

---

# Why Lottery Ticket Wins? A Theoretical Perspective of Sample Complexity on Pruned Neural Networks

---

**Shuai Zhang**

Rensselaer Polytechnic Institute  
Troy, NY, USA 12180  
zhangs21@rpi.edu

**Meng Wang**

Rensselaer Polytechnic Institute  
Troy, NY, USA 12180  
wangm7@rpi.edu

**Sijia Liu**

Michigan State University  
East Lansing, MI, USA 48824  
MIT-IBM Watson AI Lab, IBM Research  
liusiji5@msu.edu

**Pin-Yu Chen**

IBM Research  
Yorktown Heights, NY, USA 10562  
Pin-Yu.Chen@ibm.com

**Jinjun Xiong**

University at Buffalo  
Buffalo NY, USA 14260  
jinjun@buffalo.edu

## Abstract

The *lottery ticket hypothesis* (LTH) [20] states that learning on a properly pruned network (the *winning ticket*) improves test accuracy over the original unpruned network. Although LTH has been justified empirically in a broad range of deep neural network (DNN) involved applications like computer vision and natural language processing, the theoretical validation of the improved generalization of a winning ticket remains elusive. To the best of our knowledge, our work, for the first time, characterizes the performance of training a pruned neural network by analyzing the geometric structure of the objective function and the sample complexity to achieve zero generalization error. We show that the convex region near a desirable model with guaranteed generalization enlarges as the neural network model is pruned, indicating the structural importance of a winning ticket. Moreover, when the algorithm for training a pruned neural network is specified as an (accelerated) stochastic gradient descent algorithm, we theoretically show that the number of samples required for achieving zero generalization error is proportional to the number of the non-pruned weights in the hidden layer. With a fixed number of samples, training a pruned neural network enjoys a faster convergence rate to the desired model than training the original unpruned one, providing a formal justification of the improved generalization of the winning ticket. Our theoretical results are acquired from learning a pruned neural network of one hidden layer, while experimental results are further provided to justify the implications in pruning multi-layer neural networks.

## 1 Introduction

Neural network pruning can reduce the computational cost of model training and inference significantly and potentially lessen the chance of overfitting [33, 26, 15, 25, 28, 51, 58, 41]. The recent *Lottery Ticket Hypothesis* (LTH) [20] claims that a randomly initialized dense neural network al-

ways contains a so-called “winning ticket,” which is a sub-network bundled with the corresponding initialization, such that when trained in isolation, this winning ticket can achieve at least the same testing accuracy as that of the original network by running at most the same amount of training time. This so-called “improved generalization of winning tickets” is verified empirically in [20]. LTH has attracted a significant amount of recent research interests [45, 70, 39]. Despite the empirical success [19, 63, 55, 11], the theoretical justification of winning tickets remains elusive except for a few recent works. [39] provides the first theoretical evidence that within a randomly initialized neural network, there exists a good sub-network that can achieve the same test performance as the original network. Meanwhile, recent work [42] trains neural network by adding the  $\ell_1$  regularization term to obtain a relatively sparse neural network, which has a better performance numerically.

However, the theoretical foundation of network pruning is limited. The existing theoretical works usually focus on finding a sub-network that achieves a tolerable loss in either expressive power or training accuracy, compared with the original dense network [2, 71, 61, 43, 4, 3, 35, 5, 59]. To the best of our knowledge, there exists no theoretical support for the *improved* generalization achieved by winning tickets, i.e., pruned networks with faster convergence and better test accuracy.

**Contributions:** This paper provides the *first* systematic analysis of learning pruned neural networks with a finite number of training samples in the oracle-learner setup, where the training data are generated by a unknown neural network, the *oracle*, and another network, the *learner*, is trained on the dataset. Our analytical results also provide a justification of the LTH from the perspective of the sample complexity. In particular, we provide the *first* theoretical justification of the improved generalization of winning tickets. Specific contributions include:

1. **Pruned neural network learning via accelerated gradient descent (AGD):** We propose an AGD algorithm with tensor initialization to learn the pruned model from training samples. Our algorithm converges to the oracle model linearly, which has guaranteed generalization.
2. **First sample complexity analysis for pruned networks:** We characterize the required number of samples for successful convergence, termed as the *sample complexity*. Our sample complexity bound depends linearly on the number of the non-pruned weights and is a significant reduction from directly applying conventional complexity bounds in [69, 66, 67].
3. **Characterization of the benign optimization landscape of pruned networks:** We show analytically that the empirical risk function has an *enlarged* convex region for a pruned network, justifying the importance of a good sub-network (i.e., the winning ticket).
4. **Characterization of the improved generalization of winning tickets:** We show that gradient-descent methods converge faster to the oracle model when the neural network is properly pruned, or equivalently, learning on a pruned network returns a model closer to the oracle model with the same number of iterations, indicating the improved generalization of winning tickets.

**Notations.** Vectors are bold lowercase, matrices and tensors are bold uppercase. Scalars are in normal font, and sets are in calligraphy and blackboard bold font.  $\mathbf{I}$  denote the identity matrix.  $\mathbb{N}$  and  $\mathbb{R}$  denote the sets of nature number and real number, respectively.  $\|z\|$  denotes the  $\ell_2$ -norm of a vector  $z$ , and  $\|Z\|_2$ ,  $\|Z\|_F$  and  $\|Z\|_\infty$  denote the spectral norm, Frobenius norm and the maximum value of matrix  $Z$ , respectively.  $[Z]$  stands for the set of  $\{1, 2, \dots, Z\}$  for any number  $Z \in \mathbb{N}$ . In addition,  $f(r) = \mathcal{O}(g(r))$  ( or  $f(r) = \Omega(g(r))$  ) if  $f \leq C \cdot g$  ( or  $f \geq C \cdot g$  ) for some constant  $C > 0$  when  $r$  is large enough.  $f(r) = \Theta(g(r))$  if both  $f(r) = \mathcal{O}(g(r))$  and  $f(r) = \Omega(g(r))$  holds, where  $c \cdot g \leq f \leq C \cdot g$  for some constant  $0 \leq c \leq C$  when  $r$  is large enough.

## 1.1 Related Work

**Network pruning.** Network pruning methods seek a compressed model while maintaining the expressive power. Numerical experiments have shown that over 90% of the parameters can be pruned without a significant performance loss [10]. Examples of pruning methods include irregular weight pruning [25], structured weight pruning [57], neuron-based pruning [28], and projecting the weights to a low-rank subspace [13].

**Winning tickets.** [20] employs an *Iterative Magnitude Pruning* (IMP) algorithm to obtain the proper sub-network and initialization. IMP and its variations [22, 46] succeed in deeper networks like Residual Networks (Resnet)-50 and Bidirectional Encoder Representations from Transformers (BERT) network [11]. [21] shows that IMP succeeds in finding the “winning ticket” if the ticket is

stable to stochastic gradient descent noise. In parallel, [36] shows numerically that the “winning ticket” initialization does not improve over a random initialization once the correct sub-networks are found, suggesting that the benefit of “winning ticket” mainly comes from the sub-network structures. [18] analyzes the sample complexity of IMP from the perspective of recovering a sparse vector in a linear model rather than learning neural networks.

**Feature sparsity.** High-dimensional data often contains redundant features, and only a subset of the features is used in training [6, 14, 27, 60, 68]. Conventional approaches like wrapper and filter methods score the importance of each feature in a certain way and select the ones with highest scores [24]. Optimization-based methods add variants of the  $\ell_0$  norm as a regularization to promote feature sparsity [68]. Different from network pruning where the feature dimension still remains high during training, the feature dimension is significantly reduced in training when promoting feature sparsity.

**Over-parameterized model.** When the number of weights in a neural network is much larger than the number of training samples, the landscape of the objective function of the learning problem has no spurious local minima, and first-order algorithms converge to one of the global optima [37, 44, 64, 50, 9, 49, 38]. However, the global optima is not guaranteed to generalize well on testing data [62, 64].

**Generalization analyses.** The existing generalization analyses mostly fall within three categories. One line of research employs the Mean Field approach to model the training process by a differential equation assuming infinite network width and infinitesimal training step size [12, 40, 56]. Another approach is the neural tangent kernel (NTK) [30], which requires strong and probably unpractical over-parameterization such that the nonlinear neural network model behaves as its linearization around the initialization [1, 17, 72, 73]. The third line of works follow the oracle-learner setup, where the data are generated by an unknown oracle model, and the learning objective is to estimate the oracle model, which has a generalization guarantee on testing data. However, the objective function has intractably many spurious local minima even for one-hidden-layer neural networks [48, 47, 64]. Assuming an infinite number of training samples, [8, 16, 52] develop learning methods to estimate the oracle model. [23, 69, 66, 67] extend to the practical case of a finite number of samples and characterize the sample complexity for recovering the oracle model. Because the analysis complexity explodes when the number of hidden layers increases, all the analytical results about estimating the oracle model are limited to one-hidden-layer neural networks, and the input distribution is often assumed to be the standard Gaussian distribution.

## 2 Problem Formulation

In an oracle-learner model, given any input  $\mathbf{x} \in \mathbb{R}^d$ , the corresponding output  $y$  is generated by a pruned one-hidden-layer neural network, called **oracle**, as shown in Figure 1. The oracle network is equipped with  $K$  neurons where the  $j$ -th neuron is connected to any arbitrary  $r_j^*$  ( $r_j^* \leq d$ ) input features. Let  $\mathbf{W}^* = [\mathbf{w}_1^*, \dots, \mathbf{w}_K^*] \in \mathbb{R}^{d \times K}$  denotes all the weights (pruned ones are represented by zero). The number of non-zero entries in  $\mathbf{w}_j^*$  is at most  $r_j^*$ . The oracle network is not unique because permuting neurons together with the corresponding weights does not change the output. Therefore, the output label  $y$  obtained by the oracle network satisfies <sup>1</sup>

$$y = \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_j^{*T} \mathbf{x}) + \xi := g(\mathbf{x}; \mathbf{W}^*) + \xi = g(\mathbf{x}; \mathbf{W}^* \mathbf{P}) + \xi, \quad (1)$$

where  $\xi$  is arbitrary unknown additive noise bounded by some constant  $|\xi|$ ,  $\phi$  is the rectified linear unit (ReLU) activation function with  $\phi(z) = \max\{z, 0\}$ , and  $\mathbf{P} \in \{0, 1\}^{K \times K}$  is any permutation matrix.  $\mathbf{M}^*$  is a **mask matrix** for the oracle network, such that  $M_{j,i}^*$  equals to 1 if the weight  $w_{j,i}^*$  is not pruned, and 0 otherwise. Then,  $\mathbf{M}^*$  is an indicator matrix for the non-zero entries of  $\mathbf{W}^*$  with  $\mathbf{M}^* \odot \mathbf{W}^* = \mathbf{W}^*$ , where  $\odot$  is entry-wise multiplication.

Based on  $N$  pairs of training samples  $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$  generated by the oracle, we train on a **learner** network equipped with the same number of neurons in the oracle network. However, the  $j$ -th neuron in the learner network is connected to  $r_j$  input features rather than  $r_j^*$ . Let  $r_{\min}$ ,  $r_{\max}$ , and  $r_{\text{ave}}$  denote the minimum, maximum, and average value of  $\{r_j\}_{j=1}^K$ , respectively. Let  $\mathbf{M}$  denote the

<sup>1</sup>It is extendable to binary classification, and the output is generated by  $\text{Prob}(y_n = 1 | \mathbf{x}_n) = g(\mathbf{x}_n; \mathbf{W}^*)$ .

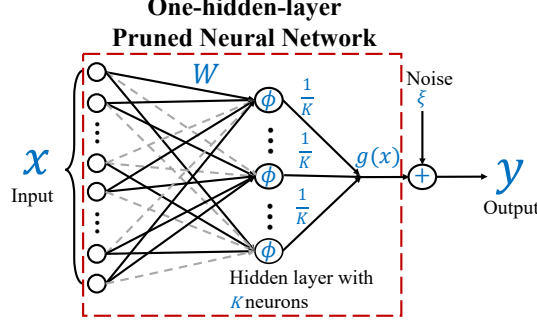


Figure 1: Illustration of the model

mask matrix with respect to the learner network, and  $w_j$  is the  $j$ -th column of  $\mathbf{W}$ . The empirical risk function is defined as

$$\hat{f}_{\mathcal{D}}(\mathbf{W}) = \frac{1}{2N} \sum_{n=1}^N \left( \frac{1}{K} \sum_{j=1}^K \phi(w_j^T \mathbf{x}_n) - y_n \right)^2. \quad (2)$$

When the mask  $\mathbf{M}$  is given, the learning objective is to estimate a proper weight matrix  $\mathbf{W}$  for the learner network from the training samples  $\mathcal{D}$  via solving

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times K}} \hat{f}_{\mathcal{D}}(\mathbf{W}) \quad \text{s.t.} \quad \mathbf{M} \odot \mathbf{W} = \mathbf{W}. \quad (3)$$

$\mathbf{M}$  is called an **accurate mask** if the support of  $\mathbf{M}$  covers the support of a permutation of  $\mathbf{M}^*$ , i.e., there exists a permutation matrix  $\mathbf{P}$  such that  $(\mathbf{M}^* \mathbf{P}) \odot \mathbf{M} = \mathbf{M}^*$ . When  $\mathbf{M}$  is accurate, and  $\xi = 0$ , there exists a permutation matrix  $\mathbf{P}$  such that  $\mathbf{W}^* \mathbf{P}$  is a global optimizer to (3). Hence, if  $\mathbf{W}^* \mathbf{P}$  can be estimated by solving (3), one can learn the oracle network accurately, which has guaranteed generalization performance on the testing data.

We assume  $\mathbf{x}_n$  is independent and identically distributed from the standard Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I}_{d \times d})$ . The Gaussian assumption is motivated by the data whitening [34] and batch normalization techniques [29] that are commonly used in practice to improve learning performance. Moreover, training one-hidden-layer neural network with multiple neurons has intractable many fake minima [47] without any input distribution assumption. In addition, the theoretical results in Section 3 assume an accurate mask, and inaccurate mask is evaluated empirically in Section 4.

The questions that this paper addresses include: 1. **what algorithm** to solve (3)? 2. what is the **sample complexity** for the accurate estimate of the weights in the oracle network? 3. what is the **impact of the network pruning** on the difficulty of the learning problem and the performance of the learned model?

### 3 Algorithm and Theoretical Results

Section 3.1 studies the geometric structure of (3), and the main results are in Section 3.2. Section 3.3 briefly introduces the proof sketch and technical novelty, and the limitations are in Section 3.4.

#### 3.1 Local Geometric Structure

Theorem 1 characterizes the local convexity of  $\hat{f}_{\mathcal{D}}$  in (3). It has two important implications.

1. **Strictly locally convex near ground truth:**  $\hat{f}_{\mathcal{D}}$  is strictly convex near  $\mathbf{W}^* \mathbf{P}$  for some permutation matrix  $\mathbf{P}$ , and the radius of the convex ball is negatively correlated with  $\sqrt{\tilde{r}}$ , where  $\tilde{r}$  is in the order of  $r_{\text{ave}}$ . Thus, the convex ball enlarges as any  $r_j$  decreases.

2. **Importance of the winning ticket architecture:** Compared with training on the dense network directly, training on a properly pruned sub-network has a larger local convex region near  $\mathbf{W}^* \mathbf{P}$ , which may lead to easier estimation of  $\mathbf{W}^* \mathbf{P}$ . To some extent, this result can be viewed as a theoretical validation of the importance of the winning architecture (a good sub-network) in [20]. Formally, we have

**Theorem 1 (Local Convexity).** Assume the mask  $\mathbf{M}$  of the learner network is accurate. Suppose constants  $\varepsilon_0, \varepsilon_1 \in (0, 1)$  and the number of samples satisfies

$$N = \Omega(\varepsilon_1^{-2} K^4 \tilde{r} \log q), \quad (4)$$

for some large constant  $q > 0$ , where

$$\tilde{r} = \frac{1}{8K^4} \left( \sum_{k=1}^K \sum_{j=1}^K (1 + \delta_{j,k})(r_j + r_k)^{\frac{1}{2}} \right)^2, \quad (5)$$

$\delta_{j,k}$  is 1 if the indices of non-pruned weights in the  $j$ -th and  $k$ -th neurons overlap and 0 otherwise. Then, there exists a permutation matrix  $\mathbf{P}$  such that for any  $\mathbf{W}$  that satisfies

$$\|\mathbf{W} - \mathbf{W}^* \mathbf{P}\|_F = \mathcal{O}\left(\frac{\varepsilon_0}{K^2}\right), \quad \text{and } \mathbf{M} \odot \mathbf{W} = \mathbf{W}, \quad (6)$$

its Hessian of  $\hat{f}_{\mathcal{D}}$ , with probability at least  $1 - K \cdot q^{-r_{\min}}$ , is bounded as:

$$\Theta\left(\frac{1 - \varepsilon_0 - \varepsilon_1}{K^2}\right) \mathbf{I} \preceq \nabla^2 \hat{f}_{\mathcal{D}}(\mathbf{W}) \preceq \Theta\left(\frac{1}{K}\right) \mathbf{I}. \quad (7)$$

**Remark 1.1 (Parameter  $\tilde{r}$ ):** Clearly  $\tilde{r}$  is a monotonically increasing function of any  $r_j$  from (5). Moreover, one can check that  $\frac{1}{8} r_{\text{ave}} \leq \tilde{r} \leq r_{\text{ave}}$ . Hence,  $\tilde{r}$  is in the order of  $r_{\text{ave}}$ .

**Remark 1.2 (Local landscape):** Theorem 1 shows that with enough samples as shown in (4), in a local region of  $\mathbf{W}^* \mathbf{P}$  as shown in (6), all the eigenvalues of the Hessian matrix of the empirical risk function are lower and upper bounded by two positive constants. This property is useful in designing efficient algorithms to recover  $\mathbf{W}^* \mathbf{P}$ , as shown in Section 3.2.

**Remark 1.3 (Size of the convex region):** When the number of samples  $N$  is fixed and  $r$  changes,  $\varepsilon_1$  can be  $\Theta(\sqrt{\tilde{r}/N})$  while (4) is still met.  $\varepsilon_0$  in (7) can be arbitrarily close to but smaller than  $1 - \varepsilon_1$  so that the Hessian matrix is still positive definite. Then from (6), the radius of the convex ball is  $\Theta(1) - \Theta(\sqrt{\tilde{r}/N})$ , indicating an enlarged region when  $\tilde{r}$  decreases. The enlarged convex region serves as an important component in proving the faster convergence rate, summarized in Theorem 2. Besides this, as Figure 1 shown in [20], the authors claim that the learning is stable if the linear interpolation of the learned models with SGD noises still remain similar in performance, which is summarized as the concept “linearly connected region.” Intuitively, we conjecture that the winning ticket shows a better performance in the stability analysis because it has a larger convex region. In the other words, a larger convex region indicates that the learning is more likely to be stable in the linearly connected region.

### 3.2 Convergence Analysis with Accelerated Gradient Descent

We propose to solve the non-convex problem (3) via the accelerated gradient descent (AGD) algorithm, summarized in Algorithm 1. Compared with the vanilla gradient descent (GD) algorithm, AGD has an additional momentum term, denoted by  $\beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)})$ , in each iteration. AGD enjoys a faster convergence rate than vanilla GD in solving optimization problems, including learning neural networks [65]. Vanilla GD can be viewed as a special case of AGD by letting  $\beta = 0$ . The initial point  $\mathbf{W}^{(0)}$  can be obtained through a tensor method, and the details are provided in Appendix B.

---

#### Algorithm 1 Accelerated Gradient Descent (AGD) Algorithm

---

- 1: **Input:** training data  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , gradient step size  $\eta$ , momentum parameter  $\beta$ , and an initialization  $\mathbf{W}^{(0)}$  by the tensor initialization method;
  - 2: Partition  $\mathcal{D}$  into  $T = \log(1/\varepsilon)$  disjoint subsets, denoted as  $\{\mathcal{D}_t\}_{t=1}^T$ ;
  - 3: **for**  $t = 1, 2, \dots, T$  **do**
  - 4:      $\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \cdot \mathbf{M} \odot \nabla \hat{f}_{\mathcal{D}_t}(\mathbf{W}^{(t)}) + \beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)})$
  - 5: **end for**
  - 6: **Return:**  $\mathbf{W}^{(T)}$
- 

The theoretical analyses of our algorithm are summarized in Theorem 2 (convergence) and Lemma 1 (Initialization). The significance of these results can be interpreted from the following aspects.

1. **Linear convergence to the oracle model:** Theorem 2 implies that if initialized in the local convex region, the iterates generated by AGD converge linearly to  $\mathbf{W}^* \mathbf{P}$  for some  $\mathbf{P}$  when noiseless. When there is noise, they converge to a point  $\mathbf{W}^{(T)}$ . The distance between  $\mathbf{W}^{(T)}$  and  $\mathbf{W}^* \mathbf{P}$  is proportional to the noise level and scales in terms of  $\mathcal{O}(\sqrt{\tilde{r}/N})$ . Moreover, when  $N$  is fixed, the convergence rate of AGD is  $\Theta(\sqrt{\tilde{r}/K})$ . Recall that Algorithm 1 reduces to the vanilla GD by setting  $\beta = 0$ . The rate for the vanilla GD algorithm here is  $\Theta(\sqrt{\tilde{r}/K})$  by setting  $\beta = 0$  by Theorem 2, indicating a slower convergence than AGD. Lemma 1 shows the tensor initialization method indeed returns an initial point in the convex region.

