

REVISED NTK ANALYSIS OF OPTIMIZATION AND GENERALIZATION WITH ITS EXTENSIONS TO ARBITRARY INITIALIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent theoretical works based on the neural tangent kernel (NTK) have shed light on the optimization and generalization of over-parameterized neural networks, and partially bridge the gap between their practical success and classical learning theory. However, the existing NTK-based analysis has a limitation that the scaling of the initial parameter should decrease with respect to the sample size which is contradictory to the practical initialization scheme. To address this issue, in this paper, we present the revised NTK analysis of optimization and generalization of overparametrized neural networks, which successfully remove the dependency on the sample size of the initialization. Based on our revised analysis, we further extend our theory that allow for arbitrary initialization, not limited to Gaussian initialization. Under our initialization-independent analysis, we propose NTK-based regularizer that can improve the model generalization, thereby illustrating the potential to bridge the theory and practice while also supporting our theory. Our numerical simulations demonstrate that the revised theory indeed can achieve the significantly lower generalization error bound compared to existing error bound. Also importantly, the proposed regularizer also corroborate our theory on the arbitrary initialization with fine-tuning scenario, which takes the first step for NTK theory to be promisingly applied to real-world applications.

1 INTRODUCTION

Though neural networks (NNs) have achieved great success in practice, it remains a well-known mystery that over-parameterized NNs generalize well and do not suffer from overfitting even with a simple first-order optimization (Neyshabur et al., 2014; Livni et al., 2014; Zhang et al., 2016; Arora et al., 2018), seemingly contradicting the traditional learning theory. To theoretically explain this phenomenon, extensive research has been conducted, and one of the main directions is based on the neural tangent kernel (NTK). Given a NN $f_{\theta}(\cdot)$ parametrized by θ and n training inputs $\{\mathbf{x}_i\}_{i=1}^n$, the NTK is defined as a Gram matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$ induced by the structure of target prediction function whose (i, j) -th entry is given by

$$\mathbf{H}_{ij} := \left\langle \frac{\partial f_{\theta}(\mathbf{x}_i)}{\partial \theta}, \frac{\partial f_{\theta}(\mathbf{x}_j)}{\partial \theta} \right\rangle. \quad (1)$$

The NTK was introduced in Jacot et al. (2018) to control the dynamics of learning NNs. In the over-parameterized regime, the trained parameter of NN is close to its initialization, which also makes the NTK almost unchanged throughout the training process. This stability of NTK allows the learning dynamics of NN to be easily analyzed throughout the training process, thus making it possible to derive training and generalization error bounds by using existing learning theory.

As a representative study using NTK, Arora et al. (2019a;b) showed the following important results for over-parameterized NNs:

- (a) (Training error bound) A training error bound reflecting a tighter characterization of training speed than that of studies was proposed in Arora et al. (2019a). This bound implies that not only is a network able to represent any finite sample perfectly (as shown in Zhang et al. (2016)), but

the speed at which a network learns training samples varies depending on a complexity measure reflecting how well the data is ordered.

- (b) (Generalization error bound) A generalization error bound, referred to as *complexity measure of data* (CMD) was proposed in Arora et al. (2019a). CMD has no stringent conditions on certain properties of the trained NN and the true model; it only depends on input \mathbf{x} and label y of training data and the initial parameter scale κ of NN, hence we can compute the bound *before* actually training the network. CMD is considered as one of the most important achievements in the topic of generalizability of over-parameterized NNs and has been the cornerstone of many follow-up studies in recent years (Arora et al., 2019b; Su & Yang, 2019; Oymak et al., 2019; Xu et al., 2019; Zhang et al., 2019; Hu et al., 2019; Du et al., 2018).

The above results (a)~(b) derived in Arora et al. (2019a;b) provide the upper bounds on the training/generalization error that are *uniformly* available over all network scaling (e.g., initialization) parameter κ . Note that Arora et al. (2019a;b) focus only on the case where

$$\kappa = o(1) \text{ with respect to } n \quad (2)$$

in order to have meaningful bounds that can converge to zero as n increases.

Surprisingly, however, we prove in this paper that, contrary to the theories of Arora et al. (2019a;b), the training/generalization errors in (a)-(b) *do not* hold when κ decreases with respect to n as in Eq. 2. The high-level reason for this is as follows. As trained parameter is known to be close to its initial one in the over-parameterized regime (Jacot et al., 2018; Arora et al., 2019a;b; Du et al., 2018; 2019; Ji & Telgarsky, 2019; Cao & Gu, 2019; Ma et al., 2019; Li & Liang, 2018; Hu et al., 2019; Oymak et al., 2019), the output value of trained NN $|f_{\theta}(\mathbf{x})|$ is close to that of its initial NN $|f_{\theta(0)}(\mathbf{x})|$. Meanwhile, the output value of its initial NN decreases to zero as n increases if the scale κ of initialization decreases with respect to n . Thus, it cannot guarantee zero training error under the condition Eq. 2, as the output value of trained NN decreases to zero but the target label does not.

We further resolve the above issue and revise the analyses of Arora et al. (2019a;b) without major modifications of the original statements. Hence, our revision makes it possible for results (a)–(b) to maintain their original meanings and implications without any issue on decreasing κ . Our revised analyses provide tighter results on training/generalization error bounds, and with these improvements, we can guarantee the bounds to converge to zero even when κ is a constant w.r.t. n .

Building on the proof technique in the revised theory, we further extend our analysis based on the Gram matrix and network initial parameters drawn from a Gaussian distribution to allow for arbitrary initialization. Grounded on our initialization-independent analysis, we propose the NTK regularizer that can boost the model generalization, which is unavailable in previous studies. Note that our initialization-independent analysis enables NTK theory to be applied to pre-trained networks, which is expected to provide the connections to various practical scenarios such as fine-tuning. Toward this, we empirically verify that our revised theory indeed achieves lower generalization error bounds than the baseline theory and present the potential of the proposed NTK regularizer in fine-tuning scenario.

Notations. The sets $\{1, 2, \dots, i\}$ and $\{i, i + 1, \dots, j\}$ are denoted by $\{i\}$ and $\{i : j\}$, respectively. The Frobenius norm is denoted by $\|\cdot\|$. For a matrix \mathbf{H} , $\lambda_{\min}(\mathbf{H})$ denotes the smallest eigenvalue of \mathbf{H} . Training samples are given as n input-label pairs $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ generated from a data distribution $\mathcal{D}(\mathbf{x}, y)$, i.i.d. For simplicity, we assume that $\|\mathbf{x}\| = 1$ for $\mathbf{x} \sim \mathcal{D}$. We denote inputs and labels in the training dataset by $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$ and $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$, respectively.

2 NTK-BASED ANALYSIS FOR TRAINING/GENERALIZATION ERROR BOUNDS

We first review the training and test error bounds of Arora et al. (2019a) in Section 2.1, and disprove and revise them in Sections 2.2 and 2.3, respectively.

2.1 PRELIMINARY: TRAINING/GENERALIZATION ERROR BOUNDS OF ARORA ET AL. (2019A)

Consider a two-layer ReLU network $f_{\mathbf{W}, \mathbf{a}}(\mathbf{x})$ with scalar outputs as in Arora et al. (2019a):

$$f_{\mathbf{W}, \mathbf{a}}(\mathbf{x}) := \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\mathbf{w}_r^\top \mathbf{x}). \quad (3)$$

Here $\mathbf{x} \in \mathbb{R}^d$ is a given input datapoint, $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m) \in \mathbb{R}^{d \times m}$ is the weight parameter in the first layer, $\mathbf{a} = (a_1, \dots, a_m)^\top \in \mathbb{R}^m$ is the weight parameter in the second layer, and $\sigma(\cdot)$ is the ReLU activation. The setting indicates that there are m hidden neurons.

Using n samples (\mathbf{X}, \mathbf{y}) , we train the neural network Eq. 3 so that its prediction function $f_{\mathbf{W}, \mathbf{a}}(\cdot)$ minimizes the following squared error

$$L(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (y_i - f_{\mathbf{W}, \mathbf{a}}(\mathbf{x}_i))^2 \quad (4)$$

by updating the network parameter \mathbf{W} via the discrete time optimization of gradient descent (GD) as

$$\mathbf{W}(k+1) := \mathbf{W}(k) - \eta \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} \Big|_{\mathbf{W}=\mathbf{W}(k)}.$$

We denote by $\mathbf{u}(k) = (u_1(k), \dots, u_n(k))^\top = (f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x}_1), \dots, f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x}_n))^\top \in \mathbb{R}^n$ the network output with trained parameter $\mathbf{W}(k)$ at k -step. The parameter \mathbf{W} is assumed to be randomly initialized as $\mathbf{w}_r \sim \mathcal{N}(0, \kappa^2 \mathbf{I}_d)$ using standard deviation κ for $r \in \{m\}$ as in Arora et al. (2019a). Each element of \mathbf{a} is independently initialized (and fixed) as following $\mathcal{U}(\{-1, 1\})$.

By setting the network $f_{\theta}(\mathbf{x})$ and its parameter θ of NTK in Eq. 1 as $f_{\mathbf{W}, \mathbf{a}}(\mathbf{x})$ and $\mathbf{W}(0)$, respectively, Arora et al. (2019a) derived a specific NTK (with $m = \infty$) as Gram matrix $\mathbf{H}^\infty \in \mathbb{R}^{n \times n}$ as follows: given data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ of n input training samples, (i, j) -th entry of \mathbf{H}^∞ is given by

$$\mathbf{H}_{ij}^\infty := \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_d)} [\mathbf{x}_i^\top \mathbf{x}_j \mathbb{I}\{\mathbf{w}^\top \mathbf{x}_i \geq 0, \mathbf{w}^\top \mathbf{x}_j \geq 0\}] = \frac{\mathbf{x}_i^\top \mathbf{x}_j (\pi - \arccos(\mathbf{x}_i^\top \mathbf{x}_j))}{2\pi}$$

where $\mathbb{I}\{\cdot\}$ is the indicator function. We use λ_0 to denote $\lambda_{\min}(\mathbf{H}^\infty)$. Then, all NTKs obtained by updated parameters $\{\mathbf{W}(k)\}_{k=0}^\infty$ are close to \mathbf{H}^∞ in the over-parameterized regime. Using this fact and extending Du et al. (2018) to hold for arbitrary κ , Arora et al. (2019a) provided the following theorem, which guarantees zero training error with a convergence rate depending on λ_0 .

Theorem 2.1 (Theorem 3.1 in Arora et al. (2019a)). *Fix a failure probability $\delta \in (0, 1)$. Suppose that $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^2 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, $\lambda_0 > 0$, and $\eta = O(\frac{\lambda_0}{n^2})$. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for any κ and all $k \geq 0$,*

$$\|\mathbf{y} - \mathbf{u}(k+1)\|^2 \leq \left(1 - \frac{\eta \lambda_0}{2}\right) \|\mathbf{y} - \mathbf{u}(k)\|^2. \quad (5)$$

As a corollary of Theorem 2.1, Arora et al. (2019a) showed a new training error bound reflecting a tighter characterization of training speed such that its convergence rate is mainly affected by the training data belonging to the top eigenspaces of \mathbf{H}^∞ . This bound is given as follows.

Corollary 2.1 (The training error bound, Theorem 4.1 in Arora et al. (2019a)). *Suppose all conditions in Theorem 2.1 hold. Then, with probability at least $1 - \delta$ for $\delta \in (0, 1)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, for all $k \geq 0$,*

$$\frac{1}{\sqrt{n}} \|\mathbf{y} - \mathbf{u}(k)\| = \sqrt{\frac{1}{n} \sum_{i=1}^n (1 - \eta \lambda_i)^{2k} (\mathbf{v}_i^\top \mathbf{y})^2} \pm O\left(\frac{\kappa}{\delta} + \frac{n^3}{\sqrt{m} \lambda_0^2 \kappa \delta^2}\right), \quad (6)$$

where $\{\mathbf{v}_i\}_{i=1}^n$ are orthonormal eigenvectors of \mathbf{H}^∞ and $\{\lambda_i\}_{i=1}^n$ are the corresponding eigenvalues.

This bound, given as the right-hand side in Eq. 6, reflects the convergence rate in more details by using all the spectral information of \mathbf{H}^∞ (i.e., $\{\lambda_i\}_{i=1}^n$), but the training error bound in Eq. 5 reflects only the least influential part (i.e., $\lambda_0 = \lambda_n$) among these information. This improvement over Theorem 2.1 allows to demonstrate that true labels yield faster learning speeds than random labels (Arora et al., 2019a; Zhang et al., 2016). Meanwhile, for this bound in Eq. 6 to converge to zero, its second term must also decrease to zero and the corresponding condition is given as follows.

Remark. For the error term $\frac{\kappa}{\delta}$ in Eq. 6 to decrease to 0 w.r.t. n , it should hold that $\kappa = o(1)$ w.r.t. n .

Using Theorem 2.1, Arora et al. (2019a) also derived the following generalization error bound, named complexity measure of data (CMD).

Theorem 2.2 (The generalization error bound, Theorem 5.1 in Arora et al. (2019a)). Suppose that all conditions except $\lambda_0 > 0$ in Theorem 2.1 hold and we fix a failure probability $\delta \in (0, 1)$. Suppose also that $m = \tilde{\Omega}(\kappa^{-2} \text{poly}(n, \lambda_0^{-1}, \delta^{-1}))$. Suppose further that $\lambda_0 > 0$ holds with probability at least $1 - \delta/3$ for n i.i.d. training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ from true model distribution \mathcal{D} . Consider any loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$ that is 1-Lipschitz in the first argument. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ and the training samples, the neural network $f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x})$ trained by GD for $k \geq \Omega(\frac{1}{\eta \lambda_0} \log \frac{n}{\delta})$ iterations has population loss² $L_{\mathcal{D}}(f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x})) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell(f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x}), y)]$ bounded as

$$L_{\mathcal{D}}(f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x})) \leq \underbrace{\sqrt{\frac{2\mathbf{y}^\top (\mathbf{H}^\infty)^{-1} \mathbf{y}}{n}}}_{\text{CMD}} + \underbrace{O\left(\frac{\sqrt{n}\kappa}{\lambda_0 \delta}\right)}_{\text{Error term } \mathcal{E}} + O\left(\sqrt{\frac{\log \frac{n}{\lambda_0 \delta}}{n}}\right). \quad (7)$$

The CMD bound, given in the right-hand side of Eq. 7, only depends on the training samples (e.g., $\mathbf{y}, \mathbf{H}^\infty, \lambda_0$) and κ . This makes it possible to know whether a NN can generalize without actually training the NN, as mentioned above.

Remark. For the error term $\mathcal{E} = \frac{\sqrt{n}\kappa}{\lambda_0 \delta}$ in Eq. 7 to decrease to 0 with respect to the sample size n , it should hold that $\kappa = o\left(\frac{\lambda_0}{\sqrt{n}}\right)$.

2.2 DISPROOF OF EXISTING NTK-BASED TRAINING/GENERALIZATION ERROR BOUNDS

From the remarks above, κ should follow $o(1)$ and $o(\lambda_0/\sqrt{n})$ for the training and test errors in Eq. 6 and Eq. 7 to approach zero, respectively. In fact, we have shown in Figure 1 that λ_0 does not increase with n in standard benchmark datasets, thus $o(\lambda_0/\sqrt{n})$ implies $o(1)$. Hence, κ should follow $o(1)$ for both training and generalization errors in Eq. 6 and Eq. 7 to approach zero. These conditions on κ can be allowed only if the original Theorem 2.1 is valid for such κ , as Theorem 2.1 says.

