a deterministic upper bound. Hence, we exploit the property of the sigmoid activation to show that the gradient and the Hessian of the cross-entropy loss are upper bounded with high probability in order to establish the uniform concentration property.

To the best of our knowledge, combining the analysis of gradient descent and initialization, this work provides the first globally convergent algorithm for the recovery of one-hidden-layer neural networks using the *cross-entropy* loss function.

C. Related Work

Due to the scope, we focus on the most relevant literature on theoretical and algorithmic aspects of learning shallow neural networks via nonconvex optimization. The parameter recovery viewpoint is relevant to the success of nonconvex learning in signal processing problems such as matrix completion, phase retrieval, blind deconvolution, dictionary learning and tensor decomposition [21]–[28], to name a few; see also the overview article [29]. The statistical model for data generation effectively removes worst-case instances and allows us to focus on average-case performance, which often possess much benign geometric properties that enable global convergence of simple local search algorithms.

The studies of one-hidden-layer network model can be further categorized into two classes, landscape analysis and model recovery. In the landscape analysis, it is known that if the network size is large enough compared to the data input, then there are no spurious local minima in the optimization landscape, and all local minima are global [30], [31], [32], [33]. For the case with multiple neurons ($2 \leq K \leq d$) in the under-parameterized setting, the work of Tian [34] studied the landscape of the population squared loss surface with ReLU activations. In particular, there exist spurious bad local minima in the optimization landscape [35], [36] even at the population level. Zhong et. al. [15] provided several important geometric characterizations for the regression problem using a variety of activation functions and the squared loss.

In the model recovery problem, the number of neurons is smaller than the input dimension, and all the existing works discussed below assumed the squared loss and (sub-)Gaussian inputs. In the case with a single neuron ($K = 1$), [14] showed that gradient descent converges linearly when the activation function is ReLU, with a zero initialization, as long as the sample complexity is $O(d)$ for the regression problem. When the activation function is quadratic, [37] shows that randomly initialized gradient descent converges fast to the global optimum at a near-optimal sample complexity. On the other hand, [17] showed that when $\phi(\cdot)$ has bounded first, second and third derivatives, there is no other critical points than the unique global minimum (within a constrained region of interest), and (projected) gradient descent converges linearly with an arbitrary initialization, as long as the sample complexity is $O(d \log^2 d)$ for the classification problem. Moreover, in the case with multiple neurons, [19] showed that projected gradient descent with a local initialization converges linearly for smooth activations with bounded second derivatives for the regression problem, [38] showed that gradient descent with tensor initialization converges linearly to a neighborhood of the ground truth using ReLU activations, and [39] showed the linear convergence of gradient descent with the spectral initialization using quadratic activations. For CNN with ReLU activations, [18] shows that gradient descent converges to the ground truth with random initialization for the population risk function based on the squared loss under Gaussian inputs. Moreover, [20] shows that gradient descent successfully learns a two-

layer convolutional neural network despite the existence of bad local minima. From a technical perspective, our study differs from all the aforementioned work in that the cross-entropy loss function we analyze has a very different form. Furthermore, we study the model recovery classification problem under the multi-neuron case, which has not been studied before.

Finally, we note that several papers study one-hidden-layer or two-layer neural networks with different structures under Gaussian input. For example, [40] studied the overlapping convolutional neural network, [16] studied a two-layer feedforward networks with ReLU activations and identity mapping, and [41] introduced the Porcupine Neural Network.

D. Paper Organization and Notations

The rest of the paper is organized as follows. Section II presents the main results on local geometry and local linear convergence of gradient descent. Section III discusses the initialization based on the tensor method. Numerical examples are demonstrated in Section IV, and finally, conclusions are drawn in Section V. Details of the technical proofs are delayed in the supplemental materials.

Throughout this paper, we use boldface letters to denote vectors and matrices, e.g. \mathbf{w} and \mathbf{W} . The transpose of \mathbf{W} is denoted by \mathbf{W}^\top , and $\|\mathbf{W}\|$, $\|\mathbf{W}\|_F$ denote the spectral norm and the Frobenius norm. For a positive semidefinite (PSD) matrix \mathbf{A} , we write $\mathbf{A} \succeq 0$. The identity matrix is denoted by \mathbf{I} . The gradient and the Hessian of a function $f(\mathbf{W})$ is denoted by $\nabla f(\mathbf{W})$ and $\nabla^2 f(\mathbf{W})$, respectively.

Denote $\|\cdot\|_{\psi_1}$ as the sub-exponential norm of a random variable. We use c, C, C_1, \dots to denote constants whose values may vary from place to place. For nonnegative functions $f(x)$ and $g(x)$, $f(x) = O(g(x))$ means there exist positive constants c and a such that $f(x) \leq cg(x)$ for all $x \geq a$; $f(x) = \Omega(g(x))$ means there exist positive constants c and a such that $f(x) \geq cg(x)$ for all $x \geq a$.

II. GRADIENT DESCENT AND ITS PERFORMANCE GUARANTEE

To estimate the network parameter \mathbf{W}^* , since (4) is a highly nonconvex function, vanilla gradient descent with an arbitrary initialization may get stuck at local minima. Therefore, we implement gradient descent (GD) with a well-designed initialization scheme that is described in details in Section III. In this section, we focus on the performance of the local update rule

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla f_n(\mathbf{W}_t),$$

where η is the constant step size. The algorithm is summarized in Algorithm 1.

Note that throughout the execution of GD, the same set of training samples is used which is the standard implementation of gradient descent. Consequently the analysis is challenging due to the statistical dependence of the iterates with the data.

A. Uniform local strong convexity

We first characterize the local strong convexity of $f_n(\cdot)$ in a neighborhood of the ground truth. We use the Euclidean ball

Algorithm 1 Gradient Descent (GD)

Input: Training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, step size η , iteration T
Initialization: $\mathbf{W}_0 \leftarrow \text{INITIALIZATION}(\{(\mathbf{x}_i, y_i)\}_{i=1}^n)$
Gradient Descent: for $t = 0, 1, \dots, T-1$

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla f_n(\mathbf{W}_t).$$

Output: \mathbf{W}_T

to denote the local neighborhood of \mathbf{W}^* for FCN or of \mathbf{w}^* for CNN.

$$\mathbb{B}(\mathbf{W}^*, r) = \{\mathbf{W} \in \mathbb{R}^{d \times K} : \|\mathbf{W} - \mathbf{W}^*\|_F \leq r\}, \quad (6a)$$

$$\mathbb{B}(\mathbf{w}^*, r) = \{\mathbf{w} \in \mathbb{R}^m : \|\mathbf{w} - \mathbf{w}^*\|_2 \leq r\}, \quad (6b)$$

where r is the radius of the ball. With slight abuse of notations, we will drop the subscript FCN or CNN for simplicity, whenever it is clear from the context that the result is for FCN when the argument is $\mathbf{W} \in \mathbb{R}^{d \times K}$ and for CNN when the argument is $\mathbf{w} \in \mathbb{R}^m$. Further, $\sigma_i(\mathbf{W})$ denotes the i -th largest singular value of \mathbf{W}^* . Let the condition number be $\kappa = \sigma_1/\sigma_K$, and $\lambda = \prod_{i=1}^K (\sigma_i/\sigma_K)$. Moreover, we introduce an important quantity $\rho(\sigma)$ regarding $\phi(z)$, the sigmoid activation function, that captures the geometric properties of the loss function for neural networks (1) and (2).

Definition 1 (Key quantity for FCN). *Let $z \sim \mathcal{N}(0, 1)$ and define $\alpha_q(\sigma) = \mathbb{E}[\phi'(\sigma \cdot z)z^q]$, $\forall q \in \{0, 1, 2\}$, and $\beta_q(\sigma) = \mathbb{E}[\phi'(\sigma \cdot z)^2 z^q]$, $\forall q \in \{0, 2\}$. Define $\rho_{\text{FCN}}(\sigma)$ as*

$$\rho_{\text{FCN}}(\sigma) = \min \{\beta_0(\sigma) - \alpha_0^2(\sigma), \beta_2(\sigma) - \alpha_2^2(\sigma)\} - \alpha_1^2(\sigma).$$

Definition 2 (Key quantity for CNN). *Let $z \sim \mathcal{N}(0, \sigma^2)$ and define $\rho_{\text{CNN}}(\sigma)$ as*

$$\rho_{\text{CNN}}(\sigma) = \min \left\{ \mathbb{E}[(\phi'(z)z)^2], \mathbb{E}[\phi'(z)^2] \right\}.$$

Note that Definition 1 for FCN is different from that in [15, Property 3.2] but consistent with [15, Lemma D.4] which removes the third term in [15, Property 3.2]. For the activation function considered in this paper, the first two terms suffice. Definition 2 for CNN is a newly distilled quantity in this paper tailored to the special structure of CNN.

The quantity $\rho(\sigma)$ plays an important role in the following theorem which guarantees the Hessian of the empirical risk function in the local neighborhood of the ground truth is positive definite with high probability for both FCN and CNN.

Theorem 1 (Local Strong Convexity). *Consider the classification model with FCN (1) or CNN (2) and the sigmoid activation function.*

- For FCN, assume $\|\mathbf{w}_k^*\|_2 \leq 1$ for all k . There exist constants c_1 and c_2 such that as soon as sample size

$$n_{\text{FCN}} \geq c_1 \cdot dK^5 \log^2 d \cdot \left(\frac{\kappa^2 \lambda}{\rho_{\text{FCN}}(\sigma_K)} \right)^2,$$

with probability at least $1 - d^{-10}$, we have for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$,

$$\Omega \left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \right) \cdot \mathbf{I} \preceq \nabla^2 f_n(\mathbf{W}) \preceq \Omega(1) \cdot \mathbf{I},$$

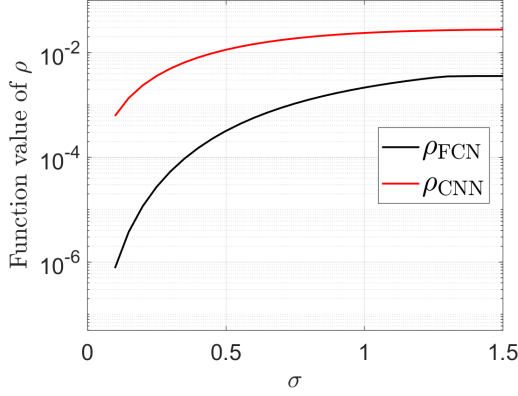


Fig. 2. Illustration $\rho(\sigma)$ for both FCN and CNN with the sigmoid activation.

where $r_{\text{FCN}} := \frac{c_2}{\sqrt{K}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}$.

- For CNN, assume $\|\mathbf{w}^*\|_2 \leq 1$. There exist constants c_3 and c_4 such that as soon as sample size

$$n_{\text{CNN}} \geq c_3 \cdot dK^5 \log^2 d \cdot \left(\frac{1}{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)} \right)^2,$$

with probability at least $1 - d^{-10}$, we have for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})$,

$$\Omega\left(\frac{1}{K} \cdot \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)\right) \cdot \mathbf{I} \preceq \nabla^2 f_n(\mathbf{w}) \preceq \Omega(K) \cdot \mathbf{I},$$

where $r_{\text{CNN}} := \frac{c_4}{K^2} \cdot \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)$.

We note that for FCN (1), all column permutations of \mathbf{W}^* are equivalent global minimum of the loss function, and Theorem 1 applies to all such permutation matrices of \mathbf{W}^* . The proof of Theorem 1 is outlined in Appendix B.

A pivot observation from the lower bound of the Hessian is that the sign of $\rho(\cdot)$ will determine whether the Hessian is positive definite or not, since K, κ, λ are all positive. We depict $\rho(\sigma)$ as a function of σ in a certain range for the sigmoid activation in Fig. 2. It can be seen from Fig. 2 that $\rho(\sigma)$ is monotonic increasing when σ increases, and we have $\rho(\sigma) > 0$ as long as $\sigma > 0$. When \mathbf{W}^* is orthogonal, κ and λ are both 1, $\rho(\sigma)$ is a constant, hence the lower bound of Hessian is on the order of $\frac{1}{K^2}$ for FCN. However, in the worst case where the columns of \mathbf{W}^* is linear dependent, then $\kappa, \lambda, \rho(\sigma)$ are infinite, and the local strong convexity doesn't hold for FCN case. Furthermore, the value of $\rho_{\text{CNN}}(\sigma)$ is much larger than $\rho_{\text{FCN}}(\sigma)$ for the same input.

Theorem 1 guarantees that for both FCN (1) and CNN (2) the Hessian of the empirical cross-entropy loss function $f_n(\mathbf{W})$ is positive definite in a neighborhood of the ground truth \mathbf{W}^* , as long as the sample size n is sufficiently large and the columns of \mathbf{W}^* are linearly independent. The bounds in Theorem 1 depend on the dimension parameters of the network (n and K), as well as the ground truth ($\rho_{\text{FCN}}(\sigma_K), \lambda, \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)$).

B. Performance Guarantees of GD

For the classification problem, due to the nature of quantized labels, \mathbf{W}^* is no longer a critical point of $f_n(\mathbf{W})$. By the

strong convexity of the empirical risk function $f_n(\mathbf{W})$ in the local neighborhood of \mathbf{W}^* , there can exist at most one critical point in $\mathbb{B}(\mathbf{W}^*, r)$, which is the unique local minimizer in $\mathbb{B}(\mathbf{W}^*, r)$ if it exists. The following theorem shows that there indeed exists such a critical point $\widehat{\mathbf{W}}_n$, which is provably close to the ground truth \mathbf{W}^* , and gradient descent converges linearly to $\widehat{\mathbf{W}}_n$.

Theorem 2 (Performance Guarantees of Gradient Descent). *Assume the assumptions in Theorem 1 hold. Under the event that local strong convexity holds,*

- for FCN, there exists a critical point in $\mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$ such that

$$\|\widehat{\mathbf{W}}_n - \mathbf{W}^*\|_{\text{F}} \leq c_1 \frac{K^{9/4} \kappa^2 \lambda}{\rho_{\text{FCN}}(\sigma_K)} \sqrt{\frac{d \log n}{n}},$$

and if the initial point $\mathbf{W}_0 \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$, GD converges linearly to $\widehat{\mathbf{W}}_n$, i.e.

$$\|\mathbf{W}_t - \widehat{\mathbf{W}}_n\|_{\text{F}} \leq \left(1 - \frac{c_2 \eta \rho_{\text{FCN}}(\sigma_K)}{K^2 \kappa^2 \lambda}\right)^t \|\mathbf{W}_0 - \widehat{\mathbf{W}}_n\|_{\text{F}},$$

for $\eta \leq c_3$, where c_1, c_2, c_3 are constants;

- for CNN, there exists a critical point in $\mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})$ such that

$$\|\widehat{\mathbf{w}}_n - \mathbf{w}^*\|_2 \leq c_4 \frac{K}{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)} \cdot \sqrt{\frac{d \log n}{n}},$$

and if the initial point $\mathbf{w}_0 \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})$, GD converges linearly to $\widehat{\mathbf{w}}_n$, i.e.

$$\|\mathbf{w}_t - \widehat{\mathbf{w}}_n\|_2 \leq \left(1 - \frac{c_5 \eta \rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K}\right)^t \|\mathbf{w}_0 - \widehat{\mathbf{w}}_n\|_2,$$

for $\eta \leq c_6/K$, where c_4, c_5, c_6 are constants.

Similarly to Theorem 1, for FCN (1) Theorem 2 also holds for all column permutations of \mathbf{W}^* . The proof can be found in Appendix C. Theorem 2 guarantees that the existence of critical points in the local neighborhood of the ground truth, which GD converges to, and also shows that the critical points converge to the ground truth \mathbf{W}^* at the rate of $O(K^{9/4} \sqrt{d \log n}/n)$ for FCN (1) and $O(K \sqrt{d \log n}/n)$ for CNN(2) with respect to increasing the sample size n . Therefore, \mathbf{W}^* can be recovered consistently as n goes to infinity. Moreover, for both FCN (1) and CNN (2) gradient descent converges linearly to $\widehat{\mathbf{W}}_n$ (or resp. $\widehat{\mathbf{w}}_n$) at a linear rate, as long as it is initialized in the basin of attraction. To achieve ϵ -accuracy, i.e. $\|\mathbf{W}_t - \widehat{\mathbf{W}}_n\|_{\text{F}} \leq \epsilon$ (or resp. $\|\mathbf{w}_t - \widehat{\mathbf{w}}_n\|_2 \leq \epsilon$), it requires a computational complexity of $O(ndK^4 \log(1/\epsilon))$ (or resp. $O(ndK^2 \log(1/\epsilon))$), which is linear in n, d and $\log(1/\epsilon)$.

