

Improved Linear Convergence of Training CNNs with Generalizability Guarantees: A One-hidden-layer Case

Shuai Zhang, *Student Member, IEEE*, Meng Wang, *Member, IEEE*,
Jinjun Xiong, *Member, IEEE*, Sijia Liu, *Member, IEEE*, Pin-Yu Chen, *Member, IEEE*

Abstract—We analyze the learning problem of one-hidden-layer non-overlapping convolutional neural networks with the rectified linear unit (ReLU) activation function from the perspective of model estimation. The training outputs are assumed to be generated by the neural network with the unknown ground-truth parameters plus some additive noise, and the objective is to estimate the model parameters by minimizing a non-convex squared loss function of the training data. Assuming that the training set contains a finite number of samples generated from the Gaussian distribution, we prove that the accelerated gradient descent algorithm with a proper initialization converges to the ground-truth parameters (up to the noise level) with a linear rate even though the learning problem is non-convex. Moreover, the convergence rate is proved to be faster than the vanilla gradient descent. The initialization can be achieved by the existing tensor initialization method. In contrast to the existing works that assume an infinite number of samples, we theoretically establish the sample complexity of the required number of training samples. Although the neural network considered here is not deep, this is the first work to show that accelerated gradient descent algorithms can find the global optimizer of the non-convex learning problem of neural networks. This is also the first work that characterizes the sample complexity of gradient-based methods in learning convolutional neural networks with the non-smooth ReLU activation function. This work also provides the tightest bound so far of the estimation error with respect to the output noise.

Index Terms—convolutional neural networks, generalizability, global optimality, accelerated gradient descent, linear convergence

I. INTRODUCTION

Neural networks, especially convolutional neural networks (CNNs), have demonstrated superior performance in machine learning for image classification [16] and recognition [19], natural language processing [5], and strategic game program [33]. Compared with fully connected neural networks, CNNs require fewer coefficients and can better capture local features [20], and thus perform well in applications like image and video processing.

Learning a neural network needs to find appropriate parameters for the hidden layers using the training data and is achieved by minimizing a non-convex empirical loss function over the choices of the model parameters. The non-convex learning problem is usually solved by a first-order

gradient descent (GD) algorithm. The convergence to the global optimal, however, is not guaranteed naturally due to the existence of spurious local minima. Another major hurdle to the widespread acceptance of deep learning is the lack of analytical performance guarantees about whether the parameters learned from the training data perform well on the testing data, i.e., the generalizability of the learned model. A learned model generalizes well to the testing data provided that it is a global minimizer of the population loss function, which takes the expectation over the distribution of testing samples. Since the distribution is unknown, one minimizes the empirical loss function of the training data assuming that the training data are drawn from the same distribution. Moreover, a large number of training samples are required to obtain a network model with powerful feature representation capability [6], while the method may perform poorly when the number of training samples is small [4]. The theoretical characterization of the required size of the training data for a given network architecture is vastly unavailable.

To analyze the learning performance, one line of research focuses on the over-parameterized case that the number of parameters in the neural network is larger than the number of training samples [1], [2], [14], [15], [24], [27], [30], [34]. In particular, the optimization problem has no spurious local minima [24], [34], [43], and GD methods can indeed find the global minimum of the empirical loss function. Nevertheless, the over-parameterized models may experience overfitting issues in practice [42], [43]. Moreover, when over-parameterized, there is no guarantee by VC-dimension learning theory that the empirical loss function is close to the population loss and thus, the generalizability of the learned model to the testing data is unknown. Ref. [1] develops a new analysis tool to explore the generalizability under over-parameterization assumption. The convergence rate provided by [1] is sub-linear, and the sizes of neural networks increase as a polynomial function of the inverse of the desired testing error, which implies a high computational cost. Moreover, the training error and the generalization error are analyzed separately, and it is not clear if both a small training error and a small generalization error can be achieved simultaneously.

Refs. [18], [39] study the convergence to the global optimal for shallow neural networks when the data is linearly separable. Assuming the Rectified Linear Unit (ReLU) activation function and the hinge loss function, ref. [39] can detect all the spurious local minima and saddle points, and the

The first two authors are with the Dept. of Electrical, Computer, and Systems Engineering, Rensselaer Polytechnic Institute, 110 8th Street, Troy, NY, 12180. Email: {zhangs21, wangm7}@rpi.edu. The other authors are with IBM Thomas J. Watson Research Center, jinjun@us.ibm.com, {Sijia.Liu, Pin-Yu.Chen}@ibm.com.

generalization error of the learned model approaches zero when the number of samples goes to infinite. However, if the data are linearly separable, simple algorithms, such as Perceptron [29], can find a classifier in finite steps. Moreover, the detection method of the spurious local minima and saddle points in [39] only apply the ReLU activation function and hinge loss function, and the method does not extend to other activation functions and loss functions.

One recent line of research assumes the existence of a ground-truth model that maps the input data to the output data. Then the set of the ground-truth model parameters is a global minimizer of both the population and the empirical risk functions. The learning problem can be viewed as a model estimation problem. If the parameters are accurately estimated, the generalizability to the testing data is guaranteed. This paper follows this line of research.

To simplify the analysis, one standard trick in this line of research is to assume that the number of input data is infinite so that the empirical loss function is simplified to the population loss function that is easier to analyze. Most existing theoretical results are centered on one-hidden-layer shallow neural networks as the analyses quickly become intractable when the number of layers increases. The input data are usually assumed to follow the Gaussian distribution [32] or some distributions that are rationally invariant [7]. Refs. [3], [8], [36] analyze the landscape of the population loss function of a simple one-hidden-layer neural network with only one or two nodes and show that there exists a considerably large convex region near the global optimum. Then a random initial point lies in this region with a constant probability, and gradient descent algorithms converge to the global minimum. This result does not easily generalize to general neural networks as spurious local minima are fairly common for neural networks with even one hidden layer but multiple nodes [31]. Some works [10], [22], [23] seek to obtain a good optimization landscape through changing the neural network structure. Ref. [22] adds an identity mapping after the hidden layer to improve the convergence of GD algorithm. An additional regularization term is added to the loss function in [10] such that the ground-truth parameters are still close to the global minimum, and spurious local minima are excluded. An exponential node is added in each layer of an arbitrary neural network such that all local minima are global minima [23]. Another work [12] developed a new iterative algorithm named Convotron, which applied a modified gradient descent update in each iteration and does not require initialization.

In the practical case of a finite number of samples, the nice properties of the population loss function do not directly generalize to the empirical loss function. Some recent works study the training performance with a finite number of samples [9], [40], [44]–[46]. If the number of samples is greater than a certain threshold, referred to the sample complexity, ref. [40] shows that the iterates converge to the ground-truth parameters for one-hidden-layer neural networks. However, the sample complexity is sub-optimal as it is a high order polynomial with respect to the dimension of the input data. With the tensor initialization method [46], GD algorithms are proved to converge to the ground-truth parameters linearly in one-

hidden-layer neural networks, and the sample complexity is nearly linear in the dimension of the input data [9], [44]–[46]. However, the analyses in [9], [45], [46] are limited to smooth activation functions and exclude the widely used non-smooth activation function, ReLU. Among them, only ref. [44] studies the ReLU activation function, but focuses on fully connected neural networks. Ref. [44] can only guarantee the convergence to the ground truth up to some nonzero estimation error, even when the data are noiseless.

The majority of the above works assume that data are noiseless, which may not be realistic in practice. Only [10] and [44] consider the cases that the output data contain additive noise that is independent of the input. The noise is assumed to be zero mean in [10], and the authors analyze the stochastic gradient descent through expectation. Thus, the noise does not affect their analyses and results. The result in [44] guarantees the convergence of GD provided that the initialization is sufficiently close to the ground-truth parameters, but no discussion is provided about whether the initialization in [44] satisfies this assumption or not.

All the aforementioned works analyze standard GD algorithms. It is well known that Accelerated Gradient Descent (AGD) methods such as Nesterov’s accelerated gradient (NAG) method [26] and Heavy ball method [28] converge faster than vanilla GD. However, the analyses for GD do not generalize directly to AGD because of the additional momentum term introduced in AGD. Only refs. [35] and [41] explore the numerical performance of AGD in neural networks. No theoretical analysis of AGD is reported in [35]. Ref. [41] analyzes AGD from a general optimization perspective, and it is not clear whether the neural network learning problem satisfies the assumptions in [41].

This paper provides novel contributions to the theoretical analyses of neural networks in three aspects. First, this paper provides the first theoretical analysis of AGD methods in learning neural networks. We prove analytically that the AGD method can converge to the ground-truth parameters linearly, and its convergence rate is faster than vanilla GD. Second, it is the first work that explicitly proves the convergence of the proposed learning algorithm to the ground-truth parameters (or nearby) when the data contain noise. We characterize the relationship between the learning accuracy and the noise level quantitatively. Our error bound is much tighter than that in [44], and [44] makes assumptions about the initialization without any justification. In the special case of noiseless data, our parameter estimation is exact, while the method in [44] is not. Third, it provides the first tight generalizability analysis of the widely used convolutional neural networks with the nonsmooth ReLU activation functions. Specifically, we prove that for one-hidden-layer non-overlapping convolutional neural networks, if initialized using the tensor method, and the number of samples exceeds our characterized sample complexity, both GD and ADG converge to a global minimum linearly up to the noise level. Our sample complexity is order-wise optimal with respect to the dimension of the node parameters. Our estimation error bound of the ground-truth parameters is much tighter than a direct application of the existing results for fully connected neural networks such as [44] to CNN.

The rest of this paper is organized as follows. Section II introduces the problem formulation. The algorithm and major theorems are presented in Section III. Section IV shows the simulation results, and Section V concludes the paper. All the proofs are in the Appendix.

Notation: Vectors are bold lowercase, matrices and tensors are bold uppercase, and scalars are in normal font. For instance, \mathbf{Z} is a matrix, and \mathbf{z} is a vector. z_i denotes the i -th entry of \mathbf{z} , and Z_{ij} denotes the (i, j) -th entry of \mathbf{Z} . \mathbf{I} and \mathbf{e}_i denote the identity matrix and the i -th standard basis vector. \mathbf{Z}^T denotes the transpose of \mathbf{Z} , similarly for \mathbf{z}^T . $\|\mathbf{z}\|$ denotes the ℓ_2 -norm of a vector \mathbf{z} , and $\|\mathbf{Z}\|_2$ and $\|\mathbf{Z}\|_F$ denotes the spectral norm and Frobenius norm of a matrix \mathbf{Z} , respectively. We use $\sigma_i(\mathbf{Z})$ to denote the i -th largest singular value of \mathbf{Z} . The outer product of a group of vectors $\mathbf{z}_i \in \mathbb{R}^{n_i}$, $1 \leq i \leq l$ and $l \in \mathbb{N}^+$, is defined as $\mathbf{T} = \mathbf{z}_1 \otimes \cdots \otimes \mathbf{z}_l \in \mathbb{R}^{n_1 \times \cdots \times n_l}$ with $T_{j_1, \dots, j_l} = (z_1)_{j_1} \cdots (z_l)_{j_l}$. Let \mathcal{L}_i be a linear operator from \mathbb{R}^{n_i} to \mathbb{R}^{d_i} with $1 \leq i \leq l$, then $\mathbf{T}(\mathcal{L}_1, \dots, \mathcal{L}_l) = \mathcal{L}_1(\mathbf{z}_1) \otimes \cdots \otimes \mathcal{L}_l(\mathbf{z}_l) \in \mathbb{R}^{d_1 \times \cdots \times d_l}$. Moreover, $f(d) = O(g(d))$ means that if for some constant $C > 0$, $f(d) \leq Cg(d)$ holds when d is sufficiently large. $f(d) = \Theta(g(d))$ means that for some constants $c > 0$ and $C > 0$, $cg(d) \leq f(d) \leq Cg(d)$ holds when d is sufficiently large. In the Appendix, we use $f(d) \gtrsim (\lesssim)g(d)$ to denote there exists some positive constant C such that $f(d) \geq (\leq)C \cdot g(d)$ when d is sufficiently large.

II. PROBLEM FORMULATION

Following [45], we consider the regression setup in this paper as follows. Given N input data $\mathbf{x}_n \in \mathbb{R}^p$, $n = 1, 2, \dots, N$, that are independent and identically distributed (i.i.d.) from the standard Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I}_{p \times p})$, the resulting outputs $y_n \in \mathbb{R}$, $n = 1, 2, \dots, N$, are generated from $\{\mathbf{x}_n\}_{n=1}^N$ by a one-hidden-layer non-overlapping convolutional neural network shown in Fig. 1. The hidden layer has K nodes. We use the vector $\mathbf{w}_j^* \in \mathbb{R}^d$ to denote the weight parameters for the j -th node in the hidden layer and define the weight matrix $\mathbf{W}^* = [\mathbf{w}_1^*, \mathbf{w}_2^*, \dots, \mathbf{w}_K^*] \in \mathbb{R}^{d \times K}$. Followed by the hidden layer, there is a pooling layer with ground-truth parameters $\mathbf{v}^* \in \mathbb{R}^d$. We assume $K < d$ throughout the paper because K is the constant, while d increases as the dimension of the input data increases. $\sigma_i = \sigma_i(\mathbf{W}^*)$ denotes the i -th largest singular value of \mathbf{W}^* . We define $\kappa = \sigma_1(\mathbf{W}^*)/\sigma_K(\mathbf{W}^*)$ as the conditional number of \mathbf{W}^* and $\gamma = \prod_{j=1}^K (\sigma_j(\mathbf{W}^*)/\sigma_K(\mathbf{W}^*))$.

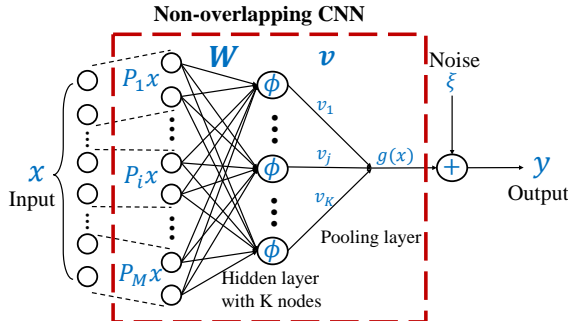


Fig. 1: One-hidden-layer non-overlapping CNN

Each input data \mathbf{x}_n is partitioned into M non-overlapping patches, denoted by $\mathbf{P}_i \mathbf{x}_n \in \mathbb{R}^d$, $i = 1, \dots, M$. $\mathbf{P}_i \in \mathbb{R}^{d \times p}$, $i = 1, \dots, M$, are a series of matrices that satisfy the following properties: (1) there exists one and only one non-zero entry with value 1 in each row of \mathbf{P}_i ; (2) $\langle \mathbf{P}_{i_1}, \mathbf{P}_{i_2} \rangle = 0$ for $i_1 \neq i_2$.¹ A simple example of $\{\mathbf{P}_i\}_{i=1}^M$ is

$$\mathbf{P}_i = \left[\underbrace{\mathbf{0}_{d \times d} \cdots \mathbf{0}_{d \times d}}_{(i-1) \text{ submatrices}} \quad \mathbf{I}_{d \times d} \quad \underbrace{\mathbf{0}_{d \times d} \cdots \mathbf{0}_{d \times d}}_{(M-i) \text{ submatrices}} \right].$$

The output y_n can be written as

$$y_n = g(\mathbf{x}_n) + \xi_n = \sum_{j=1}^K \sum_{i=1}^M v_j^* \phi(\mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x}_n) + \xi_n \quad (1)$$

for $1 \leq n \leq N$, where ξ_n is the additive stochastic noise.

Throughout this paper, we assume bounded noise with zero mean and use $|\xi|$ to denote the upper bound such that $|\xi_n| \leq |\xi|$ for all n . In practice, the mapping from the input to output data may not be modeled exactly by a neural network due to the random fluctuations or measurement errors in the data. The additive noise better characterizes the relations in real datasets.

The activation function $\phi(z) = \max\{z, 0\}$ is the ReLU function, which is widely used in various applications [11], [13], [21], [25]. Note that if the activation function is homogeneous, such as ReLU, one can assume v_j^* to be either $+1$ or -1 without loss of generality. That is because $v_j^* \phi(\mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x}_n) = \text{sign}(v_j^*) \phi(|v_j^*| \mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x}_n)$ for a homogeneous ϕ . We can just let $\tilde{\mathbf{w}}_j^* = |v_j^*| \mathbf{w}_j^*$ and $\tilde{v}_j^* = \text{sign}(v_j^*)$ and use $\{\tilde{\mathbf{w}}_j^*\}_{j=1}^K$ and $\{\tilde{v}_j^*\}_{j=1}^K$ as ground-truth parameters equivalently. Therefore, we assume $v_j^* \in \{+1, -1\}$ for any $1 \leq j \leq K$ throughout the paper.

Given any estimated $\mathbf{W} \in \mathbb{R}^{d \times K}$ and $\mathbf{v} \in \mathbb{R}^K$ of the weight matrix \mathbf{W}^* and \mathbf{v}^* , the empirical squared loss function² of the training set $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$ is defined as

$$\hat{f}_{\mathcal{D}}(\mathbf{W}, \mathbf{v}) = \frac{1}{2N} \sum_{n=1}^N \left(\sum_{j=1}^K v_j \sum_{i=1}^M \phi(\mathbf{w}_j^T \mathbf{P}_i \mathbf{x}_n) - y_n \right)^2. \quad (2)$$

Our goal is to estimate the ground-truth weight matrix \mathbf{W}^* and \mathbf{v}^* via solving the following problem:

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times K}, \mathbf{v} \in \mathbb{R}^K} \hat{f}_{\mathcal{D}}(\mathbf{W}, \mathbf{v}). \quad (3)$$

Clearly $(\mathbf{W}^*, \mathbf{v}^*)$ is a global minimizer to (3) when measurements are noiseless, i.e., $\xi_n = 0$ for all n . However, (3) is a non-convex optimization problem and is not easy to solve.

III. ALGORITHM AND THEORETICAL RESULTS

We propose to solve the non-convex problem (3) via the Heavy Ball method [28]. The algorithm is initialized via the tensor method [46]. Although the tensor initialization is designed for fully connected neural networks in [46], we can extend it to non-overlapping convolutional neural networks

¹Such requirement on \mathbf{P}_i guarantees the independence of each patches and will be used in the proofs.

²Besides the mean squared error, another choice of the loss function, especially in classification problems, is the cross entropy, see. e.g., [9].

with minor changes. $\hat{\mathbf{v}}$ is estimated through the tensor initialization. During each iteration, we update \mathbf{W} through the AGD algorithm. Compared with the vanilla gradient descent, in the $(t+1)$ -th iteration, an additional momentum term, denoted by $\beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)})$, is added to the update, where $\mathbf{W}^{(t)}$ is the estimation in iteration t . The momentum represents the direction of the previous iterations. Hence, besides moving along the gradient descent direction with a step size of η , $\mathbf{W}^{(t)}$ is further moved along the direction of previous steps with a parameter of β . During each iteration, a fresh subset of data is applied to estimate the gradient descent. Such disjoint subsets guarantee the independence of $\hat{f}_{\mathcal{D}_t}$ over the iterations. This is a standard analysis technique [45], [46] and not necessary in numerical experiments. The initialization algorithm is summarized in Section III-A, and Algorithm 1 summarizes our proposed algorithm to solve (3).