However, in this section, we show that Theorem 2.1 actually does not hold under these conditions on κ (i.e., decreasing κ). Toward this, we visit the case where λ_0 satisfies the following mild condition

$$\lambda_0 = \mathcal{O}(n^\gamma) > 0 \text{ for some constant } \gamma \leq 1. \quad (8)$$

Under Eq. 8, we claim that an additional condition for κ (i.e., non-decreasing κ) is needed for the statements in Theorem 2.1 to hold.

Theorem 2.3. Suppose the condition Eq. 8 holds for a constant $\gamma \leq 1$ and $m = \Omega(n^{3-2\gamma})$. Suppose further $\kappa = o(1)$ for n . Then, for any η satisfying $0 < \lambda_0 \eta < 2$, there exists a finite integer k (and n) with probability at least $1 - \delta$ for $\delta \in (0, 1)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ such that

$$\|\mathbf{y} - \mathbf{u}(k+1)\|^2 > \left(1 - \frac{\eta \lambda_0}{2}\right) \|\mathbf{y} - \mathbf{u}(k)\|^2. \quad (9)$$

It can be clearly seen that Eq. 5 and Eq. 9 are contradictory and hence the following corollaries hold.

Corollary 2.2. Theorem 2.1 does not hold if the condition $\kappa = o(1)$ w.r.t. n and Eq. 8 hold.

Corollary 2.3. Corollary 2.1 fails to guarantee that NNs attain zero training loss if Eq. 8 holds.

Corollary 2.4. Theorem 2.2 fails to guarantee that NNs attain zero gen. err. if $\lambda_0 = O(\sqrt{n}) > 0$.

The question that naturally arises at this point is how easily the condition Eq. 8 is satisfied in practice. In addition to the observation that λ_0 does not increase with n in practice as shown in Figure 1, we also find a simple sufficient condition for Eq. 8 to hold, provided in the following proposition:

Proposition 2.1. Suppose that n input samples are not parallel, i.e., $\mathbf{x}_i \neq c\mathbf{x}_j$ for any $c \in \mathbb{R}$ and different $i, j \in \{n\}^2$. Then, $\lambda_0 = O(\sqrt{n}) > 0$ holds.

Proposition 2.1 confirms that the condition $\lambda_0 = O(\sqrt{n}) > 0$ for Corollary 2.4 holds (i.e., Theorem 2.2 fails) easily in the practically common case where the training data is not parallel.

²Arora et al. (2019a) claimed that Eq. 7 holds for a general loss, but in fact they implicitly assumed squared loss and did not consider a general loss in the proof.

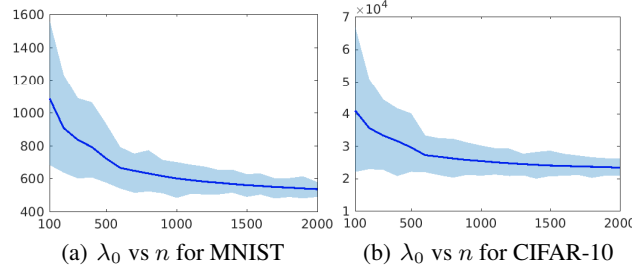


Figure 1: (a) and (b) show the value of λ_0 w.r.t. sample size n for the standard benchmark image datasets, MNIST and CIFAR-10, respectively. This result shows $\lambda_0 = O(1)$ holds easily in practice.

2.3 REVISING NTK-BASED TRAINING/GENERALIZATION ERROR BOUNDS

By Theorem 2.3, κ should not decrease w.r.t. n (i.e., $\kappa = o(1)$) in order for the statements in Theorem 2.1 to hold. Thus, we derive tighter bounds so that we avoid the case of setting decreasing κ (i.e., $\kappa = \Theta(1)$). In fact, Du et al. (2018) already showed that this is possible for Theorem 2.1 with $\kappa = \Theta(1)$. Here, we revise training and generalization bounds in Corollaries 2.1 and 2.2.

Theorem 2.4 (Revision of Corollary 2.1). *Suppose all conditions in Theorem 2.1 hold and $\kappa = \Theta(1)$. Then, with probability at least $1 - \delta$ for $\delta \in (0, 1)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for all $k \geq 0$,*

$$\frac{1}{\sqrt{n}} \|\mathbf{y} - \mathbf{u}(k)\| = \sqrt{\frac{1}{n} \sum_{i=1}^n (1 - \eta \lambda_i)^{2k} \left(\mathbf{v}_i^\top (\mathbf{y} - \mathbf{u}(0)) \right)^2} + O\left(\frac{n^3}{\sqrt{m} \lambda_0^2 \delta^2}\right), \quad (10)$$

where $\{\mathbf{v}_i\}_{i=1}^n$ are orthonormal eigenvectors of \mathbf{H}^∞ and $\{\lambda_i\}_{i=1}^n$ are the corresponding eigenvalues.

Compared to Corollary 2.1, the training error bound in Theorem 2.4 does not have the term κ/δ . Accordingly, this tighter bound can converge to zero as the iteration number k increases even in the case for $\kappa = \Theta(1)$. This is formally stated in the following:

Proposition 2.2. *Suppose all conditions in Theorem 2.1 hold, $\kappa = \Theta(1)$, and $m = \Omega\left(\frac{n^{6+\alpha}}{\lambda_0^4}\right)$ with any $\alpha > 0$. Then, with probability at least $1 - \delta$ for $\delta \in (0, 1)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, the right hand side of Eq. 10 converges to zero as k and n increase.*

We also revise the CMD bound in Theorem 2.2 as follows, under the condition $\kappa = \Theta(1)$.

Theorem 2.5 (Revision of Theorem 2.2). *Suppose that all conditions except $\lambda_0 > 0$ in Theorem 2.1 hold and we fix a failure probability $\delta \in (0, 1)$. Suppose also that $\lambda_0 > 0$ holds with probability at least $1 - \delta/3$ for n i.i.d. training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ from data distribution \mathcal{D} , and that $\kappa = \Theta(1)$ and $m = \tilde{\Omega}(\text{poly}(n, \lambda_0^{-1}, \delta^{-1}))$. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ and the training samples, it follows that for any $k \geq \Omega\left(\frac{1}{\eta \lambda_0} \log \frac{n}{\delta}\right)$,*

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \frac{1}{2} |y - f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x})|^2 = \underbrace{\sqrt{\frac{2(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^\infty)^{-1} (\mathbf{y} - \mathbf{u}(0))}{n}}}_{\text{Revised CMD}} + O\left(\sqrt{\frac{\log \frac{n}{\lambda_0 \delta}}{n}}\right). \quad (11)$$

Compared to the original CMD bound in Eq. 7, the revised version Eq. 11 does not have the error term in Eq. 7, $(\sqrt{n\kappa})/(\lambda_0 \delta)$, which is the culprit for generalization error bound to blow up. By applying Corollary 6.2 in Arora et al. (2019a) to our setting, we can also bound the first term in Eq. 11 even with the introduction of $\mathbf{u}(0)$ exactly as in the original CMD bound:

Proposition 2.3. *Suppose $y_i - u_i(0) = g(\mathbf{x}_i) := \sum_j \alpha_j (\beta_j^\top \mathbf{x}_i)^{p_j}$ for all $i \in \{n\}$, where for each j , $p_j \in \{1, 2, 4, 6, \dots\}$ and $\alpha_j \in \mathbb{R}$ and $\beta_j \in \mathbb{R}^d$ are any constants w.r.t. n . Then,*

$$\sqrt{\frac{2(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^\infty)^{-1} (\mathbf{y} - \mathbf{u}(0))}{n}} \leq \frac{6 \sum_j p_j |\alpha_j| \|\beta_j\|^{p_j}}{\sqrt{n}} = O\left(\frac{1}{\sqrt{n}}\right). \quad (12)$$

In that sense, the revised bound directly improves the CMD (in addition to fixing it) by only removing its error term without sacrificing any additional assumption. Also, in Section 4, we will demonstrate that our proposed generalization error bound based on revised CMD could be much smaller than the original bound by showing that revised CMD does not significantly differ in value from the existing CMD for benchmark datasets.

3 TOWARDS GENERALIZED ANALYSIS ALLOWING ARBITRARY INITIALIZATION

Although we present the revised NTK analysis in Theorem 2.5, our theory depends on the Gram matrix \mathbf{H}^∞ and the initial network parameter $\mathbf{W}(0)$ (due to the initial network output $\mathbf{u}(0)$), both of which depend on the Gaussian distribution. In this section, we extend our theory that can allow for arbitrary initialization. Based on our extended initialization-independent analysis, we introduce promising applications that can bridge our theory to practice.

3.1 PROBLEM SETUP

Similar to our revised NTK theory in Section 2, we could consider 2-layer ReLU networks, but for our theory and the application in the subsequent section (Section 3.3), we consider more generalized settings. Toward this, we consider a multi-layer neural network $f(\cdot)$ with an arbitrary depth, consisting of two components: a feature extractor $\phi_{\mathbf{V}}(\cdot)$ and a linear classifier $f_{\mathbf{W}}(\cdot)$, as follows:

$$h_{\mathbf{V}, \mathbf{W}}(\mathbf{x}) := f_{\mathbf{W}}(\phi_{\mathbf{V}}(\mathbf{z})) = f_{\mathbf{W}}(\mathbf{x}). \quad (13)$$

where $\mathbf{z} \in \mathbb{R}^{d_i}$ is an input datapoint, $\mathbf{x} \in \mathbb{R}^{d_h}$ here denotes the representation encoded by the feature extractor $\phi_{\mathbf{V}}(\cdot) : \mathbb{R}^{d_i} \rightarrow \mathbb{R}^{d_h}$, and $f_{\mathbf{W}}(\cdot) : \mathbb{R}^{d_h} \rightarrow \mathbb{R}^{d_o}$ denotes a (linear) classifier. Each component $\phi(\cdot)$ and $f(\cdot)$ is characterized by the trainable parameters \mathbf{V} and \mathbf{W} respectively. Note that the feature extractor (or also called an encoder) $\phi_{\mathbf{V}}(\cdot)$ can be any neural network that yields d_h -dimensional representations, e.g., the feature vectors after global average pooling layer in popular CNNs such as VGG (Simonyan & Zisserman, 2014), ResNet (He et al., 2016), EfficientNet (Tan, 2019), and many other architectures. Throughout the paper, we assume that the parameter of feature extractor is fixed by $\mathbf{V} = \mathbf{V}^*$.

Allowing theoretical analysis in subsequent sections, we replace the linear classifier $f(\cdot)$ with a two-layer ReLU classifier with the vector-valued outputs, i.e., $f_{\mathbf{W}}(\mathbf{x}) = (f_{\mathbf{W}}(\mathbf{x})[1], \dots, f_{\mathbf{W}}(\mathbf{x})[d_o])^\top \in \mathbb{R}^{d_o}$ whose i -th output in this setting is computed by

$$f_{\mathbf{W}}(\mathbf{x})[i] := \frac{\sqrt{d_o}}{\sqrt{m}} \sum_{r=1}^m \mathbf{a}_i[r] \sigma(\mathbf{w}_r^\top \mathbf{x}), \quad (14)$$

where $\mathbf{x} = \phi_{\mathbf{V}}(\mathbf{z}) \in \mathbb{R}^{d_h}$ is a latent feature, $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_m] \in \mathbb{R}^{d_h \times m}$ is the trainable parameter in the first layer of the replaced classifier, $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_{d_o}] \in \mathbb{R}^{m \times d_o}$ is the (randomly fixed) weight matrix in the second layer, and $\sigma(\cdot)$ is the ReLU activation.

Similar to the network considered in Section 2, to initialize the parameter \mathbf{A} , we directly extend the setting of Arora et al. (2019a) so that only diagonal blocks are randomly initialized as follows: $\mathbf{a}_i[r] \sim \mathcal{U}(\{-1, 1\})$ for $i \in \{d_o\}$ and $r \in \{\bar{m} \cdot i - \bar{m} + 1 : \bar{m} \cdot i\}$, otherwise $\mathbf{a}_i[r] = 0$, where $\bar{m} = m/d_o$ and we assume \bar{m} is an integer throughout the paper for simplicity. In Eq. 14, we use the scaling factor of $\sqrt{d_o/m}$, which is comparable to the scaling of $1/\sqrt{m}$ used in related studies (Bai & Lee, 2019; Nitanda & Suzuki, 2019; Zhang et al., 2019; Du et al., 2018; Arora et al., 2019a; Du et al., 2019). While we have the additional $\sqrt{d_o}$ term, the vector-valued network we consider in Eq. 14 can be divided into d_o scalar-valued networks, and hence the effective number of hidden units for each component is equal to m/d_o . We provide the detailed proof on the equivalence between the vector-valued network and the multiple number of scalar-valued networks in Appendix C.

3.2 GENERALIZATION ERROR WITH ARBITRARY INITIALIZATION

Let $\mathbf{X} := [\mathbf{x}_1 := \phi_{\mathbf{V}^*}(\mathbf{z}_1), \dots, \mathbf{x}_n := \phi_{\mathbf{V}^*}(\mathbf{z}_n)]^\top \in \mathbb{R}^{n \times d_h}$ be a set of d_h -dimensional latent feature vectors obtained by the feature extractor $\phi_{\mathbf{V}^*}(\mathbf{Z})$ with the input dataset $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n]$.

With the transformed dataset \mathbf{X} , we specify the training process by using gradient descent (GD) optimization and assuming the training loss $\mathcal{L}(\cdot)$ as the mean square error (MSE), specified as

$$\mathcal{L}(h_{\mathbf{V}^*, \mathbf{W}}(\mathbf{Z}), \mathbf{Y}) = L(\mathbf{W}) := \frac{1}{2} \|\mathbf{Y} - F(\mathbf{W})\|^2 \quad (15)$$

where $\mathbf{Y} \in \mathbb{R}^{d_o \times n}$ is a multi-class labels and $F(\mathbf{W}) := [f_{\mathbf{W}}(\mathbf{x}_1), \dots, f_{\mathbf{W}}(\mathbf{x}_n)] \in \mathbb{R}^{d_o \times n}$ denotes an another representation of prediction function $h_{\mathbf{V}^*, \mathbf{W}}(\mathbf{z})$. Then, the network parameter \mathbf{W} is assumed to be updated via GD on the loss $L(\mathbf{W})$. For our analysis, we define the Gram matrix for each class $i \in \{1 : d_o\}$ computed on the model parameter $\mathbf{W}(k)$ as

$$[\mathbf{H}_i(k)]_{pq} := [\mathbf{H}_i(\mathbf{W}(k))]_{pq} := \frac{d_o}{m} \mathbf{x}_p^\top \mathbf{x}_q \sum_{r \in \mathcal{M}_i} [\mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_p \geq 0, \mathbf{w}_r(k)^\top \mathbf{x}_q \geq 0\}]. \quad (16)$$

Note that we use λ_0^i to denote $\lambda_{\min}(\mathbf{H}_i(0))$ and λ_0 to denote $\min(\{\lambda_0^i\}_{i=1}^{d_o})$.