III. INITIALIZATION VIA TENSOR METHOD

Our initialization adopts the tensor method proposed in [15]. The initialization method works for the FCN model directly, and works for the CNN model with slight modification as presented in [42]. To avoid unnecessary repetitions from the previous work, we focus on the FCN case to outline the algorithm and remark the difference. We recommend the readers refer to [15], [42] for more details.

A. Preliminary and Algorithm

We start with introducing the necessary definitions which can be found in [15]. We first define a product $\tilde{\otimes}$ as follows. If $\mathbf{v} \in \mathbb{R}^d$ is a vector and \mathbf{I} is the identity matrix, then $\mathbf{v} \tilde{\otimes} \mathbf{I} = \sum_{j=1}^d [\mathbf{v} \otimes \mathbf{e}_j \otimes \mathbf{e}_j + \mathbf{e}_j \otimes \mathbf{v} \otimes \mathbf{e}_j + \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}]$. If \mathbf{M} is a symmetric rank- r matrix factorized as $\mathbf{M} = \sum_{i=1}^r \mathbf{s}_i \mathbf{v}_i \mathbf{v}_i^\top$ and \mathbf{I} is the identity matrix, then

$$\mathbf{M} \tilde{\otimes} \mathbf{I} = \sum_{i=1}^r \mathbf{s}_i \sum_{j=1}^d \sum_{l=1}^6 \mathbf{A}_{l,i,j}, \quad (7)$$

where $\mathbf{A}_{1,i,j} = \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_j$, $\mathbf{A}_{2,i,j} = \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{e}_j$, $\mathbf{A}_{3,i,j} = \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{e}_j$, $\mathbf{A}_{4,i,j} = \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}_i$, $\mathbf{A}_{5,i,j} = \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{e}_j \otimes \mathbf{v}_i$ and $\mathbf{A}_{6,i,j} = \mathbf{e}_j \otimes \mathbf{e}_j \otimes \mathbf{v}_i \otimes \mathbf{v}_i$. This allows us to introduce the following quantities.

Definition 3. Define M_1, M_2, M_3, M_4 and $m_{1,i}, m_{2,i}, m_{3,i}, m_{4,i}$ as follows:

$$\begin{aligned} M_1 &= \mathbb{E}[y \cdot \mathbf{x}], \\ M_2 &= \mathbb{E}[y \cdot (\mathbf{x} \otimes \mathbf{x} - \mathbf{I})], \\ M_3 &= \mathbb{E}[y \cdot (\mathbf{x}^{\otimes 3} - \mathbf{x} \tilde{\otimes} \mathbf{I})], \\ M_4 &= \mathbb{E}[y \cdot (\mathbf{x}^{\otimes 4} - (\mathbf{x} \otimes \mathbf{x}) \tilde{\otimes} \mathbf{I} + \mathbf{I} \tilde{\otimes} \mathbf{I})], \end{aligned}$$

$m_{l,i} = g_{l,i}(\|\mathbf{w}_i^*\|)$, $\forall l = 0, 1, 2, 3, 4$, where $g_{1,i}(\sigma) = \gamma_1(\sigma)$, $g_{2,i}(\sigma) = \gamma_2(\sigma) - \gamma_0(\sigma)$, $g_{3,i}(\sigma) = \gamma_3(\sigma) - 3\gamma_1(\sigma)$, $g_{4,i}(\sigma) = \gamma_4(\sigma) + 3\gamma_0(\sigma) - 6\gamma_2(\sigma)$, and $\gamma_j(\sigma) = \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\phi(\sigma \cdot z) z^j]$, $\forall j = 0, 1, 2, 3, 4$.

We further define a tensor operation as follows. For a tensor $\mathbf{T} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ and three matrices $\mathbf{A} \in \mathbb{R}^{n_1 \times d_1}$, $\mathbf{B} \in \mathbb{R}^{n_2 \times d_2}$, $\mathbf{C} \in \mathbb{R}^{n_3 \times d_3}$, the (i, j, k) -th entry of the tensor $\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C})$ is given by

$$\sum_{i'}^{n_1} \sum_{j'}^{n_2} \sum_{k'}^{n_3} \mathbf{T}_{i',j',k'} \mathbf{A}_{i',i} \mathbf{B}_{j',j} \mathbf{C}_{k',k}. \quad (8)$$

Armed this with definition, we define the following useful quantities.

Definition 4. Let $\alpha \in \mathbb{R}^d$ denote a randomly picked vector. We define \mathbf{P}_2 and \mathbf{P}_3 as follows: $\mathbf{P}_2 = \mathbf{M}_{j_2}(\mathbf{I}, \mathbf{I}, \alpha, \dots, \alpha)$, where $j_2 = \min\{j \geq 2 | M_j \neq 0\}$, and $\mathbf{P}_3 = \mathbf{M}_{j_3}(\mathbf{I}, \mathbf{I}, \mathbf{I}, \alpha, \dots, \alpha)$, where $j_3 = \min\{j \geq 3 | M_j \neq 0\}$.

We further denote $\bar{\mathbf{w}} = \mathbf{w}/\|\mathbf{w}\|$. An important implication of Definition 3 and 4 is that the non-zero matrix \mathbf{P}_2 and non-zero tensor \mathbf{P}_3 is in the form of $\sum_{i=1}^K m_{j_2,i} (\alpha^\top \bar{\mathbf{w}}_i^*)^{j_2-2} \bar{\mathbf{w}}_i^{*\otimes 2}$, $\sum_{i=1}^K m_{j_3,i} (\alpha^\top \bar{\mathbf{w}}_i^*)^{j_3-3} \bar{\mathbf{w}}_i^{*\otimes 3}$, see [15, Claim 5.5]. The basic strategy is to extract the direction, magnitude information from the empirical version of \mathbf{P}_2 and \mathbf{P}_3 . Hence estimating \mathbf{W}^* can be decomposed as the following two steps.

- Step 1 Estimate the direction of each column of \mathbf{W}^* by decomposing \mathbf{P}_2 to approximate the subspace spanned by $\{\bar{\mathbf{w}}_1^*, \bar{\mathbf{w}}_2^*, \dots, \bar{\mathbf{w}}_K^*\}$ (denoted by \mathbf{V}), then reduce the third-order tensor \mathbf{P}_3 to a lower-dimension tensor $\mathbf{R}_3 = \mathbf{P}_3(\mathbf{V}, \mathbf{V}, \mathbf{V}) \in \mathbb{R}^{K \times K \times K}$, and apply non-orthogonal tensor decomposition on \mathbf{R}_3 to output the estimate $s_i \mathbf{V}^\top \bar{\mathbf{w}}_i^*$, where $s_i \in \{1, -1\}$ is a random sign.
- Step 2 Approximate the magnitude of \mathbf{w}_i^* and the sign s_i by solving a linear system of equations.

The initialization algorithm based on the tensor method is outlined in Algorithm 2. For more implementation details about Algorithm 2, e.g., power method, we refer to [15].

Algorithm 2 Initialization via Tensor Method

Input: Partition n pairs of data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ into three subsets $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3$.

Output:

- 1: Estimate $\hat{\mathbf{P}}_2$ of \mathbf{P}_2 from data set \mathcal{D}_1 .
 - 2: $\mathbf{V} \leftarrow \text{POWERMETHOD}(\hat{\mathbf{P}}_2, K)$.
 - 3: Estimate $\hat{\mathbf{R}}_3$ of $\mathbf{P}_3(\mathbf{V}, \mathbf{V}, \mathbf{V})$ from data set \mathcal{D}_2 .
 - 4: $\{\hat{\mathbf{u}}_i\}_{i \in [K]} \leftarrow \text{KCL}(\hat{\mathbf{R}}_3)$.
 - 5: $\{\mathbf{w}_i^{(0)}\}_{i \in [K]} \leftarrow \text{RECMAG}(\mathbf{V}, \{\hat{\mathbf{u}}_i\}_{i \in [K]}, \mathcal{D}_3)$.
-

B. Performance Guarantee of Initialization

For the classification problem, we make the following technical assumptions, similarly to [15, Assumption 5.3] for the regression problem.

Assumption 1. The activation function $\phi(z)$ satisfies the following conditions:

- 1) If $M_j \neq 0$, then

$$\sum_{i=1}^K m_{j,i} (\mathbf{w}_i^{*\top} \alpha)^{j-2} \bar{\mathbf{w}}_i^* \bar{\mathbf{w}}_i^{*\top} \neq \mathbf{0},$$

$$\sum_{i=1}^K m_{j,i} (\bar{\mathbf{w}}_i^{*\top} \alpha)^{j-3} (\mathbf{V}^\top \bar{\mathbf{w}}_i^*) \text{vec}((\mathbf{V}^\top \bar{\mathbf{w}}_i^*)(\mathbf{V}^\top \bar{\mathbf{w}}_i^*)^\top)^\top \neq \mathbf{0},$$

for $j \geq 3$.

- 2) At least one of M_3 and M_4 is non-zero.

Assumption 1 is to guarantee that the key terms still contain the magnitude information about \mathbf{w}_j^* . It can be verified that for sigmoid activation $m_{3,i}$ is non-zero for $\sigma > 0$, hence it will satisfy Assumption 1. Furthermore, we do not require the homogeneous assumption (i.e., $\phi(az) = a^p z$ for an integer p) required in [15], which can be restrictive. Instead, we assume the following condition on the curvature of the activation function around the ground truth, which holds for a larger class of activation functions such as sigmoid and tanh.

Assumption 2. Let l_1 be the index of the first nonzero M_i where $i = 1, \dots, 4$. For the activation function $\phi(\cdot)$, there exists a positive constant δ such that $g_{l_1,i}(\cdot)$ is strictly monotone over the interval $(\|\mathbf{w}_i^*\| - \delta, \|\mathbf{w}_i^*\| + \delta)$, and the derivative of $g_{l_1,i}(\cdot)$ is lower bounded by some constant for all i .

It can be numerically verified that sigmoid activation will also satisfy Assumption 2. We next present the performance guarantee for the initialization algorithm in the following theorem.

Theorem 3. For the classification model (1), under Assumptions 1 and 2, for any $0 < \epsilon < 1$ and $\zeta > 1$, if the sample size $n \geq d \cdot \text{poly}(K, \kappa, \zeta, \log d, 1/\epsilon)$, then the output $\mathbf{W}_0 \in \mathbb{R}^{d \times K}$ of Algorithm 2 satisfies

$$\|\mathbf{W}_0 - \mathbf{W}^*\|_F \leq \epsilon \text{poly}(K, \kappa) \|\mathbf{W}^*\|_F, \quad (9)$$

with probability at least $1 - d^{-\Omega(\zeta)}$.

The proof of Theorem 3 consists of (a) showing the estimation of the direction of \mathbf{W}^* is sufficiently accurate and (b) showing the approximation of the norm of \mathbf{W}^* is accurate enough. The proof of part (a) is the same as that in [15], but our argument in part (b) is different, where we relax the homogeneous assumption on activation functions. More details can be found in the supplementary materials in Appendix E.

IV. NUMERICAL EXPERIMENTS

For FCN, we first implement gradient descent to verify that the empirical risk function is strongly convex in the local region around \mathbf{W}^* . If we initialize multiple times in such a local region, it is expected that gradient descent converges to the same critical point $\widehat{\mathbf{W}}_n$, with the same set of training samples. Given a set of training samples, we randomly initialize multiple times, and then calculate the variance of the output of gradient descent. Denote the output of the ℓ th run as $\widehat{\mathbf{w}}_n^{(\ell)} = \text{vec}(\widehat{\mathbf{W}}_n^{(\ell)})$ and the mean of the runs as $\bar{\mathbf{w}}$. The error is calculated as $\text{SD}_n = \sqrt{\frac{1}{L} \sum_{\ell=1}^L \|\widehat{\mathbf{w}}_n^{(\ell)} - \bar{\mathbf{w}}\|^2}$, where $L = 20$ is the total number of random initializations. Adopted from [17], it quantifies the standard deviation of the estimator $\widehat{\mathbf{W}}_n$ under different initializations with the same set of training samples. We say an experiment is successful, if $\text{SD}_n \leq 10^{-4}$. We generate the ground truth \mathbf{W}^* from Gaussian matrices, and the training samples are generated using the FCN (1). Figure 3 (a) shows the successful rate of gradient descent by averaging over 50 sets of training samples for each pair of n and d , where $K = 3$ and $d = 15, 20, 25$ respectively. The maximum iterations for gradient descent is set as $\text{iter}_{\max} = 3500$. It can be seen that as long as the sample complexity is large enough, gradient descent converges to the same local minima with high probability.

We next show that the statistical accuracy of the local minimizer for gradient descent if it is initialized close enough to the ground truth. Suppose we initialize around the ground truth such that $\|\mathbf{W}_0 - \mathbf{W}^*\|_F \leq 0.1 \cdot \|\mathbf{W}^*\|_F$. We calculate the average estimation error as $\sum_{\ell=1}^L \|\widehat{\mathbf{W}}_n^{(\ell)} - \mathbf{W}^*\|_F^2 / (L \|\mathbf{W}^*\|_F^2)$ over $L = 100$ Monte Carlo simulations with random initializations. Fig. 3 (b) shows the average estimation error with respect to the sample complexity when $K = 3$ and $d = 20, 35, 50$ respectively. It can be seen that the estimation error decreases gracefully as we increase the sample size and matches with the theoretical prediction of error rates reasonably well.

Similarly, for CNN, we first verify that the empirical risk function is locally strongly convex using the same method as before. We generate the entries of true weights \mathbf{w}^* from standard Gaussian distribution, and generate the training samples using the CNN model (2). In Fig. 4 (a), we say an experiment is successful if $\text{SD}_n \leq 10^{-14}$, and the successful rate is calculated over 100 sets of training samples with $K = 3$ and $d = 15, 24, 30$ respectively. Then we verify the performance of gradient descent in Fig. 4 (b). Suppose we initialized in the neighborhood of \mathbf{w}^* , i.e., $\|\mathbf{w}_0 - \mathbf{w}^*\|_2 \leq 0.9 \cdot \|\mathbf{w}^*\|_2$, for fixed d, K, n , the average error is calculated over $L = 100$ Monte Carlo simulations. It can be seen that the error decreases as we increase the number of samples.

V. CONCLUSIONS

In this paper, we have studied the model recovery problem of a one-hidden-layer neural network using the cross-entropy loss in a multi-neuron classification problem. In particular, we have characterized the sample complexity to guarantee local strong convexity in a neighborhood (whose size we have characterized as well) of the ground truth when the training data are generated from a classification model for two types of neural network models: fully-connected network and non-overlapping convolutional network. This guarantees that with high probability, gradient descent converges linearly to the ground truth if initialized properly. In the future, it will be interesting to extend the analysis in this paper to more general class of activation functions, particularly ReLU-like activations.

APPENDIX A

GRADIENT AND HESSIAN OF POPULATION LOSS

For the convenience of analysis, we first provide the gradient and the Hessian formula for the cross-entropy loss using FCN and CNN here.