Algorithm 1 Accelerated Gradient Descent Algorithm with Tensor Initialization

- 1: **Input:** training data $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$, gradient step size η , momentum parameter β , and thresholding error parameter ε ;
 - 2: **Initialization:** $\mathbf{W}^{(0)}$, $\hat{\mathbf{v}}$ through Tensor Initialization via Subroutine 1;
 - 3: Partition \mathcal{D} into $T = \log(1/\varepsilon)$ disjoint subsets, denoted as $\{\mathcal{D}_i\}_{i=1}^T$;
 - 4: **for** $t = 1, 2, \dots, T$ **do**
 - 5: $\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)}, \hat{\mathbf{v}}) + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)})$
 - 6: **end for**
 - 7: **Return:** $\mathbf{W}^{(T)}$ and $\hat{\mathbf{v}}$.
-

A. Initialization via tensor method

In this section, we first briefly introduce the tensor initialization method that is built upon Algorithm 1 in [46]. We then provide the first theoretical performance guarantee of the tensor initialization method when the output contains noise in Lemma 1, while the result in [46] only applies to noiseless measurements.

The tensor initialization method in [46] is designed for the fully connected neural networks. To handle the convolutional neural networks, the definitions of the high-order moments (see (5)-(7)) are modified by replacing \mathbf{x} in Definition 5.1 in [46] with $\mathbf{P}_i \mathbf{x}$. All the other steps mainly follow [46].

Following [46], we define a special outer product, denoted by $\tilde{\otimes}$. For any vector $\mathbf{v} \in \mathbb{R}^{d_1}$ and $\mathbf{Z} \in \mathbb{R}^{d_1 \times d_2}$,

$$\mathbf{v} \tilde{\otimes} \mathbf{Z} = \sum_{i=1}^{d_2} (\mathbf{v} \otimes \mathbf{z}_i \otimes \mathbf{z}_i + \mathbf{z}_i \otimes \mathbf{v} \otimes \mathbf{z}_i + \mathbf{z}_i \otimes \mathbf{z}_i \otimes \mathbf{v}), \quad (4)$$

where \otimes is the outer product and \mathbf{z}_i is the i -th column of \mathbf{Z} . Next, we pick any $i \in \{1, 2, \dots, K\}$ and define

$$\mathbf{M}_{i,1} = \mathbb{E}_{\mathbf{x}} \{y \mathbf{x}\} \in \mathbb{R}^d, \quad (5)$$

$$\mathbf{M}_{i,2} = \mathbb{E}_{\mathbf{x}} \left\{ y [(\mathbf{P}_i \mathbf{x}) \otimes (\mathbf{P}_i \mathbf{x}) - \mathbf{I}] \right\} \in \mathbb{R}^{d \times d}, \quad (6)$$

$$\mathbf{M}_{i,3} = \mathbb{E}_{\mathbf{x}} \left\{ y [(\mathbf{P}_i \mathbf{x})^{\otimes 3} - (\mathbf{P}_i \mathbf{x}) \tilde{\otimes} \mathbf{I}] \right\} \in \mathbb{R}^{d \times d \times d}, \quad (7)$$

where $\mathbf{z}^{\otimes 3} := \mathbf{z} \otimes \mathbf{z} \otimes \mathbf{z}$, and $\mathbb{E}_{\mathbf{x}}$ is the expectation over \mathbf{x} .

From Claim 5.2 in [46], there exist some known constants $\psi_i, i = 1, 2, 3$, such that

$$\mathbf{M}_{i,1} = \sum_{j=1}^K \psi_1 \cdot v_j^* \|\mathbf{w}_j^*\| \cdot \bar{\mathbf{w}}_j^*, \quad (8)$$

$$\mathbf{M}_{i,2} = \sum_{j=1}^K \psi_2 \cdot v_j^* \|\mathbf{w}_j^*\| \cdot \widehat{\bar{\mathbf{w}}_j^*}^* \widehat{\bar{\mathbf{w}}_j^*}^T, \quad (9)$$

$$\mathbf{M}_{i,3} = \sum_{j=1}^K \psi_3 \cdot v_j^* \|\mathbf{w}_j^*\| \cdot \bar{\mathbf{w}}_j^{*\otimes 3}, \quad (10)$$

where $\bar{\mathbf{w}}_j^* = \mathbf{w}_j^* / \|\mathbf{w}_j^*\|_2$ in (5)-(7) is the normalization of \mathbf{w}_j^* .

$\mathbf{M}_{i,1}$, $\mathbf{M}_{i,2}$ and $\mathbf{M}_{i,3}$ can be estimated through the samples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, and let $\widehat{\mathbf{M}}_{i,1}$, $\widehat{\mathbf{M}}_{i,2}$, $\widehat{\mathbf{M}}_{i,3}$ denote the corresponding estimates. First, we will decompose the rank- k tensor $\mathbf{M}_{i,3}$ and obtain the $\{\bar{\mathbf{w}}_j^*\}_{j=1}^K$. By applying the tensor decomposition method [17] to $\widehat{\mathbf{M}}_{i,3}$, the outputs, denoted by $\widehat{\bar{\mathbf{w}}_j^*}$, are the estimations of $\{s_j \bar{\mathbf{w}}_j^*\}_{j=1}^K$, where s_j is an unknown sign. Second, we will estimate s_j , v_j^* and $\|\mathbf{w}_j^*\|_2$ through $\mathbf{M}_{i,1}$ and $\mathbf{M}_{i,2}$. Note that $\mathbf{M}_{i,2}$ does not contain the information of s_j because s_j^2 is always 1. Then, through solving the following two optimization problem:

$$\hat{\alpha}_1 = \arg \min_{\alpha_1 \in \mathbb{R}^K} : \left| \widehat{\mathbf{M}}_{i,1} - \sum_{j=1}^K \psi_1 \alpha_{1,j} \widehat{\bar{\mathbf{w}}_j^*} \right|, \quad (11)$$

$$\hat{\alpha}_2 = \arg \min_{\alpha_2 \in \mathbb{R}^K} : \left| \widehat{\mathbf{M}}_{i,2} - \sum_{j=1}^K \psi_2 \alpha_{2,j} \widehat{\bar{\mathbf{w}}_j^*} \widehat{\bar{\mathbf{w}}_j^*}^T \right|,$$

The estimation of s_j can be given as $\hat{s}_j = \text{sign}(\hat{\alpha}_{1,j} / \hat{\alpha}_{2,j})$. Also, we know that $|\hat{\alpha}_{1,j}|$ is the estimation of $\|\mathbf{w}_j^*\|$ and $\hat{v}_j = \text{sign}(\hat{\alpha}_{1,j} / s_j)$. Thus, $\mathbf{W}^{(0)}$ is given as $[\text{sign}(\hat{\alpha}_{2,1}) \hat{\alpha}_{1,1} \widehat{\bar{\mathbf{w}}_1^*}, \dots, \text{sign}(\hat{\alpha}_{2,K}) \hat{\alpha}_{1,K} \widehat{\bar{\mathbf{w}}_K^*}]$.

To reduce the computational complexity of tensor decomposition, one can project $\widehat{\mathbf{M}}_{i,3}$ to a lower-dimensional tensor [46]. The idea is to first estimate the subspace spanned by $\{\mathbf{w}_j^*\}_{j=1}^K$, and let $\widehat{\mathbf{V}}$ denote the estimated subspace. Then, from (7) and (10), we know that $\mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) \in \mathbb{R}^{K \times K \times K}$ is represented by

$$\begin{aligned} & \mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) \\ &= \mathbb{E}_{\mathbf{x}} \left\{ y [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x})^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}) \tilde{\otimes} \mathbf{I}] \right\} \\ &= \sum_{j=1}^K \psi_3 (\widehat{\mathbf{V}}^T \mathbf{w}_j^*) \cdot (\widehat{\mathbf{V}}^T \bar{\mathbf{w}}_j^*)^{\otimes 3} \end{aligned} \quad (12)$$

and can be estimated by training samples as well. Next, one can decompose the estimate $\widehat{\mathbf{M}}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$ to obtain unit vectors $\{\hat{\mathbf{u}}_j\}_{j=1}^K \in \mathbb{R}^K$. Since $\bar{\mathbf{w}}_j^*$ lies in the subspace \mathbf{V} , we have $\mathbf{V} \mathbf{V}^T \bar{\mathbf{w}}_j^* = \bar{\mathbf{w}}_j^*$. Then, $\widehat{\mathbf{V}} \hat{\mathbf{u}}_j$ is an estimate of $s_j \bar{\mathbf{w}}_j^*$. The initialization process is summarized in Subroutine 1.

Subroutine 1 Tensor Initialization Method

- 1: **Input:** training data $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$;
 - 2: Partition \mathcal{D} into three disjoint subsets $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3$;
 - 3: Calculate $\widehat{M}_{i,1}, \widehat{M}_{i,2}$ following (5), (6) using $\mathcal{D}_1, \mathcal{D}_2$, respectively;
 - 4: Obtain the estimate subspace \widehat{V} of $\widehat{M}_{i,2}$;
 - 5: Calculate $\widehat{M}_{i,3}(\widehat{V}, \widehat{V}, \widehat{V})$ using (12) through \mathcal{D}_3 ;
 - 6: Obtain $\{\widehat{\mathbf{u}}_j\}_{j=1}^K$ via tensor decomposition method [17];
 - 7: Obtain $\widehat{\alpha}_1, \widehat{\alpha}_2$ by solving optimization problem (11);
 - 8: **Return:** $\mathbf{w}_j^{(0)} = \text{sign}(\widehat{\alpha}_{2,j})\widehat{\alpha}_{1,j}\widehat{V}\widehat{\mathbf{u}}_j$ and $\widehat{\mathbf{v}} = \text{sign}(\widehat{\alpha}_2)$, $j = 1, \dots, K$.
-

B. Parameter estimation through accelerated gradient descent

In this part, we provide the major theoretical results. Lemma 1 provides the first error bound of the initialization using the tensor initialization method in the presence of noise. Based on the tensor initialization method, Theorem 1 summarizes the recovery accuracy of \mathbf{W}^* using Algorithm 1.

Lemma 1. Assume the noise level $|\xi| \leq KM\sigma_1$ and the number of samples $N \geq C_1\kappa^8 M^2 K d \log^4 d$ for some large positive constant C_1 , the tensor initialization method in Subroutine 1 outputs $\widehat{\mathbf{v}}, \mathbf{W}^{(0)}$ such that

$$\widehat{\mathbf{v}} = \mathbf{v}^*, \quad (13)$$

and

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_2 \leq C_2\kappa^6 \sqrt{\frac{K^4 d \log d}{N}} (KM\sigma_1 + |\xi|) \quad (14)$$

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

Theorem 1. Let $\{\mathbf{W}^{(t)}\}_{t=1}^T$ be the sequence generated in Algorithm 1 with $\eta = \frac{1}{12M^2K}$. Suppose the noise level $|\xi| \leq KM\sigma_1$ and the number of samples satisfies

$$N \geq C_3\varepsilon_0^{-2}\kappa^9\gamma^3 M^3 K^8 d \log^4 d \log(1/\varepsilon) \quad (15)$$

for some constants $C_3 > 0$ and $\varepsilon_0 \in (0, \frac{1}{2})$. Then $\{\mathbf{W}^{(t)}\}_{t=1}^T$ converges linearly to \mathbf{W}^* with probability at least $1 - K^2 M^2 T \cdot d^{-10}$ as

$$\begin{aligned} \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2 &\leq \nu(\beta)^t \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_2 \\ &\quad + C_4 \sqrt{\frac{\kappa^2 \gamma M K^2 d \log d}{N}} \cdot |\xi|, \end{aligned} \quad (16)$$

and

$$\|\mathbf{W}^{(T)} - \mathbf{W}^*\|_2 \leq \varepsilon \|\mathbf{W}^*\|_2 + C_4 \sqrt{\frac{\kappa^2 \gamma M K^2 d \log d}{N}} \cdot |\xi|, \quad (17)$$

where $\nu(\beta)$ is the convergence rate that depends on β , and C_4 is some positive constant. Moreover, we have

$$\nu(\beta) < \nu(0) \quad \text{for some small nonzero } \beta, \quad (18)$$

Specifically, let $\beta^* = \left(1 - \sqrt{\frac{1-\varepsilon_0}{132\kappa^2\gamma KM}}\right)^2$, we have

$$\begin{aligned} 1 - \frac{1-\varepsilon_0}{132\kappa^2\gamma KM} &\leq \nu(0) \leq 1 - \frac{1-2\varepsilon_0}{132\kappa^2\gamma KM}, \\ \nu(\beta^*) &\leq 1 - \frac{1-\varepsilon_0}{\sqrt{132\kappa^2\gamma KM}}. \end{aligned} \quad (19)$$

Remark 1 (Zero generalization error of learned model):

Lemma 1 shows that the weight vector \mathbf{v}^* of the second layer can be exactly recovered when the noise is bounded, and there exist enough samples. Theorem 1 shows that the iterates returned by Algorithm 1 converge to \mathbf{W}^* exactly in the noiseless case or approximately in noisy case. For the convenience of presentation, we refer to the second term on the right-hand side of (16) and (17) as the noise error term. Specifically, when the relation of input \mathbf{x} and the output y can be exactly described by the CNN model, i.e., the noise $\xi = 0$, then the noise error term vanishes, and the ground-truth \mathbf{W}^* can be estimated exactly with a finite number of samples. When the noise is not zero, the noise error term decreases as the number of samples N increases in the order of $\sqrt{1/N}$. With a sufficiently large sample size, the iterates can approach \mathbf{W}^* for an arbitrarily small error. With the number of samples satisfies (15), the second error term on the right-hand side of (16) is proportional to the noise magnitude $|\xi|$. From the definition of $g(\cdot)$, one can check that $\kappa KM\sigma_1 \leq \mathbb{E}_{\mathbf{x}}|g(\mathbf{x})| \leq KM\sigma_1$ when \mathbf{x} follows $\mathcal{N}(0, 1)$. Then the condition in Lemma 1 and Theorem 1 that $|\xi| \leq KM\sigma_1$ means that the noise can be as high as the order of the average energy of the noiseless output $g(\mathbf{x})$.

Remark 2 (Faster linear convergence rate than GD in learning neural networks):

Theorem 1 indicates that the Heavy Ball step can accelerate the rate of convergence as shown in (18). Without the second momentum term, i.e., $\beta = 0$, the rate of convergence is $1 - \Theta(\frac{1}{KM})$ for the vanilla GD. If β is selected appropriately, the rate of convergence is improved and upper bounded by $1 - \Theta(\frac{1}{\sqrt{KM}})$. This is the first paper to provide theoretical guarantees for the convergence of AGD methods in learning neural networks.

Remark 3 (Sample complexity analysis): Theorem 1 requires $O(M^3 K^8 d \log^4 d \log(\frac{1}{\varepsilon}))$ number of samples for the successful estimation. K is the number of nodes in the hidden layer and usually a fixed constant for a given neural network. d is the dimension of patches and scales with the size of input data. ε is the estimation error of \mathbf{W}^* . Note that the degree of freedom of \mathbf{W}^* is Kd . The required number of samples in Theorem 1 depends on $d \log^4 d$ and thus is nearly optimal with respect to d .

C. Comparisons with related works

We compare our results with all the exiting works to the best of our knowledge that provide generalizability guarantees. We focus on the following three aspects.

(1) Tensor initialization method and AGD algorithm:

Tensor initialization method is first introduced and analyzed in [46] for fully connected neural networks with homogeneous activation functions. Ref. [9] extends the analysis to the non-homogeneous sigmoid activation. However, both works only consider noiseless settings. When reduced to the case of fully connected neural networks without noise, i.e., $\xi = 0$ and $M = 1$, the bound in (14) is as tight as that in [46].

Existing works only consider the convergence of GD instead of AGD in neural networks. Due to the additional momentum term, the analysis of GD does not directly generalize to AGD.

Specifically, the convergence of GD is based on establishing $\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 \leq \nu \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2$ for some $|\nu| < 1$, so this analysis does not directly apply to AGD. Instead, our analysis of AGD is based on the augmented iteration as $\begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix}$, and the convergence rate is calculated as a function of β . Note our analysis also applies to the special case that $\beta = 0$, i.e., the GD algorithm.

(2) Noisy outputs: Refs. [10], [44] consider noisy outputs in fully connected neural networks. In [10], the authors analyze stochastic gradient descent through expectation, and the noise is assumed to be zero mean. Thus, the noise level does not appear in the theoretical bounds. In [44], the authors assume the existence of a proper initialization, but there is no theoretical guarantee in [44] about whether their proposed initialization method in the noisy setting can return a desirable initialization. Moreover, our error bound (16) is tighter than that in [44]. Specifically, the second term on the right-hand side of (16) only depends on noise factor ξ . In contrast, eqn. (4.1) in [44] shows that the GD algorithm converges to \mathbf{W}^* up to an estimation error that depends on both $\|\mathbf{W}^*\|_F$ and the noise level. Even when there is no noise, the additional error term in eqn. (4.1) of [44] is nonzero.

(3) Theoretical guarantees: As most existing works only focus on GD algorithm with noiseless outputs, we compare with these works by reducing to $\beta = 0$ and $\xi = 0$ in Theorem 1. Refs. [3], [8], [9], [45] consider one-hidden-layer non-overlapping convolutional neural networks. Refs. [3] and [8] show that the GD algorithm converges to the ground-truth with a constant probability from one random initialization, but the result only applies to the case of one node in the hidden layer, i.e., $K = 1$. Moreover, the analyses assume an infinite number of input samples and do not consider the sample complexity. Based on the tensor initialization method [46], refs. [9] and [45] show that the GD algorithm converges to the ground-truth with a linear convergence rate, but the result only applies to smooth activation functions, like sigmoid functions, and excludes ReLU functions. Refs. [10], [44] provide the sample complexity analysis with ReLU activation function but focus on one-hidden-layer fully connected neural networks, which can be viewed as a special case of the convolutional neural network studied in this paper by selecting $M = 1$. The sample complexity in [10] with respect to d is $\text{poly}(d)$, but the power of d is not provided explicitly. Moreover, the convergence rate in [10] is sub-linear, while our theorem shows that both GD and AGD enjoy linear convergence rates.