2. **Sample complexity for accurate estimation:** We show that the required number of samples for successful estimation of the oracle model is  $\Theta(\tilde{r} \log q \log(1/\varepsilon))$  for some large constant  $q$  and estimation accuracy  $\varepsilon$ . Our sample complexity is much less than the conventional bound of  $\Theta(d \log q \log(1/\varepsilon))$  for one-hidden-layer networks [69, 66, 67]. This is the first theoretical characterization of learning a pruned network from the perspective of sample complexity.

3. **Improved generalization of winning tickets:** We prove that with a fixed number of training samples, training on a properly pruned sub-network converges faster to  $\mathbf{W}^* \mathbf{P}$  than training on the original dense network. Our theoretical analysis justifies that training on the winning ticket can meet or exceed the same test accuracy within the same number of iterations. To the best of our knowledge, our result here provides the first theoretical justification for this intriguing empirical finding of “improved generalization of winning tickets” by [20].

**Theorem 2** (Convergence). *Assume the mask  $\mathbf{M}$  of the learner network is accurate. Suppose  $\mathbf{W}^{(0)}$  satisfies (6) and the number of samples satisfies*

$$N = \Omega(\varepsilon_0^{-2} K^6 \tilde{r} \log q \log(1/\varepsilon)) \quad (8)$$

*for some  $\varepsilon_0 \in (0, 1/2)$ . Let  $\eta = K/14$  in Algorithm 1. Then, the iterates  $\{\mathbf{W}^{(t)}\}_{t=1}^T$  returned by Algorithm 1 converges linearly to  $\mathbf{W}^* \mathbf{P}$  up to the noise level with probability at least  $1 - K^2 T \cdot q^{-r_{\min}}$*

$$\|\mathbf{W}^{(t)} - \mathbf{W}^* \mathbf{P}\|_F \leq \nu(\beta)^t \|\mathbf{W}^{(0)} - \mathbf{W}^* \mathbf{P}\|_F + \mathcal{O}\left(\sum_j \sqrt{\frac{r_j \log q}{N}}\right) \cdot |\xi|, \quad (9)$$

$$\text{and } \|\mathbf{W}^{(T)} - \mathbf{W}^* \mathbf{P}\|_F \leq \varepsilon \|\mathbf{W}^*\|_F + \mathcal{O}\left(\sum_j \sqrt{\frac{r_j \log q}{N}}\right) \cdot |\xi|, \quad (10)$$

*for a fixed permutation matrix  $\mathbf{P}$ , where  $\nu(\beta)$  is the rate of convergence that depends on  $\beta$  with  $\nu(\beta^*) = 1 - \Theta(\frac{1-\varepsilon_0}{\sqrt{K}})$  for some non-zero  $\beta^*$  and  $\nu(0) = 1 - \Theta(\frac{1-\varepsilon_0}{K})$ .*

**Lemma 1** (Initialization). *Assume the noise  $|\xi| \leq \|\mathbf{W}^*\|_2$  and the number of samples  $N = \Omega(\varepsilon_0^{-2} K^5 r_{\max} \log q)$  for  $\varepsilon_0 > 0$  and large constant  $q$ , the tensor initialization method outputs  $\mathbf{W}^{(0)}$  such that (6) holds, i.e.,  $\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F = \mathcal{O}(\frac{\varepsilon_0 \sigma_K}{K^2})$  with probability at least  $1 - q^{-r_{\max}}$ .*

**Remark 2.1** (Faster convergence on pruned network): With a fixed number of samples, when  $\tilde{r}$  decreases,  $\varepsilon_0$  can increase as  $\Theta(\sqrt{\tilde{r}})$  while (8) is still met. Then  $\nu(0) = \Theta(\sqrt{\tilde{r}/K})$  and  $\nu(\beta^*) = \Theta(\sqrt{\tilde{r}/K})$ . Therefore, when  $\tilde{r}$  decreases, both the stochastic and accelerated gradient descent converge faster. Note that as long as  $\mathbf{W}^{(0)}$  is initialized in the local convex region, not necessarily by the tensor method, Theorem 2 guarantees the accurate recovery. [66, 67] analyze AGD on convolutional neural networks, while this paper focuses on network pruning.

**Remark 2.2** (Sample complexity of initialization): From Lemma 1, the required number of samples for a proper initialization is  $\Omega(\varepsilon_0^{-2} K^5 r_{\max} \log q)$ . Because  $r_{\max} \leq K r_{\text{ave}}$  and  $\tilde{r} = \Omega(r_{\text{ave}})$ , this number is no greater than the sample complexity in (8). Thus, provided that (8) is met, Algorithm 1 can estimate the oracle network model accurately.

**Remark 2.3** (Inaccurate mask): The above analyses are based on the assumption that the mask of the learner network is accurate. In practice, a mask can be obtained by an iterative pruning method such as [20] or a one-shot pruning method such as [55]. In Appendix E, we prove that the magnitude pruning method can obtain an accurate mask with enough training samples. Moreover, empirical experiments in Section 4.2 and 4.3 suggest that even if the mask is not accurate, the three properties (linear convergence, sample complexity with respect to the network size, and improved generalization of winning tickets) can still hold. Therefore, our theoretical results provide some insight into the empirical success of network pruning.

### 3.3 The Sketch of Proofs and Technical Novelty

Our proof outline is inspired by [69] on fully connected neural networks, however, major technical changes are made in this paper to generalize the analysis to an arbitrarily pruned network. To characterize the local convex region of  $\hat{f}_{\mathcal{D}}$  (Theorem 1), the idea is to bound the Hessian matrix of the population risk function, which is the expectation of the empirical risk function, locally and then characterize the distance between the empirical and population risk functions through the concentration bounds. Then, the convergence of AGD (Theorem 2) is established based on the desired local curvature, which in turn determines the sample complexity. Finally, to initialize in the local convex region (Lemma 1), we construct tensors that contain the weights information and apply a decomposition method to estimate the weights.

Our technical novelties are as follows. First, a direct application of the results in [69] leads to a sample complexity bound that is linear in the feature dimension  $d$ . We develop new techniques to tighten the sample complexity bound to be linear in  $\tilde{r}$ , which can be significantly smaller than  $d$  for a sufficiently pruned network. Specifically, we develop new concentration bounds (Lemmas 4 and 5 in Appendix) to bound the distance between the population and empirical risk functions rather than using the bound in [69]. Second, instead of restricting the activation to be smooth for convergence analysis, we study the case of ReLU function which is non-smooth. Third, new tensors are constructed for pruned networks (see (21)-(23) in Appendix) in computing the initialization, and our new concentration bounds are employed to reduce the required number of samples for a proper initialization. Last, Algorithm 1 employs AGD and is proved to converge faster than the GD algorithm in [69].

### 3.4 Limitations

Like most theoretical works based on the oracle-learner setup, limitations of this work include (1) one hidden layer only; and (2) the input follows the Gaussian distribution. Extension to multi-layer might be possible if the following technical challenges are addressed. First, when characterizing the local convex region, one needs to show that the Hessian matrix is positive definite. In multi-layer networks, the Hessian matrix is more complicated to compute. Second, new concentration bounds need to be developed because the input feature distributions to the second and third layers depend on the weights in previous layers. Third, the initialization approach needs to be revised. The team is also investigating the other input distributions such as Gaussian mixture models.

## 4 Numerical Experiments

The theoretical results are first verified on synthetic data, and we then analyze the pruning performance on both synthetic and real datasets. In Section 4.1, Algorithm 1 is implemented with minor modification, such that, the initial point is randomly selected as  $\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F / \|\mathbf{W}^*\|_F < \lambda$  for some  $\lambda > 0$  to reduce the computation. Algorithm 1 terminates when  $\|\mathbf{W}^{(t+1)} - \mathbf{W}^{(t)}\|_F / \|\mathbf{W}^{(t)}\|_F$  is smaller than  $10^{-8}$  or reaching 10000 iterations. In Sections 4.2 and 4.3, the Gradient Signal Preservation (GraSP) algorithm [55] and IMP algorithm [10, 20]<sup>2</sup> are implemented to prune the neural networks. As many works like [11, 10, 20] have already verified the faster convergence and better generalization accuracy of the winning tickets empirically, we only include the results of some representative experiments, such as training MNIST and CIFAR-10 on Lenet-5 [32] and Resnet-50 [27] networks, to verify our theoretical findings.

The synthetic data are generated using a oracle model in Figure 1. The input  $\mathbf{x}_n$ 's are randomly generated from Gaussian distribution  $\mathcal{N}(0, \mathbf{I}_{d \times d})$  independently, and indices of non-pruned weights of the  $j$ -th neuron are obtained by randomly selecting  $r_j$  numbers without replacement from  $[d]$ . For the convenience of generating specific  $\tilde{r}$ , the indices of non-pruned weights are almost overlapped ( $\sum_j \sum_k \delta_j \delta_k > 0.95K^2$ ) except for Figure 5. In Figures 2 and 4,  $r_j$  is selected uniformly from  $[0.9\tilde{r}, 1.1\tilde{r}]$  for a given  $\tilde{r}$ , and  $r_j$  are the same in value for all  $j$  in other figures. Each non-zero entry of  $\mathbf{W}^*$  is randomly selected from  $[-0.5, 0.5]$  independently. The noise  $\xi_n$ 's are i.i.d. from  $\mathcal{N}(0, \sigma^2)$ , and the noise level is measured by  $\sigma/E_y$ , where  $E_y$  is the root mean square of the noiseless outputs.

<sup>2</sup>The source codes used are downloaded from [https://github.com/VITA-Group/CV\\_LTH\\_Pre-training](https://github.com/VITA-Group/CV_LTH_Pre-training).

#### 4.1 Evaluation of theoretical findings on synthetic data

**Local convexity near  $\mathbf{W}^*$ .** We set the number of neurons  $K = 10$ , the dimension of the data  $d = 500$  and the sample size  $N = 5000$ . Figure 2 illustrates the success rate of Algorithm 1 when  $\tilde{r}$  changes. The  $y$ -axis is the relative distance of the initialization  $\mathbf{W}^{(0)}$  to the ground-truth. For each pair of  $\tilde{r}$  and the initial distance, 100 trails are constructed with the network weights, training data and the initialization  $\mathbf{W}^{(0)}$  are all generated independently in each trail. Each trail is called successful if the relative error of the solution  $\mathbf{W}$  returned by Algorithm 1, measured by  $\|\mathbf{W} - \mathbf{W}^*\|_F / \|\mathbf{W}^*\|_F$ , is less than  $10^{-4}$ . A black block means Algorithm 1 fails in estimating  $\mathbf{W}^*$  in all trails while a white block indicates all success. As Algorithm 1 succeeds if  $\mathbf{W}^{(0)}$  is in the local convex region near  $\mathbf{W}^*$ , we can see that the radius of convex region is indeed linear in  $-\tilde{r}^{\frac{1}{2}}$ , as predicted by Theorem 1.

**Convergence rate.** Figure 3 shows the convergence rate of Algorithm 1 when  $\tilde{r}$  changes.  $N = 5000$ ,  $d = 300$ ,  $K = 10$ ,  $\eta = 0.5$ , and  $\beta = 0.2$ . Figure 3(a) shows that the relative error decreases exponentially as the number of iterations increases, indicating the linear convergence of Algorithm 1. As shown in Figure 3(b), the results are averaged over 20 trials with different initial points, and the areas in low transparency represent the standard deviation errors. We can see that the convergence rate is almost linear in  $\sqrt{\tilde{r}}$ , as predicted by Theorem 2. We also compare with GD by setting  $\beta$  as 0. One can see that AGD has a smaller convergence rate than GD, indicating faster convergence.

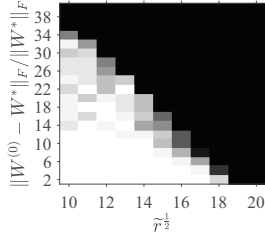


Figure 2: The radius of the local convex region against  $\tilde{r}^{\frac{1}{2}}$

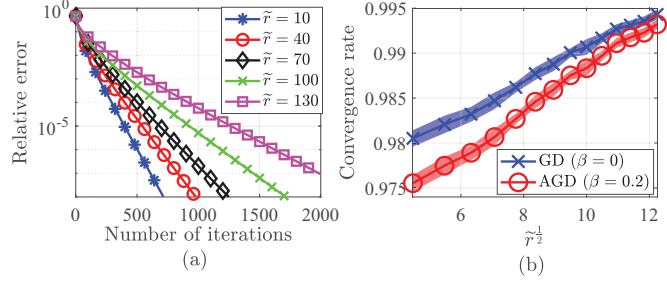


Figure 3: Convergence rate when  $\tilde{r}$  changes

**Sample complexity.** Figures 4 and 5 show the success rate of Algorithm 1 when varying  $N$  and  $\tilde{r}$ .  $d$  is fixed as 100. In Figure 4, we construct 100 independent trails for each pair of  $N$  and  $\tilde{r}$ , where the ground-truth model and training data are generated independently in each trail. One can see that the required number of samples for successful estimation is linear in  $\tilde{r}$ , as predicted by (8). In Figure 5,  $r_j$  is fixed as 20 for all neurons, but different network architectures after pruning are considered. One can see that although the number of remaining weights is the same,  $\tilde{r}$  can be different in different architectures, and the sample complexity increases as  $\tilde{r}$  increases, as predicted by (8).

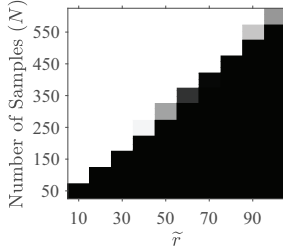


Figure 4: Sample complexity when  $\tilde{r}$  changes

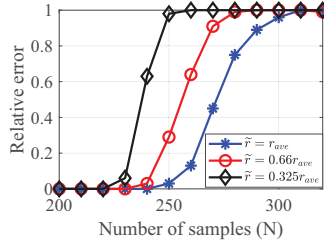


Figure 5: Relative error against  $\tilde{r}$

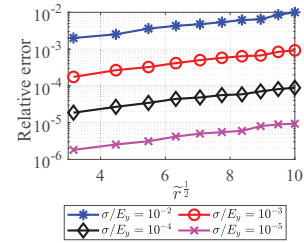


Figure 6: Relative error against  $\tilde{r}^{\frac{1}{2}}$  at different noise level

**Performance in noisy case.** Figure 6 shows the relative error of the learned model by Algorithm 1 from noisy measurements when  $\tilde{r}$  changes.  $N = 1000$ ,  $K = 10$ , and  $d = 300$ . The results are averaged over 100 independent trials, and standard deviation is around 2% to 8% of the corresponding relative errors. The relative error is linear in  $\tilde{r}^{\frac{1}{2}}$ , as predicted by (9). Moreover, the relative error is proportional to the noise level  $|\xi|$ .

#### 4.2 Performance with inaccurate mask on synthetic data

The performance of Algorithm 1 is evaluated when the mask  $\mathbf{M}$  of the learner network is inaccurate. The number of neurons  $K$  is 5. The dimension of inputs  $d$  is 100.  $r_j^*$  of the oracle model is 20 for



all  $j \in [K]$ . GraSP algorithm [55] is employed to find masks based only on early-trained weights in 20 iterations of AGD. The mask accuracy is measured by  $\|\mathbf{M}^* \odot \mathbf{M}\|_0 / \|\mathbf{M}^*\|_0$ , where  $\mathbf{M}^*$  is the mask of the oracle model. The pruning ratio is defined as  $(1 - r_{\text{ave}}/d) \times 100\%$ . The number of training samples  $N$  is 200. The model returned by Algorithm 1 is evaluated on  $N_{\text{test}} = 10^5$  samples, and the test error is measured by  $\sqrt{\sum_n |y_n - \hat{y}_n|^2 / N_{\text{test}}}$ , where  $\hat{y}_n$  is the output of the learned model with the input  $\mathbf{x}_n$ , and  $(\mathbf{x}_n, y_n)$  is the  $n$ -th testing sample generated by the oracle network.

**Improved generalization by GraSP.** Figure 7 shows the test error with different pruning ratios. For each pruning ratio, we randomly generate 1000 independent trials. Because the mask of the learner network in each trail is generated independently, we compute the average test error of the learned models in all the trails with same mask accuracy. If there are less than 10 trails for certain mask accuracy, the result of that mask accuracy is not reported as it is statistically meaningless. The test error decreases as the mask accuracy increases. More importantly, at fixed mask accuracy, the test error decreases as the pruning ratio increases. That means the generalization performance improves when  $\tilde{r}$  decreases, even if the mask is not accurate.

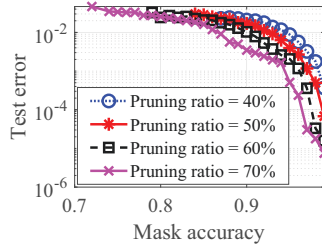


Figure 7: Test error against mask accuracy with different pruning ratios

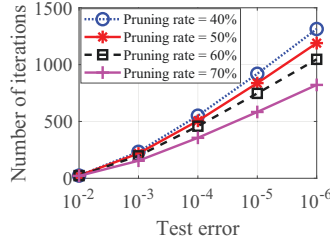


Figure 8: Convergence rate with mask accuracy in [0.85, 0.9]

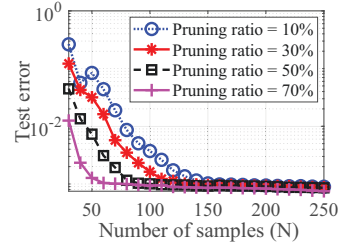


Figure 9: Test error against the number of samples with mask accuracy in [0.85, 0.9]

**Linear convergence.** Figure 8 shows the convergence rate of Algorithm 1 with different pruning ratios. We show the smallest number of iterations required to achieve a certain test error of the learned model, and the results are averaged over the independent trials with mask accuracy between 0.85 and 0.90. Even with inaccurate mask, the test error converges linearly. Moreover, as the pruning ratio increases, Algorithm 1 converges faster.

**Sample complexity with respect to the pruning ratio.** Figure 9 shows the test error when the number of training samples  $N$  changes. All the other parameters except  $N$  remain the same. The results are averaged over the trials with mask accuracy between 0.85 and 0.90. We can see the test error decreases when  $N$  increases. More importantly, as the pruning ratio increases, the required number of samples to achieve the same test error (no less than  $10^{-3}$ ) decreases dramatically. That means the sample complexity decreases as  $\tilde{r}$  decreases even if the mask is inaccurate.

### 4.3 Performance of IMP on synthetic, MNIST and CIFAR-10 datasets

We implement the IMP algorithm to obtain pruned networks on synthetic, MNIST and CIFAR-10 datasets. Figure 10 shows the test performance of a pruned network on synthetic data with different sample sizes. Here in the oracle network model,  $K = 5$ ,  $d = 100$ , and  $r_j^* = 20$  for all  $j \in [K]$ . The noise level  $\sigma/E_y = 10^{-3}$ . One observation is that for a fixed sample size  $N$  greater than 100, the test error decreases as the pruning ratio increases. This verifies that the IMP algorithm indeed prunes the network properly. It also shows that the learned model improves as the pruning progresses, verifying our theoretical result in Theorem 2 that the difference of the learned model from the oracle model decreases as  $r_j$  decreases. The second observation is that the test error decreases as  $N$  increases for any fixed pruning ratio. This verifies our result in Theorem 2 that the difference of the learned model from the oracle model decreases as the number of training samples increases. When the pruning ratio is too large (greater than 80%), the pruned network cannot explain the data properly, and thus the test error is large for all  $N$ . When the number of samples is too small, like  $N = 100$ , the test error is always large, because it does not meet the sample complexity requirement for estimating the oracle model even though the network is properly pruned.

Figures 11 and 12 show the test performance of learned models by implementing the IMP algorithm on MNIST and CIFAR-10 using Lenet-5 [32] and Resnet-50 [27] architecture, respectively. The

experiments follow the standard setup in [10] except for the size of the training sets. To demonstrate the effect of sample complexity, we randomly selected  $N$  samples from the original training set without replacement. As we can see, a properly pruned network (i.e., winning ticket) helps reduce the sample complexity required to reach the test accuracy of the original dense model. For example, training on a pruned network returns a model (e.g.,  $P_1$  and  $P_3$  in Figures 11 and 12) that has better testing performance than a dense model (e.g.,  $P_2$  and  $P_4$  in Figures 11 and 12) trained on a larger data set. Given the number of samples, we consistently find the characteristic behavior of winning tickets: That is, the test accuracy could increase when the pruning ratio increases, indicating the effectiveness of pruning. The test accuracy then drops when the network is overly pruned. The results show that our theoretical characterization of sample complexity is well aligned with the empirical performance of pruned neural networks and explains the improved generalization observed in LTH.

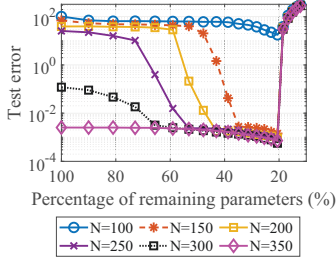


Figure 10: Test error of pruned models on the synthetic dataset

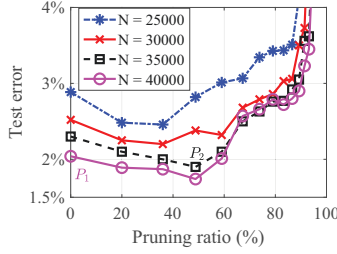


Figure 11: Test accuracy of pruned LeNet-5 on Mnist

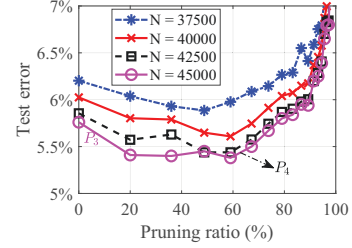


Figure 12: Test accuracy of pruned Resnet-50 on Cifar-10

## 5 Conclusions

This paper provides the first theoretical analysis of learning one-hidden-layer pruned neural networks, which offers formal justification of the improved generalization of winning ticket observed from empirical findings in LTH. We characterize analytically the impact of the number of remaining weights in a pruned network on the required number of samples for training, the convergence rate of the learning algorithm, and the accuracy of the learned model. We also provide extensive numerical validations of our theoretical findings.