A. The FCN case

Consider the population loss function $f(\mathbf{W}) = \mathbb{E}[f_n(\mathbf{W})] = \mathbb{E}[\ell(\mathbf{W}; \mathbf{x})]$, where $\ell(\mathbf{W}; \mathbf{x})$ is associated with network $H_{\text{FCN}}(\mathbf{W}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}_k^\top \mathbf{x})$. Hiding the dependence on \mathbf{x} for notational simplicity, we can calculate the gradient and the Hessian as

$$\mathbb{E} \left[\frac{\partial \ell(\mathbf{W})}{\partial \mathbf{w}_j} \right] = \mathbb{E} \left[-\frac{1}{K} \frac{(y - H(\mathbf{W}))}{H(\mathbf{W})(1 - H(\mathbf{W}))} \phi'(\mathbf{w}_j^\top \mathbf{x}) \mathbf{x} \right], \quad (10)$$

$$\mathbb{E} \left[\frac{\nabla^2 \ell(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_l} \right] = \mathbb{E} [\xi_{j,l}(\mathbf{W}) \cdot \mathbf{x} \mathbf{x}^\top], \quad (11)$$

for $1 \leq j, l \leq K$. Here, when $j \neq l$,

$$\xi_{j,l}(\mathbf{W}) = \frac{1}{K^2} \phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x}) \cdot \frac{H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W})}{H^2(\mathbf{W})(1 - H(\mathbf{W}))^2},$$

and when $j = l$,

$$\begin{aligned} \xi_{j,j}(\mathbf{W}) &= \frac{1}{K^2} \phi'(\mathbf{w}_j^\top \mathbf{x})^2 \cdot \frac{H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W})}{H^2(\mathbf{W})(1 - H(\mathbf{W}))^2} \\ &\quad - \frac{1}{K} \phi''(\mathbf{w}_j^\top \mathbf{x}) \cdot \frac{y - H(\mathbf{W})}{H(\mathbf{W})(1 - H(\mathbf{W}))}. \end{aligned}$$

B. The CNN case

For the CNN case, i.e., $H(\mathbf{w}) := H_{\text{CNN}}(\mathbf{w}, \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \phi(\mathbf{w}^\top \mathbf{x}^{(k)})$, the corresponding gradient and Hessian of the population loss function $\ell(\mathbf{w})$ is given by

$$\mathbb{E} \left[\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}} \right] = \mathbb{E} \left[-\phi'(\mathbf{w}^\top \mathbf{x}^{(1)}) \cdot \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \cdot \mathbf{x}^{(1)} \right], \quad (12)$$

$$\mathbb{E} \left[\frac{\nabla^2 \ell(\mathbf{w})}{\partial \mathbf{w}^2} \right] = \mathbb{E} \left[\sum_{j=1}^K \sum_{l=1}^K g_{j,l}(\mathbf{w}) \mathbf{x}^{(j)} \mathbf{x}^{(l)\top} \right], \quad (13)$$

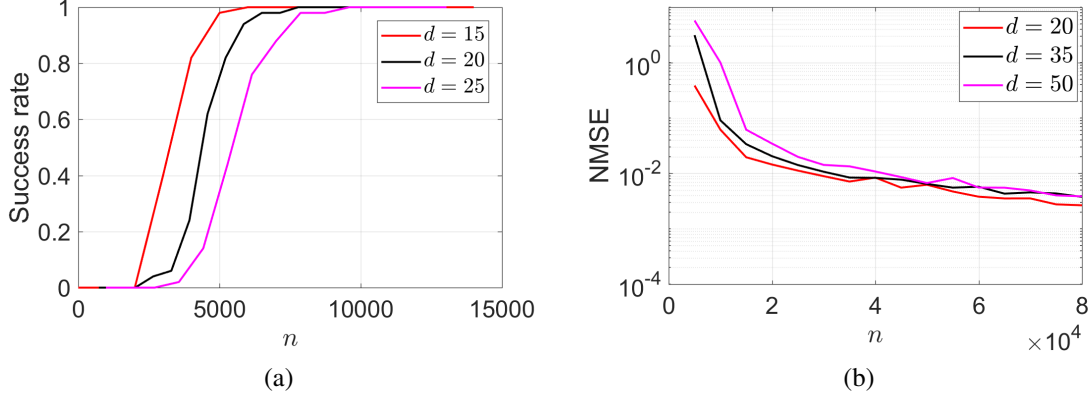


Fig. 3. For FCN (1) fix $K = 3$. (a) Success rate of converging to the same local minima with respect to the sample complexity for various d with threshold 10^{-4} ; (b) Average estimation error of gradient descent in a local neighborhood of the ground truth with respect to the sample complexity for various d . The x-axis is scaled to illuminate the correct scaling between n and d .

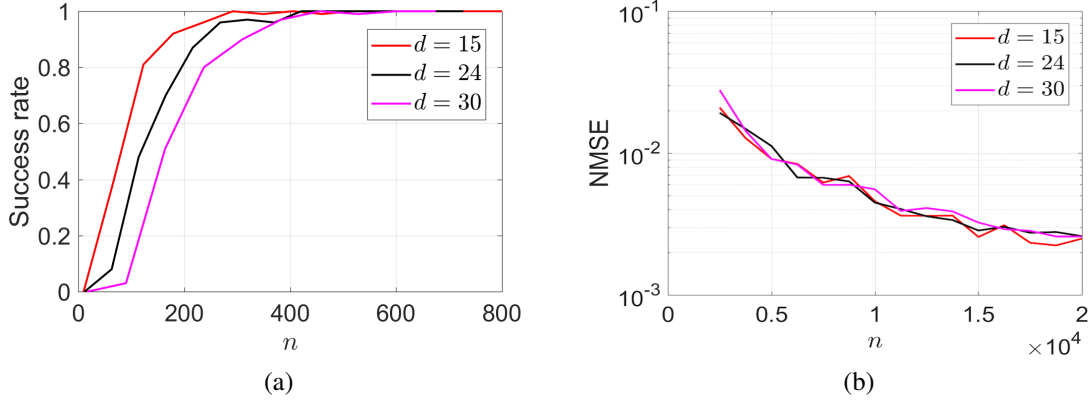


Fig. 4. For CNN (2), fix $K = 3$. (a) Success rate of converging to the same local minima with respect to the sample complexity for various d with threshold 10^{-14} ; (b) Average estimation error of gradient descent in a local neighborhood of the ground truth with respect to the sample complexity for various d . The x-axis is scaled to illuminate the correct scaling between n and d .

where when $j \neq l$,

$$g_{j,l}(\mathbf{w}) = \frac{1}{K^2} \cdot \frac{H(\mathbf{w})^2 + y - 2y \cdot H(\mathbf{w})}{(H(\mathbf{w})(1 - H(\mathbf{w})))^2} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)}),$$

and when $j = l$,

$$g_{j,j}(\mathbf{w}) = \frac{1}{K^2} \cdot \frac{H(\mathbf{w})^2 + y - 2y \cdot H(\mathbf{w})}{(H(\mathbf{w})(1 - H(\mathbf{w})))^2} \cdot \phi'(\mathbf{w}^\top \mathbf{x}^{(j)})^2 \\ - \frac{1}{K} \cdot \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \cdot \phi''(\mathbf{w}^\top \mathbf{x}^{(j)}).$$

APPENDIX B PROOF OF THEOREM 1

In order to show that the empirical loss possesses a local strong convexity, we follow the following steps:

- 1) We first show that the Hessian $\nabla^2 f(\mathbf{W})$ of the population loss function is smooth with respect to $\nabla^2 f(\mathbf{W}^*)$ (Lemma 1);
- 2) We then show that $\nabla^2 f(\mathbf{W})$ satisfies local strong convexity and smoothness in a neighborhood of \mathbf{W}^* with appropriately chosen radius, $\mathbb{B}(\mathbf{W}^*, r)$, by leveraging similar properties of $\nabla^2 f(\mathbf{W}^*)$ (Lemma 2);

3) Next, we show that the Hessian of the empirical loss function $\nabla^2 f_n(\mathbf{W})$ is close to its population counterpart $\nabla^2 f(\mathbf{W})$ uniformly in $\mathbb{B}(\mathbf{W}^*, r)$ with high probability (Lemma 3).

4) Finally, putting all the arguments together, we establish $\nabla^2 f_n(\mathbf{W})$ satisfies local strong convexity and smoothness in $\mathbb{B}(\mathbf{W}^*, r)$.

To begin, we first show that the Hessian of the population risk is smooth enough around \mathbf{W}^* in the following lemmas.

Lemma 1 (Hessian Smoothness of Population Loss). *Suppose the loss $\ell(\cdot)$ associates with FCN (1), and assume $\|\mathbf{w}_k^*\|_2 \leq 1$ for all k and $\|\mathbf{W} - \mathbf{W}^*\|_F \leq 0.7$. Then we have*

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \leq \frac{C_1}{K^{\frac{3}{2}}} \cdot \|\mathbf{W} - \mathbf{W}^*\|_F, \quad (14)$$

holds. Similarly, suppose the loss $\ell(\cdot)$ associates with CNN (2), and assume $\|\mathbf{w}^*\|_2 \leq 1$ and $\|\mathbf{w} - \mathbf{w}^*\|_2 \leq 0.7$. We have

$$\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| \leq C_2 \cdot K \cdot \|\mathbf{w} - \mathbf{w}^*\|_2, \quad (15)$$

holds. Here C_1 and C_2 denote some large constants.

The proof is provided in Appendix D-A. Together with the fact that $\nabla^2 f(\mathbf{W}^*)$ be lower and upper bounded, Lemma 1

allows us to bound $\nabla^2 f(\mathbf{W})$ in a neighborhood around ground truth, given below.

Lemma 2 (Local Strong Convexity and Smoothness of Population Loss). *If the loss $\ell(\cdot)$ associates with FCN (1), there exists some constant C_1 , such that*

$$\frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{W}) \preceq C_1 \cdot \mathbf{I},$$

holds for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$ with $r_{\text{FCN}} := \frac{C_2}{K^{\frac{1}{2}}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}$. Moreover, if loss $\ell(\cdot)$ associates with CNN (2), then we have

$$C_3 \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{w}) \preceq C_4 \cdot K \cdot \mathbf{I}, \quad (16)$$

holds for all $\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})$ with $r_{\text{CNN}} := C_5 \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K^2}$.

The proof is provided in Appendix D-B. The next step is to show the Hessian of the empirical loss function is close to the Hessian of the population loss function in a uniform sense, which can be summarized as follows.

Lemma 3. *If the loss $\ell(\cdot)$ associates with FCN (1), then there exists a constant C such that as long as $n \geq C \cdot dK \log dK$, with probability at least $1 - d^{-10}$, the following holds*

$$\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})} \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| \leq C \sqrt{\frac{dK \log n}{n}}, \quad (17)$$

where $r_{\text{FCN}} := \frac{C}{K^{\frac{1}{2}}} \cdot \frac{\rho(\sigma_K)}{\kappa^2 \lambda}$. And if the loss $\ell(\cdot)$ associates with CNN (2), then we have

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})} \|\nabla^2 f_n(\mathbf{w}) - \nabla^2 f(\mathbf{w})\| \leq CK^2 \sqrt{\frac{\frac{d}{K} \cdot \log(n)}{n}}, \quad (18)$$

holds with probability at least $1 - d^{-10}$, as long as $n \geq \frac{d}{K} \log\left(\frac{d}{K}\right)$, and $r_{\text{CNN}} := C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K^2}$.

The proof is provided in Appendix D-C. Combining the above results will give us the result. Next we assume that the loss $\ell(\cdot)$ associates with FCN, and take it as an example in the proof. Then if the loss $\ell(\cdot)$ associates with CNN, the proof follows in the same manner.

Proof of Theorem 1. With probability at least $1 - d^{-10}$,

$$\begin{aligned} & \nabla^2 f_n(\mathbf{W}) \\ & \succeq \nabla^2 f(\mathbf{W}) - \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| \cdot \mathbf{I} \\ & \succeq \Omega\left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}\right) \cdot \mathbf{I} - \Omega\left(C \cdot \sqrt{\frac{dK \log n}{n}}\right) \cdot \mathbf{I}. \end{aligned}$$

As long as the sample size n is set to satisfy

$$C \cdot \sqrt{\frac{dK \log n}{n}} \leq \frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda},$$

i.e. $n \geq C \cdot dK^5 \log^2 d \cdot \left(\frac{\kappa^2 \lambda}{\rho_{\text{FCN}}(\sigma_K)}\right)^2$, we have

$$\nabla^2 f_n(\mathbf{W}) \succeq \Omega\left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}\right) \cdot \mathbf{I}.$$

holds for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$. Similarly, we have

$$\nabla^2 f_n(\mathbf{W}) \preceq C \cdot \mathbf{I}$$

holds for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$. \square

APPENDIX C PROOF OF THEOREM 2

We have established that $f_n(\mathbf{W})$ is strongly convex in $\mathbb{B}(\mathbf{W}^*, r)$ in Theorem 1. Thus there exists at most one critical point in $\mathbb{B}(\mathbf{W}^*, r)$. The proof of Theorem 2 follows the steps below:

- 1) We first show that the gradient $\nabla f_n(\mathbf{W})$ concentrates around $\nabla f(\mathbf{W})$ in $\mathbb{B}(\mathbf{W}^*, r)$ (Lemma 4), and then invoke [17, Theorem 2] to guarantee that there indeed exists a critical point $\widehat{\mathbf{W}}_n$ in $\mathbb{B}(\mathbf{W}^*, r)$;
- 2) We next show that $\widehat{\mathbf{W}}_n$ is close to \mathbf{W}^* and gradient descent converges linearly to $\widehat{\mathbf{W}}_n$ with a properly chosen step size.

To begin, the following lemma establishes that $\nabla f_n(\mathbf{W})$ uniformly concentrates around $\nabla f(\mathbf{W})$.

Lemma 4. *If the loss $\ell(\cdot)$ associates with FCN (1) with $r_{\text{FCN}} := \frac{C}{K^{\frac{1}{2}}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}$, and $\|\mathbf{w}_k^*\|_2 \leq 1$ for all k , then*

$$\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})} \|\nabla f_n(\mathbf{W}) - \nabla f(\mathbf{W})\| \leq C \sqrt{\frac{d \sqrt{K} \log n}{n}}$$

holds with probability at least $1 - d^{-10}$, as long as $n \geq CdK \log(dK)$. If the loss $\ell(\cdot)$ associates with CNN (2), with $r_{\text{CNN}} := C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^*\|_2)}{K^2}$ and $\|\mathbf{w}^*\|_2 \leq 1$, then

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})} \|\nabla f_n(\mathbf{w}) - \nabla f(\mathbf{w})\| \leq C \cdot \sqrt{\frac{d \log n}{n}} \quad (19)$$

holds with probability at least $1 - d^{-10}$ as long as $n \geq C \frac{d}{K} \log\left(\frac{d}{K}\right)$.

The proof is provided in Appendix D-D. Notice that for the population risk function $f(\mathbf{W})$, \mathbf{W}^* is the unique critical point in $\mathbb{B}(\mathbf{W}^*, r)$ due to local strong convexity. With Lemma 3 and Lemma 4, we can invoke [17, Theorem 2], which guarantees the following.

Corollary 1. *If the loss $\ell(\cdot)$ associates with FCN or CNN, there exists one and only one critical point $\widehat{\mathbf{W}}_n \in \mathbb{B}(\mathbf{W}^*, r)$ that satisfies $\nabla f_n(\widehat{\mathbf{W}}_n) = \mathbf{0}$ correspondingly.*

Again, since the proof for the case with the loss $\ell(\cdot)$ associating with FCN is the same as that for CNN, we next take FCN as an example.