IV. SIMULATION

The input data $\{\mathbf{x}_n\}_{n=1}^N$ are randomly selected from the Gaussian distribution $\mathcal{N}(0, \mathbf{I})$. The number of patches M is selected as a factor of the signal dimension p , and all the patches have the same size d with $d = p/M$. Entries of the weight matrix \mathbf{W}^* are i.i.d generated from $\mathcal{N}(0, 1^2)$. The noise $\{\xi_n\}_{n=1}^N$ are i.i.d from $\mathcal{N}(0, \sigma^2)$, and the noise level is measured by σ/E_y , where E_y is the average energy of the noiseless outputs $\{g(\mathbf{x}_n)\}_{n=1}^N$ calculated as $E_y = \sqrt{\frac{1}{N} \sum_{n=1}^N |g(\mathbf{x}_n)|^2}$. The output data $\{y_n\}_{n=1}^N$ are generated

by (1). In the following numerical experiments, the whole dataset $\{\mathbf{x}_n, y_n\}_{n=1}^N$ instead of a fresh subset is used to calculate the gradient in each iteration. The initialization is randomly selected from $\{\mathbf{W}_0 \mid \|\mathbf{W}_0 - \mathbf{W}^*\|_F / \|\mathbf{W}^*\|_F < 0.5\}$ and $\mathbf{v}^{(0)} = \mathbf{v}^*$ to reduce the computation. As shown in [9], [44], random initialization and the tensor method have very similar numerical performance.

If not otherwise specified, we use the following parameter setup. p is chosen as 50, and M is selected as 5. Hence, $d = p/M$ is 10. The number of nodes in hidden layer K is chosen as 5. The number of samples N is chosen as 200. The step size of the gradient η is $\frac{2K}{M^2}$, and β is selected as $(1 - \frac{1}{\sqrt{KM}})^2$. All the simulations are implemented in MATLAB 2015a on a desktop with 3.4 GHz Intel Core i7.

A. Performance of AGD with different \mathbf{v}^*

Figs. 2 and 3 show the performance of AGD with different \mathbf{v}_j^* , and the results are averaged over 100 independent trials. In Fig. 2, the relative error is defined as $\|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F / \|\mathbf{W}^*\|_F$, where $\mathbf{W}^{(t)}$ is the estimate in the t -th iteration. In Fig. 3, each trial is called a success if the relative error is less than 10^{-6} . We generate two cases of \mathbf{v}^* . In Case 1, all the entries of \mathbf{v}^* are 1, while each entry is i.i.d. selected from $\{+1, -1\}$ with equal probability in Case 2. k is set as 5, and d is set as 60 with $p = 300$. In both figures, the results of Case 1 is shown by the lines marked as “ $v_j = +1$ ”, and the second group is marked as “ $v_j \in \{+1, -1\}$ ”. We can see that the performances of these two cases are almost the same. In the following experiments, we fix \mathbf{v}_j^* as 1 for all j .

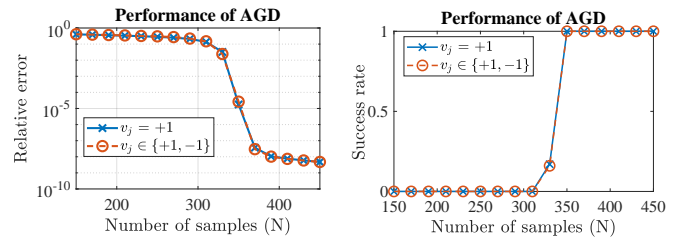


Fig. 2: Recovery error of AGD Fig. 3: Success rate of AGD under different \mathbf{v}^*

B. Performance of AGD with noiseless output

Figs. 4 and 5 show the convergence of AGD by varying K and M . In Fig. 4, η , β are calculated based the value of K , and other parameters are fixed. For each K , we conducted independent trials with random selected \mathbf{x}_n , \mathbf{W}^* and the corresponding y_n . Given K , the convergence rates of different trials vary slightly. Fig. 4 shows one example of these trials for each K . We can see that the convergence rate decreases as K increases. Similarly, Fig. 5 shows that the convergence rate decreases as M increases.

Figs. 6 and 7 show the phrase transition where the number of samples N , the dimension of patches d , and the number of nodes in the hidden layer K change. All the other parameters except N and d (or k) remain the same as the default values.

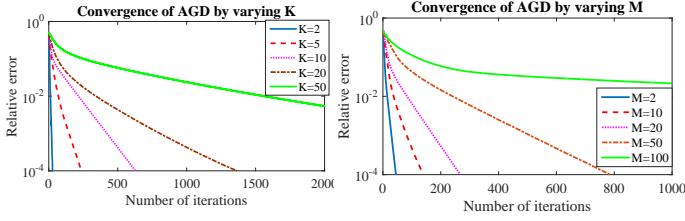


Fig. 4: Convergence of AGD with different K

with a higher noise level, the success region becomes smaller.

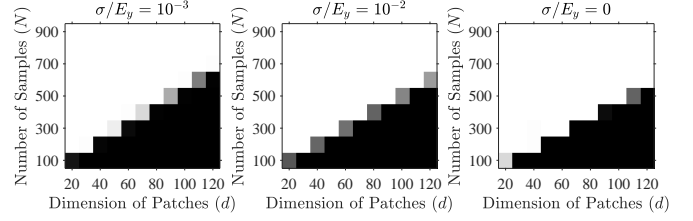


Fig. 9: The phrase transition of AGD in noisy settings

For each (N, d) or (N, K) pair, we conduct 100 independent trials. Each trial is called a success if the relative error is less than 10^{-6} . A white block means all the trials are successful, while a black one means all the trials fail.

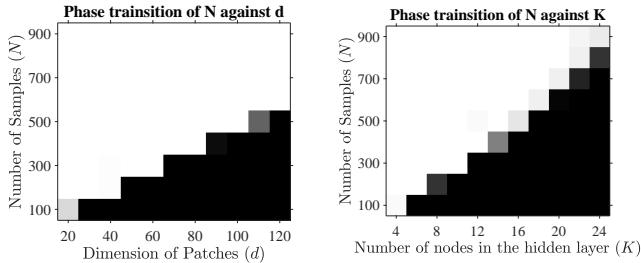


Fig. 6: Phrase transition of N against d

Fig. 7: Phrase transition of N against K

D. Comparison of GD and AGD

Fig. 10 shows the progress of both GD and AGD methods across iterations. We fix the same initialization for GD and AGD in Fig. 10(a) and (b), respectively. In both cases, β and other parameters except for η are fixed as the default values. The only difference is that the step size η is $\frac{2K}{M^2}$ in Fig. 10(a) and $\frac{3K}{M^2}$ in Fig. 10(b). One can see that starting from the same initialization, GD sometimes diverges in (b) with a large step size. By adding the heavy-ball term, the AGD method can converge to the global minimum. Moreover, when both GD and AGD converge, AGD converges faster than GD.

C. Performance of AGD with noisy output

Fig. 8 shows the relative error of AGD algorithm by varying the number of samples N in the noisy case. K is set as 5, and d is set as 60 with $p = 300$. Hence, the degree of freedom of \mathbf{W}^* is 300. Y-axis stands for the relative error, and the results are averaging over 100 independent trials. We can see that the relative errors are high when N is less than the degree of freedom as 300. Once the number of samples exceeds the degree of freedom, the relative error decreases dramatically in both noisy and noiseless settings. As N increases, the relative error in the noisy setting converges fast to the noise level.

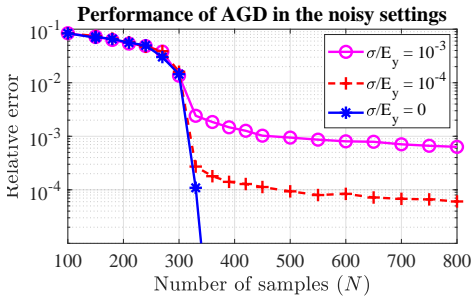


Fig. 8: The performance of Alg. 1 with noisy measurements

Fig. 9 shows the phrase transition of N against d with different noise levels. A trial is considered successful if the returned \mathbf{W} satisfies $\|\mathbf{W} - \mathbf{W}^*\|_2 / \|\mathbf{W}^*\|_2 \leq \sigma/E_y$ (or 10^{-6} in noiseless settings). As d increases, the required number of samples for all successful estimations increases as well. Also,

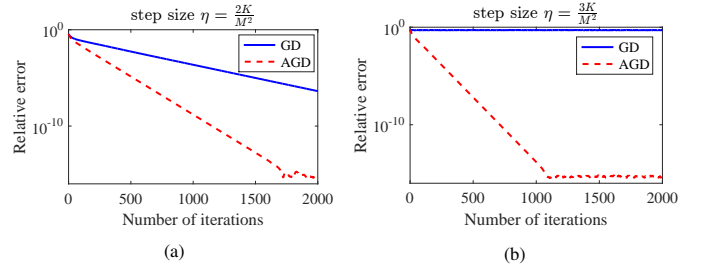


Fig. 10: Performance of AGD and GD under different η

Fig. 11 compares the convergence rates of AGD and GD. The number of samples N is set as 500, and other parameters are the default values. Each point means the smallest number of iterations needed to reach the corresponding estimation error, and the results are averaged over 100 independent trials. AGD requires a smaller number of the iterations than GD to achieve the same relative error.

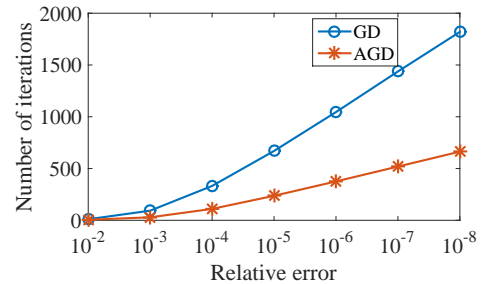


Fig. 11: Comparison of AGD and GD in number of iterations

Fig. 12 shows the phrase transition of GD and AGD by varying N and d when the output is noiseless. AGD has

a larger successful region than GD so that AGD requires a smaller number of samples to guarantee successful recovery for a given d .

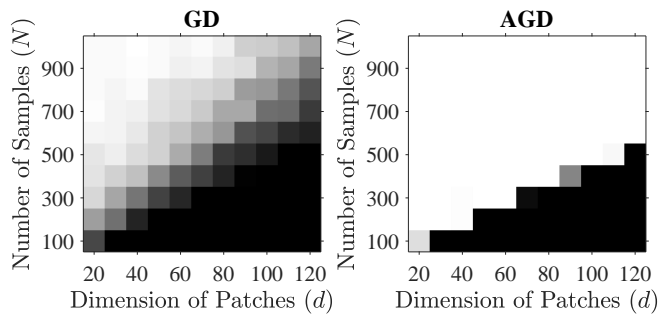


Fig. 12: The phase transition of GD and AGD

V. CONCLUSION AND FUTURE WORKS

We have analyzed the performance of (accelerated) gradient descent methods in learning one-hidden-layer non-overlapping convolutional neural networks with multiple nodes and ReLU activation function. We have shown that if the number of samples exceeds our provided sample complexity, gradient descent methods with the tensor initialization find the ground-truth parameters with a linear convergence rate. The parameters can be estimated exactly when the data are noiseless. Moreover, accelerated gradient descent is proved to converge faster than vanilla gradient descent. One future direction is to extend the analysis framework to multi-layer overlapping convolutional neural networks.

ACKNOWLEDGEMENT

This work was supported by AFOSR FA9550-20-1-0122, NSF 1932196, and the Rensselaer-IBM AI Research Collaboration (<http://airc.rpi.edu>), part of the IBM AI Horizons Network (<http://ibm.biz/AIHorizons>).

REFERENCES

- [1] Z. Allen-Zhu, Y. Li, and Y. Liang, "Learning and generalization in over-parameterized neural networks, going beyond two layers," in *Advances in neural information processing systems*, 2019, pp. 6155–6166.
- [2] O. Bousquet and A. Elisseeff, "Stability and generalization," *Journal of machine learning research*, vol. 2, no. Mar, pp. 499–526, 2002.
- [3] 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-Volume 70*. JMLR. org, 2017, pp. 605–614.
- [4] Y. Chen, H. Jiang, C. Li, X. Jia, and P. Ghamisi, "Deep feature extraction and classification of hyperspectral images based on convolutional neural networks," *IEEE Trans. Geosci. Remote.*, vol. 54, no. 10, pp. 6232–6251, July 2016.
- [5] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th International Conference on Machine Learning*, ser. ICML '08, 2008, pp. 160–167.
- [6] J. Deng, W. Dong, R. Socher, L. J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [7] S. S. Du, J. D. Lee, and Y. Tian, "When is a convolutional filter easy to learn?" *arXiv preprint, http://arxiv.org/abs/1709.06129*, 2017.
- [8] S. S. Du, J. D. Lee, Y. Tian, A. Singh, and B. Póczos, "Gradient descent learns one-hidden-layer cnn: Don't be afraid of spurious local minima," in *International Conference on Machine Learning*, 2018, pp. 1338–1347.
- [9] H. Fu, Y. Chi, and Y. Liang, "Guaranteed recovery of one-hidden-layer neural networks via cross entropy," *arXiv preprint arXiv:1802.06463*, 2018.
- [10] R. Ge, J. D. Lee, and T. Ma, "Learning one-hidden-layer neural networks with landscape design," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=BkwHObBRZ>
- [11] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, 2011, pp. 315–323.
- [12] S. Goel, A. Klivans, and R. Meka, "Learning one convolutional layer with overlapping patches," in *ICML*, 2018.
- [13] R. H. Hahnloser and H. S. Seung, "Permitted and forbidden sets in symmetric threshold-linear networks," in *Advances in Neural Information Processing Systems*, 2001, pp. 217–223.
- [14] M. Hardt, B. Recht, and Y. Singer, "Train faster, generalize better: Stability of stochastic gradient descent," in *International Conference on Machine Learning*, 2016, pp. 1225–1234.
- [15] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang, "On large-batch training for deep learning: Generalization gap and sharp minima," *arXiv preprint arXiv:1609.04836*, 2016.
- [16] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [17] V. Kuleshov, A. Chaganty, and P. Liang, "Tensor factorization via matrix factorization," in *Artificial Intelligence and Statistics*, 2015, pp. 507–516.
- [18] T. Laurent and J. Brecht, "The multilinear structure of relu networks," in *International Conference on Machine Learning*, 2018, pp. 2908–2916.
- [19] S. Lawrence, C. L. Giles, A. C. Tsoi, and A. D. Back, "Face recognition: A convolutional neural-network approach," *IEEE transactions on neural networks*, vol. 8, no. 1, pp. 98–113, 1997.
- [20] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [21] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.
- [22] 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.
- [23] S. Liang, R. Sun, J. D. Lee, and R. Srikant, "Adding one neuron can eliminate all bad local minima," in *Advances in Neural Information Processing Systems*, 2018, pp. 4355–4365.
- [24] R. Livni, S. Shalev-Shwartz, and O. Shamir, "On the computational efficiency of training neural networks," in *Advances in neural information processing systems*, 2014, pp. 855–863.
- [25] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *Proceedings of the 30th International Conference on Machine Learning*, vol. 30, no. 1, 2013.
- [26] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media, 2013, vol. 87.
- [27] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro, "Exploring generalization in deep learning," in *Advances in Neural Information Processing Systems*, 2017, pp. 5947–5956.
- [28] B. T. Polyak, "Introduction to optimization," *New York: Optimization Software, Inc*, 1987.
- [29] F. Rosenblatt, "The perceptron: a probabilistic model for information storage and organization in the brain," *Psychological review*, vol. 65, no. 6, p. 386, 1958.
- [30] D. E. Rumelhart, G. E. Hinton, R. J. Williams *et al.*, "Learning representations by back-propagating errors," *Cognitive modeling*, vol. 5, no. 3, 1988.
- [31] I. Safran and O. Shamir, "Spurious local minima are common in two-layer relu neural networks," in *International Conference on Machine Learning*, 2018, pp. 4430–4438.
- [32] O. Shamir, "Distribution-specific hardness of learning neural networks," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 1135–1163, 2018.
- [33] 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.
- [34] 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*, vol. 65, no. 2, pp. 742–769, 2018.

- [35] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *Int. Conf. Mach. Learn.*, 2013, pp. 1139–1147.
- [36] 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-Volume 70*. JMLR. org, 2017, pp. 3404–3413.
- [37] J. A. Tropp, "User-friendly tail bounds for sums of random matrices," *Foundations of computational mathematics*, vol. 12, no. 4, pp. 389–434, 2012.
- [38] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," *arXiv preprint arXiv:1011.3027*, 2010.
- [39] G. Wang, G. B. Giannakis, and J. Chen, "Learning relu networks on linearly separable data: Algorithm, optimality, and generalization," *IEEE Transactions on Signal Processing*, vol. 67, no. 9, pp. 2357–2370, 2019.
- [40] S. Wu, A. G. Dimakis, and S. Sanghavi, "Learning distributions generated by one-layer relu networks," in *Advances in Neural Information Processing Systems*, 2019, pp. 8105–8115.
- [41] T. Yang, Q. Lin, and Z. Li, "Unified convergence analysis of stochastic momentum methods for convex and non-convex optimization," *arXiv preprint arXiv:1604.03257*, 2016.
- [42] G. Yehudai and O. Shamir, "On the power and limitations of random features for understanding neural networks," in *Advances in Neural Information Processing Systems*, 2019, pp. 6594–6604.
- [43] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *arXiv preprint arXiv:1611.03530*, 2016.
- [44] X. Zhang, Y. Yu, L. Wang, and Q. Gu, "Learning one-hidden-layer relu networks via gradient descent," in *The 22nd International Conference on Artificial Intelligence and Statistics*, 2019, pp. 1524–1534.
- [45] K. Zhong, Z. Song, and I. S. Dhillon, "Learning non-overlapping convolutional neural networks with multiple kernels," *arXiv preprint arXiv:1711.03440*, 2017.
- [46] 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-Volume 70*. JMLR. org, <https://arxiv.org/abs/1706.03175>, 2017, pp. 4140–4149.

APPENDIX

A. Proof of Theorem 1

We first summarize the high-level ideas in proving Theorem 1 before presenting the technical proof. Following the recent line of research such as [45], [46], the idea is to initialize the weights \mathbf{W} near the ground-truth \mathbf{W}^* and then gradually converge to it. Our initialization is similar to [46], as discussed in Section III-A. However, our proof is more involved than that of [46] to handle the additional noise item, the non-smooth ReLU functions, the additional momentum term in accelerated gradient descent, and different neural network structures.

As for the convergence analysis, refs. [45], [46] apply the intermediate value theorem over $\nabla \hat{f}_{\mathcal{D}_t}$ at each iterate \mathbf{W}_t as

$$\nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)}) \simeq \langle \nabla^2 \hat{f}_{\mathcal{D}_t}(\widehat{\mathbf{W}}^{(t)}), \mathbf{W}^{(t)} - \mathbf{W}^* \rangle$$

for some $\widehat{\mathbf{W}}^{(t)}$ between $\mathbf{W}^{(t)}$ and \mathbf{W}^* and analyze $\nabla^2 \hat{f}_{\mathcal{D}_t}$ to obtain a recursive inequality of $\mathbf{W}^{(t)} - \mathbf{W}^*$ over t . The intermediate value theorem only applies to the continuous functions, and their analyses do not extend to our setup because with the ReLU activation function, the resulting $\nabla \hat{f}_{\mathcal{D}}$ is non-continuous. Instead, we will first prove that the population loss function f , which is defined as

$$\begin{aligned} f(\mathbf{W}) &:= \mathbb{E}_{\mathcal{D}_t} \hat{f}_{\mathcal{D}_t}(\mathbf{W}) \\ &= \mathbb{E}_{\mathbf{x}} \left(\frac{1}{K} \sum_{j=1}^K \sum_{i=1}^M \phi(\mathbf{w}_j^T \mathbf{P}_i \mathbf{x}) - y \right)^2, \end{aligned} \quad (20)$$

is locally convex near \mathbf{W}^* , and the gradient of $\hat{f}_{\mathcal{D}_t}$ is close enough to ∇f . We will then show that the iterates based on $\nabla \hat{f}_{\mathcal{D}_t}$ converge to \mathbf{W}^* .