Using the above notations, we present mild conditions required for our theorem to be satisfied.

Condition 1 (Variable R decreases fast enough for increasing n). *Given $\mathbf{W}(0)$, there exists a R satisfying the following condition for each $c \in \{1 : d_o\}$*

$$\frac{1}{n\bar{m}} \sum_{p \in \{1:n\}} \sum_{r \in \mathcal{M}_c} \mathbb{1}\{|\mathbf{w}_r(0)^\top \mathbf{x}_p| \leq R\} = \mathcal{O}\left(\frac{\lambda_0}{n^2}\right), \quad (17)$$

where $\bar{m} = m/d_o$ and $\mathcal{M}_c = \{\bar{m} \cdot c - \bar{m} + 1 : \bar{m} \cdot c\}$.

Condition 2 ($\mathbf{W}(0)$ is bounded and m is sufficiently large). *The initial weight $\mathbf{W}(0)$ satisfies the following two conditions*

$$\begin{aligned} \frac{1}{n\bar{m}} \sum_{p \in \{1:n\}} \sum_{r \in \mathcal{M}_i} \mathbb{1}\left\{|\mathbf{w}_r(0)^\top \mathbf{x}_p| = \mathcal{O}\left(\frac{n}{\sqrt{m}\lambda_0}\right)\right\} &= \mathcal{O}\left(\frac{\min(\lambda_0^2, \lambda_0^3)}{n^4}\right), \\ \frac{1}{n\bar{m}} \sum_{p \in \{1:n\}} \sum_{r \in \mathcal{M}_i} |\mathbf{w}_r(0)^\top \mathbf{x}_p|^2 &= \mathcal{O}(1). \end{aligned} \quad (18)$$

Condition 3 (Elements of $h_{\mathbf{V}^*, \mathbf{W}(k)}$ are bounded). *For an input sample $\mathbf{z} \in \mathbb{R}^{d_i}$ obtained from \mathcal{D} and for every $k \geq 0$, it follows that with probability at least $1 - \delta$ over the random configuration of \mathbf{A} and input sample \mathbf{z} , the following holds $|h_{\mathbf{V}^*, \mathbf{W}(k)}(\mathbf{z})[i]| = \mathcal{O}(1)$ for each $i \in \{1 : d_o\}$.*

Note that the condition 3 is easily satisfied as each element of network output has a bounded magnitude invariant of n in practice. Further, the conditions 1 and 2 also easily hold by a practical assumption that correlation between a target training sample and a weight column follows the Gaussian variable (i.e., they are independent of each other) if they have no deterministic relation. We formally state it as follows:

Proposition 3.1. (a) Suppose that $|\mathbf{w}_r(0)|$ is invariant of n for any $r \in \{1 : m\}$ (i.e., $|\mathbf{w}_r(0)| = \mathcal{O}(1)$). (b) Suppose that given some positive constant ϵ , $|\mathbf{w}_r(0)^\top \mathbf{x}_p| \geq \epsilon$ satisfies for any $r \in \{1 : m\}$ and $p \in \{1 : n\}$ without having randomness. (c) Suppose also that for any $r \in \{1 : m\}$ and $p \in \{1 : n\}$ with having randomness, $\mathbb{P}[|\mathbf{w}_r(0)^\top \mathbf{x}_p| \leq x] = \mathcal{O}(x)$ satisfies for any $x > 0$ (e.g., $\mathbf{w}_r(0)^\top \mathbf{x}_p$ follows the Gaussian distribution). Then, both conditions 1 and 2 hold with $R = \mathcal{O}(\frac{\lambda_0}{n^2})$ and $m = \mathcal{O}(n^\alpha)$ for some sufficiently large α .

Proposition 3.1 indicates that the conditions 1 and 2 hold even when $\mathbf{W}(0)$ is partially random (i.e., only some columns of $\mathbf{W}(0)$ are random and the others have the deterministic relation with training dataset). In addition, we show that the proposed conditions 1 and 2 even hold in the case where $\mathbf{W}(0)$ is completely random, proving its global mildness, as specified in the following remark.

Remark. For simplicity, we let $\|\mathbf{x}_j\| = 1$ for all $j \in \{1 : n\}$. Suppose that each element of $\mathbf{W}(0)$ is i.i.d. given as the normal distribution. Then, as Proposition 3.1 holds, with probability at least $1 - \delta$ over $\mathbf{W}(0)$, both the conditions 1 and 2 hold with $R = \mathcal{O}(\frac{\lambda_0}{n^2})$ and $m = \mathcal{O}(n^\alpha)$ for some sufficiently large constant α .

Under the above setup and conditions, then we are now ready to present our theorem, which shows the generalization bound with arbitrary initialization as follows.

Theorem 3.1 (Generalization Error Bound with Arbitrary Initialization). *Suppose that the conditions 1 ~ 3 hold, $\|\mathbf{Y}_i\| = O(\sqrt{n})$ for all $i \in \{1 : d_o\}$, $m = \Omega\left(\frac{n^2}{\lambda_0^2 R^2 \delta}\right)$, $m = \Omega\left(\frac{d_o \cdot n^4}{\min(\lambda_0^2, \lambda_0^4)}\right)$, the set $\{\mathbf{x}_j := \phi_{\mathbf{V}^*}(\mathbf{z}_j)\}_{j=1}^n$ of n training samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$, and $\eta = O(\frac{\lambda_0}{n^2})$. Suppose also that $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$ with probability at least $1 - \delta/3$ for n i.i.d. training samples $\{(\mathbf{z}_i, \mathbf{y}_i)\}_{i=1}^n$ from data distribution \mathcal{D} . Then, with probability at least $1 - \delta$ over the random initialization of \mathbf{A} and the training samples, it follows that for any $k \geq \Omega(\frac{1}{\eta \lambda_0} \log \frac{n}{\delta})$,*

$$\mathbb{E}_{\mathcal{D}} \left[\frac{1}{2} \|\mathbf{Y} - F(\mathbf{W})\|^2 \right] \leq \underbrace{\sum_{i=1}^{d_o} \frac{\left\| \mathbf{H}_i(0)^{-\frac{1}{2}} (\mathbf{Y}_i - h_{\mathbf{V}^*, \mathbf{W}(0)}(\mathbf{Z})_i) \right\|}{\sqrt{n}}}_{\text{Multi-class Revised CMD}} + \mathcal{O} \left(d_o \sqrt{\frac{\log \frac{n}{\lambda_0^2 \delta}}{n}} \right), \quad (19)$$

where $\mathbf{Y}_i, h_{\mathbf{V}^*, \mathbf{W}(0)}(\mathbf{z})_i \in \mathbb{R}^n$ denote all the collection of labels/outputs of the class i , resp.

Note that the upper bound 1 of condition $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$ in Theorem 3.1 can be easily extended to the case of any constant other than 1, thereby being satisfied in practice. As the second term in Eq. 19 trivially converges to 0 as n increases, so it can be interpreted that the first term in Eq. 19 represents the generalization error bound. The multi-class revised CMD term in Eq. 19 does not rely on the Gram matrix \mathbf{H}^∞ but on the Gram matrix $\mathbf{H}(0)$, which allows arbitrary initialization. Furthermore, Theorem 3.1 can be thought of as a generalization of Theorem 2.5 in the absence of a feature extractor $\phi_{\mathbf{V}^*}(\cdot)$ since we only have a 2-layer ReLU classifier in this case.

3.3 OPENING THE DOOR TO PRACTICE: NTK REGULARIZER IN FINE-TUNING

In this section, we discuss the potential applications based on our revised analysis on generalization error, and as an example, we will introduce how our bound can be utilized in practice. As an observation, we can regard the k' -th step parameter $\mathbf{W}(k')$ for any $k' \in \mathbb{N}$ as a new initial parameter $\widetilde{\mathbf{W}}(0)$ so that the parameter $\mathbf{W}(k' + 1)$ updated at the next step from k' -th step can be viewed as the parameter $\widetilde{\mathbf{W}}(1)$ updated only once from the new initial parameter $\widetilde{\mathbf{W}}(0)$. Since Theorem 3.1 allows arbitrary initialization, this observation provides the following remark.

Remark. Suppose that all conditions in Theorem 3.1 for $\mathbf{W}(0)$ hold if $\mathbf{W}(0)$ is replaced with $\mathbf{W}(k)$ at any step k . Then, the generalization error bound can be again characterized as

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{d_o} \left\| \mathbf{H}_i(k)^{-1/2} (\mathbf{Y}_i - h_{\mathbf{V}^*, \mathbf{W}(k)}(\mathbf{Z})_i) \right\|. \quad (20)$$

where $\mathbf{H}_i(k)$ is defined in Eq. 16. By the above observation and remark, the multi-class revised CMD in Eq. 20 could further be thought of as the generalization error bound of some fine-tuned networks. This fact motivates the following NTK regularizer,

$$\mathcal{R}_{\text{NTK}}(\mathbf{W}) = \frac{1}{\sqrt{n}} \sum_{i=1}^{d_o} \left\| \mathbf{H}_i(\mathbf{W})^{-1/2} (\mathbf{Y}_i - h_{\mathbf{V}^*, \mathbf{W}}(\mathbf{Z})_i) \right\|, \quad (21)$$

and training with the regularization loss can play a crucial role in directly reducing the generalization error. Therefore, we propose to solve the following optimization problem in fine-tuning: $\min_{\mathbf{W}} \{ \mathcal{L}(h_{\mathbf{V}^*, \mathbf{W}}(\mathbf{z}), \mathbf{y}) + \mu \mathcal{R}_{\text{NTK}}(\mathbf{W}) \}$ where μ represents the strength of the regularizer.

In fact, the proposed NTK regularizer in Eq. 21 requires the computations of inverse Gram matrix $\mathbf{H}_i(k)^{-1/2}$ for each class $i \in \{1 : d_o\}$, which makes network training computationally heavy in practice. To bypass this issue, we suggest to use the fixed Gram matrix $\mathbf{H}_i(0)$ computed on the initial parameter $\mathbf{W}(0)$ for all i (in fine-tuning regime, the initial parameter $\mathbf{W}(0)$ boils down to some pre-trained parameter \mathbf{W}^*). The intuition behind this is as follows: the generalization error bound depends on the multi-class revised CMD as in Theorem 2.5 for arbitrary initialization. In the over-parametrization regime, since the model parameter will remain close to its initial point during training, the Gram matrix $\mathbf{H}_i(k)$, depending on the k -th model parameter $\mathbf{W}(k)$, will also stay close to its initial Gram matrix $\mathbf{H}_i(0)$ for all i . Though the Gram matrix $\mathbf{H}_i(k)$ is fixed to $\mathbf{H}_i(0)$ in NTK regularizer in Eq. 21, the gradient-based training on regularized loss is still possible since the gradient with respect to \mathbf{W} will be backpropagated through $h_{\mathbf{V}^*, \mathbf{W}}$ in $\mathcal{R}_{\text{NTK}}(\mathbf{W})$.

Table 1: Comparisons of generalization error bound: Theorem 2.2 (baseline) vs. Theorem 2.5 (ours). Note that *the error term \mathcal{E} is absent in Theorem 2.5 of our revised analysis*. Our revised theory removing the error term achieves significantly lower generalization error bound.

Dataset	Error Term \mathcal{E} in Eq. 7	Original CMD	Revised CMD (Ours)	Original Gen. Err. Bound	Revised Gen. Err. Bound (Ours)
MNIST	7×10^3	0.5998	0.5997	$\mathcal{O}(10^3)$	$\mathcal{O}(10^{-1})$
FashionMNIST	1×10^5	0.2617	0.2618	$\mathcal{O}(10^5)$	$\mathcal{O}(10^{-1})$
CIFAR-10	4×10^3	2.0605	2.0604	$\mathcal{O}(10^3)$	$\mathcal{O}(1)$

Lastly, note that our NTK regularizer is expected to exhibit its greatest effect in boosting generalization given a considerably lack of data, rather than in cases where a sufficient number of samples are available to achieve plausible performance. Also, the inverse Gram matrix involves $\mathcal{O}(n^3)$ computations w.r.t. sample size n , we mainly focus on the lack-of-data scenario for evaluating NTK regularizer.

4 NUMERICAL SIMULATIONS

The primary goal in experiments is to verify (i) the tighter generalization error bound of our revised analysis in Section 2 and (ii) whether the generalization is indeed improved via NTK regularizer, which will corroborate our theory in Section 3.

4.1 THEORY VALIDATION: COMPARISONS OF GENERALIZATION ERROR BOUNDS

In order to compare the generalization error bounds between Theorem 2.2 and Theorem 2.5, we consider 2-layer ReLU networks with the width $m = 10000$. Since the baseline theory includes an error term \mathcal{E} which is absent in our revised bound, the key points in comparing the generalization error bound is (i) how much the baseline CMD term differs from the revised CMD term, and (ii) the scale of the error term. Toward this, we consider three benchmark datasets: (i) MNIST, (ii) FashionMNIST, and (iii) CIFAR-10. While our theory can allow the multi-dimensional outputs as in Theorem 3.1, the baseline error bound in Theorem 2.2 could only guarantee the regression or binary classification (refer to Corollary 5.2 in Arora et al. (2019a)), i.e., single output case. Thus, we randomly pick two classes for each dataset. The revised CMD term depends on the initial network output $\mathbf{u}(0)$, thus we initialize $\mathbf{W}(0)$ with the practical Kaiming normal distribution (He et al., 2015) whose scaling does not rely on the sample size n at all, which violates the conditions of Theorem 2.2.

Table 1 illustrates the direct comparison of generalization error bound. Note first that the revised CMD does not significantly differ from the original CMD. Although the revised CMD is slightly larger than the original CMD in FashionMNIST dataset, the difference is only on the scale of 10^{-3} . To compute the scale of the error term \mathcal{E} in Eq. 7, which is the second most important factor in the comparison of generalization errors, we set the failure probability $\delta = 0.01$ (larger value is also fine). Note that the error terms \mathcal{E} have the scale of $10^3 \sim 10^5$ while the CMD terms have values only about 0.6, 0.26, and 2.06 for each dataset as can be seen in Table 1. Hence, our revised analysis removing the error terms could significantly improve the existing generalization error bound. It is important to note that the results in Table 1 are not limited to solely to the width $m = 10000$, since our findings hold true across a broad range of width m from 10^2 to 10^4 , and we provide the results in Appendix A.

4.2 VERIFICATION OF MULTI-CLASS REVISED CMD VIA NTK REGULARIZER

In order to verify our theory (Theorem 3.1) on the arbitrary initialization, we use the NTK regularizer proposed in Eq. 21. Toward this, we fine-tune pre-trained models given a limited number of samples, which closely mirrors the typical scenario in medical applications. Thus, we consider the skin cancer classification for our experiments. The details on experimental settings are provided in Appendix.