We first show that $\widehat{\mathbf{W}}_n$ is close to \mathbf{W}^* . By the mean value theorem, there exists \mathbf{W}' on the straight line connecting \mathbf{W}^* and $\widehat{\mathbf{W}}_n$ such that

$$\begin{aligned} f_n(\widehat{\mathbf{W}}_n) &= f_n(\mathbf{W}^*) + \langle \nabla f_n(\mathbf{W}^*), \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*) \rangle \\ &\quad + \frac{1}{2} \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*)^\top \nabla^2 f_n(\mathbf{W}') \text{vec}(\widehat{\mathbf{W}}_n - \mathbf{W}^*) \\ &\leq f_n(\mathbf{W}^*), \end{aligned} \quad (20)$$

where the last inequality follows from the optimality of $\widehat{\mathbf{W}}_n$. By Theorem 1, we have

$$\begin{aligned} & \frac{1}{2} \text{vec} \left(\widehat{\mathbf{W}}_n - \mathbf{W}^* \right)^\top \nabla^2 f_n(\mathbf{W}') \text{vec} \left(\widehat{\mathbf{W}}_n - \mathbf{W}^* \right) \\ & \geq \Omega \left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \right) \left\| \widehat{\mathbf{W}}_n - \mathbf{W}^* \right\|_{\text{F}}^2. \end{aligned} \quad (21)$$

On the other hand, by the Cauchy-Schwarz inequality, we have

$$\begin{aligned} & \left| \left\langle \nabla f_n(\mathbf{W}^*), \text{vec} \left(\widehat{\mathbf{W}}_n - \mathbf{W}^* \right) \right\rangle \right| \\ & \leq \left\| \nabla f_n(\mathbf{W}^*) \right\|_2 \left\| \widehat{\mathbf{W}}_n - \mathbf{W}^* \right\|_{\text{F}} \\ & \leq \Omega \left(\sqrt{\frac{dK^{1/2} \log n}{n}} \right) \left\| \widehat{\mathbf{W}}_n - \mathbf{W}^* \right\|_{\text{F}}, \end{aligned} \quad (22)$$

where the last line follows from Lemma 4. Plugging (21) and (22) into (20), we have

$$\left\| \widehat{\mathbf{W}}_n - \mathbf{W}^* \right\|_{\text{F}} \leq \Omega \left(\frac{K^{\frac{9}{4}} \kappa^2 \lambda}{\rho_{\text{FCN}}(\sigma_K)} \sqrt{\frac{d \log n}{n}} \right). \quad (23)$$

Now we have established that there indeed exists a critical point in $\mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$. We can then establish the local linear convergence of gradient descent as below. Let \mathbf{W}_t be the estimate at the t -th iteration. Due to the update rule, we have

$$\begin{aligned} \mathbf{W}_{t+1} - \widehat{\mathbf{W}}_n &= \mathbf{W}_t - \eta \nabla f_n(\mathbf{W}_t) - \left(\widehat{\mathbf{W}}_n - \eta \nabla f_n(\widehat{\mathbf{W}}_n) \right) \\ &= \left(\mathbf{I} - \eta \int_0^1 \nabla^2 f_n(\mathbf{W}(\gamma)) \right) (\mathbf{W}_t - \widehat{\mathbf{W}}_n), \end{aligned}$$

where $\mathbf{W}(\gamma) = \widehat{\mathbf{W}}_n + \gamma(\mathbf{W}_t - \widehat{\mathbf{W}}_n)$ for $\gamma \in [0, 1]$. If $\mathbf{W}_t \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$, it is obvious that $\mathbf{W}(\gamma) \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})$, and by Theorem 1, we have

$$H_{\min} \cdot \mathbf{I} \preceq \nabla^2 f_n(\mathbf{W}(\gamma)) \preceq H_{\max} \cdot \mathbf{I},$$

where $H_{\min} = \Omega \left(\frac{1}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \right)$ and $H_{\max} = C$. Therefore, we have

$$\begin{aligned} \left\| \mathbf{W}_{t+1} - \widehat{\mathbf{W}}_n \right\|_{\text{F}} &\leq \left\| \mathbf{I} - \eta \int_0^1 \nabla^2 f_n(\mathbf{W}(\gamma)) \right\| \left\| \mathbf{W}_t - \widehat{\mathbf{W}}_n \right\|_{\text{F}} \\ &\leq (1 - \eta H_{\min}) \left\| \mathbf{W}_t - \widehat{\mathbf{W}}_n \right\|_{\text{F}}. \end{aligned} \quad (24)$$

Hence, by setting $\eta = \frac{1}{H_{\max}} := \Omega(C)$, we obtain

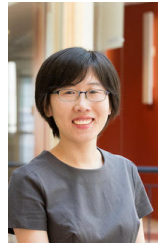
$$\left\| \mathbf{W}_{t+1} - \widehat{\mathbf{W}}_n \right\|_{\text{F}} \leq \left(1 - \frac{H_{\min}}{H_{\max}} \right) \left\| \mathbf{W}_t - \widehat{\mathbf{W}}_n \right\|_{\text{F}}, \quad (25)$$

which implies that gradient descent converges linearly to the local minimizer $\widehat{\mathbf{W}}_n$.

REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [2] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [4] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [5] G. Cybenko, "Approximation by superpositions of a sigmoidal function," *Mathematics of control, signals and systems*, vol. 2, no. 4, pp. 303–314, 1989.
- [6] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [7] A. R. Barron, "Universal approximation bounds for superpositions of a sigmoidal function," *IEEE Transactions on Information theory*, vol. 39, no. 3, pp. 930–945, 1993.
- [8] M. Telgarsky, "benefits of depth in neural networks," in *Conference on Learning Theory*, 2016, pp. 1517–1539.
- [9] Y. N. 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 *Advances in neural information processing systems*, 2014, pp. 2933–2941.
- [10] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *arXiv preprint arXiv:1611.03530*, 2016.
- [11] P. L. Bartlett, D. J. Foster, and M. J. Telgarsky, "Spectrally-normalized margin bounds for neural networks," in *Advances in Neural Information Processing Systems*, 2017, pp. 6241–6250.
- [12] A. Brutzkus, A. Globerson, E. Malach, and S. Shalev-Shwartz, "SGD learns over-parameterized networks that provably generalize on linearly separable data," in *International Conference on Learning Representations*, 2018.
- [13] M. Mondelli and A. Montanari, "On the connection between learning two-layer neural networks and tensor decomposition," in *Proceedings of Machine Learning Research*, 2019, pp. 1051–1060.
- [14] M. Soltanolkotabi, "Learning relu via gradient descent," in *Advances in Neural Information Processing Systems*, 2017, pp. 2007–2017.
- [15] K. Zhong, Z. Song, P. Jain, P. L. Bartlett, and I. S. Dhillon, "Recovery guarantees for one-hidden-layer neural networks," in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70, 2017, pp. 4140–4149.
- [16] Y. Li and Y. Yuan, "Convergence analysis of two-layer neural networks with relu activation," in *Advances in Neural Information Processing Systems*, 2017, pp. 597–607.
- [17] S. Mei, Y. Bai, and A. Montanari, "The landscape of empirical risk for nonconvex losses," *Ann. Statist.*, vol. 46, no. 6A, pp. 2747–2774, 12 2018.
- [18] A. Brutzkus and A. Globerson, "Globally optimal gradient descent for a ConvNet with Gaussian inputs," in *Proceedings of the 34th International Conference on Machine Learning*, 2017, pp. 605–614.
- [19] S. Oymak, "Learning compact neural networks with regularization," in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 3966–3975.
- [20] S. S. Du, J. D. Lee, and Y. Tian, "When is a convolutional filter easy to learn?" in *International Conference on Learning Representations*, 2018.
- [21] R. Sun and Z.-Q. Luo, "Guaranteed matrix completion via non-convex factorization," *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 6535–6579, 2016.
- [22] Y. Chen and Y. Chi, "Harnessing structures in big data via guaranteed low-rank matrix estimation," *IEEE Signal Processing Magazine*, 2018.
- [23] E. J. Candès, X. Li, and M. Soltanolkotabi, "Phase retrieval via Wirtinger flow: Theory and algorithms," *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1985–2007, April 2015.
- [24] R. Ge and T. Ma, "On the optimization landscape of tensor decompositions," in *Advances in Neural Information Processing Systems 30*, 2017, pp. 3653–3663.
- [25] R. Ge, J. D. Lee, and T. Ma, "Matrix completion has no spurious local minimum," in *Advances in Neural Information Processing Systems*, 2016, pp. 2973–2981.
- [26] J. Sun, Q. Qu, and J. Wright, "Complete dictionary recovery using nonconvex optimization," *International Conference on Machine Learning*, pp. 2351–2360, 2015.
- [27] S. Bhojanapalli, B. Neyshabur, and N. Srebro, "Global optimality of local search for low rank matrix recovery," in *Advances in Neural Information Processing Systems*, 2016, pp. 3873–3881.
- [28] C. Ma, K. Wang, Y. Chi, and Y. Chen, "Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion," in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 3345–3354.

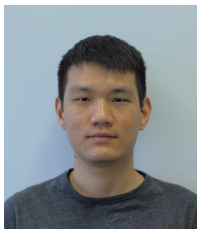
- [29] Y. Chi, Y. M. Lu, and Y. Chen, "Nonconvex optimization meets low-rank matrix factorization: An overview," *IEEE Transactions on Signal Processing*, vol. 67, no. 20, pp. 5239–5269, 2019.
- [30] M. Soltanolkotabi, A. Javanmard, and J. D. Lee, "Theoretical insights into the optimization landscape of over-parameterized shallow neural networks," *IEEE Transactions on Information Theory*, 2018.
- [31] D. Boob and G. Lan, "Theoretical properties of the global optimizer of two layer neural network," *arXiv preprint arXiv:1710.11241*, 2017.
- [32] I. Safran and O. Shamir, "On the quality of the initial basin in overspecified neural networks," in *International Conference on Machine Learning*, 2016, pp. 774–782.
- [33] Q. Nguyen and M. Hein, "The loss surface of deep and wide neural networks," in *International Conference on Machine Learning*, 2017, pp. 2603–2612.
- [34] Y. Tian, "An analytical formula of population gradient for two-layered ReLU network and its applications in convergence and critical point analysis," in *Proceedings of the 34th International Conference on Machine Learning*, 2017, pp. 3404–3413.
- [35] R. Ge, J. D. Lee, and T. Ma, "Learning one-hidden-layer neural networks with landscape design," in *International Conference on Learning Representations*, 2018.
- [36] I. Safran and O. Shamir, "Spurious local minima are common in two-layer ReLU neural networks," in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 4433–4441.
- [37] Y. Chen, Y. Chi, J. Fan, and C. Ma, "Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval," *Math. Program.* 176, 5–37, 2019.
- [38] X. Zhang, Y. Yu, L. Wang, and Q. Gu, "Learning one-hidden-layer relu networks via gradient descent," *International Conference on Artificial Intelligence and Statistics*, 2019.
- [39] Y. Li, C. Ma, Y. Chen, and Y. Chi, "Nonconvex matrix factorization from rank-one measurements," *arXiv preprint arXiv:1802.06286*, 2018.
- [40] S. Goel, A. Klivans, and R. Meka, "Learning one convolutional layer with overlapping patches," in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 1783–1791.
- [41] S. Feizi, H. Javadi, J. Zhang, and D. Tse, "Porcupine neural networks: Approximating neural network landscapes," in *Advances in Neural Information Processing Systems 31*, 2018, pp. 4836–4846.
- [42] K. Zhong, Z. Song, and I. S. Dhillon, "Learning non-overlapping convolutional neural networks with multiple kernels," *arXiv preprint arXiv:1711.03440*, 2017.
- [43] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," *Compressed Sensing, Theory and Applications*, pp. 210 – 268, 2012.



Yuejie Chi (S'09-M'12-SM'17) received the Ph.D. degree in Electrical Engineering from Princeton University in 2012, and the B.E. (Hon.) degree in Electrical Engineering from Tsinghua University, Beijing, China, in 2007. She was with The Ohio State University from 2012 to 2017. Since 2018, she is an Associate Professor with the department of Electrical and Computer Engineering at Carnegie Mellon University, where she holds the Robert E. Doherty Early Career Development Professorship. Her research interests include signal processing, statistical inference, machine learning, large-scale optimization, and their applications in data science, inverse problems, imaging, and sensing systems. She is a recipient of the PECASE Award, NSF CAREER Award, AFOSR and ONR Young Investigator Program Awards, Ralph E. Powe Junior Faculty Enhancement Award, and Google Faculty Research Award. She received IEEE SPS Early Career Technical Achievement Award and Young Author Best Paper Award from IEEE Signal Processing Society, and Best Paper Award at the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP). She has served as an Elected Member of the SPTM, SAM and MLSP Technical Committees of the IEEE Signal Processing Society. She currently serves as an Associate Editor of IEEE Trans. on Signal Processing.



Dr. Yingbin Liang (S'01-M'05-SM'16) is currently a Professor at the Department of Electrical and Computer Engineering at the Ohio State University (OSU). She received the Ph.D. degree in Electrical Engineering from the University of Illinois at Urbana-Champaign in 2005, and served on the faculty of University of Hawaii and Syracuse University before she joined OSU. Dr. Liang's research interests include machine learning, optimization, information theory and statistical signal processing. Dr. Liang received the National Science Foundation CAREER Award and the State of Hawaii Governor Innovation Award in 2009. She also received EURASIP Best Paper Award for the EURASIP Journal on Wireless Communications and Networking in 2014. She served as an Associate Editor for the Shannon Theory of the IEEE Transactions on Information Theory during 2013-2015.



Haoyu Fu received the Ph.D. degree in Electrical and Computer Engineering from The Ohio State University, Columbus, OH, USA, in 2019. His research interests include compressed sensing, statistical signal processing, optimization, and machine learning.

Supplementary Materials: Additional Proofs

APPENDIX D PROOF OF AUXILIARY LEMMAS

A. Proof of Lemma 1.

We prove the two claims for FCN and CNN separately as below.

- **The FCN case:** Let $\Delta = \nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)$. For each $(j, l) \in [K] \times [K]$, let $\Delta_{j,l} \in \mathbb{R}^{d \times d}$ denote the (j, l) -th block of Δ . Let $\mathbf{a} = [\mathbf{a}_1^\top, \dots, \mathbf{a}_K^\top]^\top \in \mathbb{R}^{dK}$. By definition,

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| = \max_{\|\mathbf{a}\|=1} \mathbf{a}^\top (\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)) \mathbf{a} = \max_{\|\mathbf{a}\|=1} \sum_{j=1}^K \sum_{l=1}^K \mathbf{a}_j^\top \Delta_{j,l} \mathbf{a}_l. \quad (26)$$

From (11) we know that

$$\Delta_{j,l} = \frac{\partial^2 f(\mathbf{W})}{\partial \mathbf{w}_j \partial \mathbf{w}_l} - \frac{\partial^2 f(\mathbf{W}^*)}{\partial \mathbf{w}_j^* \partial \mathbf{w}_l^*} = \mathbb{E} [(\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}^*)) \cdot \mathbf{x} \mathbf{x}^\top], \quad (27)$$

and then by the mean value theorem, we can further expand $\xi_{j,l}(\mathbf{W})$ as

$$\xi_{j,l}(\mathbf{W}) = \xi_{j,l}(\mathbf{W}^*) + \sum_{k=1}^K \left\langle \frac{\partial \xi_{j,l}(\widetilde{\mathbf{W}})}{\partial \widetilde{\mathbf{w}}_k}, \mathbf{w}_k - \mathbf{w}_k^* \right\rangle, \quad (28)$$

where $\widetilde{\mathbf{W}} = \eta \cdot \mathbf{W} + (1 - \eta) \mathbf{W}^*$ for some $\eta \in (0, 1)$. Thus we can write $\Delta_{j,l}$ as

$$\Delta_{j,l} = \mathbb{E} \left[\left(\sum_{k=1}^K \left\langle \frac{\partial \xi_{j,l}(\widetilde{\mathbf{W}})}{\partial \widetilde{\mathbf{w}}_k}, \mathbf{w}_k - \mathbf{w}_k^* \right\rangle \right) \cdot \mathbf{x} \mathbf{x}^\top \right], \quad (29)$$

which can be further simplified as

$$\Delta_{j,l} = \mathbb{E} \left[\left(\sum_{k=1}^K T_{j,l,k} \langle \mathbf{x}, \mathbf{w}_k - \mathbf{w}_k^* \rangle \right) \cdot \mathbf{x} \mathbf{x}^\top \right], \quad (30)$$

by the fact that $\frac{\partial \xi_{j,l}(\widetilde{\mathbf{W}})}{\partial \widetilde{\mathbf{w}}_k}$ can be written as $T_{j,l,k} \cdot \mathbf{x}$, where $T_{j,l,k} \in \mathbb{R}$ is a scalar depending on \mathbf{x} . When $j = l$, we calculate $\frac{\partial \xi_{j,l}(\widetilde{\mathbf{W}})}{\partial \widetilde{\mathbf{w}}_k}$ for illustration,