The following two lemmas are important for our proof. We leave their proofs to Appendix-C and D.

Lemma 2. *For any \mathbf{W} that satisfies*

$$\|\mathbf{W} - \mathbf{W}^*\|_2 \leq \frac{\varepsilon_0 \sigma_K}{44\kappa^2 \gamma M}, \quad (21)$$

we have

$$\frac{(1 - \varepsilon_0)M}{11\kappa^2 \gamma} \mathbf{I} \leq \nabla^2 f(\mathbf{W}) \leq 6M^2 K \mathbf{I}. \quad (22)$$

Lemma 3. *Suppose a fixed point \mathbf{W} satisfies (21). Then, for a training set \mathcal{D} with $N > d \log d$ samples, we have*

$$\begin{aligned} &\left\| \nabla f(\mathbf{W}) - \nabla \hat{f}_{\mathcal{D}}(\mathbf{W}) \right\|_2 \\ &\lesssim MK \sqrt{\frac{d \log d}{N}} \left(MK \|\mathbf{W} - \mathbf{W}^*\|_2 + |\xi| \right), \end{aligned} \quad (23)$$

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

Lemma 2 shows that the population loss function $f(\mathbf{W})$ is locally convex near \mathbf{W}^* . Then, the analysis of AGD algorithm over the empirical loss function $\hat{f}_{\mathcal{D}}(\mathbf{W})$ is based on the analysis over $f(\mathbf{W})$ and the error bound between $\nabla \hat{f}_{\mathcal{D}}(\mathbf{W})$ and $\nabla f(\mathbf{W})$ as shown in (26).

Lemma 3 describes the error bound between $\nabla f(\mathbf{W})$ and $\nabla \hat{f}_{\mathcal{D}}(\mathbf{W})$, and (23) shows that $\nabla \hat{f}_{\mathcal{D}}(\mathbf{W})$ converges to $\nabla f(\mathbf{W})$ in a small neighborhood of \mathbf{W}^* when N is large enough. A similar result is stated in Lemma 5.3 of [44] for fully connected neural networks with ReLU activation function. Fully connected neural networks can be viewed as a special kind of convolutional neural networks with $M = 1$. Moreover, even when reducing our model to the case $M = 1$, the error bound presented in (23) is much tighter than that in Lemma 5.3 of [44].

Combining Lemmas 1, 2 and 3, we will show the convergence of GD in solving (3) by mathematical induction. Conditioned on the assumption that $\mathbf{W}^{(t)}$ satisfies (21), we show that $\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2$ is related to $\|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2$ by (38). The acceleration of Heavy-Ball steps is analyzed through (32), and the result is summarized in (33). The next step is to show (38) holds for all $0 \leq t \leq T - 1$. By Lemma 1, we can choose N to be large enough so that $\mathbf{W}^{(0)}$ satisfies (21). Then in the induction step, with a large enough N and a bounded ξ , we will show that $\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 < \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2$. Then $\mathbf{W}^{(t)}$ satisfies (21) naturally. The details are as follows.

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

$$\begin{aligned} &\mathbf{W}^{(t+1)} \\ &= \mathbf{W}^{(t)} - \eta \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)}) + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) \\ &= \mathbf{W}^{(t)} - \eta \nabla f(\mathbf{W}^{(t)}) + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) \\ &\quad + \eta(\nabla f(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)})). \end{aligned} \quad (24)$$

Since $\nabla^2 f$ is a smooth function, by the intermediate value theorem, we have

$$\begin{aligned} \mathbf{W}^{(t+1)} &= \mathbf{W}^{(t)} - \eta \nabla^2 f(\widehat{\mathbf{W}}^{(t)})(\mathbf{W}^{(t)} - \mathbf{W}^*) \\ &\quad + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)}) \\ &\quad + \eta(\nabla f(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)})), \end{aligned} \quad (25)$$

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

Next, we have

$$\begin{aligned} &\begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} - \eta \nabla^2 f(\widehat{\mathbf{W}}^{(t)}) + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{W}^{(t)} - \mathbf{W}^* \\ \mathbf{W}^{(t-1)} - \mathbf{W}^* \end{bmatrix} \\ &\quad + \eta \begin{bmatrix} \nabla f(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)}) \\ \mathbf{0} \end{bmatrix}. \end{aligned} \quad (26)$$

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

$$\begin{aligned} \left\| \begin{bmatrix} \mathbf{W}^{(t+1)} - \mathbf{W}^* \\ \mathbf{W}^{(t)} - \mathbf{W}^* \end{bmatrix} \right\|_2 &= \|\mathbf{A}(\beta)\|_2 \left\| \begin{bmatrix} \mathbf{W}^{(t)} - \mathbf{W}^* \\ \mathbf{W}^{(t-1)} - \mathbf{W}^* \end{bmatrix} \right\|_2 \\ &\quad + \eta \left\| \begin{bmatrix} \nabla f(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)}) \\ \mathbf{0} \end{bmatrix} \right\|_2. \end{aligned}$$

From Lemma 3, we know that

$$\begin{aligned} &\eta \left\| \nabla f(\mathbf{W}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)}) \right\|_2 \\ &\leq C_5 \eta M^2 \sqrt{\frac{d \log d}{N_t}} \left(\|\mathbf{W} - \mathbf{W}^*\|_2 + \frac{|\xi|}{M} \right) \end{aligned} \quad (27)$$

for some constant $C_5 > 0$. Then, we have

$$\begin{aligned} &\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 \\ &\leq \left(\|\mathbf{A}(\beta)\|_2 + C_5 \eta M^2 \sqrt{\frac{d \log d}{N_t}} \right) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2 \\ &\quad + C_5 \eta M \sqrt{\frac{d \log d}{N_t}} |\xi| \\ &:= \nu(\beta) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2 + C_5 \eta M \sqrt{\frac{d \log d}{N_t}} |\xi|. \end{aligned} \quad (28)$$

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

$$\begin{aligned} \tilde{\mathbf{A}}(\beta) &:= \begin{bmatrix} \mathbf{S}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^T \end{bmatrix} \mathbf{A}(\beta) \begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{I} - \eta \mathbf{\Lambda} + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{bmatrix}. \end{aligned} \quad (29)$$

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

$$\delta_i^2 - (1 - \eta \lambda_i + \beta) \delta_i + \beta = 0. \quad (30)$$

Then, we have

$$\delta_i(\beta) = \frac{(1 - \eta \lambda_i + \beta) + \sqrt{(1 - \eta \lambda_i + \beta)^2 - 4\beta}}{2}, \quad (31)$$

and

$$|\delta_i(\beta)| = \begin{cases} \sqrt{\beta}, & \text{if } \beta \geq (1 - \sqrt{\eta \lambda_i})^2, \\ \frac{1}{2} \left| (1 - \eta \lambda_i + \beta) + \sqrt{(1 - \eta \lambda_i + \beta)^2 - 4\beta} \right|, & \text{otherwise.} \end{cases} \quad (32)$$

Note that the other root of (30) is abandoned because the root in (31) is always no less than the other root with $|1 - \eta \lambda_i| < 1$. By simple calculations, we have

$$\delta_i(0) > \delta_i(\beta), \quad \text{for } \forall \beta \in (0, (1 - \eta \lambda_i)^2). \quad (33)$$

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

Let us first assume $\mathbf{W}^{(t)}$ satisfies (21), then from Lemma 2, we know that

$$0 < \frac{(1 - \varepsilon_0)M}{11\kappa^2\gamma} \leq \lambda_i \leq 6M^2K.$$

Let $\gamma_1 = \frac{(1 - \varepsilon_0)M}{11\kappa^2\gamma}$ and $\gamma_2 = 6KM^2$. If we choose β such that

$$\beta^* = \max \{ (1 - \sqrt{\eta \gamma_1})^2, (1 - \sqrt{\eta \gamma_2})^2 \}, \quad (34)$$

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

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

$$\begin{aligned} \|\mathbf{A}(\beta^*)\|_2 &= \max_i \delta_i(\beta^*) = 1 - \sqrt{\frac{\gamma_1}{2\gamma_2}} \\ &= 1 - \sqrt{\frac{1 - \varepsilon_0}{132\kappa^2\gamma KM}} \\ &\leq 1 - \frac{1 - (3/4) \cdot \varepsilon_0}{\sqrt{132\kappa^2\gamma KM}}. \end{aligned} \quad (35)$$

Then, let

$$C_5 \eta M^2 \sqrt{\frac{d \log d}{N_t}} \leq \frac{\varepsilon_0}{4\sqrt{132\kappa^2\gamma KM}}, \quad (36)$$

we need $N_t \gtrsim \varepsilon_0^{-2} \kappa^2 \gamma M K^3 d \log d$. Combining (35) and (36), we have

$$\nu(\beta^*) \leq 1 - \frac{1 - \varepsilon_0}{\sqrt{132\kappa^2\gamma KM}}. \quad (37)$$

Let $\beta = 0$, we have

$$\nu(0) \geq \|\mathbf{A}(0)\|_2 = 1 - \frac{1 - \varepsilon_0}{132\kappa^2\gamma KM},$$

$$\nu(0) \leq \|\mathbf{A}(0)\|_2 + C_5 \eta M^2 \sqrt{\frac{d \log d}{N_t}} \leq 1 - \frac{1 - 2\varepsilon_0}{132\kappa^2\gamma KM}$$

if $N_t \gtrsim \varepsilon_0^{-2} \kappa^2 \gamma M^2 K^4 d \log d$.

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

$$\begin{aligned} \|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_2 &\leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{132\kappa^2\gamma KM}}\right) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2 \\ &\quad + 2C_5 \eta M \sqrt{\frac{d \log d}{N_t}} |\xi|, \end{aligned} \quad (38)$$

provided that $\mathbf{W}^{(t)}$ satisfies (21), and

$$N_t \gtrsim \varepsilon_0^{-2} \kappa^2 \gamma M K^3 d \log d. \quad (39)$$

Then, we can start mathematical induction of (38) over t .

Base case: According to Lemma 1, we know that (21) holds for $\mathbf{W}^{(0)}$ if

$$N \gtrsim \varepsilon_0^{-2} \kappa^9 \gamma^2 K^8 M^2 d \log^4 d. \quad (40)$$

According to (15) in Theorem 1, it is clear that the number of samples N satisfies (40), then (21) indeed holds for $t = 0$. Since (21) holds for $t = 0$ and N in (15) satisfies (39) as well, we have (38) holds for $t = 0$.

Induction step: Assuming (38) holds for $\mathbf{W}^{(t)}$, we need to show that (38) holds for $\mathbf{W}^{(t+1)}$. That is to say, we need (i) N satisfies (39); (ii) (21) holds for $\mathbf{W}^{(t+1)}$. The requirement (i) holds naturally from (15). To guarantee (ii) holds, we need

$$\eta M \sqrt{\frac{d \log d}{N_t}} \lesssim \frac{1 - \varepsilon_0}{\sqrt{132 \kappa^2 \gamma K M}} \cdot \frac{\varepsilon_0 \sigma_K}{44 \kappa^2 \gamma K^2 M}. \quad (41)$$

That requires

$$N_t \gtrsim \varepsilon_0^{-2} \kappa^8 \gamma^3 M^3 K^6 d \log d. \quad (42)$$

Therefore, when $N_t \gtrsim \varepsilon_0^{-2} \kappa^9 \gamma^3 M^3 K^8 d \log^4 d$, we know that (38) holds for all $0 \leq t \leq T - 1$ with probability at least $1 - K^2 M^2 T \cdot d^{-10}$. By simple calculations, we can obtain

$$\begin{aligned} \|\mathbf{W}^{(T)} - \mathbf{W}^*\|_2 &\leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{132 \kappa^2 \gamma K M}}\right)^T \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_2 \\ &\quad + C_4 \sqrt{\frac{\kappa^2 \gamma M K^2 d \log d}{N_t}} \cdot |\xi| \end{aligned} \quad (43)$$

for some constant $C_4 > 0$. \square

B. Proof of Lemma 1

The proof of Lemma 1 is divided into three major parts to bound I_1 , I_2 and I_3 in (50). Lemmas 4, 5 and 6 provide the error bounds for I_1 , I_2 and I_3 , respectively. Compared with the proof of Theorem 5.6 in [46] which considers noiseless measurements, we need to handle additional items corresponding with noise, and the error bounds for these items are obtained by applying matrix concentration inequalities shown in Lemma 7. The detailed proofs of Lemmas 4-6 can be found in the supplementary materials.

Lemma 4. Suppose $\mathbf{M}_{i,2}$ is defined as in (6) and $\widehat{\mathbf{M}}_{i,2}$ is the estimation of $\mathbf{M}_{i,2}$ by samples $\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$. Then, with probability $1 - d^{-10}$, we have

$$\|\widehat{\mathbf{M}}_{i,2} - \mathbf{M}_{i,2}\| \lesssim \sqrt{\frac{d \log d}{N}} (KM\sigma_1 + |\xi|), \quad (44)$$

provided that $N \gtrsim d \log^4 d$.

Lemma 5. Let $\widehat{\mathbf{V}}$ be generated by step 4 in Subroutine 1. Suppose $\mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$ is defined as in (12) and $\widehat{\mathbf{M}}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$ is the estimation of $\mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$ by samples $\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$. Further, we assume $\mathbf{V} \in \mathbb{R}^{d \times K}$ is an orthogonal basis of \mathbf{W}^* and satisfies $\|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\| \leq 1/4$.

Then, provided that $N \gtrsim K^5 \log^6 d$, with probability at least $1 - d^{-10}$, we have

$$\begin{aligned} &\|\widehat{\mathbf{M}}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})\| \\ &\lesssim (KM\sigma_1 + |\xi|) \sqrt{\frac{K^3 \log d}{N}}. \end{aligned} \quad (45)$$

Lemma 6. Suppose $\mathbf{M}_{i,1}$ is defined as in (5) and $\widehat{\mathbf{M}}_{i,1}$ is the estimation of $\mathbf{M}_{i,1}$ by samples $\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$. Then, with probability $1 - d^{-10}$, we have

$$\|\widehat{\mathbf{M}}_{i,1} - \mathbf{M}_{i,1}\| \lesssim (KM\sigma_1 + |\xi|) \sqrt{\frac{d \log d}{N}} \quad (46)$$

provided that $N \gtrsim d \log^4 d$.

Lemma 7 ([37], Theorem 1.6). Consider a finite sequence $\{\mathbf{Z}_k\}$ of independent, random matrices with dimensions $d_1 \times d_2$. Assume that such random matrix satisfies

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

Define

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

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

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

Lemma 8 ([46], Lemma E.6). Let $\mathbf{V} \in \mathbb{R}^{d \times K}$ be an orthogonal basis of \mathbf{W}^* and $\widehat{\mathbf{V}}$ be generated by step 4 in Subroutine 1. Assume $\|\widehat{\mathbf{M}}_{i,2} - \mathbf{M}_{i,2}\|_2 \leq \sigma_K(\mathbf{M}_{i,2})/10$. Then, for some small ε_0 , we have

$$\|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\|_2 \leq \frac{\|\mathbf{M}_{i,2} - \widehat{\mathbf{M}}_{i,2}\|}{\sigma_K(\mathbf{M}_{i,2})}. \quad (47)$$

Lemma 9 ([46], Lemmas E.13 and E.14). Let $\mathbf{V} \in \mathbb{R}^{d \times K}$ be an orthogonal basis of \mathbf{W}^* and $\widehat{\mathbf{V}}$ be generated by step 4 in Subroutine 1. Assume $\mathbf{M}_{i,1}$ can be written in the form of (8) with some constant ϕ_1 , and let $\widehat{\mathbf{M}}_{i,1}$ be the estimation of $\mathbf{M}_{i,1}$ by samples $\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$. Let $\widehat{\boldsymbol{\alpha}}_1$ and $\widehat{\boldsymbol{\alpha}}_2$ be the optimal solutions of (11) with $\widehat{\mathbf{w}}_j = \widehat{\mathbf{V}}\widehat{\mathbf{u}}_j$. Then, for each $j \in \{1, 2, \dots, K\}$, if

$$\begin{aligned} T_1 &:= \|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\|_2 \leq \frac{1}{\kappa^2 \sqrt{K}}, \\ T_2 &:= \|\widehat{\mathbf{u}}_j - s_j \widehat{\mathbf{V}}^T \widehat{\mathbf{w}}_j\|_2 \leq \frac{1}{\kappa^2 \sqrt{K}}, \\ T_3 &:= \|\widehat{\mathbf{M}}_{i,1} - \mathbf{M}_{i,1}\|_2 \leq \frac{1}{4} \|\mathbf{M}_{i,1}\|_2, \end{aligned} \quad (48)$$

then we have

$$\begin{aligned} |\alpha_{1,j}^* - \widehat{\alpha}_{1,j}| &\leq \left(\kappa^4 K^{\frac{3}{2}} (T_1 + T_2) + \kappa^2 K^{\frac{1}{2}} T_3 \right) |\alpha_{1,j}^*|, \\ |\alpha_{2,j}^* - \widehat{\alpha}_{2,j}| &\leq \left(\kappa^8 K^3 T_2 + \kappa^2 K^2 T_3 \right) |\alpha_{2,j}^*|, \end{aligned} \quad (49)$$

where $\alpha_{1,j}^* = s_j v_j^* \|\mathbf{w}_j^*\|_2$ and $\alpha_{2,j}^* = v_j^* \|\mathbf{w}_j^*\|_2$.