## Broader impacts

We see no ethical or immediate societal consequence of our work. This paper contributes to the theoretical foundation of both network pruning and generalization guarantee. The former encourages the development of learning method to reduce the computational cost. The latter increases the public trust in incorporating AI technology in critical domains.

## Acknowledgement

This work was supported by AFOSR FA9550-20-1-0122, ARO W911NF-21-1-0255, NSF 1932196 and the Rensselaer-IBM AI Research Collaboration (<http://airc.rpi.edu>), part of the IBM AI Horizons Network (<http://ibm.biz/AIHorizons>). We thank Tianlong Chen at University of Texas at Austin, Haolin Xiong at Rensselaer Polytechnic Institute and Yihua Zhang at Michigan State University for the help in formulating numerical experiments. We thank all anonymous reviewers for their constructive comments.

## References

- [1] Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019.
- [2] S. Arora, R. Ge, B. Neyshabur, and Y. Zhang. Stronger generalization bounds for deep nets via a compression approach. In J. Dy and A. Krause, editors, *Proceedings of the 35th International*

- Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 254–263, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [3] C. Baykal, L. Liebenwein, I. Gilitschenski, D. Feldman, and D. Rus. Data-dependent coresets for compressing neural networks with applications to generalization bounds. In *International Conference on Learning Representations*, 2018.
  - [4] C. Baykal, L. Liebenwein, I. Gilitschenski, D. Feldman, and D. Rus. Sipping neural networks: Sensitivity-informed provable pruning of neural networks. *arXiv preprint arXiv:1910.05422*, 2019.
  - [5] M. Ben, M. Osadchy, V. Braverman, S. Zhou, and D. Feldman. Data-independent neural pruning via coresets. In *International Conference on Learning Representations (ICLR)*, 2020.
  - [6] M. L. Bermingham, R. Pong-Wong, A. Spiliopoulou, C. Hayward, I. Rudan, H. Campbell, A. F. Wright, J. F. Wilson, F. Agakov, P. Navarro, et al. Application of high-dimensional feature selection: evaluation for genomic prediction in man. *Scientific reports*, 5(1):1–12, 2015.
  - [7] R. Bhatia. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.
  - [8] 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*, pages 605–614. JMLR. org, 2017.
  - [9] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud. Neural ordinary differential equations. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 6572–6583, 2018.
  - [10] T. Chen, J. Frankle, S. Chang, S. Liu, Y. Zhang, M. Carbin, and Z. Wang. The lottery tickets hypothesis for supervised and self-supervised pre-training in computer vision models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16306–16316, 2021.
  - [11] T. Chen, J. Frankle, S. Chang, S. Liu, Y. Zhang, Z. Wang, and M. Carbin. The lottery ticket hypothesis for pre-trained bert networks. *arXiv preprint arXiv:2007.12223*, 2020.
  - [12] L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS’18*, page 3040–3050, 2018.
  - [13] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, and N. de Freitas. Predicting parameters in deep learning. In *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 2*, pages 2148–2156, 2013.
  - [14] V. C. Dinh and L. S. Ho. Consistent feature selection for analytic deep neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 2420–2431. Curran Associates, Inc., 2020.
  - [15] X. Dong, S. Chen, and S. Pan. Learning to prune deep neural networks via layer-wise optimal brain surgeon. In *Advances in Neural Information Processing Systems*, pages 4857–4867, 2017.
  - [16] 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*, pages 1338–1347, 2018.
  - [17] S. S. Du, X. Zhai, B. Póczos, and A. Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019.
  - [18] B. Elesedy, V. Kanade, and Y. W. Teh. Lottery tickets in linear models: An analysis of iterative magnitude pruning. *arXiv preprint arXiv:2007.08243*, 2020.
  - [19] U. Evci, T. Gale, J. Menick, P. S. Castro, and E. Elsen. Rigging the lottery: Making all tickets winners. *International Conference on Machine Learning*, 2020.

- [20] J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019.
- [21] J. Frankle, G. K. Dziugaite, D. Roy, and M. Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pages 3259–3269. PMLR, 2020.
- [22] J. Frankle, G. K. Dziugaite, D. M. Roy, and M. Carbin. Stabilizing the lottery ticket hypothesis. *arXiv preprint arXiv:1903.01611*, 2019.
- [23] H. Fu, Y. Chi, and Y. Liang. Guaranteed recovery of one-hidden-layer neural networks via cross entropy. *IEEE transactions on signal processing*, 68:3225–3235, 2020.
- [24] I. Guyon and A. Elisseeff. An introduction to variable and feature selection. *Journal of machine learning research*, 3(Mar):1157–1182, 2003.
- [25] S. Han, J. Pool, J. Tran, and W. Dally. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pages 1135–1143, 2015.
- [26] B. Hassibi and D. G. Stork. Second order derivatives for network pruning: Optimal brain surgeon. In *Advances in neural information processing systems*, pages 164–171, 1993.
- [27] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [28] H. Hu, R. Peng, Y.-W. Tai, and C.-K. Tang. Network trimming: A data-driven neuron pruning approach towards efficient deep architectures. *arXiv preprint arXiv:1607.03250*, 2016.
- [29] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015.
- [30] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [31] V. Kuleshov, A. Chaganty, and P. Liang. Tensor factorization via matrix factorization. In *Artificial Intelligence and Statistics*, pages 507–516, 2015.
- [32] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [33] Y. LeCun, J. S. Denker, and S. A. Solla. Optimal brain damage. In *Advances in neural information processing systems*, pages 598–605, 1990.
- [34] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–48. Springer, 2012.
- [35] L. Liebenwein, C. Baykal, H. Lang, D. Feldman, and D. Rus. Provable filter pruning for efficient neural networks. In *International Conference on Learning Representations*, 2019.
- [36] Z. Liu, M. Sun, T. Zhou, G. Huang, and T. Darrell. Rethinking the value of network pruning. In *International Conference on Learning Representations*, 2018.
- [37] R. Livni, S. Shalev-Shwartz, and O. Shamir. On the computational efficiency of training neural networks. In *Advances in neural information processing systems*, pages 855–863, 2014.
- [38] Y. Lu, C. Ma, Y. Lu, J. Lu, and L. Ying. A mean field analysis of deep resnet and beyond: Towards provably optimization via overparameterization from depth. In *International Conference on Machine Learning*, pages 6426–6436. PMLR, 2020.

- [39] E. Malach, G. Yehudai, S. Shalev-Schwartz, and O. Shamir. Proving the lottery ticket hypothesis: Pruning is all you need. In *International Conference on Machine Learning*, pages 6682–6691. PMLR, 2020.
- [40] S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- [41] D. Molchanov, A. Ashukha, and D. Vetrov. Variational dropout sparsifies deep neural networks. In *International Conference on Machine Learning*, pages 2498–2507, 2017.
- [42] B. Neyshabur. Towards learning convolutions from scratch. *Advances in Neural Information Processing Systems*, 33, 2020.
- [43] L. Orseau, M. Hutter, and O. Rivasplata. Logarithmic pruning is all you need. *Advances in Neural Information Processing Systems*, 33, 2020.
- [44] S. Oymak and M. Soltanolkotabi. Toward moderate overparameterization: Global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 1(1), 2020.
- [45] V. Ramanujan, M. Wortsman, A. Kembhavi, A. Farhadi, and M. Rastegari. What’s hidden in a randomly weighted neural network? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11893–11902, 2020.
- [46] A. Renda, J. Frankle, and M. Carbin. Comparing rewinding and fine-tuning in neural network pruning. In *International Conference on Learning Representations*, 2019.
- [47] I. Safran and O. Shamir. Spurious local minima are common in two-layer relu neural networks. In *International Conference on Machine Learning*, pages 4430–4438, 2018.
- [48] O. Shamir. Distribution-specific hardness of learning neural networks. *The Journal of Machine Learning Research*, 19(1):1135–1163, 2018.
- [49] S. Singh and A. Majumdar. Deep sparse coding for non-intrusive load monitoring. 9(5):4669–4678, Feb. 2017.
- [50] 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*, 65(2):742–769, 2018.
- [51] S. Srinivas and R. V. Babu. Data-free parameter pruning for deep neural networks. *arXiv preprint arXiv:1507.06149*, 2015.
- [52] 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*, pages 3404–3413. JMLR. org, 2017.
- [53] J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12(4):389–434, 2012.
- [54] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [55] C. Wang, G. Zhang, and R. Grosse. Picking winning tickets before training by preserving gradient flow. In *International Conference on Learning Representations*, 2019.
- [56] F. Wang, K. Li, C. Liu, Z. Mi, M. Shafie-Khah, and J. P. S. Catalão. Synchronous pattern matching principle-based residential demand response baseline estimation: Mechanism analysis and approach description. *IEEE Transactions on Smart Grid*, 9(6):6972–6985, 2018.
- [57] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li. Learning structured sparsity in deep neural networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 2082–2090, 2016.

- [58] T.-J. Yang, Y.-H. Chen, and V. Sze. Designing energy-efficient convolutional neural networks using energy-aware pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5687–5695, 2017.
- [59] M. Ye, C. Gong, L. Nie, D. Zhou, A. Klivans, and Q. Liu. Good subnetworks provably exist: Pruning via greedy forward selection. In *International Conference on Machine Learning*, pages 10820–10830. PMLR, 2020.
- [60] M. Ye and Y. Sun. Variable selection via penalized neural network: a drop-out-one loss approach. In *International Conference on Machine Learning*, pages 5620–5629. PMLR, 2018.
- [61] M. Ye, L. Wu, and Q. Liu. Greedy optimization provably wins the lottery: Logarithmic number of winning tickets is enough. *Advances in Neural Information Processing Systems*, 33, 2020.
- [62] G. Yehudai and O. Shamir. On the power and limitations of random features for understanding neural networks. In *Advances in Neural Information Processing Systems*, pages 6594–6604, 2019.
- [63] H. You, C. Li, P. Xu, Y. Fu, Y. Wang, X. Chen, R. G. Baraniuk, Z. Wang, and Y. Lin. Drawing early-bird tickets: Toward more efficient training of deep networks. In *International Conference on Learning Representations*, 2019.
- [64] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- [65] S. Zhang, M. Wang, S. Liu, P.-Y. Chen, and J. Xiong. Fast learning of graph neural networks with guaranteed generalizability: one-hidden-layer case. In *2020 International Conference on Machine Learning (ICML)*, 2020.
- [66] S. Zhang, M. Wang, S. Liu, P.-Y. Chen, and J. Xiong. Guaranteed convergence of training convolutional neural networks via accelerated gradient descent. In *2020 54th Annual Conference on Information Sciences and Systems (CISS)*, 2020.
- [67] S. Zhang, M. Wang, J. Xiong, S. Liu, and P.-Y. Chen. Improved linear convergence of training cnns with generalizability guarantees: A one-hidden-layer case. *IEEE Transactions on Neural Networks and Learning Systems*, 32(6):2622–2635, 2020.
- [68] L. Zhao, Q. Hu, and W. Wang. Heterogeneous feature selection with multi-modal deep neural networks and sparse group lasso. *IEEE Transactions on Multimedia*, 17(11):1936–1948, 2015.
- [69] 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*, pages 4140–4149. JMLR. org, <https://arxiv.org/abs/1706.03175>, 2017.
- [70] H. Zhou, J. Lan, R. Liu, and J. Yosinski. Deconstructing lottery tickets: Zeros, signs, and the supermask. In *Advances in Neural Information Processing Systems 32*, pages 3597–3607. 2019.
- [71] W. Zhou, V. Veitch, M. Austern, R. P. Adams, and P. Orbanz. Non-vacuous generalization bounds at the imagenet scale: a pac-bayesian compression approach. In *International Conference on Learning Representations*, 2018.
- [72] D. Zou, Y. Cao, D. Zhou, and Q. Gu. Gradient descent optimizes over-parameterized deep relu networks. *Machine Learning*, 109(3):467–492, 2020.
- [73] D. Zou and Q. Gu. An improved analysis of training over-parameterized deep neural networks. In *Advances in Neural Information Processing Systems*, pages 2055–2064, 2019.

## Supplementary Materials for:

### Why Lottery Ticket Wins? A Theoretical Perspective of Sample Complexity on Pruned Neural Networks

We first provide an overview about techniques used in proving the landscape (Theorem 1), linear convergence to the ground truth (Theorem 2) and tensor initialization (Lemma 1).

1. **Sample complexity scales in  $\{r_j\}_{j=1}^K$ :** To guarantee the theoretical bounds depend on  $\{r_j\}_{j=1}^K$  instead of  $d$ , we define an equivalent empirical risk function as shown in (12) in Appendix A, from  $\mathbb{R}^{\sum_j r_j}$  to  $\mathbb{R}$ . Existing concentration theorems and landscape analysis built upon (2) can no longer be used here, and thus we revised or updated the corresponding lemmas, which can be found in Appendix G to I. In the initialization methods, for estimating a proper weights that match new empirical risk function, the construction of high-momenta in Appendix B and corresponding proofs in Appendix J are updated accordingly as well;

2. **Local convex region:** In proving Theorem 1 (Appendix C), we first bound the Hessian of the expectation of the new empirical risk function and then obtain the distance of the Hessian of the new empirical risk function to its expectation by concentration theorem. By triangle inequality, the Hessian of the new empirical risk function is characterized in terms of sample size  $N$ ;

3. **Linear Convergence:** In proving Theorem 2 (Appendix D), we first characterize the gradient descent term by *Intermediate Value Theorem* (IVT). However, since the empirical risk function is non-smooth due to the ReLU activation function, IVT is applied in the expectation of the empirical risk function instead, and we later show the gradient generated by finite number of samples is close to its expectation. Therefore, the iterates still converge to the ground truth with enough samples. Further, the linear convergence rate are determined by  $\|\mathbf{W}^{(t+1)} - \mathbf{W}^* \mathbf{P}\| / \|\mathbf{W}^{(t)} - \mathbf{W}^* \mathbf{P}\|$ , which turns out to be dependent on  $\beta$ ;

4. **Initialization via Tensor Method:** The major challenge for tensor initialization is to construct the proper high dimensional momenta. As we mentioned above, if one directly applies the method in [69], the sample complexity is in  $\Theta(d)$ . In this paper, we select  $\tilde{\mathbf{x}}$  (see (20) in Appendix B), which is the sum of the augmented  $\mathbf{x}_{\Omega_j}$ . In proving Lemma 1, the major idea to bound the estimations of the directions and magnitudes of  $\mathbf{w}_{j,\Omega_j}$  to the ground values, respectively (see in Appendix F).

## A Notations

In this section, we first introduce some important notations that will be used in the following proofs, and the notations are summarized in Table 1.

First, for the convenience of proofs, some notations in main contexts, namely,  $\Omega_j^*$ ,  $r_j^*$  and  $\hat{f}_{\mathcal{D}}$  will be re-defined. We emphasize here that the re-definition of these notations will not affect the presentation of theoretical results in Section 3, and the explanations can be found in the following paragraphs.

Next, given a permutation matrix  $\mathbf{P}$ , we define a group of sets  $\{\Omega_j^*\}_{j=1}^K$  with  $|\Omega_j^*| = r_j^*$ , and  $\Omega_j^*$  denotes the indices of non-zero entries in  $\mathbf{M}^* \mathbf{P}$ , which is also the non-pruned weights of the  $j$ -th neuron in the oracle model with respect to ground truth weights  $\mathbf{M}^* \mathbf{P}$ , instead of  $\mathbf{M}^*$ . Please note that the sets  $\{\Omega_j^*\}_{j=1}^K$  and  $\{r_j^*\}_{j=1}^K$  here are just a permutation of these in the main context. Since the permutation of  $\{r_j\}_{j=1}^K$  will not change the results in Section 3, we abuse the notations for the convenience of proofs. Correspondingly, for the learner model, the indices of non-pruned weights of the  $j$ -th neuron is denoted as  $\Omega_j$ , and  $|\Omega_j| = r_j$ . Therefore, we have

$$\mathbf{w}_j^T \mathbf{x} = \mathbf{w}_{j,\Omega_j}^T \mathbf{x}_{\Omega_j}, \quad (11)$$

where  $\mathbf{z}_{\Omega_j} \in \mathbb{R}^{r_j}$  is the subvector of  $\mathbf{z}$  with respect to indices  $\Omega_j$  for any vector  $\mathbf{z} \in \mathbb{R}^d$ .

Then, recall the *empirical risk function* defined in (2), it can be re-written as

$$\hat{f}_{\mathcal{D}}(\tilde{\mathbf{w}}) := \frac{1}{2N} \sum_{n=1}^N \left( \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_{j,\Omega_j}^T \mathbf{x}_{n,\Omega_j}) - y_n \right)^2, \quad (12)$$

Table 1: Table of Notations

Notation	Description
$N$	The number of training samples; a scalar in $\mathbb{Z}$
$K$	The number of neurons in the neural network; a scalar in $\mathbb{R}$
$d$	The dimension of input data; a scalar in $\mathbb{R}$
$\mathbf{x}$	The input data/features; a vector in $\mathbb{R}^d$
$y$	The output label; a scalar in $\mathbb{R}$
$\hat{f}_{\mathcal{D}}$	The empirical risk function defined in (12); a mapping from $\mathbb{R}^{\sum_j r_j}$ to $\mathbb{R}$
$f$	The population risk function defined as $f = \mathbb{E}_{\mathcal{D}} \hat{f}_{\mathcal{D}}$ ; a mapping from $\mathbb{R}^{\sum_j r_j}$ to $\mathbb{R}$
$\mathbf{P}$	The permutation matrix; a binary matrix in $\{0, 1\}^{K \times K}$
$\mathbf{W}^*$	The ground truth weights of oracle network; a matrix in $\mathbb{R}^{d \times K}$
$\mathbf{M}^*$	The mask matrix of the oracle network; a binary matrix in $\{0, 1\}^{d \times K}$
$r_j^*$	The number of non-pruned weights in the $j$ -th neuron of oracle network
$\mathbf{W}$	The ground truth weights of learner network; a matrix in $\mathbb{R}^{d \times K}$
$\mathbf{M}$	The mask matrix of the learner network; a binary matrix in $\{0, 1\}^{d \times K}$
$r_j$	The number of non-pruned weights in the $j$ -th neuron of learner network
$r_{\min}$	The minimal value in $\{r_j\}_{j=1}^K$
$r_{\max}$	The maximal value in $\{r_j\}_{j=1}^K$
$\Omega_j^*$	The indices of non-pruned weights in teacher network; a set with size of $r_j^*$
$\Omega_j$	The indices of non-pruned weights in learner network; a set with size of $r_j$
$\tilde{\mathbf{w}}$	Contains the non-pruned weights of $\mathbf{W}$ and equals to $[\mathbf{w}_{1,\Omega_1}^T, \mathbf{w}_{2,\Omega_2}^T, \dots, \mathbf{w}_{K,\Omega_K}^T]^T$ ; a vector in $\mathbb{R}^{\sum_j r_j}$
$\tilde{\mathbf{w}}^*$	Contains the non-pruned weights of the oracle model; a vector in $\mathbb{R}^{\sum_j r_j}$
$\delta_{i,j}$	A binary scalar, and the value is 1 if $\Omega_j$ and $\Omega_k$ are overlapped and 0 otherwise
$\tilde{r}$	The value of $\frac{1}{8K^4} \left( \sum_k \sum_j (1 + \delta_{j,k})(r_j + r_k)^{\frac{1}{2}} \right)^2$
$\sigma_i$	The $i$ -th largest singular value of $\mathbf{W}^* \mathbf{P}$ , and the value equals to the $i$ -th largest singular value of $\mathbf{W}^*$
$\kappa$	The value of $\sigma_1 / \sigma_K$
$\gamma$	The value of $\prod_{i=1}^K \sigma_i / \sigma_K$
$\rho$	A fixed positive constant in $\mathbb{R}^+$
$q$	Some large constant in $\mathbb{R}$

where  $\tilde{\mathbf{w}} = [\mathbf{w}_{1,\Omega_1}^T, \mathbf{w}_{2,\Omega_2}^T, \dots, \mathbf{w}_{K,\Omega_K}^T]^T \in \mathbb{R}^{\sum_j r_j}$ . Here, we abuse the notation of  $\hat{f}_{\mathcal{D}}$  to represent a mapping from  $\mathbb{R}^{\sum_j r_j}$ , instead of  $\mathbb{R}^{K \times d}$  in (2), to  $\mathbb{R}$ . In fact, under the constraint of  $\mathbf{W} = \mathbf{M} \odot \mathbf{W}$ , the degree of freedom of  $\mathbf{W}$  is actually  $\sum_j r_j$  instead of  $Kd$ , and the definition in (2) is a easier way for us to present the following proofs. Therefore, the optimization problem in (3) is equivalent as

$$\min_{\tilde{\mathbf{w}}} : \hat{f}_{\mathcal{D}}(\tilde{\mathbf{w}}). \quad (13)$$

Let us define  $\tilde{\mathbf{w}}^* = [\mathbf{w}_{1,\Omega_1}^{*T}, \mathbf{w}_{2,\Omega_2}^{*T}, \dots, \mathbf{w}_{K,\Omega_K}^{*T}]^T \in \mathbb{R}^{\sum_j r_j}$ , where  $\mathbf{w}_j^{*T}$  is the  $j$ -th column of  $\mathbf{W}^* \mathbf{P}$ , and it is clear that  $\tilde{\mathbf{w}}^*$  is the global optimal to (13). Additionally, the population risk function, which is the expectation of the empirical risk function over the data  $\mathcal{D}$ , is defined as

$$\begin{aligned} f(\tilde{\mathbf{w}}) &= \mathbb{E}_{\mathcal{D}} \hat{f}_{\mathcal{D}}(\tilde{\mathbf{w}}) = \mathbb{E}_{\mathcal{D}} \frac{1}{2N} \sum_{n=1}^N \left( \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_{j,\Omega_j}^T \mathbf{x}_{n,\Omega_j}) - y_n \right)^2 \\ &= \mathbb{E}_{\mathbf{x}} \frac{1}{2} \left( \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_{j,\Omega_j}^T \mathbf{x}_{\Omega_j}) - y \right)^2, \end{aligned} \quad (14)$$

where  $\mathbf{x} \in \mathbb{R}^d$  belongs to standard Gaussian distribution, and  $y = g(\mathbf{W}^* \mathbf{P}^*; \mathbf{x})$ .