$$\frac{\partial \xi_{j,j}(\mathbf{W})}{\partial \mathbf{w}_k} = \begin{cases} \left(-\frac{2}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})^2}{H(\mathbf{W})^3} + \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})^2} \right) \frac{1}{K} \phi(\mathbf{w}_k^\top \mathbf{x}) \mathbf{x} & k \neq j \\ \left(\frac{2}{K^2} \left(\frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})^2} - \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})^2}{H(\mathbf{W})^3} \right) + \frac{1}{K} \left(\frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})^2} - \frac{\phi'''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})} \right) \right) \frac{1}{K} \phi(\mathbf{w}_k^\top \mathbf{x}) \mathbf{x} & k = j \end{cases}, \quad (31)$$

where we have simplified the presentation by setting $y = 1$, since y is a binary random variable, and we will show that in either case $|T_{j,j,k}|$ is upper bounded, i.e., in this case

$$|T_{j,j,k}| \leq \begin{cases} \max \left\{ \frac{2}{K^3} \frac{1}{H(\widetilde{\mathbf{W}})^3}, \frac{1}{K^2} \frac{1}{H(\widetilde{\mathbf{W}})^2} \right\} & y = 1 \\ \max \left\{ \frac{2}{K^3} \frac{1}{(1-H(\widetilde{\mathbf{W}}))^3}, \frac{1}{K^2} \frac{1}{(1-H(\widetilde{\mathbf{W}}))^2} \right\} & y = 0 \end{cases},$$

since $\phi(\cdot), \phi'(\cdot), \phi''(\cdot), \phi'''(\cdot)$ are bounded. More generally, by calculating the other case we can claim that

$$|T_{j,l,k}| \leq \max \left\{ \frac{2}{K^3} \frac{1}{H(\widetilde{\mathbf{W}})^3}, \frac{1}{K^2} \frac{1}{H(\widetilde{\mathbf{W}})^2}, \frac{2}{K^3} \frac{1}{(1-H(\widetilde{\mathbf{W}}))^3}, \frac{1}{K^2} \frac{1}{(1-H(\widetilde{\mathbf{W}}))^2} \right\}, \quad (32)$$

holds for all j, l, k . Then, we can upper bound $\mathbf{a}_j^\top \Delta_{j,l} \mathbf{a}_l$ using Cauchy-Schwarz inequality,

$$\begin{aligned} \mathbf{a}_j^\top \Delta_{j,l} \mathbf{a}_l &= \mathbb{E} \left[\left(\sum_{k=1}^K T_{j,l,k} \langle \mathbf{x}, \mathbf{w}_k - \mathbf{w}_k^* \rangle \right) \cdot (\mathbf{a}_j^\top \mathbf{x}) (\mathbf{a}_l^\top \mathbf{x}) \right] \\ &\leq \sqrt{\mathbb{E} \left[\sum_{k=1}^K T_{j,l,k}^2 \right]} \cdot \sqrt{\mathbb{E} \left[\sum_{k=1}^K (\langle \mathbf{x}, \mathbf{w}_k - \mathbf{w}_k^* \rangle (\mathbf{a}_j^\top \mathbf{x}) (\mathbf{a}_l^\top \mathbf{x}))^2 \right]} \\ &\leq \sqrt{\sum_{k=1}^K \mathbb{E} [T_{j,l,k}^2]} \cdot \sqrt{\sum_{k=1}^K \|\mathbf{w}_k - \mathbf{w}_k^*\|_2^2 \cdot \|\mathbf{a}_j\|_2^2 \cdot \|\mathbf{a}_l\|_2^2}. \end{aligned} \quad (33)$$

Plug it back to (26) we can obtain the following inequality,

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \leq \max_{\|\mathbf{a}\|=1} \sum_{j=1}^K \sum_{l=1}^K \sqrt{\sum_{k=1}^K \mathbb{E}[T_{j,l,k}^2]} \cdot \sqrt{\sum_{k=1}^K \|\mathbf{w}_k - \mathbf{w}_k^*\|_2^2 \cdot \|\mathbf{a}_j\|_2^2 \cdot \|\mathbf{a}_l\|_2^2}. \quad (34)$$

Then the problem boils down to upper bound $\mathbb{E}[T_{i,j,k}^2]$, which we can apply the following lemma, whose proof can be found in Section D-E.

Lemma 5. Let $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $t = \max\{\|\mathbf{w}_1\|_2, \dots, \|\mathbf{w}_K\|_2\}$ and $z \in \mathbb{Z}$ such that $z \geq 1$, for the sigmoid activation function $\phi(x) = \frac{1}{1+e^{-x}}$, the following

$$\mathbb{E}\left[\left(\frac{1}{\frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_j^\top \mathbf{x})}\right)^z\right] \leq C_1 \cdot e^{t^2}, \quad \mathbb{E}\left[\left(\frac{1}{\left(1 - \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_j^\top \mathbf{x})\right)}\right)^z\right] \leq C_2 \cdot e^{t^2} \quad (35)$$

holds for some large enough constants C_1, C_2 that depend on the constant z .

Setting $z = 4$ and $z = 6$ in Lemma 5, together with (32) we obtain that

$$\mathbb{E}[T_{j,l,k}^2] \leq \frac{C}{K^4} \cdot e^{\max_{1 \leq i \leq k} \|\tilde{\mathbf{w}}_i\|_2^2}, \quad (36)$$

holds for some constant C . Plugging (36) into (34), we obtain

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \leq \frac{C}{K^{\frac{3}{2}}} e^{\|\tilde{\mathbf{W}}\|_F^2} \cdot \|\mathbf{W} - \mathbf{W}^*\|_F \cdot \max_{\|\mathbf{a}\|=1} \sum_{j=1}^K \sum_{l=1}^K \|\mathbf{a}_j\|_2 \|\mathbf{a}_l\|_2 \leq \frac{C}{K^{\frac{3}{2}}} e^{\|\tilde{\mathbf{W}}\|_F^2} \cdot \|\mathbf{W} - \mathbf{W}^*\|_F. \quad (37)$$

Further since $e^{\max_{1 \leq i \leq k} \|\tilde{\mathbf{w}}_i\|_2^2} \leq C$ gives that $\|\mathbf{w}_i - \mathbf{w}_i^*\|_2 \leq 0.7$, where we have used the assumption that $\max_{1 \leq i \leq k} \|\mathbf{w}_i^*\|_2^2 \leq 1$, we conclude that

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \leq \frac{C}{K^{\frac{3}{2}}} \|\mathbf{W} - \mathbf{W}^*\|_F \quad (38)$$

holds for some constant C .

- **The CNN case:** according to (13), we can calculate the upper bound of $\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\|$ by definition as

$$\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| \leq \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \mathbb{E}\left[(g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}^*)) \cdot \mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u}\right]. \quad (39)$$

We then again apply the mean value theorem to $g_{j,l}(\mathbf{w})$, such that there exists $\tilde{\mathbf{w}} = \eta \mathbf{w} + (1 - \eta) \mathbf{w}^*$ for some $\eta \in (0, 1)$,

$$g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}^*) = \langle \nabla g_{j,l}(\tilde{\mathbf{w}}), \mathbf{w} - \mathbf{w}^* \rangle.$$

Similarly to the FCN case, we can write $\nabla g_{j,l}(\tilde{\mathbf{w}})$ in the form of

$$\nabla g_{j,l}(\tilde{\mathbf{w}}) = \sum_{k=1}^K S_{j,l,k} \cdot \mathbf{x}^{(k)},$$

where $S_{j,l,k}$ is a scalar that depends on $\tilde{\mathbf{w}}$ and $\mathbf{x}^{(k)}$, $k = 1, \dots, K$. Again we take $j \neq l$ as an example to calculate $S_{j,l,k}$, by definition, and obtain

$$\begin{aligned} K^2 \cdot \frac{\partial g_{j,l}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{(1 - H(\mathbf{w})) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi''(\mathbf{w}^\top \mathbf{x}^{(l)})}{(1 - H(\mathbf{w}))^3} \cdot \mathbf{x}^{(l)} + \frac{(1 - H(\mathbf{w})) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \phi''(\mathbf{w}^\top \mathbf{x}^{(j)})}{(1 - H(\mathbf{w}))^3} \cdot \mathbf{x}^{(j)} \\ &\quad - \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{(1 - H(\mathbf{w}))^3} \cdot \left(\frac{1}{K} \sum_{k=1}^K \mathbf{x}^{(k)}\right), \end{aligned} \quad (40)$$

where we set $y = 0$ for simplification. Then we obtain

$$S_{j,l,l} = \frac{1}{K^2} \frac{(1 - H(\mathbf{w})) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi''(\mathbf{w}^\top \mathbf{x}^{(l)})}{(1 - H(\mathbf{w}))^3} - \frac{1}{K^3} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{(1 - H(\mathbf{w}))^3}. \quad (41)$$

and

$$|S_{j,l,l}| \leq \frac{1}{K^2} \frac{1}{(1 - H(\tilde{\mathbf{w}}))^3}, \quad (42)$$

hold, where we used the fact that $0 \leq H(\mathbf{w}) \leq 1$ and $\phi'(\cdot), \phi''(\cdot)$ are bounded. Hence in the same way, we can obtain

$$|S_{j,l,k}| \leq \begin{cases} \max \left\{ \frac{1}{K^2} \frac{1}{(1-H(\tilde{\mathbf{w}}))^3}, \frac{1}{K^2} \frac{1}{(H(\tilde{\mathbf{w}}))^3} \right\} & j \neq l \\ \max \left\{ \frac{1}{K} \frac{1}{(1-H(\tilde{\mathbf{w}}))^2}, \frac{1}{K} \frac{1}{(H(\tilde{\mathbf{w}}))^2} \right\} & j = l \end{cases}. \quad (43)$$

Plug these back to (39) we obtain

$$\begin{aligned} \|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| &\leq \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \mathbb{E} \left[\sum_{k=1}^K \left\langle S_{j,l,k} \cdot \mathbf{x}^{(k)}, \mathbf{w} - \mathbf{w}^* \right\rangle \cdot \mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u} \right] \\ &= \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \mathbb{E} \left[\sum_{k=1}^K S_{j,l,k} \cdot (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{x}^{(k)} \cdot \mathbf{u}^\top \mathbf{x}^{(j)} \cdot \mathbf{x}^{(l)\top} \mathbf{u} \right] \\ &\leq \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \sqrt{\mathbb{E} \left[\sum_{k=1}^K S_{j,l,k}^2 \right] \cdot \mathbb{E} \left[\sum_{k=1}^K \left((\mathbf{w} - \mathbf{w}^*)^\top \mathbf{x}^{(k)} \right)^2 (\mathbf{u}^\top \mathbf{x}^{(j)})^2 (\mathbf{x}^{(l)\top} \mathbf{u})^2 \right]} \\ &\leq \max_{\|\mathbf{u}\|_2=1} \sum_{j=1}^K \sum_{l=1}^K \sqrt{\mathbb{E} \left[\sum_{k=1}^K S_{j,l,k}^2 \right] \cdot \sum_{k=1}^K \|\mathbf{w} - \mathbf{w}^*\|_2^2 \cdot \|\mathbf{u}\|_2^2 \cdot \|\mathbf{u}\|_2^2} \\ &\leq C \cdot K \cdot e^{\|\tilde{\mathbf{w}}\|_2^2} \cdot \|\mathbf{w} - \mathbf{w}^*\|_2, \end{aligned} \quad (44)$$

where the second inequality follows from Cauchy-Schwarz inequality, and the last inequality follows from (43) and Lemma 5. Further since $e^{\|\tilde{\mathbf{w}}\|_2^2} \leq C \cdot (1 + \|\mathbf{w} - \mathbf{w}^*\|_2^2)$ given that $\|\mathbf{w} - \mathbf{w}^*\|_2 \leq 0.7$, we conclude that

$$\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^*)\| \leq C \cdot K \cdot \|\mathbf{w} - \mathbf{w}^*\|_2 \quad (45)$$

holds for some constant C and $\|\mathbf{w} - \mathbf{w}^*\| \leq 0.7$.

B. Proof of Lemma 2

We first present upper and lower bounds on the Hessian $\nabla^2 f(\mathbf{W}^*)$ of the population risk at ground truth, and then apply Lemma 1 to obtain a uniform bound in the neighborhood of \mathbf{W}^* .

• The FCN case: Recall

$$\begin{aligned} \frac{\partial^2 f(\mathbf{W}^*)}{\partial \mathbf{w}_j^2} &= \mathbb{E} \left[\frac{1}{K^2} \cdot \left(\frac{\phi'(\mathbf{w}_j^{*\top} \mathbf{x})^2}{H(\mathbf{W}^*)(1-H(\mathbf{W}^*))} \right) \mathbf{x} \mathbf{x}^\top \right], \\ \frac{\partial^2 f(\mathbf{W}^*)}{\partial \mathbf{w}_j \partial \mathbf{w}_l} &= \mathbb{E} \left[\frac{1}{K^2} \cdot \left(\frac{\phi'(\mathbf{w}_j^{*\top} \mathbf{x}) \phi'(\mathbf{w}_l^{*\top} \mathbf{x})}{H(\mathbf{W}^*)(1-H(\mathbf{W}^*))} \right) \mathbf{x} \mathbf{x}^\top \right], \end{aligned}$$

where we have applied the fact that $\mathbb{E}[y|\mathbf{x}] = H(\mathbf{W}^*)$. Let $\mathbf{a} = [\mathbf{a}_1^\top, \dots, \mathbf{a}_K^\top]^\top \in \mathbb{R}^{dK}$. Then we can write

$$\nabla^2 f(\mathbf{W}^*) \succeq \left(\min_{\|\mathbf{a}\|_2=1} \mathbf{a}^\top \nabla^2 f(\mathbf{W}^*) \mathbf{a} \right) \cdot \mathbf{I} = \min_{\|\mathbf{a}\|_2=1} \frac{1}{K^2} \mathbb{E} \left[\frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x}) \right)^2}{H(\mathbf{W}^*)(1-H(\mathbf{W}^*))} \right] \cdot \mathbf{I}. \quad (46)$$

Since $0 \leq H(\mathbf{W}^*) \leq 1$, we have that $H(\mathbf{W}^*)(1-H(\mathbf{W}^*)) \leq \frac{1}{4}$. Hence,

$$\nabla^2 f(\mathbf{W}^*) \succeq \min_{\|\mathbf{a}\|_2=1} \frac{4}{K^2} \mathbb{E} \left[\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x}) \right)^2 \right] \cdot \mathbf{I} \succeq \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \cdot \mathbf{I}, \quad (47)$$

where the last inequality follows from [15, Lemmas D.4 and D.6]. To derive an upper bound of $\nabla^2 f(\mathbf{W}^*)$, we have

$$\nabla^2 f(\mathbf{W}^*) \preceq \left(\max_{\|\mathbf{a}\|_2=1} \mathbf{a}^\top \nabla^2 f(\mathbf{W}^*) \mathbf{a} \right) \cdot \mathbf{I} = \max_{\|\mathbf{a}\|_2=1} \frac{1}{K^2} \mathbb{E} \left[\frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x}) \right)^2}{\frac{1}{K^2} \sum_{j,l} \phi(\mathbf{w}_j^{*\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{*\top} \mathbf{x}))} \right]. \quad (48)$$

Then by Cauchy-Schwarz inequality, we have

$$\frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x}) (\mathbf{a}_j^\top \mathbf{x}) \right)^2}{\frac{1}{K^2} \sum_{j,l} \phi(\mathbf{w}_j^{*\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{*\top} \mathbf{x}))} \leq \frac{\left(\sum_{j=1}^K \phi'(\mathbf{w}_j^{*\top} \mathbf{x})^2 \right) \cdot \left(\sum_{j=1}^K (\mathbf{a}_j^\top \mathbf{x})^2 \right)}{\frac{1}{K^2} \sum_{j,l} \phi(\mathbf{w}_j^{*\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{*\top} \mathbf{x}))}. \quad (49)$$