Proof of Lemma 1. we have

$$\begin{aligned}
& \|\mathbf{w}_j^* - s_j|\hat{\alpha}_{1,j}|\hat{\mathbf{V}}\hat{\mathbf{u}}_j\|_2 \\
& \leq \left\| \mathbf{w}_j^* - s_j\|\mathbf{w}_j\|_2\hat{\mathbf{V}}\hat{\mathbf{u}}_j + s_j\|\mathbf{w}_j\|_2\hat{\mathbf{V}}\hat{\mathbf{u}}_j - s_j|\hat{\alpha}_{1,j}|\hat{\mathbf{V}}\hat{\mathbf{u}}_j \right\|_2 \\
& \leq \left\| \mathbf{w}_j^* - s_j\|\mathbf{w}_j\|_2\hat{\mathbf{V}}\hat{\mathbf{u}}_j \right\|_2 + \left\| \|\mathbf{w}_j\|_2\hat{\mathbf{V}}\hat{\mathbf{u}}_j - |\hat{\alpha}_{1,j}|\hat{\mathbf{V}}\hat{\mathbf{u}}_j \right\|_2 \\
& \leq \|\mathbf{w}_j^*\|_2\|\bar{\mathbf{w}}_j^* - s_j\hat{\mathbf{V}}\hat{\mathbf{u}}_j\|_2 + \left| \|\mathbf{w}_j\|_2 - |\hat{\alpha}_{1,j}| \right| \|\hat{\mathbf{V}}\hat{\mathbf{u}}_j\|_2 \\
& \leq \sigma_1 (\|\bar{\mathbf{w}}_j^* - \hat{\mathbf{V}}\hat{\mathbf{V}}^T\bar{\mathbf{w}}_j^*\|_2 + \|\hat{\mathbf{V}}^T\bar{\mathbf{w}}_j^* - s_j\hat{\mathbf{u}}_j\|_2) \\
& \quad + \left| \|\mathbf{w}_j\|_2 - |\hat{\alpha}_{1,j}| \right| \\
& := \sigma_1 (I_1 + I_2) + I_3.
\end{aligned} \tag{50}$$

From Lemma 8, we have

$$\begin{aligned}
I_1 &= \|\bar{\mathbf{w}}_j^* - \hat{\mathbf{V}}\hat{\mathbf{V}}^T\bar{\mathbf{w}}_j^*\|_2 \leq \|\mathbf{V}\mathbf{V}^T - \hat{\mathbf{V}}\hat{\mathbf{V}}^T\|_2 \\
& \leq \frac{\|\widehat{\mathbf{M}}_{i,2} - \mathbf{M}_{i,2}\|_2}{\sigma_K(\mathbf{M}_{i,2})},
\end{aligned} \tag{51}$$

where the last inequality comes from Lemma 4. Then, from (9), we know that

$$\sigma_K(\mathbf{M}_{i,2}) \lesssim \min_{1 \leq j \leq K} \|\mathbf{w}_j\|_2 \lesssim \sigma_K. \tag{52}$$

From Theorem 3 in [17], we have

$$\begin{aligned}
I_2 &= \|\hat{\mathbf{V}}^T\bar{\mathbf{w}}_j^* - s_j\hat{\mathbf{u}}_j\|_2 \\
& \lesssim \frac{\kappa}{\sigma_K} \|\widehat{\mathbf{M}}_{i,3}(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}}) - \mathbf{M}_{i,3}(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}})\|_2.
\end{aligned} \tag{53}$$

To guarantee the condition (48) in Lemma 9 hold, according to Lemmas 4 and 5, we need $N \gtrsim \kappa^3 M^2 K d \log d$. Then, from Lemma 9, we have

$$I_3 = \left(\kappa^4 K^{3/2} (I_1 + I_2) + \kappa^2 K^{1/2} \|\widehat{\mathbf{M}}_{i,1} - \mathbf{M}_{i,1}\| \right) \sigma_1. \tag{54}$$

Since $d \gg K$, according to Lemmas 4, 5 and 6, we have

$$\|\mathbf{w}_j^* - |\hat{\alpha}_{1,j}|\hat{\mathbf{V}}\hat{\mathbf{u}}_j\|_2 \lesssim \varepsilon_0 \kappa^6 \sqrt{\frac{K^3 d \log d}{N}} (M\sigma_1 + |\xi|) \tag{55}$$

provided that $N \gtrsim d \log^4 d$.

When $N \gtrsim \varepsilon_0^{-2} \kappa^8 K^4 M d \log d$ for $\varepsilon_0 \in (0, 1)$, we have

$$|\hat{\alpha}_{1,j} - \alpha_{1,j}^*| < \varepsilon_0 |\alpha_{1,j}^*|, \text{ and } |\hat{\alpha}_{2,j} - \alpha_{2,j}^*| < \varepsilon_0 |\alpha_{2,j}^*|. \tag{56}$$

Hence, $\hat{\alpha}_{1,j}$ and $\hat{\alpha}_{2,j}$ share the same signs of $\alpha_{1,j}^*$ and $\alpha_{2,j}^*$, and $\hat{v}_j = v_j^*$. \square

C. Proof of Lemma 2

In this section, we provide the proof of Lemma 2 which shows the local convexity of f in a small neighborhood of \mathbf{W}^* . The roadmap is to first bound the smallest eigenvalue of $\nabla^2 f$ in the ground truth as shown in (59), then show that the difference of $\nabla^2 f$ between any fixed point \mathbf{W} in this region and the ground truth \mathbf{W}^* is bounded in terms of $\|\mathbf{W} - \mathbf{W}^*\|_2$ by Lemma 10 the proof of which is in the supplementary materials.

Lemma 10. Suppose \mathbf{W} satisfies (21), with any $1 \leq j \leq K$ and $1 \leq i \leq M$, we have

$$\mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_j^T \mathbf{P}_i \mathbf{x}) - \phi'(\mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x})| \leq \frac{2 \|\mathbf{w}_j^* - \mathbf{w}_j\|}{\pi \|\mathbf{w}_j^*\|}, \tag{57}$$

$$\|\nabla^2 f(\mathbf{W}^*) - \nabla^2 f(\mathbf{W})\| \leq 4M^2 K^2 \frac{\|\mathbf{W}^* - \mathbf{W}\|_2}{\sigma_K}. \tag{58}$$

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

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

and

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

The error bound of $\|\nabla^2 f(\mathbf{W}^*) - \nabla^2 f(\mathbf{W})\|_2$ can be derived from Lemma 10, and the remaining part is to bound $\nabla^2 f(\mathbf{W}^*)$. The second order derivative of f at \mathbf{W} is written as

$$\begin{aligned}
& \frac{\partial^2 f(\mathbf{W})}{\partial \mathbf{w}_{j_1} \partial \mathbf{w}_{j_2}} \\
&= \mathbb{E}_{\mathbf{x}} \left[v_i^* v_j^* \left(\sum_{i=1}^M \phi'(\mathbf{w}_{j_1}^T \mathbf{P}_i \mathbf{x}) \mathbf{P}_i \mathbf{x} \right) \left(\sum_{i=1}^M \phi'(\mathbf{w}_{j_2}^T \mathbf{P}_i \mathbf{x}) \mathbf{P}_i \mathbf{x} \right)^T \right].
\end{aligned}$$

Then, denote $\mathbf{P}_i \mathbf{x}$ by \mathbf{x}_i . For any vector $\mathbf{a} \in \mathbb{R}^{KM}$, the lower bound of $\nabla^2 f(\mathbf{W}^*)$ is derived from

$$\begin{aligned}
& \mathbf{a}^T \nabla^2 f(\mathbf{W}^*) \mathbf{a} \\
&= \mathbb{E}_{\mathbf{x}} \left(\sum_{j=1}^K \sum_{i=1}^M v_j^* \mathbf{a}_j^T \mathbf{x}_i \phi'(\mathbf{w}_j^{*T} \mathbf{x}_i) \right)^2 := \mathbb{E}_{\mathbf{x}} \left(\sum_{i=1}^M h(\mathbf{x}_i) \right)^2 \\
&= \sum_{i=1}^M \mathbb{E}_{\mathbf{x}} h^2(\mathbf{x}_i) + \frac{1}{K^2} \sum_{i_1 \neq i_2} \mathbb{E}_{\mathbf{x}} h(\mathbf{x}_{i_1}) h(\mathbf{x}_{i_2}) \\
&= \sum_{i=1}^M \mathbb{E}_{\mathbf{x}} h^2(\mathbf{x}_i) + \sum_{i_1 \neq i_2} \mathbb{E}_{\mathbf{x}_{i_1}} h(\mathbf{x}_{i_1}) \mathbb{E}_{\mathbf{x}_{i_2}} h(\mathbf{x}_{i_2}) \\
&\stackrel{(a)}{\geq} \sum_{i=1}^M \mathbb{E}_{\mathbf{x}} h^2(\mathbf{x}_i) \geq \frac{M}{11\kappa^2 \gamma} \|\mathbf{a}\|_2^2,
\end{aligned} \tag{59}$$

where (a) holds since \mathbf{x}_{i_1} and \mathbf{x}_{i_2} share the same distribution, and the last inequality comes from Lemma D.6 [46].

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

$$\begin{aligned}
& \mathbf{a}^T \nabla^2 f(\mathbf{W}^*) \mathbf{a} \\
&= \mathbb{E}_{\mathbf{x}} \left(\sum_{j=1}^K \sum_{i=1}^M v_j^* \mathbf{a}_j^T \mathbf{x}_i \phi'(\mathbf{w}_j^T \mathbf{x}_i) \right)^2 \\
&\leq \sum_{j_1=1}^K \sum_{j_2=1}^K \sum_{i_1=1}^M \sum_{i_2=1}^M \left(\mathbb{E}_{\mathbf{x}} |\mathbf{a}_{j_1}^T \mathbf{x}_{i_1}|^4 \cdot \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_1}^T \mathbf{x}_{i_1})|^4 \right. \\
&\quad \cdot \mathbb{E}_{\mathbf{x}} |\mathbf{a}_{j_2}^T \mathbf{x}_{i_2}|^4 \cdot \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2})|^4 \left. \right)^{\frac{1}{4}} \\
&\leq \sum_{j_1=1}^K \sum_{j_2=1}^K \sum_{i_1=1}^M \sum_{i_2=1}^M \left(\mathbb{E}_{\mathbf{x}} |\mathbf{a}_{j_1}^T \mathbf{x}_{i_1}|^4 \cdot \mathbb{E}_{\mathbf{x}} |\mathbf{a}_{j_2}^T \mathbf{x}_{i_2}|^4 \right)^{\frac{1}{4}} \\
&\leq 5M^2 K \|\mathbf{a}\|_2^2.
\end{aligned} \tag{60}$$

Since both (59) and (60) hold for any $\mathbf{a} \in \mathbb{R}^{Kd}$, then

$$\frac{M}{11\kappa^2\gamma}\mathbf{I} \leq \nabla^2 f(\mathbf{W}^*) \leq 5M^2K\mathbf{I}. \quad (61)$$

From the assumption in (21) and Lemma 10, we have

$$\|\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\|_2 \leq \frac{\varepsilon_0 M}{11\kappa^2\gamma}. \quad (62)$$

Combining (61) and (62) completes the whole proof. \square

D. Proof of Lemma 3

The main steps in this proof is to bound the three items in (67). Lemma 11 provides the bound for case that when $i_1 = i_2$, where $\widetilde{\mathbf{X}}_1$ (or $\widehat{\mathbf{X}}_1$) and $\widetilde{\mathbf{X}}_2$ are correlated with each other. When $i_1 \neq i_2$, $\widetilde{\mathbf{X}}_1$ (or $\widehat{\mathbf{X}}_1$) and $\widetilde{\mathbf{X}}_2$ are independent, and the corresponding results are summarized in Lemma 12. Both Lemmas 11 and 12 use the fact that $\widetilde{\mathbf{X}}_1$, $\widetilde{\mathbf{X}}_2$ and $\widehat{\mathbf{X}}_1$ are sub-Gaussian random variables, and the definition of sub-Gaussian is summarized in Definition 1. Additionally, the sub-exponential random variable is defined in Definition 2. The multiplication of two sub-Gaussian random variables belongs to the sub-exponential distribution, and this property is used in the proofs of Lemmas 11 and 12.

Lemma 5.3 in [44] provides the error bound between ∇f and $\nabla \hat{f}_{\mathcal{D}}$ for the fully connected neural networks. However, there are two major differences from our proof. First, the error bound provided in [44] is much looser than ours. Second, ref. [44] only needs to consider the case that $i_1 = i_2 = 1$ due to the fully connected neural network structures. The error bound of Lemma 5.3 in [44] is $O\left(\sqrt{\frac{d \log N}{N}}(\|\mathbf{W}^*\|_2 + |\xi|)\right)$, while the error bound in Lemma 3 is $O\left(\sqrt{\frac{d \log d}{N}}(M\|\mathbf{W}^* - \mathbf{W}\|_2 + |\xi|)\right)$, and $M = 1$ for fully connected neural networks. Since all the analyses are based on the fact that the iterates lie in a small neighborhood of \mathbf{W}^* , that is $\|\mathbf{W}^{(t)} - \mathbf{W}^*\|_2 \ll \|\mathbf{W}^*\|_2$ especially for large t . Hence, it is obvious the error bound provided in Lemma 3 is tighter.

Definition 1 (Definition 5.7, [38]). A random variable X is called a sub-Gaussian random variable if it satisfies

$$(\mathbb{E}|X|^p)^{1/p} \leq c_1 \sqrt{p} \quad (63)$$

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

$$\mathbb{E}e^{s(X - \mathbb{E}X)} \leq e^{c_2 \|X\|_{\psi_2}^2 s^2} \quad (64)$$

for all $s \in \mathbb{R}$ and some constant $c_2 > 0$, where $\|X\|_{\phi_2}$ is the sub-Gaussian norm of X defined as $\|X\|_{\psi_2} = \sup_{p \geq 1} p^{-1/2}(\mathbb{E}|X|^p)^{1/p}$.

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

Definition 2 (Definition 5.13, [38]). A random variable X is called a sub-exponential random variable if it satisfies

$$(\mathbb{E}|X|^p)^{1/p} \leq c_3 p \quad (65)$$

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

$$\mathbb{E}e^{s(X - \mathbb{E}X)} \leq e^{c_4 \|X\|_{\psi_1}^2 s^2} \quad (66)$$

for $s \leq 1/\|X\|_{\psi_1}$ and some constant $c_4 > 0$, where $\|X\|_{\psi_1}$ is the sub-exponential norm of X defined as $\|X\|_{\psi_1} = \sup_{p \geq 1} p^{-1}(\mathbb{E}|X|^p)^{1/p}$.

Lemma 11. Assume \mathbf{X} , $\mathbf{X}h_1(\mathbf{X})$ and $\mathbf{X}h_2(\mathbf{X})$ all are sub-Gaussian random vectors in \mathbb{R}^d , where h_1 and h_2 are some fixed functions from \mathbb{R}^d to \mathbb{R} . Let $\{\mathbf{X}_n\}_{n=1}^N$ be N independent samples of \mathbf{X} . Then, the following holds with probability at least $1 - d^{-10}$:

$$\begin{aligned} & \left\| \frac{1}{N} \sum_{n=1}^N \mathbf{X}_n \mathbf{X}_n^T h_1(\mathbf{X}_n) h_2(\mathbf{X}_n) - \mathbb{E} \mathbf{X} \mathbf{X}^T h_1(\mathbf{X}) h_2(\mathbf{X}) \right\|_2 \\ & \lesssim \sqrt{\frac{d \log d}{N}} \|\mathbf{X} h_1(\mathbf{X})\|_{\psi_2} \|\mathbf{X} h_2(\mathbf{X})\|_{\psi_2}. \end{aligned}$$

Lemma 12. Assume \mathbf{X}_1 and \mathbf{X}_2 are two independent sub-Gaussian random vectors in \mathbb{R}^d . Let $\{\mathbf{X}_{1,n}\}_{n=1}^N$ and $\{\mathbf{X}_{2,n}\}_{n=1}^N$ be N independent samples of \mathbf{X}_1 and \mathbf{X}_2 , respectively. Then, provided that $N \gtrsim d \log d$, the following holds with probability at least $1 - d^{-10}$:

$$\left\| \frac{1}{N} \sum_{i=1}^N \mathbf{X}_{1,n} \mathbf{X}_{2,n} - \mathbb{E} \mathbf{X}_1 \mathbf{X}_2 \right\|_2 \lesssim \sqrt{\frac{d \log d}{N}} \|\mathbf{X}_1\|_{\psi_2} \|\mathbf{X}_2\|_{\psi_2}.$$

Proof of Lemma 3. We have

$$\begin{aligned} \left[\nabla \hat{f}_{\mathcal{D}}(\mathbf{W}) \right]_k &= \sum_{j=1}^K \sum_{i_1=1}^M \sum_{i_2=1}^M \frac{v_j^* v_k^*}{N} \sum_{n=1}^N \\ & \left((\mathbf{P}_{i_1} \mathbf{x}_n) (\mathbf{P}_{i_2} \mathbf{x}_n)^T \phi'(\mathbf{w}_j^T \mathbf{P}_{i_1} \mathbf{x}_n) \phi'(\mathbf{w}_k^T \mathbf{P}_{i_2} \mathbf{x}_n) \mathbf{w}_j \right. \\ & \quad \left. - (\mathbf{P}_{i_1} \mathbf{x}_n) (\mathbf{P}_{i_2} \mathbf{x}_n)^T \phi'(\mathbf{w}_j^{*T} \mathbf{P}_{i_1} \mathbf{x}_n) \phi'(\mathbf{w}_k^T \mathbf{P}_{i_2} \mathbf{x}_n) \mathbf{w}_j^* \right) \\ & \quad + \sum_{i=1}^M \frac{v_j^* v_k^*}{N} \sum_{n=1}^N \xi_n (\mathbf{P}_i \mathbf{x}_n)^T \phi'(\mathbf{w}_k^T \mathbf{P}_i \mathbf{x}_n) \\ &= \sum_{j=1}^K \sum_{i_1=1}^M \sum_{i_2=1}^M \frac{v_j^* v_k^*}{N} \sum_{n=1}^N \left[(\mathbf{P}_{i_1} \mathbf{x}_n) (\mathbf{P}_{i_2} \mathbf{x}_n)^T \right. \\ & \quad \cdot \phi'(\mathbf{w}_j^{*T} \mathbf{P}_{i_1} \mathbf{x}_n) \phi'(\mathbf{w}_k^T \mathbf{P}_{i_2} \mathbf{x}_n) (\mathbf{w}_j - \mathbf{w}_j^*) \\ & \quad \left. + (\mathbf{P}_{i_1} \mathbf{x}_n) (\mathbf{P}_{i_2} \mathbf{x}_n)^T \right. \\ & \quad \left. \cdot (\phi'(\mathbf{w}_j^T \mathbf{P}_{i_1} \mathbf{x}_n) - \phi'(\mathbf{w}_j^{*T} \mathbf{P}_{i_1} \mathbf{x}_n)) \phi'(\mathbf{w}_k^T \mathbf{P}_{i_2} \mathbf{x}_n) \mathbf{w}_j^* \right] \\ & \quad + \sum_{i=1}^M \frac{v_j^* v_k^*}{N} \sum_{n=1}^N \xi_n (\mathbf{P}_i \mathbf{x}_n)^T \phi'(\mathbf{w}_k^T \mathbf{P}_i \mathbf{x}_n). \end{aligned}$$