Model. We use pre-trained ResNet-18 (He et al., 2016) on the ImageNet (Deng et al., 2009), which is publicly available from popular deep learning libraries (Abadi et al., 2016; Paszke et al., 2019). As

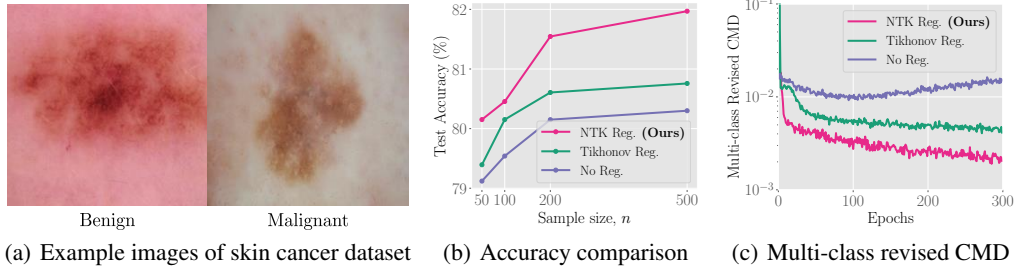


Figure 2: (a) Example images of skin cancer dataset, (b) the results on skin cancer varying the number of training samples. (c) comparison of multi-class revised CMD among different methods.

suggested in our theory, we replace the classifier of ResNet-18 with a 2-layer ReLU network, and the parameter \mathbf{A} of second layer of the classifier is initialized and fixed according to Section 3.1. The parameter \mathbf{V}^* of the feature extractor is frozen to those of the pre-trained model (one of conventional fine-tuning strategy), while only the parameter $\mathbf{W}(0)$ of the first layer of the classifier is updated.

Dataset. The skin cancer classification has been considered as one of popular medical applications in many literature (Esteva et al., 2017; Wu et al., 2022; Bello et al., 2024). The goal of this task is to classify types of skin cancer for given images into: (i) benign or (ii) malignant. In this experiment, we collect RGB images of size 224×224 of combined dataset, which consists of HAM10000 (Tschandl et al., 2018) and International Skin Imaging Collaboration (ISIC 2020). The dataset is splitted into 2077/560/660 images for train/valid/test respectively and we provide example images of the dataset in Fig. 2(a) for better understanding. The detailed information of dataset is provided in Appendix.

Baselines. To validate the efficacy of NTK regularizer, we consider two baselines: (i) no regularizer (regular fine-tuning) and (ii) Tikhonov regularizer corresponding to the case of $\mathbf{H}_i(k)$ being the identity matrix \mathbf{I} in Eq. 21, i.e., $\mathcal{R}_{\text{Tikhonov}}(\mathbf{W}) = \frac{1}{\sqrt{n}} \sum_{i=1}^{d_o} \|\mathbf{Y}_i - h_{\mathbf{V}^*, \mathbf{W}}(\mathbf{Z})_i\|$. The reason for considering the Tikhonov regularizer is to examine the role of the Gram matrix in \mathcal{R}_{NTK} .

NTK regularizer works in practice. We consider small amount of training dataset to simulate a limited-number-of-sample scenario. Toward this, we randomly choose $\{50, 100, 200, 500\}$ samples from training dataset, on which we fine-tune the pre-trained ResNet-18. As depicted in Fig. 2(b), NTK regularizer indeed improves the generalization upon regular fine-tuning. Note that the advantage of NTK regularizer over the Tikhonov regularizer clearly can be clearly observed, which indicates that the Gram matrix plays an important role in model generalization. In addition, we also compare the multi-class revised CMD term as illustrated in Fig. 2(c). An interesting observation is that the multi-class revised CMD decreases as the model generalization improves observed in Fig 2(b). This suggests that directly reducing the multi-class revised CMD can potentially enhance the generalization, which demonstrates the validity of our proposed NTK regularizer.

5 CONCLUSION

In this study, we revised the existing NTK-based theory of optimization and generalization for overparametrized neural networks. Our revised analysis successfully remove the unreasonable assumption on the initialization and provide tighter bound for generalization error. Going further, we extended our revised analysis that allow for arbitrary initialization and multi-dimensional outputs. By extending NTK theory to a network with arbitrary initialization, we were able to propose the concept of NTK regularizer, which was previously unattainable, and validate its effectiveness. *The most promising aspect of this study is that it enables the application of NTK theory to pre-trained networks.* This extension of NTK theory is expected to be applicable to various practical scenarios that require predicting the performance of pre-trained networks, such as fine-tuning, domain adaptation, out-of-distribution detection, and more. We empirically validated that our revised analysis indeed achieve significantly lower generalization error bound and also showed our NTK regularizer to be effective in fine-tuning, demonstrating that NTK theory provides a connection to real-world applications.

REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. {TensorFlow}: a system for {Large-Scale} machine learning. In *12th USENIX symposium on operating systems design and implementation (OSDI 16)*, pp. 265–283, 2016.
- Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. In *International Conference on Machine Learning*, pp. 254–263, 2018.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pp. 322–332, 2019a.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pp. 8139–8148, 2019b.
- Yu Bai and Jason D Lee. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. *arXiv preprint:1910.01619*, 2019.
- Abayomi Bello, Sin-Chun Ng, and Man-Fai Leung. Skin cancer classification using fine-tuned transfer learning of densenet-121. *Applied Sciences*, 14(17):7707, 2024.
- Yuan Cao and Quanquan Gu. Generalization error bounds of gradient descent for learning over-parameterized deep ReLU networks. *arXiv preprint:1902.01384*, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pp. 1675–1685, 2019.
- Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint:1810.02054*, 2018.
- Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Wei Hu, Zhiyuan Li, and Dingli Yu. Simple and effective regularization methods for training on noisily labeled data with generalization guarantee. *arXiv preprint:1905.11368*, 2019.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, pp. 8571–8580, 2018.
- Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow ReLU networks. *arXiv preprint:1909.12292*, 2019.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, pp. 8157–8166, 2018.

- Roi Livni, Shai Shalev-Shwartz, and Ohad Shamir. On the computational efficiency of training neural networks. In *Advances in Neural Information Processing Systems*, pp. 855–863, 2014.
- Chao Ma, Lei Wu, et al. Analysis of the gradient descent algorithm for a deep neural network model with skip-connections. *arXiv preprint:1904.05263*, 2019.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. The MIT Press, 2018.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint:1412.6614*, 2014.
- Atsushi Nitanda and Taiji Suzuki. Refined generalization analysis of gradient descent for over-parameterized two-layer neural networks with smooth activations on classification problems. *arXiv preprint:1905.09870*, 2019.
- Samet Oymak, Zalan Fabian, Mingchen Li, and Mahdi Soltanolkotabi. Generalization guarantees for neural networks via harnessing the low-rank structure of the Jacobian. *arXiv preprint:1906.05392*, 2019.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Lili Su and Pengkun Yang. On learning over-parameterized neural networks: A functional approximation prospective. In *Advances in Neural Information Processing Systems*, pp. 2637–2646, 2019.
- Mingxing Tan. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.
- Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018.
- Yinhao Wu, Bin Chen, An Zeng, Dan Pan, Ruixuan Wang, and Shen Zhao. Skin cancer classification with deep learning: a systematic review. *Frontiers in Oncology*, 12:893972, 2022.
- Zhi-Qin John Xu, Jiwei Zhang, Yaoyu Zhang, and Chengchao Zhao. A priori generalization error for two-layer ReLU neural network through minimum norm solution. *arXiv preprint:1912.03011*, 2019.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2016.
- Guodong Zhang, James Martens, and Roger B Grosse. Fast convergence of natural gradient descent for overparameterized neural networks. In *Advances in Neural Information Processing Systems*, pp. 8080–8091, 2019.

APPENDIX

A EXPERIMENTAL SETTINGS AND ADDITIONAL RESULTS

A.1 COMPARISON OF CMD FOR VARIOUS WIDTHS

Figure 3 illustrates the comparison of original CMD and the revised CMD. Note that the difference between the original CMD and revised CMD is not significant across all widths $m = 10^2 \sim 10^4$. Also, the scale of difference between the original CMD and revised CMD is only about the level of 10^{-3} . Thus, it can be clearly concluded that the scale of error term \mathcal{E} in Table 1 is dominant in generalization error bound, thus our revised theory *removing the error term* significantly improves the existing generalization error bound.

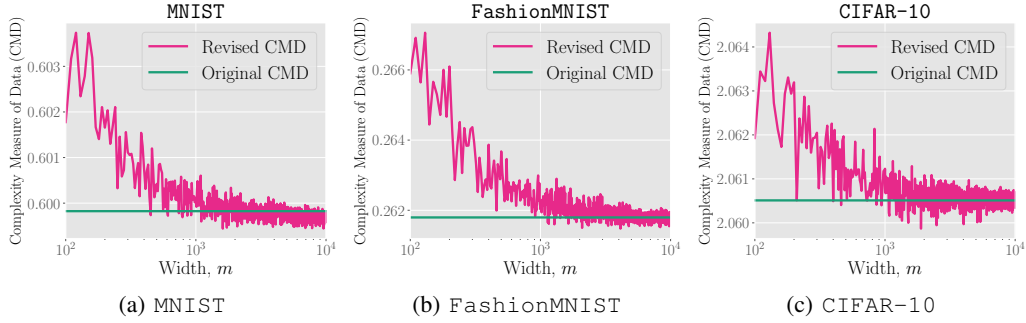


Figure 3: Comparisons of CMD term varying the width m .

A.2 ADDITIONAL RESULTS ON NTK REGULARIZER

For evaluating NTK regularizer, we conduct additional experiments and provide more results in this section. For this purpose, we consider fine-tuning scenario as in Section 4 and use ResNet-18 pre-trained on ImageNet dataset. As suggested in Section 3, we replace the linear classifier of original ResNet-18 with 2-layer ReLU networks. For 2-layer ReLU classifier, we choose the width $m = 10000$ and initialize each parameter according to Section 3.1. Also, our theorem 3.1 allows the multi-dimensional network outputs, thus we consider image classification with CIFAR-10 dataset, which corresponds to $d_o = 10$.

To simulate a limited-number-of-sample situation, we also randomly choose $n \in \{50, 100, 200, 500\}$ samples in training dataset. Note that we do not employ any data augmentation technique to solely validate the effect of NTK regularizer more accurately. As baseline methods, we consider no regularizer (regular fine-tuning) and the Tikhonov regularizer $\mathcal{R}_{\text{Tikhonov}}$ introduced in Section 4.2, which can evaluate the role of Gram matrix in \mathcal{R}_{NTK} .

Table 2 shows the results of fine-tuning ResNet-18 with varying the number of samples. Note that fine-tuning with NTK regularizer consistently outperforms baselines, which indicates that our proposed regularizer indeed can enhance the model generalization across all sample sizes. Also, directly reducing the multi-class revised CMD in 3.1 indeed helps to improve the generalization. For this experiment, we do not employ additional training technique such as weight decay, learning rate scheduling, data augmentation, and etc. In this sense, it should be emphasized that fine-tuning with NTK regularizer can achieve better performance when equipped with such training recipe.

Table 2: CIFAR-10 classification test accuracy (%) varying the training sample size n .

Methods	$n = 50$	$n = 100$	$n = 200$	$n = 500$
No Reg.	31.21	34.61	38.77	41.75
Tikhonov Reg. $\mathcal{R}_{\text{Tikhonov}}$	31.32	34.82	38.92	42.03
NTK Reg. \mathcal{R}_{NTK} (Ours)	32.71	35.47	39.41	42.65

B PROOFS OF SECTION 2

B.1 ADDITIONAL DEFINITIONS

We recall weight matrix $\mathbf{W}(k)$ at the k th step of gradient descent as

$$\mathbf{W}(k+1) := \mathbf{W}(k) - \eta \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}}|_{\mathbf{W}=\mathbf{W}(k)}. \quad (22)$$

Furthermore, we define $\mathbf{Z}(k) := \frac{\partial \mathbf{u}(k)}{\partial \text{vec}(\mathbf{W}(k))} \in \mathbb{R}^{md \times n}$. Thus, $\mathbf{Z}(k)$ is derived as

$$\mathbf{Z}(k) = \frac{1}{\sqrt{m}} \begin{pmatrix} \mathbb{I}_{1,1}(k)a_1\mathbf{x}_1 & \dots & \mathbb{I}_{1,n}(k)a_1\mathbf{x}_n \\ \dots & \dots & \dots \\ \mathbb{I}_{m,1}(k)a_m\mathbf{x}_1 & \dots & \mathbb{I}_{m,n}(k)a_m\mathbf{x}_n \end{pmatrix} \in \mathbb{R}^{md \times n}, \quad (23)$$

where $\mathbb{I}_{p,q}(k) := \mathbb{1}\{\mathbf{x}_q^\top \mathbf{w}_p(k) \geq 0\}$. Then, equation 138 can be expressed as

$$\text{vec}(\mathbf{W}(k+1)) = \text{vec}(\mathbf{W}(k)) - \eta \mathbf{Z}(k)(\mathbf{u}(k) - \mathbf{y}). \quad (24)$$

B.2 PROOF OF PROPOSITIONS

In this section, we introduce the proofs of Propositions 2.1, 2.2, and 2.3, which are given sequentially as follows.

Proposition B.1 (Proposition 2.1). *Suppose that $\mathbf{x}_i \neq a\mathbf{x}_j$ for any $a \in \mathbb{R}$ and different $i, j \in \{n\}^2$. Then, $\lambda_0 = O(\sqrt{n}) > 0$ holds.*

Proof of Proposition B.1. As $\|\mathbf{x}_i\| = 1$ for all $i \in \{n\}$, there exists a finite constant $c \leq 1$ such that $|\mathbf{x}_i^\top \mathbf{x}_j| \leq c$ for $i, j \in \{n\}$. From the definition of \mathbf{H}^* , $|[\mathbf{H}^*]_{i,j}| \leq c$ for $i, j \in \{n\}^2$. Let $\mathbf{z} \in \mathbb{R}^n$ be a vector whose elements belong to $\{-1/\sqrt{n}, 1/\sqrt{n}\}$. Then, $|(\mathbf{H}^*\mathbf{z})_i| \leq c\sqrt{n} = O(\sqrt{n})$ for $i \in \{n\}$ so that

$$\frac{\|\mathbf{H}^*\mathbf{z}\|}{\|\mathbf{z}\|} = \frac{\sqrt{\sum_{i=1}^n |(\mathbf{H}^*\mathbf{z})_i|^2}}{\sqrt{n}} = O(\sqrt{n}). \quad (25)$$

Thus, from the definition of λ_0 and equation 25, λ_0 is upper bounded as

$$\lambda_0 = \min_{\mathbf{v} \in \mathbb{R}^n \text{ s.t. } \|\mathbf{v}\|=1} \frac{\|\mathbf{H}^*\mathbf{v}\|}{\|\mathbf{v}\|} \leq \frac{\|\mathbf{H}^*\mathbf{z}\|}{\|\mathbf{z}\|} = O(\sqrt{n}). \quad (26)$$

From Theorem 3.1 in Du et al. (2018), it follows that $\lambda_0 > 0$ if $\mathbf{x}_i \neq a\mathbf{x}_j$ for any $a \in \mathbb{R}$ and different $i, j \in \{n\}^2$. Thus, the proof is completed. \square

Proposition B.2 (Proposition 2.2). *Suppose all conditions in Theorem 2.1 hold, $\kappa = \Theta(1)$, and $m = \Omega(\frac{n^{6+\alpha}}{\lambda_0^4})$ with any $\alpha > 0$. Then, with probability at least $1 - \delta$ for $\delta \in (0, 1)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, the right hand side of equation 10 converges to zero as k and n increase.*

Proof of Proposition B.2. Note that the right hand side of equation 10 is given as

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (1 - \eta\lambda_i)^{2k} \left(\mathbf{v}_i^\top (\mathbf{y} - \mathbf{u}(0)) \right)^2} + O\left(\frac{n^3}{\sqrt{m}\lambda_0^2\delta^2}\right). \quad (27)$$

Note that the first term in equation 27 is upper bounded as

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (1 - \eta\lambda_i)^{2k} \left(\mathbf{v}_i^\top (\mathbf{y} - \mathbf{u}(0)) \right)^2} \leq \frac{1}{\sqrt{n}} (1 - \eta\lambda_0)^k \|\mathbf{y} - \mathbf{u}(0)\| \stackrel{(a)}{=} O\left((1 - \eta\lambda_0)^k\right), \quad (28)$$

where (a) follows from Lemma 15. As it follows from the conditions in Theorem 2.1 that $0 < (1 - \eta\lambda_0) < 1$ holds, by applying this fact to equation 28, we obtain that the first term of the right hand side in equation 10 converges to zero as k increases.