Further since $\phi'(\mathbf{w}_j^{\star\top} \mathbf{x}) \leq \frac{1}{4}$, and

$$\sum_{j,l} \phi(\mathbf{w}_j^{\star\top} \mathbf{x}) (1 - \phi(\mathbf{w}_l^{\star\top} \mathbf{x})) \geq \sum_{j=1}^K \phi(\mathbf{w}_j^{\star\top} \mathbf{x}) (1 - \phi(\mathbf{w}_j^{\star\top} \mathbf{x})) = \sum_{j=1}^K \phi'(\mathbf{w}_j^{\star\top} \mathbf{x}) \geq 4 \sum_{j=1}^K \phi'(\mathbf{w}_j^{\star\top} \mathbf{x})^2, \quad (50)$$

we obtain

$$\mathbf{a}^\top \nabla^2 f(\mathbf{W}^*) \mathbf{a} \preceq \frac{1}{K^2} \mathbb{E} \left[\frac{CK^2}{4} \sum_{j=1}^K (\mathbf{a}_j^\top \mathbf{x})^2 \right]. \quad (51)$$

Plugging (51) back to (48), we obtain

$$\nabla^2 f(\mathbf{W}^*) \preceq C \cdot \mathbf{I}. \quad (52)$$

Thus together with the lower bound (47), we conclude that

$$\frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{W}^*) \preceq C \cdot \mathbf{I}. \quad (53)$$

From Lemma 1, we have

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \lesssim \frac{C}{K^{\frac{3}{2}}} \|\mathbf{W} - \mathbf{W}^*\|_F. \quad (54)$$

Therefore, if $\|\mathbf{W}^* - \mathbf{W}\|_F \leq 0.7$ and

$$\frac{C}{K^{\frac{3}{2}}} \cdot \|\mathbf{W} - \mathbf{W}^*\|_F \leq \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda},$$

i.e., if $\|\mathbf{W} - \mathbf{W}^*\|_F \leq \min \left\{ \frac{C}{K^{\frac{3}{2}}} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}, 0.7 \right\}$ for some constant C , we have

$$\begin{aligned} \sigma_{\min}(\nabla^2 f(\mathbf{W})) &\geq \sigma_{\min}(\nabla^2 f(\mathbf{W}^*)) - \|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| \gtrsim \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda} - \frac{C}{K^{\frac{3}{2}}} \|\mathbf{W} - \mathbf{W}^*\|_F \\ &\gtrsim \frac{4}{K^2} \cdot \frac{\rho_{\text{FCN}}(\sigma_K)}{\kappa^2 \lambda}. \end{aligned}$$

Moreover, within the same neighborhood, by the triangle inequality we have

$$\|\nabla^2 f(\mathbf{W})\| \leq \|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\| + \|\nabla^2 f(\mathbf{W}^*)\| \lesssim C.$$

- **The CNN case:** Following from (13), we have

$$\nabla^2 f(\mathbf{w}^*) = \mathbb{E} \left[\frac{\frac{1}{K^2} \sum_{j,l} \phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(l)}) \mathbf{x}^{(j)} \mathbf{x}^{(l)\top}}{H(\mathbf{w}^*) (1 - H(\mathbf{w}^*))} \right]. \quad (55)$$

By definition, we lower bound $\nabla^2 f(\mathbf{w}^*)$ by

$$\begin{aligned} &\min_{\|\mathbf{u}\|=1} \mathbb{E} \left[\frac{\frac{1}{K^2} \sum_{j,l} \phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)} \phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(l)}) \mathbf{u}^\top \mathbf{x}^{(l)}}{H(\mathbf{w}^*) (1 - H(\mathbf{w}^*))} \right] \cdot \mathbf{I} \\ &\succeq \min_{\|\mathbf{u}\|=1} \mathbb{E} \left[\frac{4}{K^2} \sum_{j,l} \phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)} \phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(l)}) \mathbf{u}^\top \mathbf{x}^{(l)} \right] \cdot \mathbf{I} \\ &= \frac{4}{K^2} \cdot \left(\min_{\|\mathbf{u}\|=1} \sum_{j \neq l} \mathbb{E} [\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)}] \cdot \mathbb{E} [\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(l)}) \mathbf{u}^\top \mathbf{x}^{(l)}] + \sum_{j=1}^K \mathbb{E} [(\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)})^2] \right) \cdot \mathbf{I}, \end{aligned}$$

where the last equality follows from the fact that $\mathbf{x}^{(j)}$ is independent from $\mathbf{x}^{(l)}$ given that $j \neq l$. Next we decompose \mathbf{u} as $\mathbf{u} = \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2^2} \cdot \mathbf{w}^* + \left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2^2} \cdot \mathbf{w}^* \right)$, and calculate the expectation as

$$\begin{aligned} \mathbb{E} [\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{u}^\top \mathbf{x}^{(j)}] &= \mathbb{E} \left[\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \left(\frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2^2} \cdot \mathbf{w}^* + \left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2^2} \cdot \mathbf{w}^* \right) \right)^\top \mathbf{x}^{(j)} \right] \\ &= \mathbb{E} \left[\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2^2} \cdot \mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right] + \mathbb{E} [\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)})] \cdot \mathbb{E} \left[\left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2^2} \cdot \mathbf{w}^* \right)^\top \mathbf{x}^{(j)} \right] \\ &= \frac{\mathbf{u}^\top \mathbf{w}^*}{\|\mathbf{w}^*\|_2^2} \mathbb{E} [\phi'(\mathbf{w}^{\star\top} \mathbf{x}^{(j)}) \mathbf{w}^{\star\top} \mathbf{x}^{(j)}], \end{aligned}$$

where the second equality follows from the independence of $\mathbf{w}^{\star\top} \mathbf{x}^{(j)}$ and $\left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^\star}{\|\mathbf{w}^\star\|_2^2} \cdot \mathbf{w}^\star\right)^\top \mathbf{x}^{(j)}$. Hence,

$$\mathbb{E} \left[\phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right) \mathbf{u}^\top \mathbf{x}^{(j)} \right] \cdot \mathbb{E} \left[\phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(l)} \right) \mathbf{u}^\top \mathbf{x}^{(l)} \right] = \left(\frac{\mathbf{u}^\top \mathbf{w}^\star}{\|\mathbf{w}^\star\|_2^2} \right)^2 (\mathbb{E} [\phi'(z) z])^2 = 0, \quad (56)$$

where $z = \mathbf{w}^{\star\top} \mathbf{x}^{(j)} \sim \mathcal{N}(0, \|\mathbf{w}^\star\|_2^2)$, and the last equality follows because $\phi'(z) z = -(\phi'(-z) \cdot (-z))$. Similarly,

$$\begin{aligned} & \mathbb{E} \left[\left(\phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right) \mathbf{u}^\top \mathbf{x}^{(j)} \right)^2 \right] \\ &= \mathbb{E} \left[\phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right)^2 \cdot \left(\left(\frac{\mathbf{u}^\top \mathbf{w}^\star}{\|\mathbf{w}^\star\|_2^2} \cdot \mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right)^2 + \left(\left(\mathbf{u} - \frac{\mathbf{u}^\top \mathbf{w}^\star}{\|\mathbf{w}^\star\|_2^2} \cdot \mathbf{w}^\star \right)^\top \mathbf{x}^{(j)} \right)^2 \right) \right] \\ &= \left(\frac{\mathbf{u}^\top \mathbf{w}^\star}{\|\mathbf{w}^\star\|_2^2} \right)^2 \cdot \mathbb{E} \left[\phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right)^2 \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right)^2 \right] + \left(\|\mathbf{u}\|_2^2 - \frac{(\mathbf{u}^\top \mathbf{w}^\star)^2}{\|\mathbf{w}^\star\|_2^2} \right) \cdot \mathbb{E} \left[\phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right)^2 \right]. \end{aligned} \quad (57)$$

Together with Definition 2, we have

$$\mathbb{E} \left[\left(\phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right) \mathbf{u}^\top \mathbf{x}^{(j)} \right)^2 \right] \geq \rho_{\text{CNN}} (\|\mathbf{w}^\star\|_2). \quad (58)$$

Hence,

$$\nabla^2 f(\mathbf{w}^\star) \succeq \frac{4}{K} \cdot \rho_{\text{CNN}} (\|\mathbf{w}^\star\|_2) \cdot \mathbf{I}. \quad (59)$$

Moreover, we apply Cauchy-Schwarz inequality and upper bound the Hessian as

$$\nabla^2 f(\mathbf{w}^\star) \leq \left(\max_{\|\mathbf{u}\|_2=1} \mathbf{u}^\top \nabla^2 f(\mathbf{w}^\star) \mathbf{u} \right) \cdot \mathbf{I} \leq \max_{\|\mathbf{u}\|_2=1} \mathbb{E} \left[\frac{\sum_{j=1}^K \left(\frac{1}{K} \phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right) \right)^2 \cdot \sum_{j=1}^K \left(\mathbf{u}^\top \mathbf{x}^{(j)} \right)^2}{H(\mathbf{w}^\star) (1 - H(\mathbf{w}^\star))} \right] \cdot \mathbf{I}. \quad (60)$$

Using (50), i.e.,

$$\frac{\frac{1}{K^2} \sum_{j=1}^K \phi' \left(\mathbf{w}^{\star\top} \mathbf{x}^{(j)} \right)^2}{H(\mathbf{w}^\star) (1 - H(\mathbf{w}^\star))} \leq \frac{1}{4}, \quad (61)$$

we upper bound the right-hand side of (60) as

$$\nabla^2 f(\mathbf{w}^\star) \preceq \max_{\|\mathbf{u}\|_2=1} \mathbb{E} \left[\frac{1}{4} \sum_{j=1}^K \left(\mathbf{u}^\top \mathbf{x}^{(j)} \right)^2 \right] \cdot \mathbf{I} = \frac{K}{4} \cdot \mathbf{I}. \quad (62)$$

Together with the lower bound, we now conclude that

$$\frac{4}{K} \cdot \rho_{\text{CNN}} (\|\mathbf{w}^\star\|_2) \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{w}^\star) \preceq \frac{K}{4} \cdot \mathbf{I}. \quad (63)$$

And following from (15) in Lemma 1, we have

$$\|\nabla^2 f(\mathbf{w}) - \nabla^2 f(\mathbf{w}^\star)\| \leq C \cdot K \cdot \|\mathbf{w} - \mathbf{w}^\star\|_2. \quad (64)$$

Thus if $\|\mathbf{w} - \mathbf{w}^\star\| \leq \min \left\{ 0.7, C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^\star\|_2)}{K^2} \right\}$, we have

$$C \cdot \frac{\rho_{\text{CNN}}(\|\mathbf{w}^\star\|_2)}{K} \cdot \mathbf{I} \preceq \nabla^2 f(\mathbf{w}) \preceq C \cdot K \cdot \mathbf{I}. \quad (65)$$

C. Proof of Lemma 3

We apply a covering type of argument to show that the Hessian of the empirical risk function concentrates around the Hessian of the population risk function uniformly, and the argument applies to both the loss associated with FCN and CNN. We first take the FCN case as an example and then we provide the necessary modifications for the proof of the CNN case.

- **The FCN case:** We adapt the analysis in [17] to our setting. Let N_ϵ be the ϵ -covering number of the Euclidean ball $\mathbb{B}(\mathbf{W}^\star, r)$. Here, we omit the subscript FCN of r for simplicity. It is known that $\log N_\epsilon \leq dK \log(3r/\epsilon)$ [43]. Let $\mathcal{W}_\epsilon = \{\mathbf{W}_1, \dots, \mathbf{W}_{N_\epsilon}\}$ be the ϵ -cover set with N_ϵ elements. For any $\mathbf{W} \in \mathbb{B}(\mathbf{W}^\star, r)$, let $j(\mathbf{W}) = \arg\min_{j \in [N_\epsilon]} \|\mathbf{W} - \mathbf{W}_{j(\mathbf{W})}\|_F \leq \epsilon$ for all $\mathbf{W} \in \mathbb{B}(\mathbf{W}^\star, r)$.

For any $\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)$, we have

$$\begin{aligned} \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| &\leq \frac{1}{n} \left\| \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i)] \right\| + \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x})] \right\| \\ &\quad + \left\| \mathbb{E}[\nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x})] - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\|. \end{aligned}$$

Hence, we have

$$\mathbb{P} \left(\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| \geq t \right) \leq \mathbb{P}(A_t) + \mathbb{P}(B_t) + \mathbb{P}(C_t),$$

where the events A_t , B_t and C_t are defined as

$$\begin{aligned} A_t &= \left\{ \sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \frac{1}{n} \left\| \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i)] \right\| \geq \frac{t}{3} \right\}, \\ B_t &= \left\{ \sup_{\mathbf{W} \in \mathcal{W}_\epsilon} \left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\| \geq \frac{t}{3} \right\}, \\ C_t &= \left\{ \sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \left\| \mathbb{E}[\nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x})] - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\| \geq \frac{t}{3} \right\}. \end{aligned}$$

In the sequel, we bound the terms $\mathbb{P}(A_t)$, $\mathbb{P}(B_t)$, and $\mathbb{P}(C_t)$, separately.

1) **Upper bound on $\mathbb{P}(B_t)$.** Before continuing, we state a useful technical lemma, whose proof can be found in [17].

Lemma 6. *Let $\mathbf{M} \in \mathbb{R}^{d \times d}$ be a symmetric $d \times d$ matrix and V_ϵ be an ϵ -cover of unit-Euclidean-norm ball $\mathbb{B}(\mathbf{0}, 1)$, then*

$$\|\mathbf{M}\| \leq \frac{1}{1 - 2\epsilon} \sup_{\mathbf{v} \in V_\epsilon} |\langle \mathbf{v}, \mathbf{M} \mathbf{v} \rangle|. \quad (66)$$

Let $V_{\frac{1}{4}}$ be a $(\frac{1}{4})$ -cover of the ball $\mathbb{B}(\mathbf{0}, 1) = \{\mathbf{W} \in \mathbb{R}^{d \times K} : \|\mathbf{W}\|_F = 1\}$, where $\log |V_{\frac{1}{4}}| \leq dK \log 12$. Following from Lemma 6, we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right\| \leq 2 \sup_{\mathbf{v} \in V_{\frac{1}{4}}} \left| \left\langle \mathbf{v}, \left(\frac{1}{n} \sum_{i=1}^n \nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})] \right) \mathbf{v} \right\rangle \right|.$$

Taking the union bound over \mathcal{W}_ϵ and $V_{\frac{1}{4}}$ yields

$$\mathbb{P}(B_t) \leq \mathbb{P} \left(\sup_{\mathbf{W} \in \mathcal{W}_\epsilon, \mathbf{v} \in V_{\frac{1}{4}}} \left| \frac{1}{n} \sum_{i=1}^n G_i \right| \geq \frac{t}{6} \right) \leq e^{dK(\log \frac{3r}{\epsilon} + \log 12)} \sup_{\mathbf{W} \in \mathcal{W}_\epsilon, \mathbf{v} \in V_{\frac{1}{4}}} \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n G_i \right| \geq \frac{t}{6} \right), \quad (67)$$

where $G_i = \langle \mathbf{v}, (\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})]) \mathbf{v} \rangle$ and $\mathbb{E}[G_i] = 0$. Let $\mathbf{a} = [\mathbf{a}_1^\top, \dots, \mathbf{a}_K^\top] \in \mathbb{R}^{dK}$. Then we can show that $\|G_i\|_{\psi_1}$ is upper bounded, which we summarize as follows, and whose proof is given in Appendix D-F.

Lemma 7. *Suppose the loss is associated with FCN. There exists some constant C such that*

$$\|G_i\|_{\psi_1} \leq C \equiv \tau^2.$$

Applying the Bernstein inequality for sub-exponential random variables [17, Theorem 9] to (67), we have that for fixed $\mathbf{W} \in \mathcal{W}_\epsilon, \mathbf{v} \in V_{\frac{1}{4}}$,

$$\mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{v}, (\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \mathbb{E}[\nabla^2 \ell(\mathbf{W}; \mathbf{x})]) \mathbf{v} \rangle \right| \geq \frac{t}{6} \right) \leq 2 \exp \left(-c \cdot n \cdot \min \left(\frac{t^2}{\tau^4}, \frac{t}{\tau^2} \right) \right), \quad (68)$$

for some universal constant c . As a result, $\mathbb{P}(B_t)$ is upper bounded by

$$2 \exp \left(-c \cdot n \cdot \min \left(\frac{t^2}{\tau^4}, \frac{t}{\tau^2} \right) + dK \log \frac{3r}{\epsilon} + dK \log 12 \right).$$

Thus as long as

$$t > C \cdot \max \left\{ \sqrt{\frac{\tau^4 (dK \log \frac{36r}{\epsilon} + \log \frac{4}{\delta})}{n}}, \frac{\tau^2 (dK \log \frac{36r}{\epsilon} + \log \frac{4}{\delta})}{n} \right\} \quad (69)$$

for some large enough constant C , we have $\mathbb{P}(B_t) \leq \frac{\delta}{2}$.