For simplification, let $\widetilde{\mathbf{X}}_{1,n} = v_j^* v_k^* (\mathbf{P}_{i_1} \mathbf{x}_n) \phi'(\mathbf{w}_j^{*T} \mathbf{P}_{i_1} \mathbf{x}_n)$ and $\widetilde{\mathbf{X}}_{2,n} = v_j^* v_k^* (\mathbf{P}_{i_2} \mathbf{x}_n) \phi'(\mathbf{w}_k^T \mathbf{P}_{i_2} \mathbf{x}_n)$. Also, let $\widehat{\mathbf{X}}_{1,n} = v_j^* v_k^* (\mathbf{P}_{i_1} \mathbf{x}_n) (\phi'(\mathbf{w}_j^T \mathbf{P}_{i_1} \mathbf{x}_n) - \phi'(\mathbf{w}_j^{*T} \mathbf{P}_{i_1} \mathbf{x}_n))$. Then, we have

$$\begin{aligned} & \left[\nabla \hat{f}_{\mathcal{D}}(\mathbf{W}) \right]_k \\ &= \frac{1}{N} \sum_{j,i_1,i_2,n} [\widetilde{\mathbf{X}}_{1,n} \widetilde{\mathbf{X}}_{2,n}^T (\mathbf{w}_j - \mathbf{w}_j^*) - \widehat{\mathbf{X}}_{1,n} \widetilde{\mathbf{X}}_{2,n}^T \mathbf{w}_j^*] \\ & \quad + \sum_{i_1=1}^M \frac{1}{N} \sum_{n=1}^N \xi_n \widetilde{\mathbf{X}}_{1,n}^T. \end{aligned}$$

Hence, we have

$$\begin{aligned}
& [\nabla f(\mathbf{W})]_k - [\nabla \hat{f}_{\mathcal{D}}(\mathbf{W})]_k \\
&= \frac{1}{N} \sum_{j, i_1, i_2, n} [(\widetilde{\mathbf{X}}_{1,n} \widetilde{\mathbf{X}}_{2,n}^T - \mathbb{E} \widetilde{\mathbf{X}}_1 \widetilde{\mathbf{X}}_2^T)(\mathbf{w}_j - \mathbf{w}_j^*) \\
&\quad - (\widehat{\mathbf{X}}_{1,n} \widehat{\mathbf{X}}_{2,n}^T - \mathbb{E} \widehat{\mathbf{X}}_1 \widehat{\mathbf{X}}_2^T) \mathbf{w}_j^*] \\
&\quad + \sum_{i_1=1}^M \frac{1}{N} \sum_{n=1}^N \xi_n \widetilde{\mathbf{X}}_{1,n}^T.
\end{aligned} \tag{67}$$

We claim that $\widetilde{\mathbf{X}}_1$ and $\widehat{\mathbf{X}}_1$ belong to the sub-Gaussian distribution. According to Definition 1, for any $\boldsymbol{\alpha} \in \mathbb{R}^d$, we have

$$\begin{aligned}
& (\mathbb{E}_{\mathbf{x}} |\boldsymbol{\alpha}^T \widetilde{\mathbf{X}}_1|^p)^{1/p} \\
&\leq (\mathbb{E}_{\mathbf{x}} |\boldsymbol{\alpha}^T \mathbf{P}_i \mathbf{x}|^p \cdot \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_j^T \mathbf{P}_i \mathbf{x})|^p)^{1/p} \\
&\leq (\mathbb{E}_{\mathbf{x}} |\boldsymbol{\alpha}^T \mathbf{P}_i \mathbf{x}|^p)^{1/p} \leq \sqrt{p},
\end{aligned} \tag{68}$$

where the last inequality holds since $\mathbf{P}_i \mathbf{x}$ is a Gaussian random vector with covariance matrix \mathbf{I}_d .

For $\widehat{\mathbf{X}}_1$, we have

$$\begin{aligned}
& (\mathbb{E}_{\mathbf{x}} |\boldsymbol{\alpha}^T \widehat{\mathbf{X}}_1|^p)^{1/p} \\
&\leq (\mathbb{E}_{\mathbf{x}} |\boldsymbol{\alpha}^T \mathbf{P}_i \mathbf{x}|^p \cdot \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_j^T \mathbf{P}_i \mathbf{x}_n) - \phi'(\mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x}_n)|^p)^{1/p} \\
&\leq \frac{2}{\pi} \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \sqrt{p},
\end{aligned}$$

where the last inequality comes from Lemma 10.

When $i_1 = i_2$, by Lemma 11, we have

$$\begin{aligned}
& \|\widetilde{\mathbf{X}}_{1,n} \widetilde{\mathbf{X}}_{2,n}^T - \mathbb{E} \widetilde{\mathbf{X}}_1 \widetilde{\mathbf{X}}_2^T\|_2 \lesssim \sqrt{\frac{d \log d}{N}}, \\
& \|\widehat{\mathbf{X}}_{1,n} \widehat{\mathbf{X}}_{2,n}^T - \mathbb{E} \widehat{\mathbf{X}}_1 \widehat{\mathbf{X}}_2^T\|_2 \lesssim \sqrt{\frac{d \log d}{N}} \cdot \frac{\|\mathbf{w}_j^* - \mathbf{w}_j\|}{\|\mathbf{w}_j^*\|}
\end{aligned} \tag{69}$$

with probability at least $1 - \frac{1}{d^{10}}$.

When $i_1 \neq i_2$, from Lemma 12, we also have

$$\begin{aligned}
& \|\widetilde{\mathbf{X}}_{1,n} \widetilde{\mathbf{X}}_{2,n}^T - \mathbb{E} \widetilde{\mathbf{X}}_1 \widetilde{\mathbf{X}}_2^T\|_2 \lesssim \sqrt{\frac{d \log d}{N}}, \\
& \|\widehat{\mathbf{X}}_{1,n} \widehat{\mathbf{X}}_{2,n}^T - \mathbb{E} \widehat{\mathbf{X}}_1 \widehat{\mathbf{X}}_2^T\|_2 \lesssim \sqrt{\frac{d \log d}{N}} \cdot \frac{\|\mathbf{w}_j^* - \mathbf{w}_j\|}{\|\mathbf{w}_j^*\|}.
\end{aligned} \tag{70}$$

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

For $\sum_{n=1}^N \xi_n \widetilde{\mathbf{X}}_{1,n}^T$, we have

$$\begin{aligned}
& \left\| \frac{1}{N} \sum_{n=1}^N \xi_n \widetilde{\mathbf{X}}_{1,n} \right\|_2 \leq |\xi| \cdot \left\| \frac{1}{N} \sum_{n=1}^N \mathbf{P}_{i_1} \mathbf{x}_n \phi'(\mathbf{w}_j^T \mathbf{P}_{i_1} \mathbf{x}_n) \right\|_2 \\
& \leq |\xi| \cdot \left\| \frac{1}{N} \sum_{n=1}^N \mathbf{P}_{i_1} \mathbf{x}_n \right\|_2 \\
& \lesssim \sqrt{\frac{d \log d}{N}} |\xi|
\end{aligned}$$

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

In conclusion, with probability at least $1 - K^2 M^2 d^{-10}$,

$$\begin{aligned}
& \|\nabla f(\mathbf{W}) - \nabla \hat{f}_{\mathcal{D}}(\mathbf{W})\|_F \\
&\leq \sum_{\substack{K, K, M, M \\ k=1, j=1 \\ i_1=1, i_2=1}} \left\| \frac{1}{N} \sum_{n=1}^N \widetilde{\mathbf{X}}_{1,n} \widetilde{\mathbf{X}}_{2,n}^T - \mathbb{E} \widetilde{\mathbf{X}}_1 \widetilde{\mathbf{X}}_2^T \right\|_2 \|\mathbf{w}_j - \mathbf{w}_j^*\|_2 \\
&\quad + \sum_{\substack{K, K, M, M \\ k=1, j=1 \\ i_1=1, i_2=1}} \left\| \frac{1}{N} \sum_{n=1}^N \widehat{\mathbf{X}}_{1,n} \widehat{\mathbf{X}}_{2,n}^T - \mathbb{E} \widehat{\mathbf{X}}_1 \widehat{\mathbf{X}}_2^T \right\|_2 \|\mathbf{w}_j^*\|_2 \\
&\quad + \sum_{k=1, i_1=1}^{K, M} |\xi| \cdot \left\| \frac{1}{N} \sum_{n=1}^N \widetilde{\mathbf{X}}_{1,n} \right\|_2 \\
&\lesssim M^2 K^2 \sqrt{\frac{d \log d}{N}} \|\mathbf{W} - \mathbf{W}^*\|_2 + MK \sqrt{\frac{d \log d}{N}} |\xi|.
\end{aligned}$$

□



Shuai Zhang received the B.E. degree from University of Science and Technology of China, Hefei, China, in 2016.

He is pursuing the Ph.D. degree in electrical engineering at Rensselaer Polytechnic Institute, Troy, NY. His research interests include signal processing and high dimensional data analysis.



Meng Wang (M'12) received the Ph.D. degree from Cornell University, Ithaca, NY, USA, in 2012.

She is an Associate Professor in the department of Electrical, Computer, and Systems Engineering at Rensselaer Polytechnic Institute. Her research interests include high dimensional data analysis and their applications in power systems monitoring and network inference.



Jinjun Xiong (M'06) received the Ph.D. degree from the University of California, Los Angeles, CA, USA, in 2006. He is a Research Staff Member and Program Director for cognitive computing systems research with the IBM T.J. Watson Research Center. He also co-directs the IBM-Illinois Center for Cognitive Computing Systems Research. His research interests include AI, Machine Learning and Systems. His research was recognized with six best paper awards and eight nominations for best paper awards at various international conferences.



Sijia Liu received the Ph.D. degree (with the All-University Doctoral Prize) in Electrical and Computer Engineering from Syracuse University, Syracuse, NY, USA, in 2016. He was a Postdoctoral Research Fellow at the University of Michigan, Ann Arbor prior to joining IBM Research, USA. He is currently a Research Staff Member at the MIT-IBM Watson AI Lab. His recent research interests include optimization for deep learning and adversarial machine learning. He received the Best Student Paper Award (Third Place) at ICASSP'17 and was among the seven finalists of the Best Student Paper Award at Asilomar'13.



Pin-Yu Chen (M'16) Dr. Pin-Yu Chen is a research staff member at IBM Thomas J. Watson Research Center, Yorktown Heights, NY, USA. He is also the chief scientist of RPI-IBM AI Research Collaboration and PI of ongoing MIT-IBM Watson AI Lab projects. Dr. Chen received his Ph.D. degree in electrical engineering and computer science from the University of Michigan, Ann Arbor, USA, in 2016. Dr. Chen's recent research is on adversarial machine learning and robustness of neural networks.

His long-term research vision is building trustworthy machine learning system. He received the NeurIPS 2017 Best Reviewer Award and the IEEE GLOBECOM 2010 GOLD Best Paper Award.

Supplementary Material

I. ADDITIONAL PROOFS OF LEMMAS IN APPENDIX B

A. Error bound for the second-order moment

Proof of Lemma 4. For $\widehat{\mathbf{M}}_{i,2} - \mathbf{M}_{i,2}$, we have

$$\begin{aligned}
& \widehat{\mathbf{M}}_{i,2} - \mathbf{M}_{i,2} \\
&= \frac{1}{N} \sum_{(\mathbf{x}_n, y_n) \in \mathcal{D}} y_n (\mathbf{P}_i \mathbf{x}_n \otimes \mathbf{P}_i \mathbf{x}_n - \mathbf{I}) - \mathbb{E}_{\mathbf{x}} y (\mathbf{P}_i \mathbf{x} \otimes \mathbf{P}_i \mathbf{x} - \mathbf{I}) \\
&= \frac{1}{N} \sum_{n=1}^N \left(\sum_{j=1}^K v_j^* \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}_n) + \xi_n \right) (\mathbf{P}_i \mathbf{x}_n \otimes \mathbf{P}_i \mathbf{x}_n - \mathbf{I}) \\
&\quad - \mathbb{E}_{\mathbf{x}} v_j^* \sum_{j=1}^K \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}) (\mathbf{P}_i \mathbf{x} \otimes \mathbf{P}_i \mathbf{x} - \mathbf{I}) \\
&= \sum_{j=1}^K v_j^* \left(\frac{1}{N} \sum_{n=1}^N \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}_n) (\mathbf{P}_i \mathbf{x}_n \otimes \mathbf{P}_i \mathbf{x}_n - \mathbf{I}) \right. \\
&\quad \left. - \mathbb{E}_{\mathbf{x}} \phi(\mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x}) (\mathbf{P}_i \mathbf{x} \otimes \mathbf{P}_i \mathbf{x} - \mathbf{I}) \right) \\
&\quad + \frac{1}{N} \sum_{n=1}^N \xi_n (\mathbf{P}_i \mathbf{x}_n \otimes \mathbf{P}_i \mathbf{x}_n - \mathbf{I}). \tag{1}
\end{aligned}$$

Following the notations in Lemma E.2 of [40], we denote

$$\mathbf{B}_{2,j}(\mathbf{x}) := \sum_{j=1}^K v_j^* \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}) (\mathbf{P}_i \mathbf{x} \otimes \mathbf{P}_i \mathbf{x} - \mathbf{I}). \tag{2}$$

Following the similar calculations of (I) - (III) in Lemma E.2 [40], we know that

$$\begin{aligned}
\|\mathbf{B}_{2,j}(\mathbf{x})\|_2 &\lesssim MK \|\mathbf{w}_j^*\|_2 d \log^{\frac{3}{2}} d, \\
\|\mathbb{E}_{\mathbf{x}} \mathbf{B}_{2,j}(\mathbf{x})\|_2 &\lesssim MK \|\mathbf{w}_j^*\|_2, \\
\|\mathbb{E}_{\mathbf{x}} \mathbf{B}_{2,j}^2(\mathbf{x})\|_2 &\lesssim M^2 K^2 \|\mathbf{w}_j^*\|^2 d
\end{aligned} \tag{3}$$

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

Define $\mathbf{Z}_{2,n} = \frac{1}{N} (\mathbf{B}_{2,j}(\mathbf{x}_n) - \mathbb{E}_{\mathbf{x}} \mathbf{B}_{2,j}(\mathbf{x}))$ for $n = 1, 2, \dots, N$, and it is obvious \mathbf{Z}_n is zero mean. Also, we have

$$\begin{aligned}
R_2 = \|\mathbf{Z}_{2,n}\|_2 &\leq \frac{1}{N} (\|\mathbf{B}_{2,j}(\mathbf{x}_n)\|_2 + \|\mathbb{E}_{\mathbf{x}} \mathbf{B}_{2,j}(\mathbf{x})\|_2) \\
&\lesssim N^{-1} MK \|\mathbf{w}_j^*\|_2 d \log^{\frac{3}{2}} d,
\end{aligned} \tag{4}$$

and

$$\begin{aligned}
\delta_2^2 &= \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{2,n}^2 \right\|_2^2 \\
&\leq \left\| \sum_{n=1}^N \frac{1}{N^2} \left(\mathbb{E} \mathbf{B}_{2,j}^2(\mathbf{x}_n) - (\mathbb{E} \mathbf{B}_{2,j}(\mathbf{x}_n))^2 \right) \right\|_2^2 \\
&\leq \frac{1}{N} \left(\|\mathbb{E} \mathbf{B}_{2,j}^2(\mathbf{x}_n)\|_2 + \|\mathbb{E} \mathbf{B}_{2,j}(\mathbf{x}_n)\|_2^2 \right) \\
&\lesssim N^{-1} M^2 K^2 \|\mathbf{w}_j^*\|^2 d.
\end{aligned} \tag{5}$$

Next, let $t = \Theta(M \|\mathbf{w}_j^*\|_2 \sqrt{\frac{d \log d}{N}})$. To make sure $\delta_2^2 \geq R_2 t / 3$, we need $N \gtrsim d \log^4 d$. Then, by Lemma 7, we have

$$\begin{aligned}
\text{Prob} \left\{ \left\| \sum_n \mathbf{Z}_{2,n} \right\|_2 \geq t \right\} &\leq 2d \exp \left(\frac{-t^2/2}{\delta^2 + Rt/3} \right) \\
&\leq 2d \exp \left(\frac{-t^2}{4\delta^2} \right).
\end{aligned} \tag{6}$$

That is

$$\left\| \sum_{n=1}^N \mathbf{Z}_{2,n} \right\|_2 \lesssim KM \|\mathbf{w}_j^*\|_2 \sqrt{\frac{d \log d}{N}} \tag{7}$$

with probability at least $1 - d^{-10}$. From Lemma 11, we know that

$$\left\| \frac{1}{N} \sum_{n=1}^N \mathbf{P}_i \mathbf{x}_n \otimes \mathbf{P}_i \mathbf{x}_n - \mathbf{I} \right\|_2 \lesssim \sqrt{\frac{d \log d}{N}} \tag{8}$$

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

In conclusion, we have

$$\|\widehat{\mathbf{M}}_{i,2} - \mathbf{M}_{i,2}\| \lesssim (KM\sigma_1 + |\nu|) \sqrt{\frac{d \log d}{N}} \tag{9}$$

with probability at least $1 - d^{-C}$ provided that $N \gtrsim d \log^4 d$. \square

B. Error bound for the third-order moment

Proof of Lemma 5. For $\widehat{\mathbf{M}}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$, we have

$$\begin{aligned}
& \widehat{\mathbf{M}}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) \\
&= \frac{1}{N} \sum_{(\mathbf{x}_n, y_n) \in \mathcal{D}} y_n [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n) \widetilde{\otimes} \mathbf{I}] \\
&\quad - \mathbb{E}_{\mathbf{x}} y [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x})^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}) \widetilde{\otimes} \mathbf{I}] \\
&= \frac{1}{N} \sum_{n=1}^N \left(\sum_{j=1}^K v_j^* \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}_n) + \xi_n \right) \\
&\quad \cdot [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n) \widetilde{\otimes} \mathbf{I}] \\
&\quad - \mathbb{E}_{\mathbf{x}} \sum_{j=1}^K \sum_{i'=1}^M v_j^* \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}) [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x})^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}) \widetilde{\otimes} \mathbf{I}] \\
&= \sum_{j=1}^K v_j^* \left(\frac{1}{N} \sum_{n=1}^N \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}_n) \right. \\
&\quad \cdot [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n) \widetilde{\otimes} \mathbf{I}] \\
&\quad \left. - \mathbb{E}_{\mathbf{x}} \phi(\mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x}) [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x})^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}) \widetilde{\otimes} \mathbf{I}] \right) \\
&\quad + \frac{1}{N} \sum_{n=1}^N \xi_n \left([(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n) \widetilde{\otimes} \mathbf{I}] \right).
\end{aligned} \tag{10}$$

Following the notations in Lemma E.8 of [40], we define

$$\mathbf{T}_j(\mathbf{x}) := \sum_{j=1}^K v_j^* \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}) [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x})^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}) \otimes \mathbf{I}]. \quad (11)$$

Then, $\mathbf{B}_{3,j}(\mathbf{x}) \in \mathbb{R}^{K \times K^2}$ is defined as flattening the tensor $\mathbf{T}_j(\mathbf{x})$ along the first dimension. Hence, we have

$$\begin{aligned} & \|\mathbf{B}_{3,j}(\mathbf{x})\|_2 \\ & \lesssim M \cdot \|\mathbf{w}_j^* \mathbf{x}\| \cdot \left(\|\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}\|_2^3 + 3K \|\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}\|_2 \right) \\ & \lesssim M \|\mathbf{w}_j^*\|_2 K^{\frac{5}{2}} \log^{\frac{5}{2}} d \end{aligned} \quad (12)$$

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

Following the similar calculations of (II) and (III) in Lemma E.8 of [40], we know that

$$\begin{aligned} & \|\mathbb{E}_{\mathbf{x}} \mathbf{B}_{3,j}(\mathbf{x})\|_2 \lesssim MK \|\mathbf{w}_j^*\|_2, \\ & \max \left\{ \left\| \mathbb{E}_{\mathbf{x}} [\mathbf{B}_{3,j}(\mathbf{x})^T \mathbf{B}_{3,j}(\mathbf{x})] \right\|_2, \left\| \mathbb{E}_{\mathbf{x}} [\mathbf{B}_{3,j}(\mathbf{x})^T \mathbf{B}_{3,j}(\mathbf{x})] \right\|_2 \right\} \\ & \lesssim M^2 K^4 \|\mathbf{w}_j^*\|_2^2. \end{aligned} \quad (13)$$

Define $\mathbf{Z}_{3,n} = \frac{1}{N} (\mathbf{B}_{3,j}(\mathbf{x}_n) - \mathbb{E}_{\mathbf{x}} \mathbf{B}_{3,j}(\mathbf{x}))$ for $n = 1, 2, \dots, N$, and it is obvious $\mathbf{Z}_{3,n}$ is zero mean. Also, we have

$$\begin{aligned} R_3 = \|\mathbf{Z}_{3,n}\|_2 & \leq \frac{1}{N} (\|\mathbf{B}_{3,j}(\mathbf{x}_n)\|_2 + \|\mathbb{E}_{\mathbf{x}} \mathbf{B}_{3,j}(\mathbf{x})\|_2) \\ & \lesssim N^{-1} M \|\mathbf{w}_j^*\|_2 K^{\frac{5}{2}} \log^{\frac{5}{2}} d, \end{aligned} \quad (14)$$

and

$$\begin{aligned} \delta_3^2 & = \left\{ \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{3,n} \mathbf{Z}_{3,n}^T \right\|_2, \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{3,n} \mathbf{Z}_{3,n}^T \right\|_2 \right\} \\ & \leq \frac{1}{N} (\|\mathbb{E} \mathbf{B}_{3,j}^2(\mathbf{x}_n)\|_2 + \|\mathbb{E} \mathbf{B}_{3,j}(\mathbf{x}_n)\|_2^2) \\ & \lesssim N^{-1} M^2 K^4 \|\mathbf{w}_j^*\|_2^2. \end{aligned} \quad (15)$$

Similar to (6), by applying Lemma 7, we have

$$\left\| \sum_{n=1}^N \mathbf{Z}_{3,n} \right\|_2 \lesssim KM \|\mathbf{w}_j^*\|_2 \sqrt{\frac{K^2 \log d}{N}} \quad (16)$$

with probability at least $1 - d^{-10}$ provided that $N \gtrsim K^5 \log^6 d$.