If $m = \Omega(\frac{n^{6+\alpha}}{\lambda_0^4})$, the second term in equation 27 is given as

$$O\left(\frac{n^3}{\sqrt{m}\lambda_0^2\delta^2}\right) = O(n^{-\frac{\alpha}{2}}). \quad (29)$$

That is, the second term of the right hand side in equation 10 also converges to zero as n increases, thereby completing the proof by applying equation 28 and equation 29 to equation 27. \square

Proposition B.3 (Proposition 2.3, variant of Corollary 6.2 in Arora et al. (2019a)). Suppose $y_i - u_i(0) = g(\mathbf{x}_i) := \sum_j \alpha_j (\beta_j^\top \mathbf{x}_i)^{p_j}$ for all $i \in \{n\}$, where for each j , $p_j \in \{1, 2, 4, 6, \dots\}$ and $\alpha_j \in \mathbb{R}$ and $\beta_j \in \mathbb{R}^d$ are any constants w.r.t. n . Then,

$$\sqrt{\frac{2(\mathbf{y} - \mathbf{u}(0))^\top \mathbf{H}^{*-1}(\mathbf{y} - \mathbf{u}(0))}{n}} \leq \frac{6 \sum_j p_j |\alpha_j| \|\beta_j\|^{p_j}}{\sqrt{n}} = O\left(\frac{1}{\sqrt{n}}\right). \quad (30)$$

Proof of Proposition B.3. The proof is completed by replacing y_i in Corollary 6.2 in Arora et al. (2019a) with $y_i - u_i(0)$ for all $i \in \{n\}$. \square

B.3 PROOF OF THEOREM 2.3

In this section, we prove Theorem 2.3. We first show some technical lemmas.

The following lemma provides an upper bound of the magnitude of the initial NN output.

Lemma 1. Suppose that set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$. Then, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\|\mathbf{u}(0)\|^2 = O\left(\frac{n\kappa^2}{\delta}\right). \quad (31)$$

Proof of Lemma 14. It follows that

$$\begin{aligned} \mathbb{E}_{\mathbf{a}}[\|\mathbf{u}(0)\|^2] &= \mathbb{E}_{\mathbf{a}}[\|(f_{\mathbf{W}(0)}(\mathbf{x}_1), \dots, f_{\mathbf{W}(0)}(\mathbf{x}_n))\|^2] \\ &= \mathbb{E}_{\mathbf{a}}\left[\sum_{j=1}^n |f_{\mathbf{W}(0)}(\mathbf{x}_j)|^2\right] \\ &= \mathbb{E}_{\mathbf{a}}\left[\frac{1}{m} \sum_{j=1}^n \left|\sum_{r \in \{1:m\}} a_r \sigma(\mathbf{w}_r^\top \mathbf{x}_j)\right|^2\right] \\ &= \frac{1}{m} \sum_{j=1}^n \sum_{r \in \{1:m\}} |\sigma(\mathbf{w}_r^\top \mathbf{x}_j)|^2. \end{aligned} \quad (32)$$

Furthermore, it follows that

$$\begin{aligned} \frac{1}{m} \sum_{j=1}^n \sum_{r \in \{1:m\}} |\sigma(\mathbf{w}_r^\top \mathbf{x}_j)|^2 &\leq \frac{1}{m} \sum_{j=1}^n \sum_{r \in \{1:m\}} |\mathbf{w}_r^\top \mathbf{x}_j|^2 \\ &\stackrel{(a)}{\leq} \frac{n}{m} \sum_{r=1}^m v_r^2, \end{aligned} \quad (33)$$

where v_r for $r \in \{m\}$ is independently sampled from $\mathcal{N}(0, \kappa^2)$ and (a) follows from the rotational invariance of Gaussian random vector and the fact that \mathbf{w}_r follows $\mathcal{N}(0, \kappa^2 \mathbf{I}_d)$ for $r \in \{1 : m\}$.

By combining equation 147 and equation 148 and using the fact that $\mathbb{E}[\sum_{r=1}^m v_r^2] = m\kappa^2$, we obtain

$$\mathbb{E}_{\mathbf{W}(0), \mathbf{a}}[\|\mathbf{u}(0)\|^2] = O(n\kappa^2).$$

Therefore, by using Markov's inequality, $\|\mathbf{u}(0)\|^2 = O(n\kappa^2/\delta)$ is satisfied with probability at least $1 - \delta$. \square

Then, by using Lemma 14, we can also obtain an upper bound of gap between the initial NN output and label as Lemma 15.

Lemma 2. Suppose that set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$. If $\|\mathbf{y}\| = O(\sqrt{n})$ is satisfied, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\|\mathbf{y} - \mathbf{u}(0)\|^2 = O\left(\frac{\max(\kappa^2, 1)n}{\delta}\right).$$

Proof of Lemma 15. It follows from Lemma 14 that with probability at least $1 - \delta$,

$$2 \max(\|\mathbf{y}\|^2, \|\mathbf{u}(0)\|^2) = O\left(\max\left(1, \frac{\kappa^2}{\delta}\right)n\right) = O\left(\frac{\max(\kappa^2, 1)n}{\delta}\right). \quad (34)$$

Then, the proof is completed by applying the following inequality to equation 149.

$$\|\mathbf{y} - \mathbf{u}(0)\|^2 \leq \|\mathbf{y}\|^2 + \|\mathbf{u}(0)\|^2 \leq 2 \max(\|\mathbf{y}\|^2, \|\mathbf{u}(0)\|^2)$$

□

The following lemma (i.e., Lemma 16) gives an upper bound of the gap between each trained weight vector and its initialization, when the training loss is reduced by the GD optimization. This lemma is the result of extending the condition for κ to arbitrary $\kappa > 0$ from $\kappa = 1$, which is given in Corollary 4.1 in Du et al. (2018).

Lemma 3 (Variant of Corollary 4.1 in Du et al. (2018)). *We are given arbitrary $\kappa > 0$. Suppose that set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$. If the following condition holds for $k' \in \{0, 1, \dots, k-1\}$,*

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq (1 - \frac{\eta\lambda_0}{2})^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2, \quad (35)$$

then for every $r \in \{m\}$,

$$\|\mathbf{w}_r(k) - \mathbf{w}_r(0)\| \leq \frac{4\sqrt{n}\|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m}\lambda_0} \quad (=: R'), \quad (36)$$

where $\mathbf{w}_j(k)$ is the column of $\mathbf{W}(k) =: [\mathbf{w}_1(k), \dots, \mathbf{w}_m(k)]$ at the k -th step of GD.

Proof of Lemma 16. Since

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{w}_r} = \frac{1}{\sqrt{m}} \sum_{q=1}^n (u_q - y_q) a_r \mathbf{x}_q \mathbf{1}(\mathbf{w}_r^\top \mathbf{x}_q \geq 0),$$

we get

$$\left\| \frac{\partial L(\mathbf{W}(k'))}{\partial \mathbf{w}_r(k')} \right\| \leq \frac{\sqrt{n}}{\sqrt{m}} \|\mathbf{y} - \mathbf{u}(k')\|.$$

Thus, we have

$$\begin{aligned} \|\mathbf{w}_r(k) - \mathbf{w}_r(0)\| &\leq \eta \sum_{k'=0}^{k-1} \left\| \frac{\partial L(\mathbf{W}(k'))}{\partial \mathbf{w}_r(k')} \right\| \\ &\leq \eta \sum_{k'=0}^{k-1} \frac{\sqrt{n}}{\sqrt{m}} \|\mathbf{y} - \mathbf{u}(k')\| \\ &\stackrel{(a)}{\leq} \frac{\sqrt{n}}{\sqrt{m}} \sum_{k'=0}^{k-1} \eta (1 - \frac{\eta\lambda_0}{2})^{k'/2} \|\mathbf{y} - \mathbf{u}(0)\| \\ &\leq \frac{\sqrt{n}}{\sqrt{m}} \sum_{k'=0}^{k-1} \eta (1 - \frac{\eta\lambda_0}{4})^{k'} \|\mathbf{y} - \mathbf{u}(0)\| \\ &\leq \frac{4\sqrt{n}\|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m}\lambda_0}, \end{aligned}$$

where (a) follows from equation 150. □

As a result of Lemma 16, from the following lemma (i.e., Lemma 4), we can obtain an upper bound of the magnitude of trained NN output, when the training loss is reduced by the GD optimization.

Lemma 4. *We are given arbitrary $\kappa > 0$. Suppose that set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$. If the following condition holds for $k' \in \{0, 1, \dots, k-1\}$,*

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq (1 - \frac{\eta\lambda_0}{2})^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2, \quad (37)$$

then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\|\mathbf{u}(k)\|^2 = O\left(\frac{n^3 \max(\kappa, 1)^2}{m\lambda_0^2 \delta^2} + \frac{n\kappa^2}{\delta^2}\right). \quad (38)$$

Proof of Lemma 4. Define a set of all weights whose distance from $\mathbf{W}(0)$ is smaller than R' as

$$\Gamma(\mathbf{W}(0), R') := \left\{ \tilde{\mathbf{W}} = [\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_m] \in \mathbb{R}^{m \times d} \mid \max_{r \in \{m\}} \|\tilde{\mathbf{w}}_r - \mathbf{w}_r(0)\| \leq R' \right\}. \quad (39)$$

Then, for any matrix $\tilde{\mathbf{W}} = [\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_m] \in \mathbb{R}^{m \times d}$ belonging to $\Gamma(\mathbf{W}(0), R')$ and any $j \in \{n\}$, it follows that

$$\begin{aligned} \mathbb{E}_{\mathbf{a}}[f_{\tilde{\mathbf{W}}}(\mathbf{x}_j)^2] &= \mathbb{E}_{\mathbf{a}} \left[\frac{1}{m} \left(\sum_{r \in \{m\}} a_r \sigma(\tilde{\mathbf{w}}_r^\top \mathbf{x}_j) \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{a}} \left[\frac{1}{m} \left(\sum_{r \in \{m\}} \sigma(\tilde{\mathbf{w}}_r^\top \mathbf{x}_j)^2 + \sum_{r, r' \in \{m\} \times \{m\}, r \neq r'} a_r a_{r'} \sigma(\tilde{\mathbf{w}}_r^\top \mathbf{x}_j) \sigma(\tilde{\mathbf{w}}_{r'}^\top \mathbf{x}_j) \right) \right] \\ &\stackrel{(a)}{=} \frac{1}{m} \left(\sum_{r \in \{m\}} \sigma(\tilde{\mathbf{w}}_r^\top \mathbf{x}_j)^2 \right) + \frac{1}{m} \left(\sum_{r, r' \in \{m\} \times \{m\}, r \neq r'} \mathbb{E}_{\mathbf{a}}[a_r a_{r'}] \sigma(\tilde{\mathbf{w}}_r^\top \mathbf{x}_j) \sigma(\tilde{\mathbf{w}}_{r'}^\top \mathbf{x}_j) \right) \\ &= \frac{1}{m} \left(\sum_{r \in \{m\}} \sigma(\tilde{\mathbf{w}}_r^\top \mathbf{x}_j)^2 \right), \end{aligned} \quad (40)$$

where (a) follows from $\tilde{\mathbf{w}}_r$ and $\tilde{\mathbf{w}}_{r'}$ are independent of the random vector \mathbf{a} (i.e., $\tilde{\mathbf{w}}_r$ and $\tilde{\mathbf{w}}_{r'}$ are only depending on $\mathbf{W}(0)$ and R' as $\tilde{\mathbf{W}}$ is an arbitrary matrix satisfying $\tilde{\mathbf{W}} \in \Gamma(\mathbf{W}(0), R')$). Thus, by using Markov's inequality, we obtain with probability at least $1 - \delta$ over the random initialization of \mathbf{a} ,

$$f_{\tilde{\mathbf{W}}}(\mathbf{x}_j)^2 \leq \frac{1}{\delta m} \sum_{r=1}^m \sigma(\tilde{\mathbf{w}}_r^\top \mathbf{x}_j)^2. \quad (41)$$

Define $\tilde{\mathbf{u}}(\tilde{\mathbf{W}}) := (f_{\tilde{\mathbf{W}}}(\mathbf{x}_1), \dots, f_{\tilde{\mathbf{W}}}(\mathbf{x}_n))^\top \in \mathbb{R}^n$. Then, applying the union bound over equation 41 for $j \in \{n\}$, the following inequalities hold with probability at least $1 - \Omega(\delta)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\begin{aligned} \|\tilde{\mathbf{u}}(\tilde{\mathbf{W}})\|^2 &\stackrel{(a)}{\leq} \frac{1}{\delta m} \sum_{j=1}^n \sum_{r=1}^m \sigma((\tilde{\mathbf{w}}_r - \mathbf{w}_r(0) + \mathbf{w}_r(0))^\top \mathbf{x}_j)^2 \\ &\leq \frac{n}{\delta m} \sum_{r=1}^m (\|\tilde{\mathbf{w}}_r - \mathbf{w}_r(0)\|^2) + \frac{1}{\delta m} \sum_{j=1}^n \sum_{r=1}^m \left| \mathbf{w}_r(0)^\top \frac{\mathbf{x}_j}{\|\mathbf{x}_j\|} \right|^2 \\ &\stackrel{(b)}{\leq} \frac{nR'^2}{\delta} + \frac{1}{\delta m} \sum_{j=1}^n \sum_{r=1}^m \left| \mathbf{w}_r(0)^\top \frac{\mathbf{x}_j}{\|\mathbf{x}_j\|} \right|^2 \\ &\stackrel{(c)}{\leq} \frac{nR'^2}{\delta} + \frac{n}{\delta m} \sum_{r=1}^m \|\mathbf{w}_r(0)\|^2 \\ &\stackrel{(d)}{\leq} \frac{nR'^2}{\delta} + \frac{n\kappa^2}{\delta^2}, \end{aligned} \quad (42)$$

where (a) follows from equation 41, (b) follows from the fact that $\tilde{\mathbf{W}}$ belongs to $\Gamma(\mathbf{W}(0), R)$ (i.e., $\|\tilde{\mathbf{w}}_r - \mathbf{w}_r(0)\| \leq R$), (c) follows from Cauchy-Schwarz inequality, and (d) follows from the fact that $\mathbb{E}[\sum_{r=1}^m \|\mathbf{w}_r(0)\|^2] = m\kappa^2$ and Markov's inequality (i.e., $\sum_{r=1}^m \|\mathbf{w}_r(0)\|^2 = m\kappa^2/\delta$ holds with probability at least δ).