- 2) **Upper bound on $\mathbb{P}(A_t)$ and $\mathbb{P}(C_t)$.** These two events will be bounded in a similar way. We first present the following useful Lemma, whose proof is provided in Appendix D-H

Lemma 8. *Suppose the loss is associated with FCN. There exists some constant C such that*

$$\mathbb{E} \left[\sup_{\mathbf{W} \neq \mathbf{W}' \in \mathbb{B}(\mathbf{W}^*, r)} \frac{\|\nabla^2 \ell(\mathbf{W}, \mathbf{x}) - \nabla^2 \ell(\mathbf{W}', \mathbf{x})\|}{\|\mathbf{W} - \mathbf{W}'\|_F} \right] \leq C \cdot d\sqrt{K}. \quad (70)$$

Consider the event C_t first. We derive

$$\begin{aligned} & \sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \|\mathbb{E} [\nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x})] - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})]\| \\ & \leq \sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \frac{\|\mathbb{E} [\nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x})] - \mathbb{E} [\nabla^2 \ell(\mathbf{W}; \mathbf{x})]\|}{\|\mathbf{W} - \mathbf{W}_{j(\mathbf{W})}\|_F} \cdot \sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \|\mathbf{W} - \mathbf{W}_{j(\mathbf{W})}\|_F \\ & \leq C \cdot d\sqrt{K} \cdot \epsilon. \end{aligned} \quad (71)$$

Therefore, C_t holds as long as

$$t \geq C \cdot d\sqrt{K} \cdot \epsilon. \quad (72)$$

We can bound the event A_t as below.

$$\begin{aligned} & \mathbb{P} \left(\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \frac{1}{n} \left\| \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i)] \right\| \geq \frac{t}{3} \right) \\ & \leq \frac{3}{t} \mathbb{E} \left[\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \left\| \frac{1}{n} \sum_{i=1}^n [\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i)] \right\| \right] \end{aligned} \quad (73)$$

$$\begin{aligned} & \leq \frac{3}{t} \mathbb{E} \left[\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \|\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i)\| \right] \\ & \leq \frac{3}{t} \mathbb{E} \left[\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \frac{\|\nabla^2 \ell(\mathbf{W}; \mathbf{x}_i) - \nabla^2 \ell(\mathbf{W}_{j(\mathbf{W})}; \mathbf{x}_i)\|}{\|\mathbf{W} - \mathbf{W}_{j(\mathbf{W})}\|_F} \right] \cdot \sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \|\mathbf{W} - \mathbf{W}_{j(\mathbf{W})}\|_F \\ & \leq \frac{C \cdot d\sqrt{K} \cdot \epsilon}{t} \end{aligned} \quad (74)$$

where (73) follows from the Markov inequality. Thus, taking

$$t \geq \frac{6\epsilon \cdot C \cdot d\sqrt{K}}{\delta} \quad (75)$$

ensures that $\mathbb{P}(A_t) \leq \frac{\delta}{2}$.

- 3) **Final step.** Let $\epsilon = \frac{\delta\tau^2}{C \cdot d\sqrt{K} \cdot ndK}$ and $\delta = d^{-10}$. Plugging ϵ and δ into (69) we need

$$t > \tau^2 \cdot \max \left\{ \frac{1}{ndK}, C \cdot \sqrt{\frac{(dK \log(36rnd^{11}K) + \log \frac{4}{\delta})}{n}}, \frac{(dK \log(36rnd^{11}K) + \log \frac{4}{\delta})}{n} \right\}.$$

The middle term can be bounded as

$$\frac{dK \log(36rnd^{11}K) + 10 \log d}{n} \leq \frac{dK \log n}{n} + \frac{dK \log 36r}{n} + \frac{11dK \log dK}{n} + \frac{10 \log d}{n}.$$

If $n \geq C \cdot dK \log(dK)$ for some large enough constant C , the first term $dK \log n$ dominates and is on the order of $dK \log(dK)$. Moreover, it decreases as n increases when $n \geq 3$. Thus we can set

$$t \geq \tau^2 \sqrt{\frac{(dK \log(36rnd^{11}K) + \log \frac{4}{\delta})}{n}} \quad (76)$$

which holds as $t \geq C' \cdot \tau^2 \sqrt{\frac{dK \log n}{n}}$ for some constant C' . By setting $t := C\tau^2 \sqrt{\frac{dK \log n}{n}}$ for sufficiently large C , as long as $n \geq C' \cdot dK \log dK$,

$$\mathbb{P} \left(\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r)} \|\nabla^2 f_n(\mathbf{W}) - \nabla^2 f(\mathbf{W})\| \geq C\tau^2 \sqrt{\frac{dK \log n}{n}} \right) \leq d^{-10}. \quad (77)$$

- **The CNN case:** If the loss is associated with CNN, we redefine G_i as $G_i = \langle \mathbf{v}, (\nabla^2 \ell(\mathbf{w}; \mathbf{x}_i) - \mathbb{E} [\nabla^2 \ell(\mathbf{w}; \mathbf{x})]) \mathbf{v} \rangle$ and we show the following Lemmas whose proof is given in Appendix D-G and Appendix D-I.

Lemma 9. Suppose the loss is associated CNN. There exists some constant C such that

$$\|G_i\|_{\psi_1} \leq C \cdot K^2 := \tau^2. \quad (78)$$

Lemma 10. Suppose the loss is associated with CNN. There exists some constant C such that

$$\mathbb{E} \left[\sup_{\mathbf{W} \neq \mathbf{W}' \in \mathbb{B}(\mathbf{W}^*, r)} \frac{\|\nabla^2 \ell(\mathbf{W}, \mathbf{x}) - \nabla^2 \ell(\mathbf{W}', \mathbf{x})\|}{\|\mathbf{W} - \mathbf{W}'\|_F} \right] \leq C \cdot d\sqrt{K}. \quad (79)$$

Following argument similar to the proof of Lemma 3, we can obtain the following concentration inequality:

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r)} \|\nabla^2 f_n(\mathbf{w}) - \nabla^2 f(\mathbf{w})\| \leq C \cdot K^2 \sqrt{\frac{\frac{d}{K} \cdot \log n}{n}}, \quad (80)$$

holds with probability at least $1 - d^{-10}$, as long as the sample complexity $n \geq C \cdot \frac{d}{K} \log\left(\frac{d}{K}\right)$.

D. Proof of Lemma 4

In order to proceed we need the following Lemma 11 whose proof is given in Appendix D-J.

Lemma 11. Suppose the loss is associated with FCN. Let \mathbf{u} be a fixed unit norm vector $\mathbf{u} = [\mathbf{u}_1^\top, \dots, \mathbf{u}_K^\top] \in \mathbb{R}^{dK}$ with $\|\mathbf{u}\|_2 = 1$. Then we have

$$\|\mathbf{u}^\top \nabla \ell(\mathbf{W}; \mathbf{x})\|_{\psi_2} \leq \sqrt{K}.$$

Suppose the loss is associated with CNN. Let \mathbf{u} be a fixed unit norm vector $\mathbf{u} \in \mathbb{R}^m$ with $\|\mathbf{u}\|_2 = 1$. Then

$$\|\langle \mathbf{u}, \nabla \ell(\mathbf{w}) \rangle\|_{\psi_2} \leq C \cdot K.$$

Following argument (details omitted) similar to the proof of Lemma 3, and applies Lemma 11, for the loss associated with FCN, we can get the following concentration inequality

$$\sup_{\mathbf{W} \in \mathbb{B}(\mathbf{W}^*, r_{\text{FCN}})} \|\nabla f_n(\mathbf{W}) - \nabla f(\mathbf{W})\|_2 \leq C \cdot \sqrt{\frac{d\sqrt{K} \log n}{n}} \quad (81)$$

with probability at least $1 - d^{-10}$, as long as the sample size $n \geq C \cdot dK \log(dK)$. For the loss associated with CNN, we obtain

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}^*, r_{\text{CNN}})} \|\nabla f_n(\mathbf{w}) - \nabla f(\mathbf{w})\| \leq C \cdot \sqrt{K} \sqrt{\frac{\frac{d}{K} \log n}{n}} = C \cdot \sqrt{\frac{d \log n}{n}}, \quad (82)$$

with probability at least $1 - d^{-10}$ as long as $n \geq C \cdot \frac{d}{K} \log\left(\frac{d}{K}\right)$.

E. Proof of Lemma 5

We take the first term in (35) as an example, since the second term follows exactly in the same way. We first derive

$$\mathbb{E} \left[\left(\frac{1}{K} \sum_{i=1}^K \phi(\mathbf{w}_i^\top \mathbf{x}) \right)^{-z} \right] \leq \mathbb{E} \left[\frac{1}{K} \sum_{i=1}^K (\phi(\mathbf{w}_i^\top \mathbf{x}))^{-z} \right], \quad (83)$$

which follows from the fact that $f(x) = x^{-z}$ is convex for $x > 0$ and $z \geq 1$. Further since $\frac{1}{\phi(x)} = 1 + e^{-x}$, and $g = \mathbf{w}_i^\top \mathbf{x} \sim \mathcal{N}(0, \sigma_i^2 = \|\mathbf{w}_i\|_2^2)$, we can exactly calculate the summands in the above equation as follows:

$$\mathbb{E} [\phi(g)^{-z}] = \mathbb{E} \left[\sum_{l=0}^z \binom{z}{l} e^{-lg} \right] = \sum_{l=0}^z \binom{z}{l} e^{\left(\frac{\sigma_i^2 l^2}{2}\right)},$$

where we use the fact that g is a Gaussian random variable. Hence, we conclude that for $t = \max(\|\mathbf{w}_1\|_2, \dots, \|\mathbf{w}_K\|_2)$ and $p \geq 1$,

$$\mathbb{E} \left[\left(\frac{1}{\frac{1}{K} \sum_{i=1}^K \phi(\mathbf{w}_i^\top \mathbf{x})} \right)^z \right] \leq C \cdot e^{t^2}, \quad (84)$$

holds for some constant C depending on z .

F. Proof of Lemma 7

The sub-exponential norm of G_i can be bounded as

$$\|G_i\|_{\psi_1} \leq \|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{W}; z) \mathbf{u} \rangle\|_{\psi_1} + \|\nabla^2 f(\mathbf{W}; z)\|,$$

where $\|\nabla^2 f(\mathbf{W}; z)\|$ is upper bounded by C due to Lemma 2. Denote the (j, l) -th block of $\nabla^2 \ell(\mathbf{W}; z)$ as $\xi_{j,l} \cdot \mathbf{x} \mathbf{x}^\top$. We can derive

$$\|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{W}; z) \mathbf{u} \rangle\|_{\psi_1} \leq \sum_{j=1}^K \sum_{l=1}^K \|\xi_{j,l} \cdot \mathbf{u}_j^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}_l\|_{\psi_1} \leq \sum_{j=1}^K \sum_{l=1}^K \sup_{t \geq 1} t^{-1} \left(\mathbb{E} |\xi_{j,l} \cdot \mathbf{u}_j^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}_l|^t \right)^{\frac{1}{t}}. \quad (85)$$

Next we show that $\xi_{j,l}$ is upper bounded by some constant for all j and l .

- For $j \neq l$,

$$|\xi_{j,l}| = \left| \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x}) \cdot (H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W}))}{(H(\mathbf{W})(1 - H(\mathbf{W})))^2} \right| = \begin{cases} \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{(1 - H(\mathbf{W}))^2} & y = 0 \\ \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{H(\mathbf{W})^2} & y = 1 \end{cases}. \quad (86)$$

Moreover,

$$\frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{(1 - H(\mathbf{W}))^2} \leq \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{(1 - \phi(\mathbf{w}_j^\top \mathbf{x}))(1 - \phi(\mathbf{w}_l^\top \mathbf{x}))} \leq \phi(\mathbf{w}_j^\top \mathbf{x}) \phi(\mathbf{w}_l^\top \mathbf{x}) \leq 1, \quad (87)$$

where the first inequality holds due to the following fact,

$$(1 - H(\mathbf{W}))^2 = \left(1 - \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_j^\top \mathbf{x}) \right)^2 \geq \frac{1}{K^2} (1 - \phi(\mathbf{w}_j^\top \mathbf{x}))(1 - \phi(\mathbf{w}_l^\top \mathbf{x})),$$

the second inequality follows because $\phi(x)(1 - \phi(x)) = \phi'(x)$. Similarly, we can show that

$$\frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_l^\top \mathbf{x})}{H(\mathbf{W})^2} \leq 1. \quad (88)$$

Thus for $j \neq l$, $|\xi_{j,l}| \leq 1$ holds.

- For $j = l$,

$$|\xi_{j,j}| \leq \left| \frac{1}{K^2} \frac{\phi'(\mathbf{w}_j^\top \mathbf{x}) \phi'(\mathbf{w}_j^\top \mathbf{x}) \cdot (H(\mathbf{W})^2 + y - 2y \cdot H(\mathbf{W}))}{(H(\mathbf{W})(1 - H(\mathbf{W})))^2} \right| + \left| \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x}) (y - H(\mathbf{W}))}{H(\mathbf{W})(1 - H(\mathbf{W}))} \right|. \quad (89)$$

For the second term in the above equation, we have

$$\left| \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x}) (y - H(\mathbf{W}))}{H(\mathbf{W})(1 - H(\mathbf{W}))} \right| = \begin{cases} \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{(1 - H(\mathbf{W}))} \leq 1 & y = 0 \\ \frac{1}{K} \frac{\phi''(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})} \leq 1 & y = 1 \end{cases},$$

which follows from the fact that the second derivative is $\phi''(x) = \phi(x)(1 - \phi(x))(1 - 2\phi(x))$, the absolute value of which can be upper bounded by $\phi(x)$ or $1 - \phi(x)$.