Similar to (12), we define $\widetilde{\mathbf{B}}$ by flattening the tensor $\sum_{n=1}^N [(\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n) \otimes \mathbf{I}]$ along the first dimension. Then, we know that

$$\begin{aligned} \|\widetilde{\mathbf{B}}\|_2 & \leq \left\| \sum_{n=1}^N \widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n \right\|_2^3 + 3K \left\| \sum_{n=1}^N \widehat{\mathbf{V}}^T \mathbf{P}_i \mathbf{x}_n \right\|_2 \\ & \lesssim \left(\frac{K \log d}{N} \right)^{\frac{3}{2}} + 3K \left(\frac{K \log d}{N} \right)^{\frac{1}{2}} \\ & \lesssim \left(\frac{K \log d}{N} \right)^{\frac{1}{2}} + \left(\frac{K^3 \log d}{N} \right)^{\frac{1}{2}} \\ & \lesssim \sqrt{\frac{K^3 \log d}{N}}, \end{aligned} \quad (17)$$

provided that $N \gtrsim K \log d$.

In conclusion, we have

$$\begin{aligned} & \left\| \widehat{\mathbf{M}}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_{i,3}(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) \right\| \\ & \lesssim (KM\sigma_1 + |\xi|) \sqrt{\frac{K^3 \log d}{N}} \end{aligned} \quad (18)$$

with probability at least $1 - d^{-C}$ provided that $N \gtrsim K^5 \log^6 d$. \square

C. Error bound for the first-order moment

Proof of Lemma 6. For $\widehat{\mathbf{M}}_{i,1} - \mathbf{M}_{i,1}$, we have

$$\begin{aligned} & \widehat{\mathbf{M}}_{i,1} - \mathbf{M}_{i,1} \\ & = \frac{1}{N} \sum_{(\mathbf{x}_n, y_n) \in \mathcal{D}} y_n (\mathbf{P}_i \mathbf{x}_n) - \mathbb{E}_{\mathbf{x}} y (\mathbf{P}_i \mathbf{x}) \\ & = \frac{1}{N} \sum_{n=1}^N \left(\sum_{j=1}^K v_j^* \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}_n) + \xi_n \right) \mathbf{P}_i \mathbf{x}_n \\ & \quad - \mathbb{E}_{\mathbf{x}} \sum_{j=1}^K v_j^* \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}) \mathbf{P}_i \mathbf{x} \\ & = \sum_{j=1}^K v_j^* \left(\frac{1}{N} \sum_{n=1}^N \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}_n) \mathbf{P}_i \mathbf{x}_n \right. \\ & \quad \left. - \mathbb{E}_{\mathbf{x}} \phi(\mathbf{w}_j^{*T} \mathbf{P}_i \mathbf{x}) \mathbf{P}_i \mathbf{x} \right) + \frac{1}{N} \sum_{n=1}^N \xi_n (\mathbf{P}_i \mathbf{x}_n). \end{aligned} \quad (19)$$

Define $\mathbf{B}_{1,j}(\mathbf{x}) := \sum_{i'=1}^M \phi(\mathbf{w}_j^{*T} \mathbf{P}_{i'} \mathbf{x}_n) \mathbf{P}_i \mathbf{x}_n$, then we have

$$\begin{aligned} & \|\mathbf{B}_{1,j}(\mathbf{x})\|_2 \lesssim M \|\mathbf{w}_j^*\|_2 d \log^{\frac{3}{2}} d; \\ & \|\mathbb{E}_{\mathbf{x}} \mathbf{B}_{1,j}(\mathbf{x})\|_2 \lesssim M \|\mathbf{w}_j^*\|_2; \\ & \left\{ \left\| \mathbb{E}_{\mathbf{x}} [\mathbf{B}_{1,j}(\mathbf{x}) \mathbf{B}_{1,j}(\mathbf{x})^T] \right\|_2, \left\| \mathbb{E}_{\mathbf{x}} [\mathbf{B}_{1,j}(\mathbf{x})^T \mathbf{B}_{1,j}(\mathbf{x})] \right\|_2 \right\} \\ & \lesssim M^2 \|\mathbf{w}_j^*\|_2^2. \end{aligned} \quad (20)$$

Next, define $\mathbf{Z}_{1,n} = \frac{1}{N} (\mathbf{B}_{1,j}(\mathbf{x}_n) - \mathbb{E}_{\mathbf{x}} \mathbf{B}_{1,j}(\mathbf{x}))$ for $n = 1, 2, \dots, N$, by calculation, we can obtain

$$R_1 = \|\mathbf{Z}_{1,n}\|_2 \lesssim N^{-1} M \|\mathbf{w}_j^*\|_2 d \log^{\frac{3}{2}} d, \quad (21)$$

and

$$\begin{aligned} \delta_2^2 & = \max \left\{ \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{1,n} \mathbf{Z}_{1,n}^T \right\|_2^2, \left| \sum_{n=1}^N \mathbf{Z}_{1,n}^T \mathbf{Z}_{1,n} \right| \right\} \\ & \lesssim N^{-1} M^2 \|\mathbf{w}_j^*\|_2^2 d. \end{aligned} \quad (22)$$

By applying Lemma 7, we have

$$\left\| \sum_{n=1}^N \mathbf{Z}_{1,n} \right\|_2 \lesssim M \|\mathbf{w}_j^*\|_2 \sqrt{\frac{d \log d}{N}} \quad (23)$$

with probability at least $1 - d^{-10}$ provided that $N \gtrsim d \log^4 d$. Since $\mathbf{P}_i \mathbf{x} \in \mathbb{R}^d$ belongs to the Gaussian distribution, we have

$$\left\| \frac{1}{N} \sum_{n=1}^N \mathbf{P}_i \mathbf{x}_n \right\|_2 \lesssim \sqrt{\frac{d \log d}{N}} \quad (24)$$

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

In conclusion, we have

$$\|\widehat{\mathbf{M}}_{i,1} - \mathbf{M}_{i,1}\| \lesssim (M\sigma_1 + |\xi|) \sqrt{\frac{d \log d}{N}} \quad (25)$$

with probability at least $1 - d^{-C}$, provided that $N \gtrsim d \log^4 d$. \square

II. ADDITIONAL PROOFS OF LEMMA IN APPENDIX C

A. Proof of Lemma 10

In this section, we present the proof of Lemma 10 which provides the error bound of the difference between $\nabla^2 f(\mathbf{W}^*)$ and $\nabla^2 f(\mathbf{W})$ for some \mathbf{W} near the ground-truth \mathbf{W}^* . The proof borrows some techniques from [46] in proving Lemma D.15 which provides the bound of $\|\nabla^2 \hat{f}_{\mathcal{D}}(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*)\|_2$ for fully connected neural networks. Our proof differs from that of Lemma D.15 in [46] mainly in two aspects. On the one hand, due to the different proof roadmaps of the main theorems, we do not need to directly bound the second derivative between the population loss function f and empirical loss function $\hat{f}_{\mathcal{D}}$. Then some steps in Lemma D.15 in [46] are simplified and modified to fit our proof. On the other hand, we need to modify steps to handle convolutional neural networks.

Proof of Lemma 10. We first bound one block of $\nabla^2 f(\mathbf{W}^*) - \nabla^2 f(\mathbf{W})$, and its mathematical expression is written as

$$\begin{aligned} & \left[\nabla^2 f(\mathbf{W}) - \nabla^2 f(\mathbf{W}^*) \right]_{j_1, j_2} \\ &= \mathbb{E}_{\mathbf{x}} \left[\sum_{i_1=1}^M \sum_{i_2=1}^M \phi'(\mathbf{w}_{j_1}^T \mathbf{x}_{i_1}) \phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) \mathbf{x}_{i_1} \mathbf{x}_{i_2}^T \right. \\ & \quad \left. - \sum_{i_1=1}^M \sum_{i_2=1}^M \phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x}_{i_1}) \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2}) \mathbf{x}_{i_1} \mathbf{x}_{i_2}^T \right]. \end{aligned} \quad (26)$$

Hence, we have

$$\begin{aligned} & \left| \mathbf{a}^T \left[\nabla^2 f(\mathbf{W}^*) - \nabla^2 f(\mathbf{W}) \right]_{j_1, j_2} \mathbf{a} \right| \\ &= \left| \mathbb{E}_{\mathbf{x}} \left[\sum_{i_1=1}^M \sum_{i_2=1}^M \left(\phi'(\mathbf{w}_{j_1}^T \mathbf{x}_{i_1}) \phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) \right. \right. \right. \\ & \quad \left. \left. \left. - \phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x}_{i_1}) \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2}) \right) \cdot (\mathbf{a}^T \mathbf{x}_{i_1}) (\mathbf{a}^T \mathbf{x}_{i_2}) \right] \right| \\ &\leq \sum_{i_1=1}^M \sum_{i_2=1}^M \mathbb{E}_{\mathbf{x}} \left[\left| \phi'(\mathbf{w}_{j_1}^T \mathbf{x}_{i_1}) \phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) \right. \right. \\ & \quad \left. \left. - \phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x}_{i_1}) \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2}) \right| \cdot |\mathbf{a}^T \mathbf{x}_{i_1}| \cdot |\mathbf{a}^T \mathbf{x}_{i_2}| \right] \\ &\leq \sum_{i_1, i_2} \mathbb{E}_{\mathbf{x}} \left[|\phi'(\mathbf{w}_{j_1}^T \mathbf{x}_{i_1})| \cdot |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})| \right. \\ & \quad \left. \cdot |\mathbf{a}^T \mathbf{x}_{i_1}| \cdot |\mathbf{a}^T \mathbf{x}_{i_2}| \right] \\ & \quad + \sum_{i_1, i_2} \mathbb{E}_{\mathbf{x}} \left[|\phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x}_{i_1}) - \phi'(\mathbf{w}_{j_1}^T \mathbf{x}_{i_1})| \right. \\ & \quad \left. \cdot |\phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})| \cdot |\mathbf{a}^T \mathbf{x}_{i_1}| \cdot |\mathbf{a}^T \mathbf{x}_{i_2}| \right] \end{aligned}$$

$$\begin{aligned} & \leq \sum_{i_1, i_2} \mathbb{E}_{\mathbf{x}} \left[|\phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})| |\mathbf{a}^T \mathbf{x}_{i_2}| \right] \\ & \quad + \sum_{i_1, i_2} \mathbb{E}_{\mathbf{x}} \left[|\phi'(\mathbf{w}_{j_1}^T \mathbf{x}_{i_1}) - \phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x}_{i_1})| |\mathbf{a}^T \mathbf{x}_{i_1}| \right] \\ & := \sum_{i_1, i_2} I(i_2, j_2) + I(i_1, j_1). \end{aligned}$$

It is easy to verify there exists an orthogonal basis such that $\mathcal{B} = \{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \dots, \mathbf{a}_d\}$ with $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3\}$ spans a subspace that contains $\mathbf{a}_1, \mathbf{w}_{j_2}$ and $\mathbf{w}_{j_2}^*$. Then, for any \mathbf{x}_{i_2} , we have a unique $\mathbf{z} = [z_1 \ z_2 \ \dots \ z_d]^T$ such that

$$\mathbf{x}_{i_2} = z_1 \mathbf{a}_1 + z_2 \mathbf{a}_2 + z_3 \mathbf{a}_3 + \dots + z_d \mathbf{a}_d.$$

Also, since $\mathbf{x}_{i_2} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, we have $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Then, we have

$$\begin{aligned} & I(i_2, j_2) \\ &= \mathbb{E}_{z_1, z_2, z_3} |\phi'(\mathbf{w}_{j_2}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_2}^{*T} \tilde{\mathbf{x}})| \cdot |\mathbf{a}^T \tilde{\mathbf{x}}| \\ &= \int |\phi'(\mathbf{w}_{j_2}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_2}^{*T} \tilde{\mathbf{x}})| \cdot |\mathbf{a}^T \tilde{\mathbf{x}}| \\ & \quad \cdot f_Z(z_1, z_2, z_3) dz_1 dz_2 dz_3, \end{aligned}$$

where $\tilde{\mathbf{x}} = z_1 \mathbf{a}_1 + z_2 \mathbf{a}_2 + z_3 \mathbf{a}_3$ and $f_Z(z_1, z_2, z_3)$ is probability density function of (z_1, z_2, z_3) . Next, we consider spherical coordinates with $z_1 = r \cos \phi_1, z_2 = r \sin \phi_1 \sin \phi_2, z_3 = r \sin \phi_1 \cos \phi_2$. Hence,

$$\begin{aligned} I(i_2, j_2) &= \int |\phi'(\mathbf{w}_{j_2}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_2}^{*T} \tilde{\mathbf{x}})| \cdot |r \cos \phi_1| \\ & \quad \cdot f_Z(r, \phi_1, \phi_2) r^2 \sin \phi_1 dr d\phi_1 d\phi_2. \end{aligned} \quad (27)$$

It is easy to verify that $\phi'(\mathbf{w}_{j_2}^T \tilde{\mathbf{x}})$ only depends on the direction of $\tilde{\mathbf{x}}$ and

$$f_Z(r, \phi_1, \phi_2) = \frac{1}{(2\pi)^{\frac{3}{2}}} e^{-\frac{r^2 + r^2 \sin^2 \phi_1 + r^2 \sin^2 \phi_1 \cos^2 \phi_2}{2}} = \frac{1}{(2\pi)^{\frac{3}{2}}} e^{-\frac{r^2}{2}}$$

only depends on r . Then, we have

$$\begin{aligned} & I(i_2, j_2) \\ &= \int |\phi'(\mathbf{w}_{j_2}^T(\tilde{\mathbf{x}}/r)) - \phi'(\mathbf{w}_{j_2}^{*T}(\tilde{\mathbf{x}}/r))| \\ & \quad \cdot |r \cos \phi_1| \cdot f_Z(r) r^2 \sin \phi_1 dr d\phi_1 d\phi_2 \\ &= \int_0^\infty r^3 f_Z(r) dr \int_0^\pi \int_0^{2\pi} |\cos \phi_1| \cdot \sin \phi_1 \\ & \quad \cdot |\phi'(\mathbf{w}_{j_2}^T(\tilde{\mathbf{x}}/r)) - \phi'(\mathbf{w}_{j_2}^{*T}(\tilde{\mathbf{x}}/r))| d\phi_1 d\phi_2 \\ &\leq \sqrt{\frac{8}{\pi}} \int_0^\infty r^2 f_Z(r) dr \int_0^\pi \int_0^{2\pi} \sin \phi_1 \\ & \quad \cdot |\phi'(\mathbf{w}_{j_2}^T(\tilde{\mathbf{x}}/r)) - \phi'(\mathbf{w}_{j_2}^{*T}(\tilde{\mathbf{x}}/r))| d\phi_1 d\phi_2 \\ &= \sqrt{\frac{8}{\pi}} \mathbb{E}_{z_1, z_2, z_3} |\phi'(\mathbf{w}_{j_2}^T \tilde{\mathbf{x}}) - \phi'(\mathbf{w}_{j_2}^{*T} \tilde{\mathbf{x}})| \\ &= \sqrt{\frac{8}{\pi}} \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})|. \end{aligned} \quad (28)$$

Define a set $\mathcal{A}_1 = \{\mathbf{x} | (\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) < 0\}$. If $\mathbf{x} \in \mathcal{A}_1$, then $\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2}$ and $\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}$ have different signs, which

means the value of $\phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2})$ and $\phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})$ are different. This is equivalent to say that

$$|\phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})| = \begin{cases} 1, & \text{if } \mathbf{x} \in \mathcal{A}_1 \\ 0, & \text{if } \mathbf{x} \in \mathcal{A}_1^c \end{cases}. \quad (29)$$

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

$$\begin{aligned} |\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2}| &\leq |\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2} - \mathbf{w}_{j_2}^T \mathbf{x}_{i_2}| \\ &\leq \|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\| \cdot \|\mathbf{x}_{i_2}\|. \end{aligned} \quad (30)$$

Define a set \mathcal{A}_2 such that

$$\begin{aligned} \mathcal{A}_2 &= \left\{ \mathbf{x} \mid \frac{|\mathbf{w}_{j_2}^{*T} \mathbf{x}|}{\|\mathbf{w}_{j_2}^*\| \|\mathbf{x}\|} \leq \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \right\} \\ &= \left\{ \theta_{\mathbf{x}, \mathbf{w}_{j_2}^*} \mid \cos \theta_{\mathbf{x}, \mathbf{w}_{j_2}^*} \leq \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \right\}. \end{aligned} \quad (31)$$

Hence, we have that

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})|^2 \\ &= \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}_{i_2}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}_{i_2})| \\ &= \text{Prob}(\mathbf{x}_{i_2} \in \mathcal{A}_1) \\ &\leq \text{Prob}(\mathbf{x}_{i_2} \in \mathcal{A}_2). \end{aligned} \quad (32)$$

Since $\mathbf{x}_{i_2} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\theta_{\mathbf{x}_{i_2}, \mathbf{w}_{j_2}^*}$ belongs to the uniform distribution on $[-\pi, \pi]$, we have

$$\begin{aligned} \text{Prob}(\mathbf{x}_{i_2} \in \mathcal{A}_2) &= \frac{\pi - \arccos \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|}}{\pi} \\ &\leq \frac{1}{\pi} \tan\left(\pi - \arccos \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|}\right) \\ &= \frac{1}{\pi} \cot\left(\arccos \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|}\right) \\ &\leq \frac{2}{\pi} \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|}. \end{aligned} \quad (33)$$

Hence, (28) and (33) suggest that

$$I_1(i_1, i_2) \leq \frac{6}{\pi} \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \|\mathbf{a}\|^2. \quad (34)$$

The same bound that shown in (34) holds for $I_2(i_1, i_2)$ as well. Therefore, we have the error bound between $\nabla^2 f(\mathbf{W}^*)$ and $\nabla^2 f(\mathbf{W})$ as

$$\begin{aligned} &\|\nabla^2 f(\mathbf{W}^*) - \nabla^2 f(\mathbf{W})\| \\ &\leq \sum_{j_1=1}^K \sum_{j_2=1}^K \max_{j_1, j_2} \left\| \left(\nabla^2 f(\mathbf{W}^*) - \nabla^2 f(\mathbf{W}) \right)_{j_1, j_2} \right\| \\ &\leq \sum_{j_1=1}^K \sum_{j_2=1}^K \frac{12M^2}{\pi} \max_{j_1, j_2} \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \\ &\leq 4M^2 K^2 \frac{\|\mathbf{W}^* - \mathbf{W}\|_2}{\sigma_K}. \end{aligned} \quad (35)$$

□

III. ADDITIONAL PROOFS OF LEMMAS IN APPENDIX D

The proofs of Lemmas 11 and 12 are based on the moment generation function (MGF) and the Chernoff bound. The major difficulty is to obtain the tight bound for the MGF. Lemmas 1 and 2 are standard techniques in evaluating the bounds for the matrix norms.