On the other hand, it follows from Lemmas 15 and 16 that $\mathbf{W}(k)$ belongs to $\Gamma(\mathbf{W}(0), R')$ with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, where $R' = \sqrt{\frac{c \max(\kappa^2, 1)n^2}{m\lambda_0^2\delta}}$ for some constant c . This implies $\tilde{\mathbf{W}}$ in equation 42 can be replaced by $\mathbf{W}(k)$ (i.e., $\tilde{\mathbf{u}}(\tilde{\mathbf{W}})$ in equation 42 can be replaced by $\tilde{\mathbf{u}}(\mathbf{W}(k)) = \mathbf{u}(k)$).

By using the union bound over the above statement (i.e., $\tilde{\mathbf{u}}(\tilde{\mathbf{W}})$ in equation 42 can be replaced by $\tilde{\mathbf{u}}(\mathbf{W}(k)) = \mathbf{u}(k)$) and the inequality in equation 42 and setting $R' = \sqrt{\frac{c \max(\kappa^2, 1)n^2}{m\lambda_0^2\delta}}$, it follows that with probability at least $1 - \Omega(\delta)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\|\mathbf{u}(k)\|^2 = O\left(\frac{n^3 \max(\kappa, 1)^2}{m\lambda_0^2\delta^2} + \frac{n\kappa^2}{\delta^2}\right). \quad (43)$$

The proof is completed by rescaling δ to a constant such that equation 43 holds with probability at least $1 - \delta$. \square

Now, by using Lemma 4, we prove Theorem 2.3 (i.e., Theorem B.1) as follows.

Theorem B.1 (Theorem 2.3). Suppose that $\|\mathbf{y}\| = O(\sqrt{n})$, $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, and $m = \Omega(n^{3-2\gamma})$. Suppose further that $\kappa = O(n^\alpha)$ holds for some constant $\alpha < 0$. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, there exists a finite integer k such that

$$\|\mathbf{y} - \mathbf{u}(k+1)\|^2 > (1 - \frac{\eta\lambda_0}{2}) \|\mathbf{y} - \mathbf{u}(k)\|^2, \quad (44)$$

where η is any constant such that $0 < \lambda_0\eta < 2$.

Proof. If the following condition equation 45 holds for $k' \in \{0, 1, \dots, \tilde{k} - 1\}$

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq (1 - \frac{\eta\lambda_0}{2})^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2, \quad (45)$$

then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that

$$\begin{aligned} \|\mathbf{u}(\tilde{k})\|^2 &\stackrel{(a)}{\leq} O\left(\frac{n^3 \max(\kappa, 1)^2}{m\lambda_0^2\delta^2} + \frac{n\kappa^2}{\delta^2}\right) \\ &\stackrel{(b)}{\leq} O\left(1 + \frac{n\kappa^2}{\delta^2}\right) \\ &= O(1 + n^{1+2\alpha}) \\ &= o(n), \end{aligned} \quad (46)$$

where (a) follows from Lemma 4 and (b) follows from $m = \Omega(n^{3-2\gamma})$.

On the other hand, if the following condition equation 47 holds for all $k \in \{0, 1, \dots\}$

$$\|\mathbf{y} - \mathbf{u}(k)\|^2 \leq (1 - \frac{\eta\lambda_0}{2})^k \|\mathbf{y} - \mathbf{u}(0)\|^2, \quad (47)$$

there exists an integer $\bar{k} \in \{0, 1, \dots\}$ (e.g., any $\bar{k} > ((\eta\lambda_0)/(2 - \eta\lambda_0))^{-1} \log(\|\mathbf{y} - \mathbf{u}(0)\|^2/\epsilon)$, as derived from equation 220) such that for arbitrary small constant ϵ invariant of n ,

$$\|\mathbf{y} - \mathbf{u}(\bar{k})\|^2 \leq \epsilon. \quad (48)$$

As ϵ is invariant of n , equation 48 implies that

$$\|\mathbf{u}(\bar{k})\|^2 = \Theta(\|\mathbf{y}\|^2). \quad (49)$$

In the case where $\tilde{k} = \bar{k}$ and $\|\mathbf{y}\|^2 = \Theta(n)$, as equation 48 implies $\|\mathbf{u}(\bar{k})\|^2 = \Theta(\|\mathbf{y}\|^2)$ in equation 49 and equation 45 also implies $\|\mathbf{u}(\bar{k})\|^2 = o(n)$ in equation 46, equation 48 and equation 45 are not satisfied at the same time. This is because $\|\mathbf{u}(\bar{k})\|^2 = \Theta(\|\mathbf{y}\|^2)$ in equation 49 is not equal to $\|\mathbf{u}(\bar{k})\|^2 = o(n)$ in equation 46 in this case ($\tilde{k} = \bar{k}$ and $\|\mathbf{y}\|^2 = \Theta(n)$).

Then, for any η satisfying that $0 < \lambda_0\eta < 2$, if equation 45 with this constant η is satisfied for all $k' \geq 0$, there should exist a integer \bar{k} satisfying equation 48, which means that equation 45 and equation 48 are satisfied at the same time. Therefore, equation 45 is not satisfied for some constant $k' \in \{0, 1, 2, \dots\}$ and for any η satisfying that $0 < \lambda_0\eta < 2$. \square

B.4 MODIFICATION OF THEOREM 4.1 IN DU ET AL. (2018)

In order to prove Theorem 2.5 stated in Section 2.3, we first prove Theorem F.1 in this section, since Theorem 2.5 is proved by using the result of Theorem F.1. Theorem F.1 is the result of extending the condition for κ to $\kappa = \Theta(1)$ from $\kappa = 1$, which is given in Theorem 4.1 in Du et al. (2018). Therefore, most of the proof processes for Theorem F.1 (and its technical lemmas) are already proved in Du et al. (2018); we provide them in this section for completeness.

To prove Theorem F.1, we first introduce some technical lemmas.

We introduce the following lemma (i.e. Lemma 5), which is Lemma 3.1 in Du et al. (2018). This result provides that the Gram matrix $\mathbf{H}(0)$ obtained in the finite NN width regime is lower bounded as λ_0 and remains near from that in the infinite NN width regime.

Lemma 5 (Lemma 3.1 in Du et al. (2018)). Define matrix $\mathbf{H}(k) \in \mathbb{R}^{n \times n}$ such that p, q -th entry of $\mathbf{H}(k)$ is given by

$$H_{pq}(k) := \frac{1}{m} \mathbf{x}_p^\top \mathbf{x}_q \sum_{r=1}^m [\mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_p \geq 0, \mathbf{w}_r(k)^\top \mathbf{x}_q \geq 0\}], \quad (50)$$

where $\mathbf{w}_j(k)$ is the j th column vector of $\mathbf{W}(k)$ such that $[\mathbf{w}_1(k), \dots, \mathbf{w}_m(k)] = \mathbf{W}(k)$. If $m = \Omega(\frac{n^2}{\lambda_0^2} \log(\frac{n}{\delta}))$, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, $\|\mathbf{H}(0) - \mathbf{H}^*\|_2 \leq \frac{\lambda_0}{4}$ and $\lambda_{\min}(\mathbf{H}(0)) \geq \frac{3}{4}\lambda_0$, where $\lambda_{\min}(\mathbf{H}(0))$ is the smallest eigenvalue of $\mathbf{H}(0)$.

The following lemma (i.e. Lemma 17) is a direct extension of Lemma 3.2 in Du et al. (2018) with respect to κ ; we further specify κ in Lemma 17 as Du et al. (2018) assume that $\kappa = 1$. This result provides that the induced Gram matrix H is lower bounded by λ_0 and remains near from the Gram matrix $\mathbf{H}(0)$.

Lemma 6 (Variant of Lemma 3.2 in Du et al. (2018)). *Suppose that $\mathbf{w}_1, \dots, \mathbf{w}_m$ are independently generated from $\mathcal{N}(0, \kappa^2 \mathbf{I})$ and $m = \Omega(\frac{n^2}{\lambda_0^2} \log(\frac{n}{\delta}))$. Then, with probability at least $1 - \delta$, the following holds. For any set of weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_m \in \mathbb{R}^d$ that satisfies $\|\mathbf{w}_r(0) - \mathbf{w}_r\| \leq \frac{c\kappa\delta\lambda_0}{n^2} := R$ for any $r \in \{m\}$, some positive constant c , then the matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$ whose p, q -th entry is defined by*

$$H_{pq} := \frac{1}{m} \mathbf{x}_p^\top \mathbf{x}_q \sum_{r=1}^m [\mathbb{1}\{\mathbf{w}_r^\top \mathbf{x}_p \geq 0, \mathbf{w}_r^\top \mathbf{x}_q \geq 0\}] \quad (51)$$

satisfies $\|\mathbf{H} - \mathbf{H}(0)\|_2 < \frac{\lambda_0}{4}$ and $\lambda_{\min}(\mathbf{H}) > \frac{\lambda_0}{2}$, where $\mathbf{H}(0)$ is defined in equation 152 and $\lambda_{\min}(\mathbf{H})$ is the smallest eigenvalue of H .

Proof of Lemma 17. The following event is defined as

$$\mathcal{E}_{qr} := \{\exists \mathbf{w} : \|\mathbf{w} - \mathbf{w}_r(0)\| \leq R, \mathbb{1}\{\mathbf{x}_q^\top \mathbf{w}_r(0) \geq 0\} \neq \mathbb{1}\{\mathbf{x}_q^\top \mathbf{w} \geq 0\}\}. \quad (52)$$

This event happens if and only if $|\mathbf{w}_r(0)^\top \mathbf{x}_q| \leq R$. Note that $\mathbf{w}_r(0)$ follows $\mathcal{N}(0, \kappa^2 \mathbf{I}_d)$. For $q \in \{n\}$, we get

$$P(\mathcal{E}_{qr}) = P_{h \sim \mathcal{N}(0,1)}(|h| \leq R) \leq \frac{2R}{\sqrt{2\pi\kappa}}. \quad (53)$$

Then, for any $(p, q) \in \{n\}^2$, it follows that

$$\begin{aligned} \mathbb{E}[|H_{pq}(0) - H_{pq}|] &= \mathbb{E}\left[\frac{1}{m} |\mathbf{x}_p^\top \mathbf{x}_q \sum_{r=1}^m (\mathbb{1}\{\mathbf{w}_r(0)^\top \mathbf{x}_p \geq 0, \mathbf{w}_r(0)^\top \mathbf{x}_q \geq 0\} - \mathbb{1}\{\mathbf{w}_r^\top \mathbf{x}_p \geq 0, \mathbf{w}_r^\top \mathbf{x}_q \geq 0\})|\right] \\ &\leq \frac{1}{m} \sum_{r=1}^m \mathbb{E}[\mathbb{1}\{\mathcal{E}_{pr} \cup \mathcal{E}_{qr}\}] \leq \frac{4R}{\sqrt{2\pi\kappa}}. \end{aligned} \quad (54)$$

By summing equation 155 over (p, q) ,

$$\mathbb{E}\left[\sum_{pq} |H_{pq}(0) - H_{pq}|\right] \leq \frac{4n^2 R}{\sqrt{2\pi\kappa}}.$$

By Markov's inequality, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\sum_{pq} |H_{pq}(0) - H_{pq}| \leq \frac{4n^2 R}{\sqrt{2\pi\kappa\delta}}.$$

Then,

$$\|\mathbf{H} - \mathbf{H}(0)\|_2 \leq \sum_{pq} |H_{pq}(0) - H_{pq}| \leq \frac{4n^2 R}{\sqrt{2\pi\kappa\delta}}. \quad (55)$$

Finally, we can obtain a lower bound of the smallest eigenvalue of \mathbf{H} ($\lambda_{\min}(\mathbf{H})$) by plugging in R and using Lemma 5 as follows.

$$\lambda_{\min}(\mathbf{H}) \geq \lambda_{\min}(\mathbf{H}(0)) - \|\mathbf{H} - \mathbf{H}(0)\|_2 \geq \lambda_{\min}(\mathbf{H}(0)) - \frac{4n^2 R}{\sqrt{2\pi\kappa\delta}} \geq \frac{\lambda_0}{2} \quad (56)$$

□

The following lemma (i.e. Lemma 18) is a direct extension of Lemma 4.1 in Du et al. (2018) with respect to κ ; we further specify κ in Lemma 18 as Du et al. (2018) assume $\kappa = 1$. We include the proof of Lemma 18 for completeness.

Lemma 7 (Variant of Lemma 4.1 in Du et al. (2018)). *Let $S_q := \{r \in \{m\} : \mathbb{1}\{\mathcal{E}_{qr}\} = 0\}$ and $(S_q)^\perp := \{m\} \setminus S_q$, where \mathcal{E}_{qr} is defined in equation 154. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, we have $\sum_{q=1}^n |(S_q)^\perp| \leq \frac{CmnR}{\delta\kappa}$ for some positive constant $C > 0$.*

Proof of Lemma 18. Note that

$$\mathbb{E}[|(S_q)^\perp|] = \sum_{r=1}^m P(\mathcal{E}_{qr}) \leq \frac{2mR}{\sqrt{2\pi\kappa}}, \quad (57)$$

where the inequality follows from (53). Then,

$$\mathbb{E}[\sum_q |(S_q)^\perp|] \leq \frac{2mnR}{\sqrt{2\pi\kappa}}, \quad (58)$$

and by Markov's inequality, with probability at least $1 - \delta$,

$$\sum_q |(S_q)^\perp| \leq \frac{CmnR}{\delta\kappa}. \quad (59)$$

□

By using Lemmas 5, 17, and 18, we prove the following theorem (i.e. Theorem F.1). Note that Theorem F.1 is a direct extension of Theorem 4.1 in Du et al. (2018) with respect to κ (from $\kappa = 1$ to $\kappa = \Theta(1)$). In the proof of Theorem F.1, we also add that there exists no contradiction in Theorem 4.1 in Du et al. (2018)

Theorem B.2. (Modification of Theorem 4.1 in Du et al. (2018)) Suppose that $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$. The DNN parameter $\mathbf{W}(k)$ is optimized via the gradient descent with the step size $\eta = O(\frac{\lambda_0}{n^2})$. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for $k \in \{0, 1, 2, \dots\}$,

$$\|\mathbf{y} - \mathbf{u}(k)\|^2 \leq \left(1 - \frac{\eta\lambda_0}{2}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|^2. \quad (60)$$

Proof of Theorem F.1. Using Theorem B.1 and the fact that the set of all conditions of Theorem F.1 implies that of Theorem B.1, we can obtain the fact that κ should not be $o(1)$ for n , in order to prove that equation 161 holds for all $k \geq 0$. Therefore, we assume that $\kappa = \Theta(1)$ and prove that equation 161 holds for all $k \geq 0$ in this case.