Moreover, for the convenience of proofs, we use  $\sigma_i$  to denote the  $i$ -th largest singular value of  $\mathbf{W}^* \mathbf{P}$ , and it is clear that  $\sigma_i(\mathbf{W}^* \mathbf{P}) = \sigma_i(\mathbf{W}^*)$  for all  $i$ . Then,  $\kappa$  is defined as  $\sigma_1 / \sigma_K$ , and



$\gamma = \prod_{i=1}^K \sigma_i / \sigma_K$ . Factor  $\rho$  is defined in Property 3.2 [69] and a fixed constant for the ReLU activation function. In addition, without special descriptions,  $\alpha = [\alpha_1^T, \alpha_2^T, \dots, \alpha_K^T]^T$  stands for any unit vector that in  $\mathbb{R}^{\sum_j r_j}$  with  $\alpha_j \in \mathbb{R}^{r_j}$ . Therefore, we have

$$\|\nabla^2 \hat{f}_{\mathcal{D}}\|_2 = \max_{\alpha} \|\alpha^T \nabla^2 \hat{f}_{\mathcal{D}} \alpha\|_2 = \max_{\alpha} \left( \sum_{j=1}^K \alpha_j^T \frac{\partial \hat{f}_{\mathcal{D}}}{\partial \mathbf{w}_j} \right)^2. \quad (15)$$

Finally, since we focus on order-wise analysis, some constant number will be ignored in the majority of the steps. In particular, we use  $h_1(z) \gtrsim$  ( or  $\lesssim, \asymp$ )  $h_2(z)$  to denote there exists some positive constant  $C$  such that  $h_1(z) \geq$  ( or  $\leq, =$ )  $C \cdot h_2(z)$  when  $z \in \mathbb{R}$  is sufficiently large.

## B Initialization via tensor method

In this section, we present the revised tensor initialization based on that in [69]. To reduce the dependency of input dimension from  $d$  to the order of  $r_{\max}$ , we need to define  $\tilde{\mathbf{x}}$  in (20) instead of directly using  $\mathbf{x}$  to generate the high order momentum as shown in (21) to (23). In addition, as  $\mathbf{w}_{j,\Omega_j}$ 's are different in dimensions, we need to define the corresponding augmented weights by inserting 0 such that augmented  $\mathbf{w}_{j,\Omega_j}$  are additive in a sense. The additional notations used in presenting are summarized in Table 2, and one can skip this part if the focus is only on the local convexity analysis (Theorem 1) and convergence analysis (Theorem 2). The intuitive reasons for selecting  $\tilde{\mathbf{x}}$  mainly lie in two aspects: first,  $\tilde{\mathbf{x}}$  is much lower dimensional vector considering  $r_j \ll d$ ; second,  $\tilde{\mathbf{x}}$  belongs to zero mean Gaussian distribution, which is rotational invariant and is correlate with  $\phi(\mathbf{w}_j^* \mathbf{x})$ . Therefore, the magnitude and direction information of  $\{\mathbf{w}_{j,\Omega_j}\}_{j=1}^K$  are separable after tensor decomposition, and the dimension of the tensors are at most in the order of  $r_{\max}$ .

Table 2: Table of Additional Notations for Tensor method

Notation	Description
$\tilde{\mathbf{x}}_{\Omega_j}^{(j)}$	The augmented vector in $\mathbb{R}^{r_{\max}}$ of $\mathbf{x}_{\Omega_j}$ by inserting 0; defined in (16)
$\mathcal{F}_j$	A linear mapping that generates a augmented vector; defined in (17)
$\mathcal{F}_j^\dagger$	The pseudo inverse of $\mathcal{F}_j$ ; a linear mapping
$\tilde{\mathbf{x}}$	The value of $\frac{1}{\sqrt{K}} \sum_j \tilde{\mathbf{x}}_{\Omega_j}^{(j)}$ ;
$\mathbf{u}_j^*$	The augmented vector in $\mathbb{R}^{r_{\max}}$ of $\mathbf{w}_{j,\Omega_j}^*$ by inserting 0; defined in (19)
$\bar{\mathbf{u}}_j^*$	The normalized vector of $\mathbf{u}_j^*$ as $\mathbf{u}_j^* / \ \mathbf{u}_j^*\ _2$
$\hat{\bar{\mathbf{u}}}_j^*$	The estimation of the normalized vector of $\mathbf{u}_j^*$
$\psi_1, \psi_2, \psi_3$	Some fixed constants depends on the distribution of $\{\mathbf{x}_{\Omega_j}\}_{j=1}^K$
$\mathbf{M}_1$	A vector in $\mathbb{R}^{r_{\max}}$ defined in (21)
$\hat{\mathbf{M}}_1$	The estimation of $\mathbf{M}_1$
$\mathbf{M}_2$	A matrix in $\mathbb{R}^{r_{\max} \times r_{\max}}$ defined in (22)
$\hat{\mathbf{M}}_2$	The estimation of $\mathbf{M}_2$
$\mathbf{M}_3$	A tensor in $\mathbb{R}^{r_{\max} \times r_{\max} \times r_{\max}}$ defined in (23)
$\hat{\mathbf{M}}_3$	The estimation of $\mathbf{M}_3$
$\mathbf{V}$	The orthogonal matrix in $\mathbb{R}^{K \times K}$ that span the sub-space of the convex hull of $\{\mathbf{u}_j\}_{j=1}^K$
$\hat{\mathbf{V}}$	The estimation of $\mathbf{V}$
$\mathbf{M}(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}})$	A tensor in $\mathbb{R}^{K \times K \times K}$ defined in (29)
$\hat{\mathbf{M}}(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}})$	The estimation of $\mathbf{M}(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}})$
$\mathbf{s}_j$	The value of $\mathbf{V} \mathbf{u}_j^*$ ; a vector in $\mathbb{R}^K$
$\hat{\mathbf{s}}_j$	The estimation of $\mathbf{s}_j$
$\alpha_j$	The value of $\ \mathbf{u}_j^*\ _2$ ; a scalar in $\mathbb{R}$
$\hat{\alpha}_j$	The estimation of $\alpha_j$

First, we define a group of augmented vectors  $\{\tilde{\mathbf{x}}_{\tilde{\Omega}_j}^{(j)}\}_{j=1}^K$  based on  $\{\mathbf{x}_{\Omega_j}\}_{j=1}^K$  such that  $\Omega_j \subseteq \tilde{\Omega}_j$  with  $|\tilde{\Omega}_j| = r_{\max}$  and

$$\tilde{x}_i^{(j)} = \begin{cases} x_i & \text{if } i \in \Omega_j \\ 0 & \text{if } i \in \tilde{\Omega}_j / \Omega_j \end{cases}. \quad (16)$$

For notation convenience, we use  $\mathcal{F}_j$  to denote the mapping from  $\mathbb{R}^{r_j}$  to  $\mathbb{R}^{r_{\max}}$  as

$$\mathcal{F}_j(\mathbf{z}) = [\mathbf{z}^T, \mathbf{0}_{(j)}^T]^T, \quad (17)$$

where  $\mathbf{0}$  is a zero vector in  $\mathbb{R}^{r_{\max}-r_j}$ . Obviously, we have

$$\tilde{\mathbf{x}}_{\tilde{\Omega}_j}^{(j)} = \mathcal{F}_j(\mathbf{x}_{\Omega_j}). \quad (18)$$

Correspondingly, the augmented weights  $\{\mathbf{u}_j^*\}_{j=1}^K$  are defined as

$$\mathbf{u}_j^* = \mathcal{F}_j(\mathbf{w}_{j,\Omega_j}^*) \quad (19)$$

for  $j \in [K]$ . The steps above guarantee the augmented weights  $\mathbf{u}_j^*$ 's are in the same dimension so that the high order momenta are able to characterize the directions of weights simultaneously. Additionally, we define

$$\tilde{\mathbf{x}} = \frac{1}{\sqrt{K}} \sum_{j=1}^K \tilde{\mathbf{x}}_{\tilde{\Omega}_j}^{(j)}, \quad (20)$$

and corresponding high order momenta are defined in the following way instead:

$$\mathbf{M}_1 = \mathbb{E}_{\mathbf{x}}\{y\tilde{\mathbf{x}}\} \in \mathbb{R}^{r_{\max}}, \quad (21)$$

$$\mathbf{M}_2 = \mathbb{E}_{\mathbf{x}}\left[y(\tilde{\mathbf{x}} \otimes \tilde{\mathbf{x}} - \mathbb{E}_{\mathbf{x}}\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T)\right] \in \mathbb{R}^{r_{\max} \times r_{\max}}, \quad (22)$$

$$\mathbf{M}_3 = \mathbb{E}_{\mathbf{x}}\left[y(\tilde{\mathbf{x}}^{\otimes 3} - \tilde{\mathbf{x}} \otimes \mathbb{E}_{\mathbf{x}}\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T)\right] \in \mathbb{R}^{r_{\max} \times r_{\max} \times r_{\max}}, \quad (23)$$

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

$$\mathbf{v} \tilde{\otimes} \mathbf{Z} = \sum_{i=1}^{d_2} (\mathbf{v} \otimes \mathbf{z}_i \otimes \mathbf{z}_i + \mathbf{z}_i \otimes \mathbf{v} \otimes \mathbf{z}_i + \mathbf{z}_i \otimes \mathbf{z}_i \otimes \mathbf{v}), \quad (24)$$

for any vector  $\mathbf{v} \in \mathbb{R}^{d_1}$  and  $\mathbf{Z} \in \mathbb{R}^{d_1 \times d_2}$ .

Following the same calculate formulas in the Claim 5.2 [69], there exist some known constants  $\psi_i, i = 1, 2, 3$ , such that

$$\mathbf{M}_1 = \sum_{j=1}^K \psi_1 \cdot \|\mathbf{u}_j^*\|_2 \cdot \bar{\mathbf{u}}_j^*, \quad (25)$$

$$\mathbf{M}_2 = \sum_{j=1}^K \psi_2 \cdot \|\mathbf{u}_j^*\|_2 \cdot \bar{\mathbf{u}}_j^* \bar{\mathbf{u}}_j^{*T}, \quad (26)$$

$$\mathbf{M}_3 = \sum_{j=1}^K \psi_3 \cdot \|\mathbf{u}_j^*\|_2 \cdot \bar{\mathbf{u}}_j^{*\otimes 3}, \quad (27)$$

where  $\bar{\mathbf{u}}_j^* = \mathbf{u}_j^* / \|\mathbf{u}_j^*\|_2$  in (21)-(23) is the normalization of  $\mathbf{u}_j^*$ .

$\mathbf{M}_1, \mathbf{M}_2$  and  $\mathbf{M}_3$  can be estimated through the samples  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , and let  $\widehat{\mathbf{M}}_1, \widehat{\mathbf{M}}_2, \widehat{\mathbf{M}}_3$  denote the corresponding estimates. First, we will decompose the rank- $k$  tensor  $\widehat{\mathbf{M}}_3$  and obtain the  $\{\bar{\mathbf{u}}_j^*\}_{j=1}^K$ . By applying the tensor decomposition method [31] to  $\widehat{\mathbf{M}}_3$ , the outputs, denoted by  $\{\widehat{\bar{\mathbf{u}}}_j^*\}_{j=1}^K$ , are the estimations of  $\{\bar{\mathbf{u}}_j^*\}_{j=1}^K$ . Next, we will estimate  $\|\mathbf{u}_j^*\|_2$  through solving the following optimization problem:

$$\widehat{\boldsymbol{\alpha}} = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^K} : \quad \left| \widehat{\mathbf{M}}_1 - \sum_{j=1}^K \psi_1 \alpha_j \widehat{\bar{\mathbf{u}}}_j^* \right|, \quad (28)$$

---

**Subroutine 1** Tensor Initialization Method

---

- 1: **Input:** training data  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ ;
  - 2: Generate augmented inputs and weights through  $\mathcal{F}_j$  as shown in (17) and (19);
  - 3: Partition  $\mathcal{D}$  into three disjoint subsets  $\mathcal{D}_1, \mathcal{D}_2, \mathcal{D}_3$ ;
  - 4: Calculate  $\widehat{\mathbf{M}}_1, \widehat{\mathbf{M}}_2$  following (21), (22) using  $\mathcal{D}_1, \mathcal{D}_2$ , respectively;
  - 5: Obtain the estimate subspace  $\widehat{\mathbf{V}}$  of  $\widehat{\mathbf{M}}_2$ ;
  - 6: Calculate  $\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  through  $\mathcal{D}_3$ ;
  - 7: Obtain  $\{\widehat{\mathbf{s}}_j\}_{j=1}^K$  via tensor decomposition method [31] on  $\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$ ;
  - 8: Obtain  $\widehat{\alpha}$  by solving optimization problem (28);
  - 9: **Return:**  $\mathbf{w}_{j,\Omega_j}^{(0)} = \mathcal{F}_j^\dagger(|\widehat{\alpha}_j|\widehat{\mathbf{V}}\widehat{\mathbf{s}}_j), j = 1, \dots, K$ .
- 

From (25) and (28), we know that  $|\widehat{\alpha}_j|$  is the estimation of  $\|\mathbf{u}_j^*\|_2$ . Thus,  $\widehat{\mathbf{U}}$  is given as  $[|\widehat{\alpha}_1|\widehat{\mathbf{u}}_1^*, \dots, |\widehat{\alpha}_j|\widehat{\mathbf{u}}_j^*, \dots, |\widehat{\alpha}_K|\widehat{\mathbf{u}}_K^*]$ .

To reduce the computational complexity of tensor decomposition, one can project  $\widehat{\mathbf{M}}_3$  to a lower-dimensional tensor [69]. 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.

Moreover, we have

$$\mathbf{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) = \mathbb{E}_{\mathbf{x}} \left[ y((\widehat{\mathbf{V}}^T \widetilde{\mathbf{x}})^{\otimes 3} - (\widehat{\mathbf{V}}^T \widetilde{\mathbf{x}}) \widetilde{\otimes} \mathbb{E}_{\mathbf{x}}(\widehat{\mathbf{V}}^T \widetilde{\mathbf{x}})(\widehat{\mathbf{V}}^T \widetilde{\mathbf{x}})^T) \right] \in \mathbb{R}^{K \times K \times K}, \quad (29)$$

Then, one can decompose the estimate  $\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  to obtain unit vectors  $\{\widehat{\mathbf{s}}_j\}_{j=1}^K \in \mathbb{R}^K$ . Since  $\overline{\mathbf{u}}^*$  lies in the subspace  $\mathbf{V}$ , we have  $\mathbf{V}\mathbf{V}^T \overline{\mathbf{u}}_j^* = \overline{\mathbf{u}}_j^*$ . Then,  $\widehat{\mathbf{V}}\widehat{\mathbf{s}}_j$  is an estimate of  $\overline{\mathbf{u}}_j^*$ . After we obtain the estimated augmented weights  $\widehat{\mathbf{u}}_j^*$ , the estimated weights can be generated through  $\widehat{\mathbf{w}}_{j,\Omega_j}^* = \mathcal{F}_j^\dagger(\widehat{\mathbf{u}}_j^*)$ , where  $\mathcal{F}_j^\dagger$  is the pseudo inverse of  $\mathcal{F}_j$ . The initialization process is summarized in Subroutine 1.

## C Proof of Theorem 1

The main idea in proving Theorem 1 is to use triangle inequality as shown in (33) by bounding the second order derivative of the population risk function and the distance between the empirical risk and population risk functions. Lemma 3 provides the lower and upper bound for the population risk function, while Lemma 4 provides the error bound between the second order derivation of empirical risk and population risk functions.

**Lemma 2** (Weyl's inequality, [7]). *Suppose  $\mathbf{B} = \mathbf{A} + \mathbf{E}$  be a matrix with dimension  $m \times m$ . Let  $\lambda_i(\mathbf{B})$  and  $\lambda_i(\mathbf{A})$  be the  $i$ -th largest eigenvalues of  $\mathbf{B}$  and  $\mathbf{A}$ , respectively. Then, we have*

$$|\lambda_i(\mathbf{B}) - \lambda_i(\mathbf{A})| \leq \|\mathbf{E}\|_2, \quad \forall i \in [m]. \quad (30)$$

**Lemma 3.** *Let  $f$  be the population risk function in (14). Assume  $\mathbf{W}$  satisfies (6), then the second-order derivative of  $f$  over  $\widetilde{\mathbf{w}}$  is bounded as*

$$\frac{(1 - \varepsilon_0)\rho}{11\kappa^2\gamma K^2} \mathbf{I} \leq \nabla^2 f(\widetilde{\mathbf{w}}) \leq \frac{7}{K} \mathbf{I}, \quad (31)$$

where  $\widetilde{\mathbf{w}}$  only contains the elements of  $\mathbf{W}$  with respect to the indices of non-pruned weights.

**Lemma 4.** *Let  $\hat{f}_{\mathcal{D}}$  and  $f$  be the empirical and population risk function in (12) and (14), respectively, then the second-order derivative of  $\hat{f}_{\mathcal{D}}$  is close to its expectation  $f$  with an upper bound as:*

$$\|\nabla^2 \hat{f}_{\mathcal{D}} - \nabla^2 f\|_2 \lesssim \frac{1}{K^2} \sum_{k=1}^K \sum_{j=1}^K (1 + \delta_{j,k}) \sqrt{\frac{(r_j + r_k) \log q}{N}} \quad (32)$$

with probability at least  $1 - q^{-r_{\min}}$ .

*Proof of Theorem 1.* Let  $\hat{\lambda}_{\max}$  and  $\hat{\lambda}_{\min}$  denote the largest and smallest eigenvalues of  $\nabla^2 \hat{f}_{\mathcal{D}}$ , respectively. Also, Let  $\lambda_{\max}$  and  $\lambda_{\min}$  denote the largest and smallest eigenvalues of  $\nabla^2 f_{\mathcal{D}}$ , respectively.

Then, from Lemma 2, we have

$$\hat{\lambda}_{\max} \leq \lambda_{\max} + \|\nabla^2 \hat{f}_{\mathcal{D}} - \nabla^2 f\|_2 \quad (33)$$

and

$$\hat{\lambda}_{\min} \geq \lambda_{\min} - \|\nabla^2 \hat{f}_{\mathcal{D}} - \nabla^2 f\|_2. \quad (34)$$

When the sample complexity satisfies  $N \gtrsim \varepsilon_1^{-2} \rho^{-2} \kappa^4 \gamma^2 K^4 \left[ \frac{1}{K^2} \sum_{k=1}^K \sum_{j=1}^K (1 + \delta_{j,k}) \sqrt{r_j + r_k} \right]^2 \log q$ , then from Lemma 4, we have

$$\|\nabla^2 \hat{f}_{\mathcal{D}} - \nabla^2 f\|_2 \leq \frac{\varepsilon_1 \rho}{11 \kappa^2 \gamma K^2}. \quad (35)$$

Then, from (33), (34) and (35), we have

$$\hat{\lambda}_{\max} \leq \frac{8}{K}, \quad (36)$$

and

$$\hat{\lambda}_{\min} \geq \frac{(1 - \varepsilon_0 - \varepsilon_1) \rho}{11 \kappa^2 \gamma K^2}, \quad (37)$$

which completes the proof.  $\square$

## D Proof of Theorem 2

The major idea in proving Theorem 2 is to first characterize the gradient descent term by intermediate value theorem. Let  $\tilde{\mathbf{w}}^{(t)}$  be the vectorized iterate  $\mathbf{W}^{(t)}$  with respect to the non-pruned weights, then we have

$$\begin{aligned} \nabla \hat{f}_{\Omega_t}(\tilde{\mathbf{w}}^{(t)}) &= f_{\Omega_t}(\tilde{\mathbf{w}}^{(t)}) + (\hat{f}_{\Omega_t}(\tilde{\mathbf{w}}^{(t)}) - f_{\Omega_t}(\tilde{\mathbf{w}}^{(t)})) \\ &= \langle \nabla^2 f_{\Omega_t}(\tilde{\mathbf{w}}^{(t)}), \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \rangle + (\hat{f}_{\Omega_t}(\tilde{\mathbf{w}}^{(t)}) - f_{\Omega_t}(\tilde{\mathbf{w}}^{(t)})), \end{aligned} \quad (38)$$

where  $\tilde{\mathbf{w}}^{(t)}$  lies in the convex hull of  $\tilde{\mathbf{w}}^{(t)}$  and  $\tilde{\mathbf{w}}^*$ . The reason that intermediate value theorem is applied on population risk function instead of empirical risk function is the non-smoothness of the empirical risk functions. Due to the non-smoothness of ReLU activation function at zero point, the empirical risk function is not smooth, either. However, the expectation of the empirical risk function over the Gaussian input  $\mathbf{x}$  is smooth. Hence, compared with smooth empirical risk function, i.e., neural networks equipped with sigmoid activation function, we have an additional lemma to bound  $\nabla \hat{f}_{\mathcal{D}_t}$  to its expectation  $\nabla f$ , which is summarized in Lemma 5.