Lemma 1 (Lemma 5.2, [38]). *Let $\mathcal{B}(0, 1) \in \{\boldsymbol{\alpha} \mid \|\boldsymbol{\alpha}\|_2 = 1, \boldsymbol{\alpha} \in \mathbb{R}^d\}$ denote a unit ball in \mathbb{R}^d . Then, a subset \mathcal{S}_ξ is called a ξ -net of $\mathcal{B}(0, 1)$ if every point $\mathbf{z} \in \mathcal{B}(0, 1)$ can be approximated to within ξ by some point $\boldsymbol{\alpha} \in \mathcal{S}_\xi$, i.e. $\|\mathbf{z} - \boldsymbol{\alpha}\|_2 \leq \xi$. Then the minimal cardinality of a ξ -net \mathcal{S}_ξ satisfies*

$$|\mathcal{S}_\xi| \leq (1 + 2/\xi)^d. \quad (36)$$

Lemma 2 (Lemma 5.3, [38]). *Let \mathbf{A} be an $N \times d$ matrix, and let \mathcal{S}_ξ be a ξ -net of $\mathcal{B}(0, 1)$ in \mathbb{R}^d for some $\xi \in (0, 1)$. Then*

$$\|\mathbf{A}\|_2 \leq (1 - \xi)^{-1} \max_{\boldsymbol{\alpha} \in \mathcal{S}_\xi} |\boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha}|. \quad (37)$$

A. Proof of Lemma 11

Proof of Lemma 11. From Lemmas 1 and 2, we know there exists a subset $\mathcal{S} \subset \{\boldsymbol{\alpha} \mid \|\boldsymbol{\alpha}\|_2 = 1, \boldsymbol{\alpha} \in \mathbb{R}^d\}$ with $|\mathcal{S}| \leq 5^d$ such that

$$\begin{aligned} &\left\| \mathbb{E} \mathbf{X} \mathbf{X} h_1(\mathbf{X}) h_2(\mathbf{X}) - \frac{1}{N} \sum_{n=1}^N \mathbf{X}_n \mathbf{X}_n h_1(\mathbf{X}) h_2(\mathbf{X}_n) \right\|_2 \\ &\leq 2 \cdot \max_{\boldsymbol{\alpha} \in \mathcal{S}} \left| \boldsymbol{\alpha}^T \left(\mathbb{E} \mathbf{X} \mathbf{X}^T h_1(\mathbf{X}) h_2(\mathbf{X}) \right. \right. \\ &\quad \left. \left. - \frac{1}{N} \sum_{n=1}^N \mathbf{X}_n \mathbf{X}_n^T h_1(\mathbf{X}) h_2(\mathbf{X}_n) \right) \boldsymbol{\alpha} \right|. \end{aligned}$$

For simplification, let $Z = (\boldsymbol{\alpha}^T \mathbf{X})^2 h_1(\mathbf{X}) h_2(\mathbf{X})$. Since $\mathbf{X} h_1(\mathbf{X})$ and $\mathbf{X} h_2(\mathbf{X})$ belong to the sub-Gaussian distribution, we have

$$\begin{aligned} \sup_{p \geq 1} (\mathbb{E} |Z|^p)^{1/p} &\leq \|\mathbf{X} h_1(\mathbf{X})\|_{\psi_2} \sqrt{p} \cdot \|\mathbf{X} h_2(\mathbf{X})\|_{\psi_2} \sqrt{p} \\ &\leq \|\mathbf{X} h_1(\mathbf{X})\|_{\psi_2} \|\mathbf{X} h_2(\mathbf{X})\|_{\psi_2} p. \end{aligned} \quad (38)$$

Hence, Z belongs to the sub-exponential distribution from Definition 2. Then, by applying the Chernoff bound, we have

$$\begin{aligned} &\text{Prob}\left(\frac{1}{N} \sum_{i=1}^N (Z_n - \mathbb{E}Z) > t\right) \\ &\leq \text{Prob}\left(\sum_{i=1}^N (Z_n - \mathbb{E}Z) > Nt\right) \\ &\leq e^{C\|Z\|_{\psi_1}^2 N s^2} / e^{s N t} = e^{-N t^2 / (C\|Z\|_{\psi_1}^2)}, \end{aligned} \quad (39)$$

where the last equality holds by choosing $s = t / (C\|Z\|_{\psi_1})$.

Then, by choosing $t = \sqrt{\frac{Cd \log d}{N}} \|Z\|_{\psi_1}$, for any $\boldsymbol{\alpha} \in \mathcal{S}$, we have

$$\begin{aligned} \left| \sum_{n=1}^N Z_n - \mathbb{E}Z \right| &\leq \sqrt{\frac{Cd \log d}{N}} \|Z\|_{\psi_1} \\ &\leq \sqrt{\frac{Cd \log d}{N}} \|\mathbf{X} h_1(\mathbf{X})\|_{\psi_2} \|\mathbf{X} h_2(\mathbf{X})\|_{\psi_2} \end{aligned} \quad (40)$$

with probability at least $1 - d^{-d}$. Since $|\mathcal{S}| \leq 5^d$, we have

$$\begin{aligned} & \left\| \mathbb{E} \mathbf{X} \mathbf{X} h_1(\mathbf{X}) h_2(\mathbf{X}) - \frac{1}{N} \sum_{n=1}^N \mathbf{X}_n \mathbf{X}_n h_1(\mathbf{X}) h_2(\mathbf{X}_n) \right\|_2 \\ & \lesssim \sqrt{\frac{d \log d}{N}} \|\mathbf{X} h_1(\mathbf{X})\|_{\psi_2} \|\mathbf{X} h_2(\mathbf{X})\|_{\psi_2} \end{aligned}$$

with probability at least $1 - (5/d)^d$. Since d is greater than any constant number, $1 - (5/d)^d$ is greater than $1 - d^{-10}$. We use $1 - d^{-10}$ as the probability to be consistent with other contents throughout this paper. \square

1) Proof of Lemma 12:

Proof of Lemma 12. From Lemmas 1 and 2, we know there exists a subset $\mathcal{S} \subset \{\boldsymbol{\alpha} \mid \|\boldsymbol{\alpha}\|_2 = 1, \boldsymbol{\alpha} \in \mathbb{R}^d\}$ with $|\mathcal{S}| \leq 5^d$ such that

$$\begin{aligned} & \left\| \mathbb{E} \mathbf{X}_1 \mathbf{X}_2 - \frac{1}{N} \sum_{n=1}^N \mathbf{X}_{1,n} \mathbf{X}_{2,n} \right\|_2 \\ & \leq 2 \cdot \max_{\boldsymbol{\alpha} \in \mathcal{S}} \left| \boldsymbol{\alpha}^T \left(\mathbb{E} \mathbf{X}_1 \mathbf{X}_2 - \sum_{n=1}^N \mathbf{X}_{1,n} \mathbf{X}_{2,n} \right) \boldsymbol{\alpha} \right|. \end{aligned} \quad (41)$$

For simplification, let \tilde{X}_1 and \tilde{X}_2 denote the random variable $\boldsymbol{\alpha}^T \mathbf{X}_1$ and $\boldsymbol{\alpha}^T \mathbf{X}_2$, respectively. Then, We have

$$\begin{aligned} & \boldsymbol{\alpha}^T \left(\mathbb{E} \mathbf{X}_1 \mathbf{X}_2 - \sum_{n=1}^N \mathbf{X}_{1,n} \mathbf{X}_{2,n} \right) \boldsymbol{\alpha} \\ & = \frac{1}{N} \sum_{n=1}^N \tilde{X}_{1,n} \tilde{X}_{2,n} - \mathbb{E}_{\tilde{X}_1, \tilde{X}_2} \tilde{X}_1 \tilde{X}_2 \\ & = \frac{1}{N} \sum_{n=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1) (\tilde{X}_{2,n} - \mathbb{E} \tilde{X}_2) \\ & \quad + \frac{1}{N} \sum_{n=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1) \mathbb{E} \tilde{X}_2 \\ & \quad + \frac{1}{N} \sum_{n=1}^N (\tilde{X}_{2,n} - \mathbb{E} \tilde{X}_2) \mathbb{E} \tilde{X}_1, \end{aligned} \quad (42)$$

where the last equality holds because \tilde{X}_1 and \tilde{X}_2 are independent. For $\frac{1}{N} \sum_{n=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1)$, since $\tilde{X}_{1,n}$ follows the sub-Gaussian distribution, we know that

$$\mathbb{E} e^{(\tilde{X}_1 - \mathbb{E} \tilde{X}_1)s} \leq e^{C \|\mathbf{X}_{1,n}\|_{\psi_2}^2 s^2} \text{ for all } s \in \mathbb{R} \quad (43)$$

with some constant $C > 0$. By applying the Chernoff bound, we have

$$\begin{aligned} & \text{Prob} \left(\frac{1}{N} \sum_{i=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1) > t \right) \\ & \leq \text{Prob} \left(\sum_{i=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1) > Nt \right) \\ & \leq e^{C \|\mathbf{X}_{1,n}\|_{\psi_2}^2 N s^2} / e^{s N t} = e^{-N t^2 / (C \|\mathbf{X}_{1,n}\|_{\psi_2}^2)}, \end{aligned} \quad (44)$$

where the last equality holds by choosing $s = t / (C \|\mathbf{X}_{1,n}\|_{\psi_2}^2)$. Then, by choosing $t = \sqrt{\frac{C d \log d}{N}} \|\mathbf{X}_{1,n}\|_{\psi_2}$ in (44), we have

$$\left| \frac{1}{N} \sum_{n=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1) \right| \lesssim \sqrt{\frac{d \log d}{N}} \|\mathbf{X}_{1,n}\|_{\psi_2} \quad (45)$$

with probability at least $1 - 2 \cdot d^{-d}$. A similar result to (45) holds for $\frac{1}{N} \sum_{n=1}^N (\tilde{X}_{2,n} - \mathbb{E} \tilde{X}_2)$ as well.

Next, we have

$$\begin{aligned} & \mathbb{E}_{\tilde{X}_1, \tilde{X}_2} e^{(\tilde{X}_1 - \mathbb{E} \tilde{X}_1)(\tilde{X}_2 - \mathbb{E} \tilde{X}_2)s} \\ & \leq \mathbb{E}_{\tilde{X}_2} \mathbb{E}_{\tilde{X}_1} e^{(\tilde{X}_1 - \mathbb{E} \tilde{X}_1)((\tilde{X}_2 - \mathbb{E} \tilde{X}_2)s)} \\ & \leq \mathbb{E}_{\tilde{X}_2} e^{C \|\mathbf{X}_{1,n}\|_{\psi_2}^2 (\tilde{X}_2 - \mathbb{E} \tilde{X}_2)^2 s^2} \\ & \stackrel{(a)}{\leq} 1 / \sqrt{1 - C^2 \|\mathbf{X}_{1,n}\|_{\psi_2}^2 \|\mathbf{X}_{2,n}\|_{\psi_2}^2 s^2}, \end{aligned} \quad (46)$$

for all $|s| \leq (C \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2})^{-1}$, where the details of (a) can be found in following contents.

Since \tilde{X}_2 belongs to sub-Gaussian distribution, we have $\mathbb{E} e^{s(\tilde{X}_2 - \mathbb{E} \tilde{X}_2)} \leq e^{C \|\mathbf{X}_{2,n}\|_{\psi_2}^2 s^2}$ for any $s \in \mathbb{R}$. Then, for any $\lambda \in (0, 1)$,

$$\mathbb{E} e^{s(\tilde{X}_2 - \mathbb{E} \tilde{X}_2) - \frac{C \|\mathbf{X}_{2,n}\|_{\psi_2}^2 s^2}{\lambda}} \leq e^{C \|\mathbf{X}_{2,n}\|_{\psi_2}^2 s^2 (1 - \frac{1}{\lambda})} \quad (47)$$

holds as well for any $s \in \mathbb{R}$. Next, we integrate over s on the both sides of (47),

$$\begin{aligned} & \int \mathbb{E} e^{s(\tilde{X}_2 - \mathbb{E} \tilde{X}_2) - \frac{C \|\mathbf{X}_{2,n}\|_{\psi_2}^2 s^2}{\lambda}} ds \\ & \leq \int e^{C \|\mathbf{X}_{2,n}\|_{\psi_2}^2 s^2 (1 - \frac{1}{\lambda})} ds. \end{aligned} \quad (48)$$

That is

$$e^{\frac{\lambda}{C \|\mathbf{X}_{2,n}\|_{\psi_2}^2} (\tilde{X}_2 - \mathbb{E} \tilde{X}_2)^2} \leq \frac{1}{\sqrt{1 - \lambda}} \text{ for any } \lambda \in (0, 1). \quad (49)$$

It is clear that we can find some $\lambda \in (0, 1)$ such that $\lambda = C_7^2 s^2$ for any $s \in (-1/C_7, 1/C_7)$. That is to say,

$$\begin{aligned} & \mathbb{E}_{\tilde{X}_2} e^{C \|\mathbf{X}_{1,n}\|_{\psi_2}^2 (\tilde{X}_2 - \mathbb{E} \tilde{X}_2)^2 s^2} \\ & \leq \frac{1}{\sqrt{1 - C^2 \|\mathbf{X}_{1,n}\|_{\psi_2}^2 \|\mathbf{X}_{2,n}\|_{\psi_2}^2 s^2}}, \end{aligned} \quad (50)$$

for all $|s| \leq (C \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2})^{-1}$.

By applying the Chernoff bound and choosing $s = (2C \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2})^{-1}$, we have

$$\begin{aligned} & \text{Prob} \left(\frac{1}{N} \sum_{i=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1) (\tilde{X}_{2,n} - \mathbb{E} \tilde{X}_2) > t \right) \\ & \leq 2e^{-Nt / (2C \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2})}. \end{aligned} \quad (51)$$

Then, let $t = \frac{2Cd \log d}{N} \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2}$, we have

$$\begin{aligned} & \left| \frac{1}{N} \sum_{i=1}^N (\tilde{X}_{1,n} - \mathbb{E} \tilde{X}_1) (\tilde{X}_{2,n} - \mathbb{E} \tilde{X}_2) \right| \\ & \leq \frac{2Cd \log d}{N} \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2} \end{aligned} \quad (52)$$

with probability at least $1 - 4 \cdot d^{-d}$.

Hence, for any $\alpha \in \mathcal{S}$, with probability at least $1 - 8 \cdot d^{-d}$,

$$\begin{aligned}
& \left| \alpha^T \left(\Sigma(\mathbb{E}\mathbf{X}_1\mathbf{X}_2 - \frac{1}{N} \sum_{n=1}^N \mathbf{X}_{1,n}\mathbf{X}_{2,n}) \right) \alpha \right| \\
& \leq \left| \frac{1}{N} \sum_{n=1}^N (\tilde{X}_{1,n} - \mathbb{E}\tilde{X}_1)(\tilde{X}_{2,n} - \mathbb{E}\tilde{X}_2) \right| \\
& \quad + \left| \frac{1}{N} \sum_{n=1}^N (\tilde{X}_{1,n} - \mathbb{E}\tilde{X}_1) \right| \cdot |\mathbb{E}\tilde{X}_2| \\
& \quad + \left| \frac{1}{N} \sum_{n=1}^N (\tilde{X}_{2,n} - \mathbb{E}\tilde{X}_2) \right| \cdot |\mathbb{E}\tilde{X}_1| \\
& \lesssim \left(\frac{d \log d}{N} + \sqrt{\frac{d \log d}{N}} + \sqrt{\frac{d \log d}{N}} \right) \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2} \\
& \lesssim \sqrt{\frac{d \log d}{N}} \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2}.
\end{aligned}$$

Since $|\mathcal{S}| \leq 5^d$, we have

$$\begin{aligned}
& \left\| \mathbb{E}\mathbf{X}_1\mathbf{X}_2 - \frac{1}{N} \sum_{n=1}^N \mathbf{X}_{1,n}\mathbf{X}_{2,n} \right\|_2 \\
& \lesssim \sqrt{\frac{d \log d}{N}} \|\mathbf{X}_{1,n}\|_{\psi_2} \|\mathbf{X}_{2,n}\|_{\psi_2}
\end{aligned} \tag{53}$$

with probability at least $1 - 8 \cdot (5/d)^d$. Since d is greater than any constant number, $1 - (5/d)^d$ is greater than $1 - 8 \cdot d^{-10}$. We use $1 - d^{-10}$ as the probability to be consistent with other contents throughout this paper. \square