We first prove that there exists no contradiction when $\kappa = \Theta(1)$ by assuming that equation 161 holds for all $k \geq 0$. Suppose that \hat{k} is any integer satisfying

$$\hat{k} \log\left(1 - \frac{\eta\lambda_0}{2}\right) \leq \log\left(\frac{\epsilon}{\|\mathbf{y} - \mathbf{u}(0)\|^2}\right), \quad (61)$$

where ϵ is arbitrary small constant invariant of n . As $-\eta\lambda_0(2 - \eta\lambda_0)^{-1} \leq \log(1 - \frac{\eta\lambda_0}{2})$, the following condition implies equation 219.

$$\hat{k} > ((\eta\lambda_0)/(2 - \eta\lambda_0))^{-1} \log(\|\mathbf{y} - \mathbf{u}(0)\|^2 / \epsilon) \quad (62)$$

From $\eta = O(\frac{\lambda_0}{n^2})$, equation 220 is implied by

$$\hat{k} = \Omega\left(\frac{n^2}{\lambda_0^2} \log\left(\frac{\|\mathbf{y} - \mathbf{u}(0)\|^2}{\epsilon}\right)\right). \quad (63)$$

Then, as we assume that equation 161 holds for all $k \geq 0$, it follows that for any \hat{k} satisfying equation 63,

$$\|\mathbf{y} - \mathbf{u}(\hat{k})\|^2 \leq \epsilon. \quad (64)$$

This implies that the value of $\|\mathbf{y} - \mathbf{u}(\check{k})\|^2$ can be arbitrarily reduced if integer \check{k} is sufficiently large. As ϵ is arbitrary small and invariant of n , there exists a pair of (ϵ, \check{k}) such that the following conditions hold at the same time:

$$\|\mathbf{y} - \mathbf{u}(\check{k})\|^2 \leq \epsilon \text{ and } \|\mathbf{u}(\check{k})\|^2 = \Theta(\|\mathbf{y}\|^2). \quad (65)$$

On the other hand, it follows from Lemma 4 that for \check{k} satisfying equation 65, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, the integer \check{k} satisfying equation 65 should also satisfy

$$\|\mathbf{u}(\check{k})\|^2 = O\left(\frac{n^3 \max(\kappa, 1)^2}{m\lambda_0^2 \delta^2} + \frac{n\kappa^2}{\delta^2}\right) \stackrel{(a)}{=} O(\kappa^2 n), \quad (66)$$

where (a) follows from $m = \Omega(\frac{n^6}{\lambda_0^4 \delta^3})$.

Since we assume $\kappa = \Theta(1)$ for n , there exists no contradiction such that equation 65 is contrary to equation 66, whereas it can happen when we assume that $\kappa = o(1)$ for n and $\|\mathbf{y}\|^2 = \Theta(n)$ (Theorem B.1).

Now we prove that equation 161 holds for all $k \geq 0$. This proof is based on that of Theorem 4.1 in Du et al. (2018). To do this, we use the induction hypothesis. We assume that $k = 0$. Then, equation 162 holds for $k' \in \{0, 1, \dots, k\} = \{0\}$.

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq (1 - \frac{\eta \lambda_0}{2})^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2 \quad (67)$$

Next, we assume that k is an integer satisfying $k > 0$. We assume that for $k' \in \{0, 1, \dots, k\}$, it holds

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq (1 - \frac{\eta \lambda_0}{2})^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2. \quad (68)$$

The gradient descent of training loss $L(\mathbf{W})$ with respect to the parameter \mathbf{w}_r for $r \in \{m\}$ can be derived as

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{w}_r} = \frac{1}{\sqrt{m}} \sum_{q=1}^n (u_q - y_q) a_r \mathbf{x}_q \mathbb{1}\{\mathbf{w}_r^\top \mathbf{x}_q \geq 0\}. \quad (69)$$

We define the event

$$\mathcal{E}_{qr} := \{\exists \mathbf{w} : \|\mathbf{w} - \mathbf{w}_r(0)\| \leq R, \mathbb{1}\{\mathbf{x}_q^\top \mathbf{w}_r(0) \geq 0\} \neq \mathbb{1}\{\mathbf{x}_q^\top \mathbf{w} \geq 0\}\}.$$

And we define $S_q := \{r \in \{m\} : \mathbb{1}\{\mathcal{E}_{qr}\} = 0\}$, $(S_q)^\perp := \{m\} \setminus S_q$, and $R := \frac{c\kappa\delta\lambda_0}{n^2}$ for some positive constant c . Then,

$$\begin{aligned} \mathbf{u}_q(k+1) - \mathbf{u}_q(k) &= \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r (\sigma(\mathbf{w}_r(k+1)^\top \mathbf{x}_q) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q)) \\ &= \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \left(\sigma\left((\mathbf{w}_r(k) - \eta \frac{\partial L(\mathbf{W}(k))}{\partial \mathbf{w}_r(k)})^\top \mathbf{x}_q\right) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q) \right) \\ &=: I_1^q + I_2^q, \end{aligned}$$

where

$$\begin{aligned} I_1^q &:= \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r \left(\sigma\left((\mathbf{w}_r(k) - \eta \frac{\partial L(\mathbf{W}(k))}{\partial \mathbf{w}_r(k)})^\top \mathbf{x}_q\right) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q) \right) \\ I_2^q &:= \frac{1}{\sqrt{m}} \sum_{r \in (S_q)^\perp} a_r \left(\sigma\left((\mathbf{w}_r(k) - \eta \frac{\partial L(\mathbf{W}(k))}{\partial \mathbf{w}_r(k)})^\top \mathbf{x}_q\right) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q) \right). \end{aligned}$$

Then, it follows that for some positive constant C , with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\begin{aligned} |I_2^q| &\leq \frac{\eta}{\sqrt{m}} \sum_{r \in (S_q)^\perp} \left| \frac{\partial L(\mathbf{W}(k))}{\partial \mathbf{w}_r(k)}^\top \mathbf{x}_q \right| \\ &\leq \frac{\eta |(S_q)^\perp|}{\sqrt{m}} \max_{r \in \{m\}} \left\| \frac{\partial L(\mathbf{W}(k))}{\partial \mathbf{w}_r(k)} \right\| \\ &\stackrel{(a)}{\leq} \frac{\eta \sqrt{n} |(S_q)^\perp| \|\mathbf{y} - \mathbf{u}(k)\|}{m} \\ &\stackrel{(b)}{\leq} \frac{C \eta m^{3/2} R \|\mathbf{y} - \mathbf{u}(k)\|}{\delta \kappa}, \end{aligned} \quad (70)$$

where (a) follows from equation 164 and (b) follows from Lemma 18.

To analyze I_1^q , by Lemma 16 and the assumption equation 163, we obtain that $\|\mathbf{w}_r(k+1) - \mathbf{w}_r(0)\| \leq R'$ and $\|\mathbf{w}_r(k) - \mathbf{w}_r(0)\| \leq R'$ for all $r \in S_q$.

Note that $R' < R$, which is equivalent to

$$R' := \frac{4\sqrt{n} \|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m} \lambda_0} < R := \frac{c\kappa\delta\lambda_0}{n^2} \Rightarrow m = \Omega\left(\frac{n^5 \|\mathbf{y} - \mathbf{u}(0)\|^2}{\lambda_0^4 \kappa^2 \delta^2}\right).$$

Note that from Lemma 15 and the assumption $\kappa = \Theta(1)$, $\|\mathbf{y} - \mathbf{u}(0)\|^2 = O(\kappa^2 n / \delta)$ with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$. Thus, it follows that

$$m = \Omega\left(\frac{n^5 \|\mathbf{y} - \mathbf{u}(0)\|^2}{\lambda_0^4 \kappa^2 \delta^2}\right) = \Omega\left(\frac{n^6}{\lambda_0^4 \delta^3}\right). \quad (71)$$

Since $R' < R$, for $r \in S_q$,

$$\mathbb{1}\{\mathbf{w}_r(k+1)^\top \mathbf{x}_q \geq 0\} = \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_q \geq 0\}.$$

Thus,

$$\begin{aligned} I_1^q &= \frac{\eta}{m} \sum_{p=1}^n \mathbf{x}_q^\top \mathbf{x}_p (u_p(k) - y_p(k)) \sum_{r \in S_q} \mathbb{1}\{\mathbf{w}_r(k+1)^\top \mathbf{x}_q \geq 0, \mathbf{w}_r(k+1)^\top \mathbf{x}_p \geq 0\} \\ &= -\eta \sum_{p=1}^n (u_p(k) - y_p(k)) (H_{qp}(k) - H_{qp}^\perp(k)), \end{aligned}$$

where

$$\begin{aligned} H_{ij}(k) &:= \frac{1}{m} \mathbf{x}_i^\top \mathbf{x}_j \sum_{r=1}^m [\mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_i \geq 0, \mathbf{w}_r(k)^\top \mathbf{x}_j \geq 0\}] \\ H_{ij}^\perp(k) &:= \frac{1}{m} \mathbf{x}_i^\top \mathbf{x}_j \sum_{r \in (S_q)^\perp} [\mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_i \geq 0, \mathbf{w}_r(k)^\top \mathbf{x}_j \geq 0\}]. \end{aligned} \quad (72)$$

By using Lemma 18, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\|\mathbf{H}^\perp(k)\|_2 \leq \sum_{(q,p)=(1,1)}^{(n,n)} |H_{qp}^\perp(k)| \leq \frac{n \sum_{q=1}^n |(S_q)^\perp|}{m} \leq \frac{Cn^2 R}{\delta \kappa}. \quad (73)$$

By using equation 164, we get

$$\|\mathbf{u}(k+1) - \mathbf{u}(k)\|_2^2 \leq \eta^2 \sum_{q=1}^n \left(\sum_{r=1}^m \left\| \frac{\partial L(\mathbf{W}(k))}{\partial \mathbf{w}_r(k)} \right\|_2 \right)^2 \leq \eta^2 n^2 \|\mathbf{y} - \mathbf{u}(k)\|_2^2. \quad (74)$$

Note that $m = \Omega(\frac{n^6}{\lambda_0^4 \delta^3})$ in equation 166 implies the condition $m = \Omega(\frac{n^2}{\lambda_0^2} \log(\frac{n}{\delta}))$ in Lemma 17 (by using $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$). Thus, we can use Lemma 17. From Lemma 17 and the fact that $R' < R$, we get

$$\lambda_{\min}(\mathbf{H}(k)) > \frac{\lambda_0}{2}, \quad (75)$$

where $\lambda_{\min}(\mathbf{H}(k))$ is the smallest eigenvalue of $\mathbf{H}(k)$. Then, by using union bound, the following inequalities hold with probability at least $1 - \Omega(\delta)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$.

$$\begin{aligned} \|\mathbf{y} - \mathbf{u}(k+1)\|_2^2 &= \|\mathbf{y} - \mathbf{u}(k)\|_2^2 - 2(\mathbf{y} - \mathbf{u}(k))^\top (\mathbf{u}(k+1) - \mathbf{u}(k)) + \|\mathbf{u}(k+1) - \mathbf{u}(k)\|_2^2 \\ &= \|\mathbf{y} - \mathbf{u}(k)\|_2^2 - 2\eta(\mathbf{y} - \mathbf{u}(k))^\top \mathbf{H}(k)(\mathbf{y} - \mathbf{u}(k)) \\ &\quad + 2\eta(\mathbf{y} - \mathbf{u}(k))^\top \mathbf{H}^\perp(k)(\mathbf{y} - \mathbf{u}(k)) - 2(\mathbf{y} - \mathbf{u}(k))^\top \mathbf{I}_2 + \|\mathbf{u}(k+1) - \mathbf{u}(k)\|_2^2 \\ &\stackrel{(a)}{\leq} (1 - \eta\lambda_0 + \frac{2C\eta n^2 R}{\delta \kappa} + \frac{2C\eta n^{3/2} R}{\delta \kappa} + \eta^2 n^2) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \\ &\stackrel{(b)}{\leq} (1 - \eta\lambda_0 + \frac{1}{5}\lambda_0\eta + O(\frac{\lambda_0\eta}{\sqrt{n}}) + \eta^2 n^2) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \\ &\stackrel{(c)}{\leq} (1 - \eta\lambda_0 + \frac{1}{5}\lambda_0\eta + O(\frac{\lambda_0\eta}{\sqrt{n}}) + \frac{1}{5}\lambda_0\eta) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \\ &\leq (1 - \frac{\eta\lambda_0}{2}) \|\mathbf{y} - \mathbf{u}(k)\|_2^2, \end{aligned} \quad (76)$$

where $\mathbf{I}_2 := (I_2^1, \dots, I_2^n)^\top$, (a) follows from equation 170, equation 165, equation 168, and equation 169, (b) follows from the fact that $R := \frac{c\kappa\delta\lambda_0}{n^2}$ can be less than $\frac{\kappa\delta\lambda_0}{10Cn^2}$ by properly setting c , and (c) follows from the definition of step size $\eta = O(\frac{\lambda_0}{n^2})$ (i.e., η can be set less than $\lambda_0/(5n^2)$).

We can rescale δ to a constant such that the following condition equation 172 holds with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$.

$$\|\mathbf{y} - \mathbf{u}(k+1)\|_2^2 \stackrel{(b)}{\leq} \left(1 - \frac{\eta\lambda_0}{2}\right) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \quad (77)$$

Therefore, by using the induction hypothesis with equation 172, with probability at least $1 - \delta$, it follows that for $k \in \{0, 1, 2, \dots\}$,

$$\|\mathbf{y} - \mathbf{u}(k)\|^2 \leq \left(1 - \frac{\eta\lambda_0}{2}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|^2. \quad (78)$$

□

B.5 PROOF OF THEOREM 2.4

In this section, we prove Theorem 2.4. We first show some technical lemmas.

The following lemma (i.e., Lemma 19) gives an upper bound of the gap between each trained weight vector and its initialization. This is the result of fixing κ in Lemma C.1 in Arora et al. (2019a) as $\kappa = \Theta(1)$.

Lemma 8 (Specific case of Lemma C.1 in Arora et al. (2019a) and Corollary of Lemma 16). *Under same setting as Theorem F.1, i.e., $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,*

$$\|\mathbf{w}_r(k) - \mathbf{w}_r(0)\| \leq \frac{4\sqrt{n} \|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m}\lambda_0} = O\left(\frac{\kappa n}{\sqrt{m}\lambda_0\sqrt{\delta}}\right) (:= R) \quad (79)$$

Proof of Lemma 19. The condition (150) is satisfied if the conditions in Theorem F.1 hold. Then, the proof is completed by combining Lemma 16 and the fact that $\|\mathbf{y} - \mathbf{u}(0)\| = O(\frac{\kappa\sqrt{n}}{\sqrt{\delta}})$ holds with probability at least $1 - \delta$ (which is obtained from Lemma 15 and the assumption $\kappa = \Theta(1)$). □

The following lemma (i.e., Lemma 20) is the result of fixing κ in Lemma C.2 in Arora et al. (2019a) as $\kappa = \Theta(1)$. Therefore, we omit the proof of Lemma 20 as Lemma 20 is a specific case of Lemma C.2 in Arora et al. (2019a).

Lemma 9 (Specific case of Lemma C.2 in Arora et al. (2019a)). *Under same setting as Theorem F.1, i.e., $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, $m = \Omega(\frac{n^6}{\lambda_0^4 \delta^3})$, set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, it follows that with probability at least $1 - 4\delta$ over the random initialization, for all $k \geq 0$ we have*

$$\|\mathbf{H}(k) - \mathbf{H}(0)\| = O\left(\frac{n^3}{\sqrt{m}\lambda_0\delta^{3/2}}\right), \quad (80)$$

$$\|\mathbf{Z}(k) - \mathbf{Z}(0)\| = O\left(\sqrt{\frac{n^2}{\sqrt{m}\lambda_0\delta^{3/2}}}\right).$$

We also introduce Lemma C.3 in Arora et al. (2019a) as follows.