Hence,

$$\begin{aligned} \|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{W}; z) \mathbf{u} \rangle\|_{\psi_1} &\leq C \cdot \sum_{j=1}^K \sum_{l=1}^K \sup_{t \geq 1} t^{-1} \left(\sqrt{\mathbb{E}[(\mathbf{u}_j^\top \mathbf{x})^{2t}]} \cdot \sqrt{\mathbb{E}[(\mathbf{u}_l^\top \mathbf{x})^{2t}]} \right)^{\frac{1}{t}} \\ &\leq C \cdot \sum_{j=1}^K \sum_{l=1}^K \|\mathbf{u}_j\|_2 \|\mathbf{u}_l\|_2 \cdot \sup_{t \geq 1} t^{-1} ((2t-1)!!)^{\frac{1}{t}} \\ &\leq C \equiv \tau^2, \end{aligned} \quad (90)$$

where the last inequality holds because

$$\begin{aligned} \sup_{t \geq 1} t^{-1} ((2t-1)!!)^{\frac{1}{t}} &\leq \sup_{t \geq 1} t^{-1} ((2t)^t)^{\frac{1}{t}} \leq 2, \\ \sum_{j=1}^K \sum_{l=1}^K \|\mathbf{u}_j\|_2 \|\mathbf{u}_l\|_2 &\leq \sum_{j=1}^K \sum_{l=1}^K \frac{\|\mathbf{u}_j\|_2^2 + \|\mathbf{u}_l\|_2^2}{2} = \frac{1}{2}. \end{aligned} \quad (91)$$

Thus, we conclude

$$\|G_i\|_{\psi_1} \leq C \equiv \tau^2.$$

G. Proof of Lemma 9

Again the sub-exponential norm of G_i can be bounded as

$$\|G_i\|_{\psi_1} \leq \|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}; z) \mathbf{u} \rangle\|_{\psi_1} + \|\nabla^2 f(\mathbf{w}; z)\|,$$

where $\|\nabla^2 f(\mathbf{W}; z)\|$ is upper bounded by $C \cdot K$ due to Lemma 2. Applying the triangle inequality, the sub-exponential norm of $\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}) \mathbf{u} \rangle$ can be bounded as

$$\|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}) \mathbf{u} \rangle\|_{\psi_1} \leq \sum_{j \neq l} \|g_{j,l}(\mathbf{w}) \mathbf{u}^\top \mathbf{x}^{(j)} \mathbf{u}^\top \mathbf{x}^{(l)}\|_{\psi_1} + \sum_{j=l} \|g_{j,l}(\mathbf{w}) \mathbf{u}^\top \mathbf{x}^{(j)} \mathbf{u}^\top \mathbf{x}^{(l)}\|_{\psi_1}. \quad (92)$$

Hence, we have

$$\left| \frac{1}{K^2} \frac{H(\mathbf{w})^2 + y - 2y \cdot H(\mathbf{w})}{(H(\mathbf{w})(1 - H(\mathbf{w})))^2} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)}) \right| = \begin{cases} \frac{1}{K^2} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)})}{H(\mathbf{w})^2} \leq 1 & y = 1 \\ \frac{1}{K^2} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \phi'(\mathbf{w}^\top \mathbf{x}^{(l)})}{(1 - H(\mathbf{w}))^2} \leq 1 & y = 0 \end{cases},$$

$$\left| \frac{1}{K} \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \right| = \begin{cases} \frac{1}{K} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{H(\mathbf{w})} \leq 1 & y = 1 \\ \frac{1}{K} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{1 - H(\mathbf{w})} \leq 1 & y = 0 \end{cases}.$$

Plugging it back to (92), we obtain

$$\|\langle \mathbf{u}, \nabla^2 \ell(\mathbf{w}) \mathbf{u} \rangle\|_{\psi_1} \leq \sum_{j \neq l} \left\| \left(\mathbf{u}^\top \mathbf{x}^{(j)} \right) \left(\mathbf{u}^\top \mathbf{x}^{(l)} \right) \right\|_{\psi_1} + \sum_{j=1}^K \left\| \left(\mathbf{u}^\top \mathbf{x}^{(j)} \right)^2 \right\|_{\psi_1} \leq C \cdot K^2. \quad (93)$$

H. Proof of Lemma 8

As noted before, we can write the (j, l) -th block of $\nabla^2 \ell(\mathbf{W}; \mathbf{z})$ as $\xi_{j,l}(\mathbf{W}) \mathbf{x} \mathbf{x}^\top$. Then we can obtain the following bound,

$$\|\nabla^2 \ell(\mathbf{W}; z) - \nabla^2 \ell(\mathbf{W}'; z)\| \leq \sum_{j=1}^K \sum_{l=1}^K |\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}')| \cdot \|\mathbf{x} \mathbf{x}^\top\|. \quad (94)$$

Using the same method as shown in the proof of Lemma 1, we can upper bound $|\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}')|$ as

$$|\xi_{j,l}(\mathbf{W}) - \xi_{j,l}(\mathbf{W}')| \leq \left(\max_k |T_{j,l,k}| \right) \cdot \|\mathbf{x}\|_2 \cdot \sqrt{K} \cdot \|\mathbf{W} - \mathbf{W}'\|_F,$$

where following from (32),

$$|T_{j,l,k}| \leq \max \left\{ \frac{2}{K^3} \frac{1}{H(\mathbf{W})^3}, \frac{1}{K^2} \frac{1}{H(\mathbf{W})^2}, \frac{2}{K^3} \frac{1}{(1 - H(\mathbf{W}))^3}, \frac{1}{K^2} \frac{1}{(1 - H(\mathbf{W}))^2} \right\}. \quad (95)$$

And thus, if $\|\mathbf{W} - \mathbf{W}'\|_F \leq 0.7$ we have

$$\mathbb{E} \left[\sup_{\mathbf{W} \neq \mathbf{W}'} \frac{\|\nabla^2 \ell(\mathbf{W}) - \nabla^2 \ell(\mathbf{W}')\|}{\|\mathbf{W} - \mathbf{W}'\|_F} \right] \leq \sqrt{K} \cdot K^2 \cdot \mathbb{E} \left[\left(\max_{j,l,k} |T_{j,l,k}| \right) \cdot \|\mathbf{x}\|_2 \cdot \|\mathbf{x} \mathbf{x}^\top\| \right] \leq C \cdot d\sqrt{K}. \quad (96)$$

Thus we only need to set $J^* \geq C \cdot d\sqrt{K}$ for some large enough C .

I. Proof of Lemma 10

Following from (13) we can write

$$\|\nabla^2 \ell(\mathbf{w}) - \nabla^2 \ell(\mathbf{w}')\| \leq \sum_{j=1}^K \sum_{l=1}^K |g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}')| \cdot \|\mathbf{x}^{(j)} \mathbf{x}^{(l)\top}\|. \quad (97)$$

Similarly, the analysis in the proof of Lemma 1 implies that

$$|g_{j,l}(\mathbf{w}) - g_{j,l}(\mathbf{w}')| \leq \left(\max_k |S_{j,l,k}| \right) \cdot \sqrt{K} \|\mathbf{x}\|_2 \cdot \|\mathbf{w} - \mathbf{w}'\|_2, \quad (98)$$

where we upper-bound $S_{j,l,k}$ in (43) as

$$|S_{j,l,k}| \leq \begin{cases} \max \left\{ \frac{1}{K^2} \frac{1}{(1 - H(\mathbf{w}))^3}, \frac{1}{K^2} \frac{1}{(H(\mathbf{w}))^3} \right\} & j \neq l \\ \max \left\{ \frac{1}{K} \frac{1}{(1 - H(\mathbf{w}))^2}, \frac{1}{K} \frac{1}{(H(\mathbf{w}))^2} \right\} & j = l \end{cases}. \quad (99)$$

Hence, if $\|\mathbf{w} - \mathbf{w}'\|_2 \leq 0.7$, we have

$$\mathbb{E} \left[\sup_{\mathbf{w} \neq \mathbf{w}'} \frac{\|\nabla^2 \ell(\mathbf{w}) - \nabla^2 \ell(\mathbf{w}')\|}{\|\mathbf{w} - \mathbf{w}'\|_F} \right] \leq \sqrt{K} \cdot \sum_{j=1}^K \sum_{l=1}^K \mathbb{E} \left[\left(\max_k |S_{j,l,k}| \right) \cdot \|\mathbf{x}\|_2 \cdot \|\mathbf{x}^{(j)} \mathbf{x}^{(l)\top}\| \right] \leq C \cdot d\sqrt{K}. \quad (100)$$

Thus, in this case we can set $J^* \geq C \cdot d\sqrt{K}$ as well.

J. Proof of Lemma 11

- **The FCN case:** Following from (10), we have

$$\langle \nabla \ell(\mathbf{W}), \mathbf{u} \rangle = \frac{1}{K} \sum_{j=1}^K \left(\frac{(y - H(\mathbf{W})) \cdot \phi'(\mathbf{w}_j^\top \mathbf{x})}{H(\mathbf{W})(1 - H(\mathbf{W}))} \right) (\mathbf{u}_j^\top \mathbf{x}),$$

and by definition, we can upper-bound the sub-Gaussian norm as

$$\|\langle \nabla \ell(\mathbf{W}), \mathbf{u} \rangle\|_{\psi_2} \leq \begin{cases} \frac{1}{K} \sum_{j=1}^K \left\| \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})}{(1 - \frac{1}{K} \sum_{l=1}^K \phi(\mathbf{w}_l^\top \mathbf{x}))} \mathbf{u}_j^\top \mathbf{x} \right\|_{\psi_2} & y = 0 \\ \frac{1}{K} \sum_{j=1}^K \left\| \frac{\phi'(\mathbf{w}_j^\top \mathbf{x})}{\frac{1}{K} \sum_{l=1}^K \phi(\mathbf{w}_l^\top \mathbf{x})} \mathbf{u}_j^\top \mathbf{x} \right\|_{\psi_2} & y = 1 \end{cases}.$$

Thus we conclude that

$$\|\langle \nabla \ell(\mathbf{W}), \mathbf{u} \rangle\|_{\psi_2} \leq \sum_{j=1}^K \|\mathbf{u}_j\|_2 \leq \sqrt{K}, \quad (101)$$

and the directional gradient is \sqrt{K} -sub-Gaussian.

- **The CNN case:** Following from (12), we have

$$\langle \nabla \ell(\mathbf{w}), \mathbf{u} \rangle = - \sum_{j=1}^K \frac{1}{K} \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \cdot (\mathbf{u}^\top \mathbf{x}^{(j)}),$$

where

$$\left| \phi'(\mathbf{w}^\top \mathbf{x}^{(j)}) \frac{y - H(\mathbf{w})}{H(\mathbf{w})(1 - H(\mathbf{w}))} \right| = \begin{cases} \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{\sum_{j=1}^K \frac{1}{K} \phi(\mathbf{w}^\top \mathbf{x}^{(j)})} \leq K & y = 1 \\ \frac{\phi'(\mathbf{w}^\top \mathbf{x}^{(j)})}{\sum_{j=1}^K \frac{1}{K} (1 - \phi(\mathbf{w}^\top \mathbf{x}^{(j)}))} \leq K & y = 0 \end{cases}.$$

Then the sub-Gaussian norm of $\langle \nabla \ell(\mathbf{w}), \mathbf{u} \rangle$ is upper bounded as

$$\|\langle \nabla \ell(\mathbf{w}), \mathbf{u} \rangle\|_{\psi_2} \leq K \cdot \frac{1}{K} \sum_{j=1}^K \|\mathbf{u}^\top \mathbf{x}^{(j)}\|_{\psi_2} \leq C \cdot K. \quad (102)$$

Hence, the directional gradient is K -sub-Gaussian.

APPENDIX E PROOF OF THEOREM 3

The proof contains two parts. Part (a) proves that the estimation of the direction of \mathbf{W}^* is sufficiently accurate, which follows from the arguments similar to those in [15] and is only briefly summarized below. Part (b) is different, where we do not require the homogeneous condition for the activation function, and instead, our proof is based on a mild condition in Assumption 2. We detail our proof in part (b).

(a) In order to estimate the direction of each \mathbf{w}_i for $i = 1, \dots, K$, [15] showed that for the regression problem, if the sample size $n \geq d \text{poly}(K, \kappa, \zeta, \log d)$, where $\zeta > 1$ is any constant, then

$$\|\bar{\mathbf{w}}_i^* - s_i \mathbf{V} \hat{\mathbf{u}}_i\| \leq \epsilon \text{poly}(K, \kappa) \quad (103)$$

holds with probability at least $1 - d^{-\Omega(\zeta)}$. Such a result also holds for the classification problem with only slight difference in the proof as we describe as follows. The main idea of the proof is to bound the estimation error of \mathbf{P}_2 and \mathbf{R}_3 via Bernstein inequality. For the regression problem, Bernstein inequality was applied to terms associated with each neuron individually, and the bounds were then put together via the triangle inequality in [15]. However, for the classification problem here, we apply Bernstein inequality to the terms associated with all neurons together. Another difference is that the label y_i of the classification

model is bounded by nature, whereas the output y_i in the regression model needs to be upper-bounded via homogeneously bounded conditions of the activation function. A reader can refer to [15] for the details of the proof for this part.

(b) In order to estimate $\|\mathbf{w}_i\|$ for $i = 1, \dots, K$, we provide a different proof from [15], which does not require the homogeneous condition on the activation function, but assumes a more relaxed condition in Assumption 2.

We define a quantity Q_1 as follows:

$$Q_1 = M_{l_1}(\mathbf{I}, \underbrace{\boldsymbol{\alpha}, \dots, \boldsymbol{\alpha}}_{(l_1-1)}), \quad (104)$$

where l_1 is the first non-zero index such that $M_{l_1} \neq 0$. For example, if $l_1 = 3$, then Q_1 takes the following form

$$Q_1 = M_3(\mathbf{I}, \boldsymbol{\alpha}, \boldsymbol{\alpha}) = \frac{1}{K} \sum_{i=1}^K m_{3,i}(\|\mathbf{w}_i^*\|) (\boldsymbol{\alpha}^\top \bar{\mathbf{w}}_i^*)^2 \bar{\mathbf{w}}_i^*, \quad (105)$$

where $\bar{\mathbf{w}} = \mathbf{w}/\|\mathbf{w}\|$ and by definition

$$m_{3,i}(\|\mathbf{w}_i^*\|) = \mathbb{E} [\phi(\|\mathbf{w}_i^*\| \cdot z) z^3] - 3\mathbb{E} [\phi(\|\mathbf{w}_i^*\| \cdot z) z]. \quad (106)$$

Clearly, Q_1 has information of $\|\mathbf{w}_i^*\|$, which can be estimated by solving the following optimization problem:

$$\beta^* = \operatorname{argmin}_{\beta \in \mathbb{R}^K} \left\| \frac{1}{K} \sum_{i=1}^K \beta_i s_i \bar{\mathbf{w}}_i^* - Q_1 \right\|, \quad (107)$$

where each entry of the solution takes the form

$$\beta_i^* = s_i^3 m_{3,i}(\|\mathbf{w}_i^*\|) (\boldsymbol{\alpha}^\top s_i \bar{\mathbf{w}}_i^*)^2. \quad (108)$$

In the initialization, we substitute \hat{Q}_1 (estimated from training data) for Q_1 , $\mathbf{V}\hat{\mathbf{u}}_i$ (estimated in part (a)) for $s_i \bar{\mathbf{w}}_i^*$ into (107), and obtain an estimate $\hat{\beta}$ of β^* . We then substitute $\hat{\beta}$ for β^* and $\mathbf{V}\hat{\mathbf{u}}_i$ for $s_i \bar{\mathbf{w}}_i^*$ into (108) to obtain an estimate \hat{a}_i of $\|\mathbf{w}_i^*\|$ via the following equation

$$\hat{\beta}_i = s_i^3 m_{3,i}(\hat{a}_i) (\boldsymbol{\alpha}^\top \mathbf{V}\hat{\mathbf{u}}_i)^2. \quad (109)$$

Furthermore, since $m_{l_1,i}(x)$ has fixed sign for $x > 0$ and for $l_1 \geq 1$, s_i can be estimated correctly from the sign of $\hat{\beta}_i$ for $i = 1, \dots, K$.

For notational simplicity, let $\beta_{1,i}^* := \frac{\beta_i^*}{s_i^3 (\boldsymbol{\alpha}^\top s_i \bar{\mathbf{w}}_i^*)^2}$ and $\hat{\beta}_{1,i} := \frac{\hat{\beta}_i}{s_i^3 (\boldsymbol{\alpha}^\top \mathbf{V}\hat{\mathbf{u}}_i)^2}$, and then (108) and (109) become

$$\hat{\beta}_{1,i} = m_{3,i}(\hat{a}_i), \quad \beta_{1,i}^* = m_{3,i}(\|\mathbf{w}_i^*\|). \quad (110)$$

By Assumption 2 and (108), there exists a constant $\delta' > 0$ such that the inverse function $g(\cdot)$ of $m_{3,1}(\cdot)$ has upper-bounded derivative in the interval $(\beta_{1,i}^* - \delta', \beta_{1,i}^* + \delta')$, i.e., $|g'(x)| < \Gamma$ for a constant Γ . By employing the result in [15], if the sample size $n \geq d \text{poly}(K, \kappa, t, \log d)$, then \hat{Q}_1 and Q_1 , $\mathbf{V}\hat{\mathbf{u}}_i$ and $s_i \bar{\mathbf{w}}_i^*$ can be arbitrarily close so that $|\beta_{1,i}^* - \hat{\beta}_{1,i}| < \min\{\delta', \frac{r}{\sqrt{K}\Gamma}\}$.

Thus, by (110) and the mean value theorem, we obtain

$$|\hat{a}_i - \|\mathbf{w}_i^*\|| = |g'(\xi)| |\beta_{1,i}^* - \hat{\beta}_{1,i}|, \quad (111)$$

where ξ is between $\beta_{1,i}^*$ and $\hat{\beta}_{1,i}$, and hence $|g'(\xi)| < \Gamma$. Therefore, $|\hat{a}_i - \|\mathbf{w}_i^*\|| \leq \frac{r}{\sqrt{K}}$, which is the desired result.