The momentum term  $\beta(\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)})$  plays an important role in determining the convergence rate, and the recursive rule is obtained in the following way:

$$\begin{bmatrix} \tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \end{bmatrix} = \mathbf{A}(\beta) \begin{bmatrix} \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t-1)} - \tilde{\mathbf{w}}^* \end{bmatrix}, \quad (39)$$

where  $\mathbf{A}(\beta)$  is a matrix with respect to the value of  $\beta$  and defined in (44). Then, we know  $\tilde{\mathbf{w}}^{(t)}$ , which is equivalent to  $\mathbf{W}^{(t)}$ , converges to the ground-truth with a linear rate which is the largest singular value of matrix  $\mathbf{A}(\beta)$ . Recall that AGD reduces to GD with  $\beta = 0$ , so our analysis applies to GD method as well. We are able to show the convergence rate of AGD is faster than GD by proving the largest singular value of  $\mathbf{A}(\beta)$  is smaller than  $\mathbf{A}(0)$  for some  $\beta > 0$ .

**Lemma 5.** Let  $\hat{f}_{\mathcal{D}}$  and  $f$  be the empirical and population risk function in (12) and (14), respectively, then the first-order derivative of  $\hat{f}_{\mathcal{D}}$  is close to its expectation  $f$  with an upper bound as:

$$\|\nabla \hat{f}_{\mathcal{D}}(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{w}})\|_2 \lesssim \frac{1}{K^2} \sum_{k=1}^K \sum_{j=1}^K (1 + \delta_{j,k}) \sqrt{\frac{r_k \log q}{N}} \|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2 + \frac{1}{K} \sum_{k=1}^K \sqrt{\frac{r_k \log q}{N}} \cdot |\xi| \quad (40)$$

with probability at least  $1 - q^{-r_{\min}}$ , where  $\tilde{\mathbf{w}}$  only contains the elements of  $\mathbf{W}$  with respect to the indices of non-pruned weights.

*Proof of Theorem 2.* Since  $\|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F = \|\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^*\|_2$ , we can explore the converges of  $\{\tilde{\mathbf{w}}^{(t)}\}_{t=1}^T$  instead. Recall that

$$\begin{aligned}\tilde{\mathbf{w}}^{(t+1)} &= \tilde{\mathbf{w}}^{(t)} - \eta \nabla \hat{f}_{\mathcal{D}_t}(\tilde{\mathbf{w}}^{(t)}) + \beta(\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^{(t-1)}) \\ &= \tilde{\mathbf{w}}^{(t)} - \eta \nabla f(\tilde{\mathbf{w}}^{(t)}) + \beta(\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^{(t-1)}) \\ &\quad + \eta(\nabla f(\tilde{\mathbf{w}}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\tilde{\mathbf{w}}^{(t)})).\end{aligned}\tag{41}$$

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

$$\begin{aligned}\tilde{\mathbf{w}}^{(t+1)} &= \tilde{\mathbf{w}}^{(t)} - \eta \nabla^2 f(\hat{\mathbf{w}}^{(t)})(\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^*) + \beta(\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^{(t-1)}) \\ &\quad + \eta(\nabla f(\tilde{\mathbf{w}}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\tilde{\mathbf{w}}^{(t)})),\end{aligned}\tag{42}$$

where  $\hat{\mathbf{w}}^{(t)}$  lies in the convex hull of  $\tilde{\mathbf{w}}^{(t)}$  and  $\tilde{\mathbf{w}}^*$ .

Next, we have

$$\begin{aligned}\begin{bmatrix} \tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \end{bmatrix} &= \begin{bmatrix} \mathbf{I} - \eta \nabla^2 f(\hat{\mathbf{w}}^{(t)}) + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t-1)} - \tilde{\mathbf{w}}^* \end{bmatrix} \\ &\quad + \eta \begin{bmatrix} \nabla f(\tilde{\mathbf{w}}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\tilde{\mathbf{w}}^{(t)}) \\ 0 \end{bmatrix}\end{aligned}\tag{43}$$

Let

$$\mathbf{A}(\beta) = \begin{bmatrix} \mathbf{I} - \eta \nabla^2 f(\hat{\mathbf{w}}^{(t)}) + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix},\tag{44}$$

so we have

$$\left\| \begin{bmatrix} \tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \end{bmatrix} \right\|_2 = \|\mathbf{A}(\beta)\|_2 \left\| \begin{bmatrix} \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t-1)} - \tilde{\mathbf{w}}^* \end{bmatrix} \right\|_2 + \eta \left\| \begin{bmatrix} \nabla f(\tilde{\mathbf{w}}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\tilde{\mathbf{w}}^{(t)}) \\ 0 \end{bmatrix} \right\|_2.\tag{45}$$

From Lemma 5, we know that

$$\begin{aligned}\eta \left\| \nabla f(\tilde{\mathbf{w}}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\tilde{\mathbf{w}}^{(t)}) \right\|_2 &\leq \frac{C_5 \eta}{K^2} \sum_{k=1}^K \sum_{j=1}^K (1 + \delta_{j,k}) \sqrt{\frac{r_k \log q}{N_t}} \|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2 \\ &\quad + \frac{C_5 \eta}{K} \sum_{k=1}^K \sqrt{\frac{r_k \log q}{N_t}} \cdot |\xi|\end{aligned}\tag{46}$$

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

$$\begin{aligned}\|\tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^*\|_2 &\leq \left( \|\mathbf{A}(\beta)\|_2 + \frac{C_5 \eta}{K^2} \sum_{k=1}^K \sum_{j=1}^K (1 + \delta_{j,k}) \sqrt{\frac{r_k \log q}{N_t}} \right) \|\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^*\|_2 \\ &\quad + \frac{C_5 \eta}{K} \sum_{k=1}^K \sqrt{\frac{r_k \log q}{N_t}} \cdot |\xi| \\ &:= \nu(\beta) \|\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^*\|_2 + \frac{C_5 \eta}{K} \sum_{k=1}^K \sqrt{\frac{r_k \log q}{N_t}} \cdot |\xi|.\end{aligned}\tag{47}$$

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

$$\mathbf{A}(\beta) := \begin{bmatrix} \mathbf{S}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{S}^T \end{bmatrix} \mathbf{A}(\beta) \begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} = \begin{bmatrix} \mathbf{I} - \eta \mathbf{\Lambda} + \beta \mathbf{I} & \beta \mathbf{I} \\ \mathbf{I} & 0 \end{bmatrix}\tag{48}$$

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  $\mathbf{A}(\beta)$  share the same eigenvalues.

Let  $\lambda_i$  be the  $i$ -th eigenvalue of  $\nabla^2 f(\hat{\mathbf{w}}^{(t)})$ , then the corresponding  $i$ -th eigenvalue of (48), denoted by  $\delta_i(\beta)$ , satisfies

$$\nu_i^2 - (1 - \eta \lambda_i + \beta) \delta_i + \beta = 0.\tag{49}$$

Then, we have

$$\delta_i(\beta) = \frac{(1 - \eta\lambda_i + \beta) + \sqrt{(1 - \eta\lambda_i + \beta)^2 - 4\beta}}{2}, \quad (50)$$

and

$$|\delta_i(\beta)| = \begin{cases} \sqrt{\beta}, & \text{if } \beta \geq (1 - \sqrt{\eta\lambda_i})^2, \\ \frac{1}{2} \left| (1 - \eta\lambda_i + \beta) + \sqrt{(1 - \eta\lambda_i + \beta)^2 - 4\beta} \right|, & \text{otherwise.} \end{cases} \quad (51)$$

Note that the other root of (49) is abandoned because the root in (50) is always larger than or at least equal to the other root with  $|1 - \eta\lambda_i| < 1$ . By simple calculation, we have

$$\delta_i(0) > \delta_i(\beta), \quad \text{for } \forall \beta \in (0, (1 - \eta\lambda_i)^2), \quad (52)$$

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

Let us first assume  $\tilde{\mathbf{w}}^{(t)}$  satisfies (6), then from Lemma 3, we know that

$$0 < \frac{(1 - \varepsilon_0)}{11\kappa^2\gamma K^2} \leq \lambda_i \leq \frac{7}{K}$$

provided that  $N_t \gtrsim \varepsilon_0^{-2}\rho^{-1}\kappa^2\gamma K^3 \left[ \frac{1}{K^2} \sum_j \sum_k (1 + \delta_{j,k}) \sqrt{r_k + r_j} \right]^2 \log q$ . Let  $\gamma_1 = \frac{\rho(1-\varepsilon_0)}{11\kappa^2\gamma K^2}$  and  $\gamma_2 = \frac{7}{K}$ . If we choose  $\beta$  such that

$$\beta^* = \max \{ (1 - \sqrt{\eta\gamma_1})^2, (1 - \sqrt{\eta\gamma_2})^2 \}, \quad (53)$$

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

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

$$\begin{aligned} \|\mathbf{A}(\beta^*)\|_2 &= \max_i \delta_i(\beta^*) = 1 - \sqrt{\frac{\gamma_1}{2\gamma_2}} = 1 - \sqrt{\frac{1 - \varepsilon_0}{154\rho^{-1}\kappa^2\gamma K}} \\ &\leq 1 - \frac{1 - 3/4 \cdot \varepsilon_0}{\sqrt{154\rho^{-1}\kappa^2\gamma K}}. \end{aligned} \quad (54)$$

Then, let

$$\frac{C_5\eta}{K^2} \sum_{k=1}^K \sum_{j=1}^K (1 + \delta_{j,k}) \sqrt{\frac{r_k \log q}{N_t}} \leq \frac{\varepsilon_0}{4\sqrt{154\rho^{-1}\kappa^2\gamma K}}, \quad (55)$$

we need  $N_t \gtrsim \varepsilon_0^{-2}\rho^{-1}\kappa^2\gamma K^3 \left[ \frac{1}{K^2} \sum_j \sum_k (1 + \delta_{j,k}) \sqrt{r_k} \right]^2 \log q$ .

Combine (54) and (55), we have

$$\nu(\beta^*) \leq 1 - \frac{1 - \varepsilon_0}{\sqrt{154\rho^{-1}\kappa^2\gamma K}}. \quad (56)$$

While let  $\beta = 0$ , we have

$$\nu(0) \geq \|\mathbf{A}(0)\|_2 = 1 - \frac{1 - \varepsilon_0}{154\rho^{-1}\kappa^2\gamma K} \quad (57)$$

and

$$\nu(0) \leq 1 - \frac{1 - 2\varepsilon_0}{154\rho^{-1}\kappa^2\gamma K} \quad (58)$$

if  $N_t \gtrsim \varepsilon_0^{-2}\rho^{-1}\kappa^2\gamma K^4 \left[ \frac{1}{K^2} \sum_j \sum_k (1 + \delta_{j,k}) \sqrt{r_k + r_j} \right]^2 \log q$ .

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

$$\|\tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^*\|_2 \leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{154\kappa^2\gamma K}}\right) \|\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^*\|_2 + \frac{C\eta}{K} \sum_{k=1}^K \sqrt{\frac{r_k \log q}{N_t}} |\xi|. \quad (59)$$

if  $\tilde{\mathbf{w}}^{(t+1)}$  satisfies (6) and  $N_t \gtrsim \varepsilon_0^{-2}\rho^{-1}\kappa^2\gamma K^4 \left[ \frac{1}{K^2} \sum_j \sum_k (1 + \delta_{j,k}) \sqrt{r_k + r_j} \right]^2 \log q$ .

Then, we can start mathematical induction of (59) over  $t$ .

**Base case:** (6) holds for  $\tilde{\mathbf{w}}^{(0)}$  naturally from the assumption in Theorem 2. Since (6) holds and the number of samples exceeds the required bound in (59), we have (59) holds for  $t = 0$ .

**Induction step:** Assume (59) holds for  $t$ , to make sure the mathematical induction of (59) holds, we need  $\tilde{\mathbf{w}}^{(t+1)}$  satisfies (6). That is

$$\sum_{k=1}^K \frac{\eta}{K} \sqrt{\frac{r_k \log q}{N_t}} \lesssim \frac{1 - \varepsilon_0}{\sqrt{132\kappa^2\gamma K}} \cdot \frac{\varepsilon_0 \sigma_K}{44\kappa^2\gamma K^2}. \quad (60)$$

Hence, we need

$$N_t \gtrsim \varepsilon_0^{-2} \kappa^8 \gamma^3 K^6 \left( \frac{1}{K} \sum_k \sqrt{r_k} \right)^2 \log q. \quad (61)$$

In addition, with (6) and (59) hold for all  $t \leq T$ , the following equation

$$\left\| \begin{bmatrix} \tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \end{bmatrix} \right\|_\infty = \|\mathbf{A}(\beta)\|_2 \left\| \begin{bmatrix} \tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^* \\ \tilde{\mathbf{w}}^{(t-1)} - \tilde{\mathbf{w}}^* \end{bmatrix} \right\|_\infty + \eta \left\| \begin{bmatrix} \nabla f(\tilde{\mathbf{w}}^{(t)}) - \nabla \hat{f}_{\mathcal{D}_t}(\tilde{\mathbf{w}}^{(t)}) \\ 0 \end{bmatrix} \right\|_\infty \quad (62)$$

holds as well, and  $\|\mathbf{A}(\beta)\|_2$  is bounded by  $\nu(\beta)$ . Hence, (59) also holds in infinity norm as

$$\|\tilde{\mathbf{w}}^{(t+1)} - \tilde{\mathbf{w}}^*\|_\infty \leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{154\kappa^2\gamma K}}\right) \|\tilde{\mathbf{w}}^{(t)} - \tilde{\mathbf{w}}^*\|_\infty + 2C\eta \sqrt{\frac{r \log q}{N_t}} |\xi|. \quad (63)$$

In conclusion, when  $N_t \gtrsim \varepsilon_0^{-2} \kappa^8 \gamma^3 K^6 \left( \frac{1}{K^2} \sum_k \sum_j (1 + \delta_{j,k}) \sqrt{r_k + r_j} \right)^2 \log d$ , we know that (59) holds for all  $1 \leq t \leq T$  with probability at least  $1 - K^2 T \cdot q^{-r_{\min}}$ . By simple calculation, we can obtain

$$\|\tilde{\mathbf{w}}^{(T)} - \tilde{\mathbf{w}}^*\|_2 \leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{132\kappa^2\gamma K}}\right)^T \|\tilde{\mathbf{w}}^{(0)} - \tilde{\mathbf{w}}^*\|_2 + \frac{C}{K} \sum_{k=1}^K \sqrt{\frac{\kappa^2 \gamma K^2 r_k \log q}{N_t}} \cdot |\xi|. \quad (64)$$

for some constant  $C > 0$ .

□

## E Obtaining a proper learner network via magnitude pruning

In this section, we show that how one can combine Algorithm 1 and magnitude pruning to find a proper learner network such that  $r_j \geq r_j^*$  and  $\Omega_j \supseteq \Omega_j^*$  from a fully-connected network under some assumptions. Suppose the number of samples is at least  $\Omega(K^6 d \log q \log(1/\varepsilon))$ , we train directly on the fully-connected dense network using Algorithm 1. The number of iteration in line 2 of Algorithm 1 is set as  $T_1 = \Theta(\log(2\widehat{W}_{\max}/\widehat{W}_{\min}))$ , where  $\widehat{W}_{\min}$  and  $\widehat{W}_{\max}$  denote the smallest and largest value of  $\mathbf{W}^*$ , respectively. From (63), after  $T_1$  iterations, the returned model, denote by  $\mathbf{W}^{(T_1)}$ , is close to the ground-truth  $\mathbf{W}^*$ . Specifically, if  $\mathbf{W}_{i,j}^* \neq 0$  and  $\mathbf{W}_{i',j'}^* = 0$ , then  $\mathbf{W}_{i,j}^{(T_1)} > \mathbf{W}_{i',j'}^{(T_1)}$  for any  $i, j, i', j'$ . Then we sort the weights based on their absolute values and prune them sequentially starting from the least absolute value. As long as the ratio of pruned weights is at most  $\left(1 - \frac{\sum_j r_j}{Kd}\right)$ , all the weights are removed correctly, leading to a proper learner network. In fact, if we remove exactly  $1 - \frac{\sum_j r_j}{Kd}$  fraction of weights, the pruned network has the same architecture as the oracle network.

Specifically, suppose  $\mathbf{M}^{(t)}$  to denote the mask matrix by truncating the smallest  $\left(1 - \frac{\sum_j r_j}{Kd}\right)$  fraction of entries in iterate  $\mathbf{W}^{(t)}$ . Let  $\mathbf{M}^*$  denote the ground-truth mask matrix for the oracle network, the following corollary holds from Theorem 2.

**Corollary 1.** Suppose the noise  $|\xi| \leq \widehat{W}_{\min}^*$  and the number of samples satisfies  $N = \Omega(K^6 d \log q \log(1/\varepsilon))$ . Let  $\{\mathbf{W}^{(t_1)}\}_{t_1=1}^{T_1}$  be the iterates generated from Algorithm 1 by setting  $r = d$ . Then, for any  $T_1 \geq \log(\widehat{W}_{\max}^*/\widehat{W}_{\min}^*)$ , we have

$$\mathbf{M}^{(T_1)} = \mathbf{M}^*. \quad (65)$$

*Proof of Corollary 1.* If we train on the dense network, from (63), we know that

$$\|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_\infty \leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{154\kappa^2\gamma K}}\right) \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_\infty + 2C\eta\sqrt{\frac{d \log q}{N_t}}|\xi|. \quad (66)$$

Hence, we have

$$\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_\infty \leq \left(1 - \frac{1 - \varepsilon_0}{\sqrt{154\kappa^2\gamma K}}\right)^{T_1} \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_\infty + 2C\eta\sqrt{\frac{d \log q}{N_t}}|\xi|. \quad (67)$$

With  $T_1 \geq \log(2\widehat{W}_{\max}^*/\widehat{W}_{\min}^*)$ , we have

$$\left(1 - \frac{1 - \varepsilon_0}{\sqrt{154\kappa^2\gamma K}}\right)^{T_1} \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_\infty \leq \frac{1}{4}\widehat{W}_{\min}^* \cdot \frac{\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_\infty}{\|\mathbf{W}^*\|_\infty} \leq \frac{1}{4}\widehat{W}_{\min}^*. \quad (68)$$

Since  $N = \Omega(K^6 d \log q \log(1/\varepsilon))$  and  $|\xi| \leq \widehat{W}_{\min}^*$ , we have

$$2C\eta\sqrt{\frac{d \log q}{N_t}}|\xi| \leq \frac{1}{4}\widehat{W}_{\min}^*. \quad (69)$$

From (68) and (69), we know that

$$\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_\infty \leq \frac{1}{2}\widehat{W}_{\min}^*. \quad (70)$$

Therefore, for any entry in  $\mathbf{W}_{i,j}^{(T_1)}$ , if the corresponding entry in augmented ground-truth weights  $\mathbf{W}^*$  is zero, we have

$$|W_{i,j}^{(T_1)}| \leq \frac{1}{2}\widehat{W}_{\min}^*; \quad (71)$$

if the corresponding entry in  $\mathbf{W}^*$  is non-zero, we have

$$|W_{i,j}^{(T_1)}| \geq |\widehat{W}_{i,j}^*| - \frac{1}{2}\widehat{W}_{\min}^* \geq \frac{1}{2}\widehat{W}_{\min}^*. \quad (72)$$

As we know that there are only  $\sum_j r_j/(Kd)$  fraction of non-zero weights in the ground-truth model,  $\mathbf{M}^{(T_1)} = \mathbf{M}^*$  holds.  $\square$

## F Proof of Lemma 1

Instead of providing the proof for Lemma 1, we turn to prove a more general bound for the performance of tensor initialization method as shown in Lemma 6. One can easily verify that Lemma 1 holds naturally from Lemma 6.

Recall that in Appendix B, the estimation of  $\mathbf{w}_{j,\Omega_j}^*$  are converted into estimate the augmented vector  $\mathbf{u}_j^*$ . Further, the estimation of  $\mathbf{u}_j^*$  are divided into estimating three parts: (1) the estimation of the magnitude of  $\mathbf{u}_j^*$ , which is denoted as  $\widehat{\alpha}_j$ ; (2) the estimation of the subspace of  $\mathbf{u}_j^*$ , which is denoted as  $\widehat{\mathbf{V}}$ ; (3) the estimation of the representation of  $\mathbf{u}_j^*$  on subspace  $\mathbf{V}$ , which is denoted as  $\widehat{\mathbf{s}}_j$ . Hence, the major idea of proving Lemma 6 is to characterize the difference of these three estimations to its ground-truth, which are summarized in Lemmas 7, 8 and 9, respectively.

**Lemma 6.** Assume the noise level  $|\xi| \leq K\sigma_1$  and the number of samples  $N \gtrsim \kappa^8 K^5 r_{\max} \log^6 q$  with some large constant  $q$ , the tensor initialization method in Subroutine 1 outputs  $\mathbf{W}^{(0)}$  such that

$$\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_2 \lesssim \kappa^6 \sqrt{\frac{K r_{\max} \log q}{N}} (\sigma_1 + |\xi|) \quad (73)$$

with probability at least  $1 - q^{-r_{\max}}$ .



### F.1 Proof of Lemma 6

**Lemma 7.** Suppose  $\mathbf{M}_2$  is defined as in (22) and  $\widehat{\mathbf{M}}_2$  is the estimation of  $\mathbf{M}_2$  by samples  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ . Then, with probability  $1 - q^{-r_{\max}}$ , we have

$$\|\widehat{\mathbf{M}}_2 - \mathbf{M}_2\| \lesssim \sqrt{\frac{r_{\max} \log q}{N}} (\sigma_1 + |\xi|), \quad (74)$$

provided that  $N \gtrsim r_{\max} \log^4 q$ .

**Lemma 8.** Let  $\widehat{\mathbf{V}}$  be generated by step 4 in Subroutine 1. Suppose  $\mathbf{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  is defined as in (29) and  $\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  is the estimation of  $\mathbf{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$  by samples  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ . Further, we assume  $\mathbf{V} \in \mathbb{R}^{r \times K}$  is an orthogonal basis of  $\{\mathbf{u}_j^*\}_{j=1}^K$  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  $1 - q^{-r_{\max}}$ , we have

$$\|\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})\| \lesssim \sqrt{\frac{\log q}{N}} (\sigma_1 + |\xi|). \quad (75)$$

**Lemma 9.** Suppose  $\mathbf{M}_1$  is defined as in (21) and  $\widehat{\mathbf{M}}_1$  is the estimation of  $\mathbf{M}_1$  by samples  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ . Then, with probability  $1 - q^{-r_{\max}}$ , we have

$$\|\widehat{\mathbf{M}}_1 - \mathbf{M}_1\| \lesssim \sqrt{\frac{r_{\max} \log q}{N}} (\sigma_1 + |\xi|) \quad (76)$$

provided that  $N \gtrsim r_{\max} \log^4 d$ .

**Lemma 10** ([53], 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 11** ([69], Lemma E.6). Let  $\mathbf{V} \in \mathbb{R}^{r \times K}$  be an orthogonal basis of  $\widetilde{\mathbf{w}}^*$  and  $\widehat{\mathbf{V}}$  be generated by step 4 in Subroutine 1. Assume  $\|\widehat{\mathbf{M}}_2 - \mathbf{M}_2\|_2 \leq \sigma_K(\mathbf{M}_2)/10$ . Then, we have

$$\|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\|_2 \leq \frac{\|\mathbf{M}_2 - \widehat{\mathbf{M}}_2\|}{\sigma_K(\mathbf{M}_2)}. \quad (77)$$

**Lemma 12** ([69], Lemmas E.13 and E.14). Let  $\mathbf{V} \in \mathbb{R}^{r \times K}$  be an orthogonal basis of  $\widetilde{\mathbf{w}}^*$  and  $\widehat{\mathbf{V}}$  be generated by step 4 in Subroutine 1. Assume  $\mathbf{M}_1$  can be written in the form of (25) with some constant  $\psi_1$ , and let  $\widehat{\mathbf{M}}_1$  be the estimation of  $\mathbf{M}_1$  by samples  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ . Let  $\widehat{\boldsymbol{\alpha}}$  be the optimal solutions of (28) with  $\widehat{\mathbf{u}}_j^* = \widehat{\mathbf{V}}\widehat{\mathbf{s}}_j$ . Then, for each  $j \in \{1, 2, \dots, K\}$ , if

$$\begin{aligned} T_1 &:= \|\mathbf{V}\mathbf{V}^T - \widehat{\mathbf{V}}\widehat{\mathbf{V}}^T\|_2 \leq \frac{1}{\kappa^2 \sqrt{K}}, \\ T_2 &:= \|\widehat{\mathbf{u}}_j^* - \widehat{\mathbf{V}}^T \widehat{\mathbf{s}}_j\|_2 \leq \frac{1}{\kappa^2 \sqrt{K}}, \\ T_3 &:= \|\widehat{\mathbf{M}}_1 - \mathbf{M}_1\|_2 \leq \frac{1}{4} \|\mathbf{M}_1\|_2, \end{aligned} \quad (78)$$

then we have

$$|\alpha_j^* - \widehat{\alpha}_j| \leq \left( \kappa^4 K^{\frac{3}{2}} (T_1 + T_2) + \kappa^2 K^{\frac{1}{2}} T_3 \right) |\alpha_j^*|, \quad (79)$$

where  $\alpha_j^* = \|\mathbf{u}_j^*\|_2$ .

*Proof of Lemma 1.* By simple calculation, we have

$$\begin{aligned}
& \| \mathbf{u}_j^* - |\hat{\alpha}_j| \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 \\
& \leq \| \mathbf{u}_j^* - \|\mathbf{u}_j^*\|_2 \hat{\mathbf{V}} \hat{\mathbf{s}}_j + \|\mathbf{u}_j^*\|_2 \hat{\mathbf{V}} \hat{\mathbf{s}}_j - |\hat{\alpha}_j| \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 \\
& \leq \| \mathbf{u}_j^* - \|\mathbf{u}_j^*\|_2 \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 + \| \|\mathbf{u}_j^*\|_2 \hat{\mathbf{V}} \hat{\mathbf{s}}_j - |\hat{\alpha}_j| \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 \\
& \leq \| \mathbf{u}_j^* \|_2 \| \bar{\mathbf{u}}_j^* - \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 + \| \|\mathbf{u}_j^*\|_2 - |\hat{\alpha}_j| \| \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 \\
& \leq \sigma_1 (\| \bar{\mathbf{u}}_j^* - \hat{\mathbf{V}} \hat{\mathbf{V}}^T \bar{\mathbf{u}}_j^* \|_2 + \| \hat{\mathbf{V}}^T \bar{\mathbf{u}}_j^* - \hat{\mathbf{s}}_j \|_2) + \| \|\mathbf{u}_j^*\|_2 - |\hat{\alpha}_j| \| \\
& := \sigma_1 (I_1 + I_2) + I_3.
\end{aligned} \tag{80}$$

From Lemma 11, we have

$$I_1 = \| \bar{\mathbf{u}}_j^* - \hat{\mathbf{V}} \hat{\mathbf{V}}^T \bar{\mathbf{u}}_j^* \|_2 \leq \| \mathbf{V} \mathbf{V}^T - \hat{\mathbf{V}} \hat{\mathbf{V}}^T \|_2 \leq \frac{\| \hat{\mathbf{M}}_2 - \mathbf{M}_2 \|_2}{\sigma_K(\mathbf{M}_2)}, \tag{81}$$

where the last inequality comes from Lemma 7. Then, from (26), we know that

$$\sigma_K(\mathbf{M}_2) \lesssim \min_{1 \leq j \leq K} \| \mathbf{u}_j^* \|_2 = \min_{1 \leq j \leq K} \| \tilde{\mathbf{w}}_{j, \Omega_j}^* \|_2 \lesssim \sigma_K. \tag{82}$$

From Theorem 3 in [31], we have

$$I_2 = \| \hat{\mathbf{V}}^T \bar{\mathbf{u}}_j^* - \hat{\mathbf{s}}_j \|_2 \lesssim \frac{\kappa}{\sigma_K} \| \hat{\mathbf{M}}_3(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}}) - \mathbf{M}_3(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}}) \|_2. \tag{83}$$

To guarantee the condition (78) in Lemma 12 hold, according to Lemmas 7 and 8, we need  $N \gtrsim \kappa^3 K r_{\max} \log q$ . Then, from Lemma 12, we have

$$I_3 = \left( \kappa^4 K^{3/2} (I_1 + I_2) + \kappa^2 K^{1/2} \| \hat{\mathbf{M}}_1 - \mathbf{M}_1 \| \right) \sigma_1. \tag{84}$$

When  $r_{\max} \gg K$ , according to Lemmas 7, 8 and 9, we have

$$\| \mathbf{u}_j^* - |\hat{\alpha}_j| \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 \lesssim \kappa^6 \sqrt{\frac{r_{\max} \log q}{N}} (\sigma_1 + |\xi|) \tag{85}$$

provided that  $N \gtrsim K^3 r_{\max} \log^4 d$ .

In conclusion, we have

$$\begin{aligned}
\| \mathbf{W}^{(0)} - \mathbf{W}^* \|_F &= \| \tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}^{(0)} \|_2 \leq \sqrt{K} \cdot \| \mathbf{w}_{j, \Omega_j}^* - \mathbf{w}_{j, \Omega_j}^{(0)} \|_2 \\
&= \sqrt{K} \cdot \| \mathcal{F}_j^\dagger(\mathbf{u}_j^* - \hat{\mathbf{u}}_j^*) \|_2 \\
&\leq \sqrt{K} \cdot \| \mathbf{u}_j^* - \hat{\mathbf{u}}_j^* \|_2 \\
&= \sqrt{K} \cdot \| \mathbf{u}_j^* - |\hat{\alpha}_j| \hat{\mathbf{V}} \hat{\mathbf{s}}_j \|_2 \\
&\lesssim \kappa^6 \sqrt{\frac{K r_{\max} \log q}{N}} (\sigma_1 + |\xi|).
\end{aligned} \tag{86}$$

□

## G Additional proof of the lemmas in Appendix C

### G.1 Proof of Lemma 3

The eigenvalues of  $\nabla^2 f$  at any fixed point  $\tilde{\mathbf{w}}$  is bounded through the ones at the ground truth  $\tilde{\mathbf{w}}^*$  by using Lemma 2. The eigenvalues of  $\nabla^2 f$  at ground truth  $\tilde{\mathbf{w}}^*$  is bounded in (89) and (90).

**Lemma 13.** *Let  $f$  be the population risk function in (14) and  $\tilde{\mathbf{w}}$  satisfy (6), then we have*

$$\| \nabla^2 f(\tilde{\mathbf{w}}) - \nabla^2 f(\tilde{\mathbf{w}}^*) \|_2 \leq \frac{4 \| \tilde{\mathbf{w}}^* - \tilde{\mathbf{w}} \|_2}{\sigma_K}. \tag{87}$$

*Proof of Lemma 3.* Let  $\lambda_{\max}(\tilde{\mathbf{w}})$  and  $\lambda_{\min}(\tilde{\mathbf{w}})$  denote the largest and smallest eigenvalues of  $\nabla^2 f_{\mathcal{D}}$  at point  $\tilde{\mathbf{w}}$ , respectively. Then, from Lemma 2, we have

$$\begin{aligned} \lambda_{\max}(\tilde{\mathbf{w}}) &\leq \lambda_{\max}(\tilde{\mathbf{w}}^*) + \|\nabla^2 f(\tilde{\mathbf{w}}) - \nabla^2 f(\tilde{\mathbf{w}}^*)\|_2, \\ \text{and } \lambda_{\min}(\tilde{\mathbf{w}}) &\geq \lambda_{\min}(\tilde{\mathbf{w}}^*) - \|\nabla^2 f(\tilde{\mathbf{w}}) - \nabla^2 f(\tilde{\mathbf{w}}^*)\|_2. \end{aligned} \quad (88)$$

Next, we provide the lower bound of Hessian of population function at ground truth  $\tilde{\mathbf{w}}^*$ . Then, we have

$$\begin{aligned} \min_{\|\alpha\|_2=1} \alpha^T \nabla^2 f(\tilde{\mathbf{w}}^*) \alpha &= \frac{1}{K^2} \min_{\|\alpha\|_2=1} \mathbb{E}_{\mathbf{x}} \left( \sum_{j=1}^K \alpha_j^T \mathbf{x}_{\Omega_j} \phi'(\mathbf{w}_{j,\Omega_j}^{*T} \mathbf{x}_{\Omega_j}) \right)^2 \\ &= \frac{1}{K^2} \min_{\|\tilde{\alpha}\|_2=1, \text{supp}(\tilde{\alpha}_j)=\text{supp}(\mathbf{w}_j^*)} \mathbb{E}_{\mathbf{x}} \left( \sum_{j=1}^K \tilde{\alpha}_j^T \mathbf{x} \phi'(\mathbf{w}_j^{*T} \mathbf{x}) \right)^2 \\ &\geq \frac{1}{K^2} \min_{\|\tilde{\alpha}\|_2=1} \mathbb{E}_{\mathbf{x}} \left( \sum_{j=1}^K \tilde{\alpha}_j^T \mathbf{x} \phi'(\mathbf{w}_j^{*T} \mathbf{x}) \right)^2 \\ &\geq \frac{\rho}{11\kappa^2 \lambda K^2}, \end{aligned} \quad (89)$$

where  $\tilde{\alpha} \in \mathbb{R}^{Kd}$  with  $\tilde{\alpha}_j \in \mathbb{R}^d$ , and the last inequality comes from Lemma D.6 [69].

Next, the upper bound of Hessian of population function at ground truth  $\tilde{\mathbf{w}}^*$  can be bounded in the following way. For any  $\alpha$ , we have

$$\begin{aligned} \alpha^T \nabla^2 f(\tilde{\mathbf{w}}^*) \alpha &= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \left( \sum_{j=1}^K \alpha_j^T \mathbf{x}_{\Omega_j} \phi'(\mathbf{w}_{j,\Omega_j}^{*T} \mathbf{x}_{\Omega_j}) \right)^2 \\ &\leq \frac{2}{K^2} \cdot \mathbb{E}_{\mathbf{x}} \sum_{j=1}^K \left( \alpha_j^T \mathbf{x}_{\Omega_j} \phi'(\mathbf{w}_{j,\Omega_j}^{*T} \mathbf{x}_{\Omega_j}) \right)^2 \\ &= \frac{2}{K^2} \sum_{j=1}^K \mathbb{E}_{\mathbf{x}} \left( \alpha_j^T \mathbf{x}_{\Omega_j} \phi'(\mathbf{w}_{j,\Omega_j}^{*T} \mathbf{x}_{\Omega_j}) \right)^2 \\ &\leq \frac{2}{K^2} \sum_{j=1}^K \left( \mathbb{E}_{\mathbf{x}} (\alpha_j^T \mathbf{x}_{\Omega_j})^4 \mathbb{E}_{\mathbf{x}} |\phi'|^4 \right)^{\frac{1}{2}} \\ &\leq \frac{2}{K^2} \cdot K \cdot 3 = \frac{6}{K}. \end{aligned} \quad (90)$$

Then, from Lemma 13, when  $\tilde{\mathbf{w}}$  satisfies (6), we have that

$$\|\nabla^2 f(\tilde{\mathbf{w}}) - \nabla^2 f(\tilde{\mathbf{w}}^*)\|_2 \leq \frac{\varepsilon_0 \rho}{11\kappa^2 \gamma}. \quad (91)$$

Hence, from (88) and (91), we have that

$$\frac{(1 - \varepsilon_0) \rho}{11\kappa^2 \gamma K^2} \mathbf{I} \leq \nabla^2 f(\tilde{\mathbf{w}}) \leq \frac{7}{K} \mathbf{I}. \quad (92)$$

□

## G.2 Proof of Lemma 4

We first show that the second order derivative of  $\hat{f}_{\mathcal{D}}$  is a sum of several random sub-exponential variables as shown in (101) and (102). Then, by concentration theory, i.e., Chernoff bound, we can show that the error bound of  $\nabla^2 \hat{f}_{\mathcal{D}}$  to its expectation.

**Definition 1** (Definition 5.7, [54]). A random variable  $X$  is called a sub-Gaussian random variable if it satisfies

$$(\mathbb{E}|X|^p)^{1/p} \leq c_1 \sqrt{p} \quad (93)$$

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

$$\mathbb{E}e^{s(X-\mathbb{E}X)} \leq e^{c_2 \|X\|_{\psi_2}^2 s^2} \quad (94)$$

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, [54]). A random variable  $X$  is called a sub-exponential random variable if it satisfies

$$(\mathbb{E}|X|^p)^{1/p} \leq c_3 p \quad (95)$$

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

$$\mathbb{E}e^{s(X-\mathbb{E}X)} \leq e^{c_4 \|X\|_{\psi_1}^2 s^2} \quad (96)$$

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 14** (Lemma 5.2, [54]). Let  $\mathcal{B}(0, 1) \in \{\boldsymbol{\alpha} \mid \|\boldsymbol{\alpha}\|_2 = 1, \boldsymbol{\alpha} \in \mathbb{R}^d\}$  denote a unit ball in  $\mathbb{R}^d$ . Then, a subset  $\mathcal{S}_\xi$  is called a  $\xi$ -net of  $\mathcal{B}(0, 1)$  if every point  $\mathbf{z} \in \mathcal{B}(0, 1)$  can be approximated to within  $\xi$  by some point  $\boldsymbol{\alpha} \in \mathcal{B}(0, 1)$ , i.e.  $\|\mathbf{z} - \boldsymbol{\alpha}\|_2 \leq \xi$ . Then the minimal cardinality of a  $\xi$ -net  $\mathcal{S}_\xi$  satisfies

$$|\mathcal{S}_\xi| \leq (1 + 2/\xi)^d. \quad (97)$$

**Lemma 15** (Lemma 5.3, [54]). Let  $\mathbf{A}$  be an  $d_1 \times d_2$  matrix, and let  $\mathcal{S}_\xi(d)$  be a  $\xi$ -net of  $\mathcal{B}(0, 1)$  in  $\mathbb{R}^d$  for some  $\xi \in (0, 1)$ . Then

$$\|\mathbf{A}\|_2 \leq (1 - \xi)^{-1} \max_{\boldsymbol{\alpha}_1 \in \mathcal{S}_\xi(d_1), \boldsymbol{\alpha}_2 \in \mathcal{S}_\xi(d_2)} |\boldsymbol{\alpha}_1^T \mathbf{A} \boldsymbol{\alpha}_2|. \quad (98)$$

*Proof of Lemma 4.* Recall the definition of  $f$  and  $\hat{f}$  in (14) and (12), we have

$$\begin{aligned} & \frac{\partial^2 f}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} - \frac{\partial^2 \hat{f}_{\mathcal{D}}}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} \\ &= \mathbb{E}_{\mathbf{x}} \left[ \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{\Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}}) \mathbf{x}_{\Omega_{j_1}} \mathbf{x}_{\Omega_{j_2}}^T \right. \\ & \quad \left. - \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{n, \Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{n, \Omega_{j_2}}) \mathbf{x}_{n, \Omega_{j_1}} \mathbf{x}_{n, \Omega_{j_2}}^T \right]. \end{aligned} \quad (99)$$

For any  $\boldsymbol{\alpha}$ , we have

$$\begin{aligned} & \|\nabla^2 f - \nabla^2 \hat{f}_{\mathcal{D}}\|_2 \\ &= \max_{\|\boldsymbol{\alpha}\|_2=1} \left| \boldsymbol{\alpha}^T (\nabla^2 f - \nabla^2 \hat{f}_{\mathcal{D}}) \boldsymbol{\alpha} \right| \\ &= \sum_{j_1=1}^K \sum_{j_2=1}^K \max_{\|\boldsymbol{\alpha}\|_2=1} \left| \boldsymbol{\alpha}_{j_1}^T \left( \frac{\partial^2 f}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} - \frac{\partial^2 \hat{f}_{\mathcal{D}}}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} \right) \boldsymbol{\alpha}_{j_2} \right| \\ &= \frac{1}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \max_{\|\boldsymbol{\alpha}\|_2=1} \mathbb{E}_{\mathbf{x}} \left[ \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{\Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}}) \boldsymbol{\alpha}_{j_1}^T \mathbf{x}_{\Omega_{j_1}} \boldsymbol{\alpha}_{j_2}^T \mathbf{x}_{\Omega_{j_2}} \right. \\ & \quad \left. - \frac{1}{N} \sum_{n=1}^N \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{n, \Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{n, \Omega_{j_2}}) \boldsymbol{\alpha}_{j_1}^T \mathbf{x}_{n, \Omega_{j_1}} \boldsymbol{\alpha}_{j_2}^T \mathbf{x}_{n, \Omega_{j_2}} \right]. \end{aligned} \quad (100)$$

Then, define  $Z_n(j_1, j_2) = \phi(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{n, \Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{n, \Omega_{j_2}}) \boldsymbol{\alpha}_{j_1}^T \mathbf{x}_{n, \Omega_{j_1}} \boldsymbol{\alpha}_{j_2}^T \mathbf{x}_{n, \Omega_{j_2}}$ , and we say  $Z$  belongs to sub-Exponential distribution by Definition 2. If  $|\Omega_{j_1} \cap \Omega_{j_2}| \neq \emptyset$ , namely,  $\Omega_{j_1}$  and  $\Omega_{j_2}$  are not disjointed, we have

$$\begin{aligned} (\mathbb{E}|Z_n|^p)^{1/p} &\leq \left( \mathbb{E} \left| (\boldsymbol{\alpha}_{j_1}^T \mathbf{x}_{n, \Omega_{j_1}}) \cdot (\boldsymbol{\alpha}_{j_2}^T \mathbf{x}_{n, \Omega_{j_2}}) \right|^p \right)^{1/p} \\ &\leq \left( \mathbb{E} \left| (\boldsymbol{\alpha}_{j_1}^T \mathbf{x}_{n, \Omega_{j_1}}) \right|^{2p} \right)^{1/(2p)} \cdot \left( \mathbb{E} \left| (\boldsymbol{\alpha}_{j_2}^T \mathbf{x}_{n, \Omega_{j_2}}) \right|^{2p} \right)^{1/(2p)} \\ &\leq C_{\mathbf{x}} \cdot \sqrt{2p} \cdot C_{\mathbf{x}} \sqrt{2p} \\ &= 2C_{\mathbf{x}}^2 \cdot p. \end{aligned} \quad (101)$$

While if  $|\Omega_{j_1} \cap \Omega_{j_2}| = \emptyset$ , namely,  $\Omega_{j_1}$  and  $\Omega_{j_2}$  are disjointed, we have

$$\begin{aligned} (\mathbb{E}|Z_n|^p)^{1/p} &\leq \left( \mathbb{E} \left| (\boldsymbol{\alpha}_{j_1}^T \mathbf{x}_{n, \Omega_{j_1}}) \cdot (\boldsymbol{\alpha}_{j_2}^T \mathbf{x}_{n, \Omega_{j_2}}) \right|^p \right)^{1/p} \\ &= \left( \mathbb{E} \left| (\boldsymbol{\alpha}_{j_1}^T \mathbf{x}_{n, \Omega_{j_1}}) \right|^p \right)^{1/p} \cdot \left( \mathbb{E} \left| (\boldsymbol{\alpha}_{j_2}^T \mathbf{x}_{n, \Omega_{j_2}}) \right|^p \right)^{1/p} \\ &\leq C_{\mathbf{x}} \cdot \sqrt{p} \cdot C_{\mathbf{x}} \sqrt{p} \\ &= C_{\mathbf{x}}^2 \cdot p. \end{aligned} \quad (102)$$

Then, we have

$$\mathbb{E}_{Z_n} e^{s(Z_n - \mathbb{E}Z_n)} \leq e^{-C \|Z_n\|_{\psi_1}^2 s^2} \quad (103)$$

for some constant  $C > 0$  and any  $s \in \mathbb{R}$ . From Chernoff bound, we have

$$\text{Prob} \left\{ \left| \frac{1}{N} \sum_{n=1}^N (Z_n - \mathbb{E}Z_n) \right| < t \right\} \leq 1 - \frac{e^{-C \|Z_n\|_{\psi_1}^2 \cdot N s^2}}{e^{N s t}}. \quad (104)$$

Let us select  $t = \|Z_n\|_{\psi_1} \sqrt{\frac{(r_{j_1} + r_{j_2}) \log q}{N}}$  and  $s = \frac{\sqrt{2}}{C \|Z_n\|_{\psi_1}^2} \cdot t$ , then we have

$$\left| \frac{1}{N} \sum_{n=1}^N (Z_n(j_1, j_2) - \mathbb{E}Z_n(j_1, j_2)) \right| \leq \|Z_n\|_{\psi_1} \sqrt{\frac{(r_{j_1} + r_{j_2}) \log q}{N}} \quad (105)$$

with probability at least  $1 - q^{-(r_{j_1} + r_{j_2})}$ .

Hence, from Lemma 15, we have

$$\begin{aligned} &\max_{\|\boldsymbol{\alpha}_{j_1}\|_2 \leq 1, \|\boldsymbol{\alpha}_{j_2}\|_2 \leq 1} \left| \boldsymbol{\alpha}_{j_1}^T \left( \frac{\partial^2 f}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} - \frac{\partial^2 \hat{f}_{\mathcal{D}}}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} \right) \boldsymbol{\alpha}_{j_2} \right| \\ &\leq 2 \left| \frac{1}{N} \sum_{n=1}^N (Z_n - \mathbb{E}Z_n) \right| \end{aligned} \quad (106)$$

with probability at least  $1 - (|\mathcal{S}_{\frac{1}{2}}(r_{j_1})| \cdot |\mathcal{S}_{\frac{1}{2}}(r_{j_2})|) \cdot q^{-(r_{j_1} + r_{j_2})}$ , where  $\mathcal{S}_{\frac{1}{2}}(r_{j_1})$  and  $\mathcal{S}_{\frac{1}{2}}(r_{j_2})$  are the covering sets defined in Lemma 14. From Lemma 14, we know that  $|\mathcal{S}_{\frac{1}{2}}(r_{j_1})| \cdot |\mathcal{S}_{\frac{1}{2}}(r_{j_2})| \leq 5^{(r_{j_1} + r_{j_2})}$ . As long as  $q$  is a constant that is larger than 5, (106) holds with the probability at least  $1 - \left(\frac{q}{5}\right)^{-(r_{j_1} + r_{j_2})}$ . For notation simplification, we use probability  $1 - q^{-(r_{j_1} + r_{j_2})}$  instead.

From (101) and (102), we know that

$$\|Z_n(j_1, j_2)\|_{\psi_1} \leq \begin{cases} 2C_{\mathbf{x}}^2, & \text{if } \Omega_{j_1} \text{ and } \Omega_{j_2} \text{ are joint sets} \\ C_{\mathbf{x}}^2, & \text{if } \Omega_{j_1} \text{ and } \Omega_{j_2} \text{ are disjoint sets} \end{cases}. \quad (107)$$

Hence, we have

$$\begin{aligned}
& \|\nabla^2 f(\tilde{\mathbf{w}}) - \nabla^2 \hat{f}_\Omega(\tilde{\mathbf{w}})\|_2 \\
& \leq \sum_{j_1=1}^K \sum_{j_2=1}^K \max_{\|\alpha\|_2=1} \left| \alpha_{j_1}^T \left( \frac{\partial^2 f}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} - \frac{\partial^2 \hat{f}_\mathcal{D}}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} \right) \alpha_{j_2} \right| \\
& \leq \frac{2}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \max_{\|\alpha\|_2=1} \left| \frac{1}{N} \sum_{n=1}^N (Z_n(j_1, j_2) - \mathbb{E} Z_n(j_1, j_2)) \right| \\
& \lesssim \frac{1}{K^2} \sum_{j_1=1}^K \sum_{j_2=1}^K \sqrt{\frac{(1 + \delta_{j_1, j_2})^2 (r_{j_1} + r_{j_2}) \log q}{N}}
\end{aligned} \tag{108}$$

with probability at least  $1 - q^{-r_{\min}}$ , where  $\delta_{j_1, j_2}$  equals to 0 if  $\Omega_{j_1}$  and  $\Omega_{j_2}$  are disjoint and 1 otherwise.  $\square$

## H Proof of Lemma 5

*Proof of Lemma 5.* The first-order derivative of the empirical risk function is written as

$$\begin{aligned}
\frac{\partial \hat{f}_\mathcal{D}}{\partial \mathbf{w}_{k, \Omega_k}} &= \frac{1}{K \cdot N} \sum_{n=1}^N \left( y_n - \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{w}_{j, \Omega_j}^T \mathbf{x}_{n, \Omega_j}) \right) \mathbf{x}_{n, \Omega_k} \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k}) \\
&= \frac{1}{K^2 \cdot N} \sum_{n=1}^N \sum_{j=1}^K (\phi(\mathbf{w}_{j, \Omega_j}^{*T} \mathbf{x}_{n, \Omega_j}) - \phi(\mathbf{w}_{j, \Omega_j}^T \mathbf{x}_{n, \Omega_j})) \mathbf{x}_{n, \Omega_k} \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k}) \\
&\quad + \frac{1}{K \cdot N} \sum_{j=1}^K \xi_n \mathbf{x}_{n, \Omega_k} \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k})
\end{aligned} \tag{109}$$

Define  $\mathbf{z}_n(j, k) = (\phi(\mathbf{w}_{j, \Omega_j}^{*T} \mathbf{x}_{n, \Omega_j}) - \phi(\mathbf{w}_{j, \Omega_j}^T \mathbf{x}_{n, \Omega_j})) \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k}) \mathbf{x}_{n, \Omega_k}$ . Then, for any  $\alpha_k \in \mathbb{R}^r$ , we have

$$\begin{aligned}
& p^{-1} \left( \mathbb{E}_{\mathbf{x}} |\alpha_k^T \mathbf{z}_n|^p \right)^{\frac{1}{p}} \\
&= p^{-1} \left( \mathbb{E}_{\mathbf{x}} |(\alpha_k^T \mathbf{x}_{n, \Omega_k}) (\phi(\mathbf{w}_{j, \Omega_j}^{*T} \mathbf{x}_{n, \Omega_j}) - \phi(\mathbf{w}_{j, \Omega_j}^T \mathbf{x}_{n, \Omega_j})) \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k})|^p \right)^{\frac{1}{p}} \\
&\leq p^{-1} \left( \mathbb{E}_{\mathbf{x}} |(\alpha_k^T \mathbf{x}_{n, \Omega_k}) (\phi(\mathbf{w}_{j, \Omega_j}^{*T} \mathbf{x}_{n, \Omega_j}) - \phi(\mathbf{w}_{j, \Omega_j}^T \mathbf{x}_{n, \Omega_j}))|^p \right)^{\frac{1}{p}}.
\end{aligned} \tag{110}$$

If  $\Omega_j$  and  $\Omega_k$  are joint, then

$$\begin{aligned}
& p^{-1} \left( \mathbb{E}_{\mathbf{x}} |\alpha_k^T \mathbf{z}_n|^p \right)^{\frac{1}{p}} \\
&\leq p^{-1} \left( \mathbb{E}_{\mathbf{x}} |\alpha_k^T \mathbf{x}_{n, \Omega_k}|^{2p} \right)^{\frac{1}{2p}} \cdot \left( \mathbb{E}_{\mathbf{x}} |\phi(\mathbf{w}_{j, \Omega_j}^{*T} \mathbf{x}_{n, \Omega_j}) - \phi(\mathbf{w}_{j, \Omega_j}^T \mathbf{x}_{n, \Omega_j})|^{2p} \right)^{\frac{1}{2p}} \\
&\leq p^{-1} \left( \mathbb{E}_{\mathbf{x}} |\alpha_k^T \mathbf{x}_{n, \Omega_k}|^{2p} \right)^{\frac{1}{2p}} \cdot \left( \mathbb{E}_{\mathbf{x}} |(\mathbf{w}_{j, \Omega_j}^* - \mathbf{w}_{j, \Omega_j})^T \mathbf{x}_{n, \Omega_j}|^{2p} \right)^{\frac{1}{2p}} \\
&\leq 2 \|\mathbf{w}_{j, \Omega_j}^* - \mathbf{w}_{j, \Omega_j}\|_2 \leq 2 \|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|_2.
\end{aligned} \tag{111}$$

If  $\Omega_j$  and  $\Omega_k$  are disjoint, then

$$\begin{aligned}
& p^{-1} \left( \mathbb{E}_{\mathbf{x}} |\alpha_k^T \mathbf{z}_n|^p \right)^{\frac{1}{p}} \\
&\leq p^{-1} \left( \mathbb{E}_{\mathbf{x}} |\alpha_j^T \mathbf{x}_{n, \Omega_j}|^p \right)^{\frac{1}{p}} \cdot \left( \mathbb{E}_{\mathbf{x}} |\phi(\mathbf{w}_{j, \Omega_j}^{*T} \mathbf{x}_{n, \Omega_j}) - \phi(\mathbf{w}_{j, \Omega_j}^T \mathbf{x}_{n, \Omega_j})|^p \right)^{\frac{1}{p}} \\
&\leq p^{-1} \left( \mathbb{E}_{\mathbf{x}} |\alpha_j^T \mathbf{x}_{n, \Omega_j}|^p \right)^{\frac{1}{p}} \cdot \left( \mathbb{E}_{\mathbf{x}} |(\mathbf{w}_{j, \Omega_j}^* - \mathbf{w}_{j, \Omega_j})^T \mathbf{x}_{n, \Omega_j}|^p \right)^{\frac{1}{p}} \\
&\leq \|\mathbf{w}_{j, \Omega_j}^* - \mathbf{w}_{j, \Omega_j}\|_2 \leq \|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|_2.
\end{aligned} \tag{112}$$

Following similar steps in (104), by Chernoff bound, we have

$$\left\| \frac{1}{N} \sum_{n=1}^N (\mathbf{z}_n - \mathbb{E}_{\mathbf{x}} \mathbf{z}_n) \right\|_2 \lesssim \|\mathbf{z}_n(j, k)\|_{\psi_1} \sqrt{\frac{r_j \log q}{N}} \cdot \|\mathbf{w}_{j, \Omega_j}^* - \mathbf{w}_{j, \Omega_j}\|_2 \quad (113)$$

with probability at least  $1 - q^{-r_j}$ , where

$$\|\mathbf{z}_n(j, k)\|_{\psi_1} = \begin{cases} 2\|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2, & \text{if } \Omega_k \text{ and } \Omega_j \text{ are joint,} \\ \|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2, & \text{if } \Omega_k \text{ and } \Omega_j \text{ are disjoint} \end{cases} \quad (114)$$

That is  $\|\mathbf{z}_n(j, k)\|_{\psi_1} = (1 + \delta_{j,k})\|\tilde{\mathbf{w}} - \tilde{\mathbf{w}}^*\|_2$ . Also, we know that  $\mathbf{x}_{n, \Omega_k} \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k})$  belongs to sub-Gaussian distribution as well. Then, by Chernoff bound, we have

$$\begin{aligned} \left\| \frac{1}{N} \sum_{n=1}^N \xi_n \mathbf{x}_{n, \Omega_k} \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k}) \right\|_2 &\lesssim |\xi| \cdot \left\| \frac{1}{N} \sum_{n=1}^N \mathbf{x}_{n, \Omega_k} \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k}) \right\|_2 \\ &\lesssim |\xi| \cdot \sqrt{\frac{r_k \log q}{N}} \end{aligned} \quad (115)$$

with probability at least  $q^{-r_k}$ .

In conclusion, we have

$$\begin{aligned} \|\nabla \hat{f}_{\mathcal{D}} - \nabla f\|_2 &\leq \sum_{k=1}^K \left\| \frac{\partial \hat{f}_{\mathcal{D}}}{\partial \mathbf{w}_k} - \frac{\partial f}{\partial \mathbf{w}_k} \right\|_2 \\ &\leq \sum_{k=1}^K \frac{1}{K^2} \sum_{j=1}^K \left\| \frac{1}{N} \sum_{n=1}^N (\mathbf{z}_n(j, k) - \mathbb{E}_{\mathbf{x}} \mathbf{z}_n(j, k)) \right\|_2 \\ &\quad + \sum_{k=1}^K \frac{1}{K} \left\| \frac{1}{N} \sum_{n=1}^N \xi_n \mathbf{x}_{n, \Omega_k} \phi'(\mathbf{w}_{k, \Omega_k}^T \mathbf{x}_{n, \Omega_k}) \right\|_2 \\ &\lesssim \frac{1}{K^2} \sum_{k=1}^K \sum_{j=1}^K \sqrt{\frac{(1 + \delta_{j,k})^2 r_k \log q}{N}} \|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|_2 + \frac{1}{K} \sum_{k=1}^K \sqrt{\frac{r_k \log q}{N}} \cdot |\xi|. \end{aligned} \quad (116)$$

□

## I Proof of Lemma 13

*Proof of Lemma 13.* Recall the definition of population risk function, we have

$$\frac{\partial^2 f(\mathbf{w}^*)}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} = \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^{*T} \mathbf{x}_{\Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^{*T} \mathbf{x}_{\Omega_{j_2}}) \mathbf{x}_{\Omega_{j_1}} \mathbf{x}_{\Omega_{j_2}}^T \quad (117)$$

and

$$\frac{\partial^2 f(\mathbf{w})}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} = \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{\Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}}) \mathbf{x}_{\Omega_{j_1}} \mathbf{x}_{\Omega_{j_2}}^T \quad (118)$$

Then, we have

$$\begin{aligned}
& \frac{\partial^2 f(\mathbf{w}^*)}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} - \frac{\partial^2 f(\mathbf{w})}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} \\
&= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \left[ \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^{*T} \mathbf{x}_{\Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^{*T} \mathbf{x}_{\Omega_{j_2}}) - \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{\Omega_{j_1}}) \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}}) \right] \mathbf{x}_{\Omega_{j_1}} \mathbf{x}_{\Omega_{j_2}}^T \\
&= \frac{1}{K^2} \mathbb{E}_{\mathbf{x}} \left[ \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^{*T} \mathbf{x}_{\Omega_{j_1}}) (\phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^{*T} \mathbf{x}_{\Omega_{j_2}}) - \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}})) \right. \\
&\quad \left. + \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}}) (\phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^{*T} \mathbf{x}_{\Omega_{j_1}}) - \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{\Omega_{j_1}})) \right] \mathbf{x}_{\Omega_{j_1}} \mathbf{x}_{\Omega_{j_2}}^T \\
&= \frac{1}{K^2} \left[ \mathbb{E}_{\mathbf{x}} \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^{*T} \mathbf{x}_{\Omega_{j_1}}) (\phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^{*T} \mathbf{x}_{\Omega_{j_2}}) - \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}})) \mathbf{x}_{\Omega_{j_1}} \mathbf{x}_{\Omega_{j_2}}^T \right. \\
&\quad \left. + \mathbb{E}_{\mathbf{x}} \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}}) (\phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^{*T} \mathbf{x}_{\Omega_{j_1}}) - \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^T \mathbf{x}_{\Omega_{j_1}})) \mathbf{x}_{\Omega_{j_1}} \mathbf{x}_{\Omega_{j_2}}^T \right] \\
&:= \frac{1}{K^2} (\mathbf{I}_1 + \mathbf{I}_2).
\end{aligned} \tag{119}$$

For any  $\alpha_{j_1} \in \mathbb{R}^{r_{j_1}}$  and  $\alpha_{j_2} \in \mathbb{R}^{r_{j_2}}$ , we have

$$\begin{aligned}
& \max_{\|\alpha_{j_1}\|_2, \|\alpha_{j_2}\|_2=1} \alpha_{j_1}^T \mathbf{I}_1 \alpha_{j_2} \\
&= \max_{\|\alpha_{j_1}\|_2, \|\alpha_{j_2}\|_2=1} \mathbb{E}_{\mathbf{x}} \phi'(\mathbf{w}_{j_1, \Omega_{j_1}}^{*T} \mathbf{x}_{\Omega_{j_1}}) (\phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^{*T} \mathbf{x}_{\Omega_{j_2}}) - \phi'(\mathbf{w}_{j_2, \Omega_{j_2}}^T \mathbf{x}_{\Omega_{j_2}})) \\
&\quad \cdot (\alpha_{j_1}^T \mathbf{x}_{\Omega_{j_1}}) \cdot (\alpha_{j_2}^T \mathbf{x}_{\Omega_{j_2}}) \\
&\leq \max_{\|\mathbf{a}\|_2=1} \mathbb{E}_{\mathbf{x}} \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}) (\phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^T \mathbf{x})) \cdot (\mathbf{a}^T \mathbf{x})^2,
\end{aligned} \tag{120}$$

where  $\mathbf{a} \in \mathbb{R}^d$ . Let  $I = \phi'(\mathbf{w}_{j_1}^{*T} \mathbf{x}) (\phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^T \mathbf{x})) \cdot (\mathbf{a}^T \mathbf{x})^2$ . It is easy to verify there exists a basis such that  $\mathcal{B} = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{a}_4^\perp, \dots, \mathbf{a}_d^\perp\}$  with  $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$  spans a subspace that contains  $\mathbf{a}, \mathbf{w}_{j_2}$  and  $\mathbf{w}_{j_2}^*$ . Then, for any  $\mathbf{x}$ , we have a unique  $\mathbf{z} = [z_1 \ z_2 \ \dots \ z_d]^T$  such that

$$\mathbf{x} = z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c} + \dots + z_d \mathbf{a}_d^\perp.$$

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

$$\begin{aligned}
I &= \mathbb{E}_{z_1, z_2, z_3} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})| \cdot |\mathbf{a}^T \mathbf{x}|^2 \\
&= \int |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})| \cdot |\mathbf{a}^T \mathbf{x}|^2 \cdot f_Z(z_1, z_2, z_3) dz_1 dz_2 dz_3,
\end{aligned}$$

where  $\mathbf{x} = z_1 \mathbf{a} + z_2 \mathbf{b} + z_3 \mathbf{c}$  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,

$$I = \int |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})| \cdot |r \cos \phi_1|^2 \cdot f_Z(R, \phi_1, \phi_2) R^2 \sin \phi_1 dR d\phi_1 d\phi_2. \tag{121}$$

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

$$f_Z(R, \phi_1, \phi_2) = \frac{1}{(2\pi)^{\frac{3}{2}}} e^{-\frac{x_1^2 + x_2^2 + x_3^2}{2}} = \frac{1}{(2\pi)^{\frac{3}{2}}} e^{-\frac{R^2}{2}}$$



only depends on  $R$ . Then, we have

$$\begin{aligned}
& I(i_2, j_2) \\
&= \int |\phi'(\mathbf{w}_{j_2}^T(\mathbf{x}/R)) - \phi'(\mathbf{w}_{j_2}^{*T}(\mathbf{x}/R))| \cdot |R \cos \phi_1|^2 \cdot f_Z(R) R^2 \sin \phi_1 dR d\phi_1 d\phi_2 \\
&= \int_0^\infty r^4 f_Z(R) dR \int_0^\pi \int_0^{2\pi} |\cos \phi_1|^2 \cdot \sin \phi_1 \cdot |\phi'(\mathbf{w}_{j_2}^T(\mathbf{x}/R)) - \phi'(\mathbf{w}_{j_2}^{*T}(\mathbf{x}/R))| d\phi_1 d\phi_2 \\
&\leq \sqrt{\frac{8}{\pi}} \int_0^\infty R^2 f_Z(R) dR \int_0^\pi \int_0^{2\pi} \sin \phi_1 \cdot |\phi'(\mathbf{w}_{j_2}^T(\mathbf{x}/R)) - \phi'(\mathbf{w}_{j_2}^{*T}(\mathbf{x}/R))| d\phi_1 d\phi_2 \\
&= \sqrt{\frac{8}{\pi}} \mathbb{E}_{z_1, z_2, z_3} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})| \\
&= \sqrt{\frac{8}{\pi}} \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})|.
\end{aligned} \tag{122}$$

Define a set  $\mathcal{A}_1 = \{\mathbf{x} | (\mathbf{w}_{j_2}^{*T} \mathbf{x})(\mathbf{w}_{j_2}^T \mathbf{x}) < 0\}$ . If  $\mathbf{x} \in \mathcal{A}_1$ , then  $\mathbf{w}_{j_2}^{*T} \mathbf{x}$  and  $\mathbf{w}_{j_2}^T \mathbf{x}$  have different signs, which means the value of  $\phi'(\mathbf{w}_{j_2}^T \mathbf{x})$  and  $\phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})$  are different. This is equivalent to say that

$$|\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})| = \begin{cases} 1, & \text{if } \mathbf{x} \in \mathcal{A}_1 \\ 0, & \text{if } \mathbf{x} \in \mathcal{A}_1^c \end{cases}. \tag{123}$$

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

$$|\mathbf{w}_{j_2}^{*T} \mathbf{x}| \leq |\mathbf{w}_{j_2}^{*T} \mathbf{x} - \mathbf{w}_{j_2}^T \mathbf{x}| \leq \|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\| \cdot \|\mathbf{x}\|. \tag{124}$$

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

$$\mathcal{A}_2 = \left\{ \mathbf{x} \mid \frac{|\mathbf{w}_{j_2}^{*T} \mathbf{x}|}{\|\mathbf{w}_{j_2}^*\| \|\mathbf{x}\|} \leq \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \right\} = \left\{ \theta_{\mathbf{x}, \mathbf{w}_{j_2}^*} \mid |\cos \theta_{\mathbf{x}, \mathbf{w}_{j_2}^*}| \leq \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \right\}. \tag{125}$$

Hence, we have that

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})|^2 &= \mathbb{E}_{\mathbf{x}} |\phi'(\mathbf{w}_{j_2}^T \mathbf{x}) - \phi'(\mathbf{w}_{j_2}^{*T} \mathbf{x})| \\
&= \text{Prob}(\mathbf{x} \in \mathcal{A}_1) \\
&\leq \text{Prob}(\mathbf{x} \in \mathcal{A}_2).
\end{aligned} \tag{126}$$

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

$$\begin{aligned}
\text{Prob}(\mathbf{x} \in \mathcal{A}_2) &= \frac{\pi - \arccos \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|}}{\pi} \leq \frac{1}{\pi} \tan(\pi - \arccos \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|}) \\
&= \frac{1}{\pi} \cot(\arccos \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|}) \\
&\leq \frac{2}{\pi} \frac{\|\mathbf{w}_{j_2}^* - \mathbf{w}_{j_2}\|}{\|\mathbf{w}_{j_2}^*\|} \\
&= \frac{2}{\pi} \frac{\|\mathbf{w}_{j_2, \Omega_{j_2}}^* - \mathbf{w}_{j_2, \Omega_{j_2}}\|}{\|\mathbf{w}_{j_2, \Omega_{j_2}}^*\|} \\
&\leq \frac{2}{\pi} \frac{\|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|}{\sigma_K}.
\end{aligned} \tag{127}$$

Hence, (122) and (127) suggest that

$$I \leq \frac{6}{\pi} \frac{\|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|}{\sigma_K}. \tag{128}$$

The same bound that shown in (128) holds for  $\mathbf{I}_2$  as well.

Therefore, we have

$$\begin{aligned}
& \|\nabla^2 f(\tilde{\mathbf{w}}) - \nabla^2 f(\tilde{\mathbf{w}}^*)\|_2 \\
& \leq \sum_{j_1=1}^K \sum_{j_2=1}^K \left\| \frac{\partial^2 f(\tilde{\mathbf{w}}^*)}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} - \frac{\partial^2 f(\tilde{\mathbf{w}})}{\partial \mathbf{w}_{j_1, \Omega_{j_1}} \partial \mathbf{w}_{j_2, \Omega_{j_2}}} \right\|_2 \\
& \leq \|\mathbf{I}_1 + \mathbf{I}_2\|_2 \leq \|\mathbf{I}_1\|_2 + \|\mathbf{I}_2\|_2 \\
& \leq \frac{12}{\pi} \frac{\|\tilde{\mathbf{w}}^* - \tilde{\mathbf{w}}\|_2}{\sigma_K}
\end{aligned} \tag{129}$$

□

## J Additional proofs of lemmas in Appendix F

### J.1 Error bound for the second-order moment

*Proof of Lemma 7.* For  $\widehat{\mathbf{M}}_2 - \mathbf{M}_2$ , we have

$$\begin{aligned}
& \widehat{\mathbf{M}}_2 - \mathbf{M}_2 \\
& = \frac{1}{N} \sum_{n=1}^N y_n (\tilde{\mathbf{x}}_n \otimes \tilde{\mathbf{x}}_n - \mathbb{E} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T) - \mathbb{E}_{\mathbf{x}} y(\tilde{\mathbf{x}} \otimes \tilde{\mathbf{x}} - \mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T) \\
& = \frac{1}{N} \sum_{n=1}^N \left( \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j} + \xi_n) (\tilde{\mathbf{x}}_n \otimes \tilde{\mathbf{x}}_n - \mathbb{E} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T) \right. \\
& \quad \left. - \mathbb{E}_{\mathbf{x}} \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{\tilde{\Omega}_j}) (\tilde{\mathbf{x}} \otimes \tilde{\mathbf{x}} - \mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T) \right) \\
& = \frac{1}{K \cdot N} \sum_{n=1}^N \sum_{j=1}^K \left( \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) (\tilde{\mathbf{x}}_n \otimes \tilde{\mathbf{x}}_n - \mathbb{E} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T) - \mathbb{E}_{\mathbf{x}} \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{\tilde{\Omega}_j}) (\tilde{\mathbf{x}} \otimes \tilde{\mathbf{x}} - \mathbb{E} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T) \right) \\
& \quad + \frac{1}{N} \sum_{n=1}^N \xi_n (\tilde{\mathbf{x}}_n \otimes \tilde{\mathbf{x}}_n - \mathbb{E} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T)
\end{aligned} \tag{130}$$

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

$$\mathbf{B}_2(\mathbf{x}_n) := \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) (\tilde{\mathbf{x}}_n \otimes \tilde{\mathbf{x}}_n - \mathbb{E} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T). \tag{131}$$

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

$$\begin{aligned}
\|\mathbf{B}_2(\mathbf{x})\|_2 & \lesssim \sigma_1 r_{\max} \log^{\frac{3}{2}} q, \\
\|\mathbb{E}_{\mathbf{x}} \mathbf{B}_2(\mathbf{x})\|_2 & \lesssim \sigma_1, \\
\|\mathbb{E}_{\mathbf{x}} \mathbf{B}_2^2(\mathbf{x})\|_2 & \lesssim \frac{1}{K} \sigma_1^2 r_{\max}
\end{aligned} \tag{132}$$

hold with probability at least  $1 - q^{-r_{\max}}$ .

Define  $\mathbf{Z}_{2,n} = \frac{1}{N} (\mathbf{B}_2(\mathbf{x}_n) - \mathbb{E}_{\mathbf{x}} \mathbf{B}_2(\mathbf{x}))$  for  $\mathbf{x}_n$  with  $n \in [N]$ , and it is obvious  $\mathbf{Z}_{2,n}$  is zero mean. Also, we have

$$R_2 = \|\mathbf{Z}_{2,n}\|_2 \leq \frac{1}{N} (\|\mathbf{B}_2(\mathbf{x}_n)\|_2 + \|\mathbb{E}_{\mathbf{x}} \mathbf{B}_2(\mathbf{x})\|_2) \lesssim \frac{1}{K} \sum_{j=1}^K N^{-1} \sigma_1 r_k \log^{\frac{3}{2}} q, \tag{133}$$

and

$$\begin{aligned}
\delta_2^2 &= \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{2,n}^2 \right\|_2^2 \leq \left\| \sum_{n=1}^N \frac{1}{N^2} \left( \mathbb{E} \mathbf{B}_2^2(\mathbf{x}_n) - (\mathbb{E} \mathbf{B}_2(\mathbf{x}_n))^2 \right) \right\|_2 \\
&\leq \frac{1}{N} \left( \|\mathbb{E} \mathbf{B}_2^2(\mathbf{x}_n)\|_2 + \|\mathbb{E} \mathbf{B}_2(\mathbf{x}_n)\|_2^2 \right) \\
&\lesssim N^{-1} \sigma_1^2 r_{\max}.
\end{aligned} \tag{134}$$

Next, let  $t = \Theta(\sigma_1 \sqrt{\frac{r_{\max} \log q}{N}})$ . To make sure  $\delta_2^2 \geq R_2 t/3$ , we need  $N \gtrsim r_{\max} \log^4 q$ . Then, by Lemma 10, we have

$$\text{Prob} \left\{ \left\| \sum_n \mathbf{Z}_{2,n} \right\|_2 \geq t \right\} \leq 2r \exp \left( \frac{-t^2/2}{\delta_2^2 + R_2 t/3} \right) \leq 2r \exp \left( \frac{-t^2}{4\delta_2^2} \right). \tag{135}$$

That is

$$\left\| \sum_{n=1}^N \mathbf{Z}_{2,n} \right\|_2 \lesssim \sigma_1 \sqrt{\frac{r_{\max} \log q}{N}} \tag{136}$$

with probability at least  $1 - q^{-r_{\max}}$ . Because  $\tilde{\mathbf{x}}_n$  belongs to the sub-Gaussian distribution, we know that

$$\left\| \frac{1}{N} \sum_{n=1}^N (\tilde{\mathbf{x}}_n \otimes \tilde{\mathbf{x}}_n - \mathbb{E} \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T) \right\|_2 \lesssim \sqrt{\frac{r_{\max} \log q}{N}} \tag{137}$$

with probability at least  $1 - q^{-r_{\max}}$ .

In conclusion, we have

$$\|\widehat{\mathbf{M}}_2 - \mathbf{M}_2\| \frac{1}{K} \sum_{k=1}^K \lesssim (\sigma_1 + |\xi|) \sqrt{\frac{r_{\max} \log q}{N}} \tag{138}$$

with probability at least  $1 - q^{-r_{\max}}$  provided that  $N \gtrsim r_{\max} \log^4 q$ .  $\square$

## J.2 Error bound for the third-order moment

*Proof of Lemma 8.* For  $\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}})$ , we have

$$\begin{aligned}
&\widehat{\mathbf{M}}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) - \mathbf{M}_3(\widehat{\mathbf{V}}, \widehat{\mathbf{V}}, \widehat{\mathbf{V}}) \\
&= \frac{1}{N} \sum_{n=1}^N y_n [(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n) \otimes (\mathbb{E}(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)^T)] \\
&\quad - \mathbb{E}_{\mathbf{x}} y [(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^{\otimes 3} - (\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}) \otimes \mathbb{E}(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^T] \\
&= \frac{1}{N} \sum_{n=1}^N \left( \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) + \xi_n \right) \cdot [(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n) \otimes (\mathbb{E}(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)^T)] \\
&\quad - \mathbb{E}_{\mathbf{x}} \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{\tilde{\Omega}_j}) [(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^{\otimes 3} - (\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}) \otimes (\mathbb{E}(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^T)] \\
&= \frac{1}{K \cdot N} \sum_{n=1}^N \sum_{j=1}^K \left[ \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) \cdot [(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)^{\otimes 3} - (\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n) \otimes (\mathbb{E}(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}_n)^T)] \right. \\
&\quad \left. - \mathbb{E}_{\mathbf{x}} \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{\tilde{\Omega}_j}) [(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^{\otimes 3} - (\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}) \otimes (\mathbb{E}(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^T)] \right] \\
&\quad + \frac{1}{N} \sum_{n=1}^N \xi_n [(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^{\otimes 3} - (\widehat{\mathbf{V}}^T \tilde{\mathbf{x}}) \otimes (\mathbb{E}(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})(\widehat{\mathbf{V}}^T \tilde{\mathbf{x}})^T)]
\end{aligned} \tag{139}$$

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

$$\mathbf{T}(\mathbf{x}) := \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) \cdot [(\hat{\mathbf{V}}^T \mathbf{x}_n)^{\otimes 3} - (\hat{\mathbf{V}}^T \mathbf{x}_n) \otimes (\mathbb{E}(\hat{\mathbf{V}}^T \mathbf{x}_n)(\hat{\mathbf{V}}^T \mathbf{x}_n)^T)]. \quad (140)$$

Then,  $\mathbf{B}_3(\mathbf{x}) \in \mathbb{R}^{K \times K^2}$  is defined as flattening the tensor  $\mathbf{T}(\mathbf{x})$  along the first dimension. Hence, we have

$$\begin{aligned} \|\mathbf{B}_3(\mathbf{x})\|_2 &\lesssim \max_j |\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}| \cdot \left( \|\hat{\mathbf{V}}^T \mathbf{x}_n\|_2^3 + 3K \|\hat{\mathbf{V}}^T \mathbf{x}_n\|_2 \right) \\ &\lesssim \sigma_1 K^{\frac{3}{2}} \log^{\frac{5}{2}} q \end{aligned} \quad (141)$$

with probability at least  $1 - q^{-K}$ .

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

$$\begin{aligned} \|\mathbb{E}_{\mathbf{x}} \mathbf{B}_3(\mathbf{x})\|_2 &\lesssim \sigma_1, \\ \max \left\{ \|\mathbb{E}_{\mathbf{x}} [\mathbf{B}_3(\mathbf{x})^T \mathbf{B}_3(\mathbf{x})]\|_2, \|\mathbb{E}_{\mathbf{x}} [\mathbf{B}_3(\mathbf{x})^T \mathbf{B}_3(\mathbf{x})]\|_2 \right\} &\lesssim K^2 \sigma_1^2. \end{aligned} \quad (142)$$

Define  $\mathbf{Z}_{3,n} = \frac{1}{N} (\mathbf{B}_3(\mathbf{x}_n) - \mathbb{E}_{\mathbf{x}} \mathbf{B}_3(\mathbf{x}))$  for  $(\mathbf{x}_n, y_n) \in \mathcal{D}$ , and it is obvious  $\mathbf{Z}_{3,n}$  is zero mean. Also, we have

$$\begin{aligned} R_3 = \|\mathbf{Z}_{3,n}\|_2 &\leq \frac{1}{N} (\|\mathbf{B}_3(\mathbf{x}_n)\|_2 + \|\mathbb{E}_{\mathbf{x}} \mathbf{B}_3(\mathbf{x})\|_2) \\ &\lesssim N^{-1} \sigma_1 K^{\frac{3}{2}} \log^{\frac{5}{2}} q, \end{aligned} \quad (143)$$

and

$$\begin{aligned} \delta_3^2 &= \left\{ \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{3,n} \mathbf{Z}_{3,n}^T \right\|_2, \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{3,n} \mathbf{Z}_{3,n}^T \right\|_2 \right\} \leq \frac{1}{N} (\|\mathbb{E} \mathbf{B}_3^2(\mathbf{x}_n)\|_2 + \|\mathbb{E} \mathbf{B}_3(\mathbf{x}_n)\|_2^2) \\ &\lesssim N^{-1} K^2 \sigma_1^2. \end{aligned} \quad (144)$$

Similar to (135), by applying Lemma 10, we have

$$\left\| \sum_{n=1}^N \mathbf{Z}_{3,n} \right\|_2 \lesssim \sigma_1 \sqrt{\frac{\log q}{N}} \quad (145)$$

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

Similar to (141), we define  $\mathbf{B}$  by flattening the tensor  $\sum_{n=1}^N [(\hat{\mathbf{V}}^T \tilde{\mathbf{x}})^{\otimes 3} - (\hat{\mathbf{V}}^T \tilde{\mathbf{x}}) \otimes (\mathbb{E}(\hat{\mathbf{V}}^T \tilde{\mathbf{x}})(\hat{\mathbf{V}}^T \tilde{\mathbf{x}})^T)]$  along the first dimension. Then, we know that

$$\begin{aligned} \|\mathbf{B}\|_2 &\leq \left\| \sum_{n=1}^N \hat{\mathbf{V}}^T \tilde{\mathbf{x}}_n \right\|_2^3 + 3K \left\| \sum_{n=1}^N \hat{\mathbf{V}}^T \tilde{\mathbf{x}}_n \right\|_2 \lesssim \left( \frac{K^{-4} \log q}{N} \right)^{\frac{3}{2}} + 3K \left( \frac{K^{-4} \log q}{N} \right)^{\frac{1}{2}} \\ &\lesssim \left( \frac{\log q}{N} \right)^{\frac{1}{2}} + \left( \frac{\log q}{N} \right)^{\frac{1}{2}} \\ &\lesssim \sqrt{\frac{\log q}{N}}, \end{aligned} \quad (146)$$

provided that  $N \gtrsim K^5 \log q$ .

In conclusion, we have

$$\left\| \widehat{\mathbf{M}}_3(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}}) - \mathbf{M}_3(\hat{\mathbf{V}}, \hat{\mathbf{V}}, \hat{\mathbf{V}}) \right\| \lesssim (\sigma_1 + |\xi|) \sqrt{\frac{\log q}{N}} \quad (147)$$

with probability at least  $1 - q^{-K}$  provided that  $N \gtrsim K^3 \log^6 q$ .

□

### J.3 Error bound for the first-order moment

*Proof of Lemma 9.* For  $\widehat{\mathbf{M}}_1 - \mathbf{M}_1$ , we have

$$\begin{aligned}
\widehat{\mathbf{M}}_1 - \mathbf{M}_1 &= \frac{1}{N} \sum_{n=1}^N y_n \tilde{\mathbf{x}}_n - \mathbb{E}_{\mathbf{x}} y \tilde{\mathbf{x}} \\
&= \frac{1}{N} \sum_{n=1}^N \left( \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) + \xi_n \right) \tilde{\mathbf{x}}_n - \mathbb{E}_{\mathbf{x}} \sum_{j=1}^K \frac{1}{K} \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{\tilde{\Omega}_j}) \tilde{\mathbf{x}} \\
&= \frac{1}{K \cdot N} \sum_{j=1}^K \sum_{n=1}^N \left( \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) \tilde{\mathbf{x}}_n - \mathbb{E}_{\mathbf{x}} \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{\tilde{\Omega}_j}) \tilde{\mathbf{x}} \right) + \frac{1}{N} \sum_{n=1}^N \xi_n \cdot \tilde{\mathbf{x}}_n.
\end{aligned} \tag{148}$$

Define  $\mathbf{B}_1(\mathbf{x}) := \frac{1}{K} \sum_{j=1}^K \phi(\mathbf{u}_j^{*T} \tilde{\mathbf{x}}_{n, \tilde{\Omega}_j}) \tilde{\mathbf{x}}_n$ , then we have

$$\begin{aligned}
\|\mathbf{B}_1(\mathbf{x})\|_2 &\lesssim \frac{1}{K} \sum_{k=1}^K \sigma_1 r_k \log^{\frac{3}{2}} q; \\
\|\mathbb{E}_{\mathbf{x}} \mathbf{B}_1(\mathbf{x})\|_2 &\lesssim \sigma_1; \\
\left\{ \|\mathbf{E}_{\mathbf{x}}[\mathbf{B}_1(\mathbf{x}) \mathbf{B}_1(\mathbf{x})^T]\|_2, \|\mathbf{E}_{\mathbf{x}}[\mathbf{B}_{1,j}(\mathbf{x})^T \mathbf{B}_1(\mathbf{x})]\|_2 \right\} &\lesssim \sigma_1^2.
\end{aligned} \tag{149}$$

Next, define  $\mathbf{Z}_{1,n} = \frac{1}{N} (\mathbf{B}_{1,j}(\mathbf{x}_n) - \mathbb{E}_{\mathbf{x}} \mathbf{B}_2(\mathbf{x}))$  for  $(\mathbf{x}_n, y_n) \in \mathcal{D}$ , by calculation, we can obtain

$$R_1 = \|\mathbf{Z}_{1,n}\|_2 \lesssim N^{-1} \sigma_1 r_{\max} \log^{\frac{3}{2}} q, \tag{150}$$

and

$$\delta_1^2 = \max \left\{ \left\| \sum_{n=1}^N \mathbb{E} \mathbf{Z}_{1,n} \mathbf{Z}_{1,n}^T \right\|_2^2, \left| \sum_{n=1}^N \mathbf{Z}_{1,n}^T \mathbf{Z}_{1,n} \right| \right\} \lesssim N^{-1} \sigma_1^2 r_{\max}. \tag{151}$$

By applying Lemma 10, we have

$$\left\| \sum_{n=1}^N \mathbf{Z}_{1,n} \right\|_2 \lesssim \sigma_1 \sqrt{\frac{r_{\max} \log q}{N}} \tag{152}$$

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

$$\left\| \frac{1}{N} \sum_{n=1}^N \tilde{\mathbf{x}}_n \right\|_2 \lesssim \sqrt{\frac{r_{\max} \log q}{N}} \tag{153}$$

with probability at least  $1 - q^{-r_{\max}}$ .

In conclusion, we have

$$\|\widehat{\mathbf{M}}_1 - \mathbf{M}_1\| \lesssim (\sigma_1 + |\xi|) \sqrt{\frac{r_{\max} \log q}{N}} \tag{154}$$

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