Lemma 10 (Lemma C.3 in Arora et al. (2019a)). *With probability at least $1 - \delta$, we have $\|\mathbf{H}^* - \mathbf{H}(0)\| = O\left(\frac{n\sqrt{\log(n/\delta)}}{\sqrt{m}}\right)$.*

Then, by using the above lemmas, we prove the following proposition. This result is a revision of Theorem 4.1 in Arora et al. (2019a) by removing a κ -affected value (i.e., $(1 - \eta\lambda_0)^k \frac{\sqrt{n\kappa}}{\delta}$) in the original bound given as in (33) in Arora et al. (2019a). Therefore, this proposition is our major contribution to prove Theorem 2.4.

Proposition B.4 (Modification/revision of Theorem 4.1 in Arora et al. (2019a)). *Under same setting as Theorem F.1, i.e., $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, it follows with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ that for all $k \in \{0, 1, \dots\}$,*

$$\mathbf{u}(k) - \mathbf{y} = -(I - \eta\mathbf{H}^*)^k \mathbf{y} + \mathbf{e}(k), \quad (81)$$

where

$$\|\mathbf{e}(k)\| = O\left(k\left(1 - \frac{\eta\lambda_0}{4}\right)^{k-1} \left(\frac{\eta n^{7/2}}{\sqrt{m}\lambda_0\delta^2}\right)\right). \quad (82)$$

Proof of Proposition F.1. We define $u_q(k) := f_{\mathbf{W}(k)}(\mathbf{x}_q)$ as the q th entry of $\mathbf{u}(k) := (f_{\mathbf{W}(k)}(\mathbf{x}_1), \dots, f_{\mathbf{W}(k)}(\mathbf{x}_n))^\top$. Then,

$$u_q(k+1) - u_q(k) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r [\sigma(\mathbf{w}_r(k+1)^\top \mathbf{x}_q) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q)]. \quad (83)$$

We define the event

$$\mathcal{A}_{qr} := \{|\mathbf{w}_r(0)^\top \mathbf{x}_q| \leq R\},$$

where $R = O(\frac{\kappa n}{\sqrt{m\lambda_0\sqrt{\delta}}})$. Let $S_q := \{r \in \{1 : m\} : \mathbb{1}\{\mathcal{A}_{qr}\} = 0\}$ and $S_q^\perp := \{1 : m\} \setminus S_q$.

Note that $\mathbf{w}_r(0)^\top \mathbf{x}_q$ has the same distribution as $\mathcal{N}(0, \kappa^2)$ so that

$$\mathbb{E}(\mathbb{1}\{\mathcal{A}_{qr}\}) = P_{h \sim \mathcal{N}(0, \kappa^2)}(|h| \leq R) = \int_{-R}^R \frac{1}{\sqrt{2\pi\kappa}} e^{-x^2/2\kappa^2} dx \leq \frac{2R}{\sqrt{2\pi\kappa}}.$$

Then,

$$\mathbb{E}(|S_q^\perp|) = \mathbb{E}(\sum_{r=1}^m \mathbb{1}\{\mathcal{A}_{qr}\}) \leq \frac{2mR}{\sqrt{2\pi\kappa}}$$

and

$$\mathbb{E}(\sum_{q=1}^n |S_q^\perp|) = \mathbb{E}(\sum_{q=1}^n \sum_{r=1}^m \mathbb{1}\{\mathcal{A}_{qr}\}) \leq \frac{2mnR}{\sqrt{2\pi\kappa}} = O(\frac{n^2\sqrt{m}}{\lambda_0\sqrt{\delta}}),$$

By Markov's inequality with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\sum_{q=1}^n |S_q^\perp| = O(\frac{n^2\sqrt{m}}{\lambda_0\delta^{3/2}}) \quad (84)$$

From (183), we get

$$\begin{aligned} u_q(k+1) - u_q(k) &= \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r [\sigma(\mathbf{w}_r(k+1)^\top \mathbf{x}_q) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q)] \\ &\quad + \frac{1}{\sqrt{m}} \sum_{r \in S_q^\perp} a_r [\sigma(\mathbf{w}_r(k+1)^\top \mathbf{x}_q) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q)] \end{aligned} \quad (85)$$

We denote the second term as $\dot{\epsilon}_q(k)$

$$\begin{aligned} |\dot{\epsilon}_q(k)| &= \left| \frac{1}{\sqrt{m}} \sum_{r \in S_q^\perp} a_r [\sigma(\mathbf{w}_r(k+1)^\top \mathbf{x}_q) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q)] \right| \\ &\leq \frac{1}{\sqrt{m}} \sum_{r \in S_q^\perp} \|\mathbf{w}_r(k+1) - \mathbf{w}_r(k)\| \\ &= \frac{\eta}{\sqrt{m}} \sum_{r \in S_q^\perp} \left\| \frac{\partial L_i(\mathbf{W}(k))}{\partial \mathbf{w}_r(k)} \right\| \\ &\leq \frac{\eta}{m} \sum_{r \in S_q^\perp} \sum_{j=1}^n |y_j - u_j(k)| \\ &\leq \frac{\eta\sqrt{n}|S_q^\perp|}{m} \|\mathbf{y} - \mathbf{u}(k)\|. \end{aligned} \quad (86)$$

For the first term in (186),

$$\begin{aligned}
& \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r [\sigma(\mathbf{w}_r(k+1)^\top \mathbf{x}_q) - \sigma(\mathbf{w}_r(k)^\top \mathbf{x}_q)] \\
&= \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_q \geq 0\} (\mathbf{w}_r(k+1) - \mathbf{w}_r(k))^\top \mathbf{x}_q \\
&= \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_q \geq 0\} \left(-\frac{\eta}{\sqrt{m}} \sum_{p=1}^n (u_p(k) - y_p) a_r \mathbf{x}_p \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_p \geq 0\} \right)^\top \mathbf{x}_q \\
&= -\frac{\eta}{m} \sum_{p=1}^n (u_p(k) - y_p) \mathbf{x}_p^\top \mathbf{x}_q \sum_{r \in S_q} \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_p \geq 0\} \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_q \geq 0\} + \bar{\epsilon}_q(k) \\
&= -\eta \sum_{p=1}^n (u_p(k) - y_p) H_{qp}(k) + \bar{\epsilon}_q(k), \tag{87}
\end{aligned}$$

where

$$\bar{\epsilon}_q(k) = \frac{\eta}{m} \sum_{p=1}^n (u_p(k) - y_p) \mathbf{x}_p^\top \mathbf{x}_q \sum_{r \in S_q^\perp} \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_p \geq 0\} \mathbb{1}\{\mathbf{w}_r(k)^\top \mathbf{x}_q \geq 0\}. \tag{88}$$

Then,

$$|\bar{\epsilon}_q(k)| \leq \frac{\eta \sqrt{n} |S_q^\perp|}{m} \|\mathbf{y} - \mathbf{u}(k)\|. \tag{89}$$

Combining (186), (187), (188) and (190),

$$u_q(k+1) - u_q(k) = -\eta \sum_{p=1}^n (u_p(k) - y_p) H_{qp}(k) + \dot{\epsilon}_q(k) + \bar{\epsilon}_q(k) \tag{90}$$

which gives

$$\mathbf{u}(k+1) - \mathbf{u}(k) = -\eta \mathbf{H}(k)(\mathbf{u}(k) - \mathbf{y}) + \boldsymbol{\epsilon}(k), \tag{91}$$

where $\boldsymbol{\epsilon}(k) = \dot{\boldsymbol{\epsilon}}(k) + \bar{\boldsymbol{\epsilon}}(k)$. Note that by using (185),

$$\|\boldsymbol{\epsilon}(k)\| \leq \|\boldsymbol{\epsilon}(k)\|_1 = \sum_{q=1}^n |\epsilon_q(k) + \bar{\epsilon}_q(k)| \leq \sum_{q=1}^n \frac{2\eta \sqrt{n} |S_q^\perp|}{m} \|\mathbf{y} - \mathbf{u}(k)\| = O\left(\frac{\eta n^{5/2}}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \|\mathbf{y} - \mathbf{u}(k)\|. \tag{92}$$

We rewrite (192) as

$$\mathbf{u}(k+1) - \mathbf{u}(k) = -\eta \mathbf{H}^*(\mathbf{u}(k) - \mathbf{y}) + \boldsymbol{\zeta}(k), \tag{93}$$

where $\boldsymbol{\zeta}(k) = \eta(\mathbf{H}^* - \mathbf{H}(k))(\mathbf{u}(k) - \mathbf{y}) + \boldsymbol{\epsilon}(k)$. Then, we get

$$\mathbf{u}(k) - \mathbf{y} = -(1 - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) + \sum_{t=0}^{k-1} (I - \eta \mathbf{H}^*)^t \boldsymbol{\zeta}(k-1-t). \tag{94}$$

From (193) and Lemmas 20 and 10, we bound $\boldsymbol{\zeta}(k)$ as

$$\begin{aligned}
\|\boldsymbol{\zeta}(k)\| &\leq \eta \|\mathbf{H}^* - \mathbf{H}(k)\|_2 \|\mathbf{y} - \mathbf{u}(k)\| + O\left(\frac{\eta n^{5/2}}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \|\mathbf{y} - \mathbf{u}(k)\| \\
&\leq \eta (\|\mathbf{H}(0) - \mathbf{H}(k)\| + \|\mathbf{H}^* - \mathbf{H}(0)\|) \|\mathbf{y} - \mathbf{u}(k)\| + O\left(\frac{\eta n^{5/2}}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \|\mathbf{y} - \mathbf{u}(k)\| \\
&= O\left(\frac{\eta n^3}{\sqrt{m} \lambda_0 \delta^{3/2}} + \frac{\eta n \sqrt{\log(n/\delta)}}{\sqrt{m}} + \frac{\eta n^{5/2}}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \|\mathbf{y} - \mathbf{u}(k)\| \\
&= O\left(\frac{\eta n^3}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \|\mathbf{y} - \mathbf{u}(k)\|, \tag{95}
\end{aligned}$$

where the last equality follows from the fact that $O(\frac{\eta n^3}{\sqrt{m} \lambda_0 \delta^{3/2}})$ implies $O(\frac{\eta n^{5/2}}{\sqrt{m} \lambda_0 \delta^{3/2}})$ and $O(\frac{\eta n \sqrt{\log(n/\delta)}}{\sqrt{m}})$.

Then,

$$\begin{aligned}
\left\| \sum_{t=0}^{k-1} (I - \eta \mathbf{H}^*)^t \zeta(k-1-t) \right\| &\leq \sum_{t=0}^{k-1} \|I - \eta \mathbf{H}^*\|^t \|\zeta(k-1-t)\| \\
&\leq \sum_{t=0}^{k-1} \|I - \eta \mathbf{H}^*\|^t \|\zeta(k-1-t)\| \\
&\leq \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t \|\zeta(k-1-t)\| \\
&\stackrel{(a)}{\leq} \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t O\left(\frac{\eta n^3}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \|\mathbf{y} - \mathbf{u}(k-1-t)\| \\
&\stackrel{(b)}{\leq} \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t O\left(\frac{\eta n^3}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \left(1 - \frac{\eta \lambda_0}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\| \\
&\stackrel{(c)}{\leq} \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t O\left(\frac{\eta n^3}{\sqrt{m} \lambda_0 \delta^{3/2}}\right) \left(1 - \frac{\eta \lambda_0}{4}\right)^k O\left(\frac{\sqrt{n}}{\sqrt{\delta}}\right) \\
&\leq k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} O\left(\frac{\eta n^{7/2}}{\sqrt{m} \lambda_0 \delta^2}\right), \tag{96}
\end{aligned}$$

where (a) follows from (196), (b) follows from Theorem F.1 such that

$$\|\mathbf{y} - \mathbf{u}(k)\| \leq \left(1 - \frac{\eta \lambda_0}{2}\right)^{k/2} \|\mathbf{y} - \mathbf{u}(0)\| \leq \left(1 - \frac{\eta \lambda_0}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|, \tag{97}$$

and (c) follows from Lemma 15 and the assumption $\kappa = \Theta(1)$ (i.e., $\|\mathbf{y} - \mathbf{u}(0)\| = O(\frac{\sqrt{n}}{\sqrt{\delta}})$).

By applying (197) to (195), it follows that under same setting as Theorem F.1, it follows that for $k \geq 0$, with probability at least $1 - \delta$,

$$\mathbf{u}(k) - \mathbf{y} = -(I - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) + \mathbf{e}(k), \tag{98}$$

where

$$\|\mathbf{e}(k)\| = O\left(k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} \left(\frac{\eta n^{7/2}}{\sqrt{m} \lambda_0 \delta^2}\right)\right). \tag{99}$$

□

As a simple corollary of Proposition F.1, now we can prove Theorem 2.4 as follows.

Theorem B.3. [Theorem 2.4, modification/revision of Theorem 4.1 in Arora et al. (2019a)] Suppose that all conditions in Theorem 2.1 hold. Suppose also that $\kappa = \Theta(1)$. Then, with probability at least $1 - \delta$ for $\delta \in (0, 1)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for all $k \geq 0$,

$$\frac{1}{\sqrt{n}} \|\mathbf{y} - \mathbf{u}(k)\| = \sqrt{\frac{1}{n} \sum_{i=1}^n (1 - \eta \lambda_i)^{2k} \left(v_i^\top (\mathbf{y} - \mathbf{u}(0))\right)^2} + O\left(\frac{n^3}{\sqrt{m} \lambda_0^2 \delta^2}\right), \tag{100}$$

where $v_1, \dots, v_n \in \mathbb{R}^n$ are orthonormal eigenvectors of \mathbf{H}^* and $\lambda_1, \dots, \lambda_n$ are the corresponding eigenvalues.

Proof of Theorem F.2. Our proof is based on that for Theorem 4.1 in Arora et al. (2019a).

From Proposition F.1, it follows that for all $k \in \{0, 1, \dots\}$,

$$\mathbf{u}(k) - \mathbf{y} = -(I - \eta \mathbf{H}^*)^k \mathbf{y} + \mathbf{e}(k), \tag{101}$$

where

$$\|\mathbf{e}(k)\| = O\left(k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} \left(\frac{\eta n^{7/2}}{\sqrt{m} \lambda_0 \delta^2}\right)\right). \tag{102}$$

Therefore, we get

$$\begin{aligned}
\|\mathbf{u}(k) - \mathbf{y}\| &\stackrel{(a)}{\leq} \left\| (\mathbf{I} - \eta \mathbf{H}^*)^k (\mathbf{u}(0) - \mathbf{y}) \right\| + \|\mathbf{e}(k)\| \\
&\stackrel{(b)}{=} \sqrt{\sum_{j=1}^n (1 - \eta \lambda_j)^{2k} (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0)))^2} + \|\mathbf{e}(k)\| \\
&\stackrel{(c)}{=} \sqrt{\sum_{j=1}^n (1 - \eta \lambda_j)^{2k} (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0)))^2} + O\left(k(1 - \frac{\eta \lambda_0}{4})^{k-1} \frac{\eta n^{7/2}}{\sqrt{m} \lambda_0 \delta^2}\right) \\
&\stackrel{(d)}{=} \sqrt{\sum_{j=1}^n (1 - \eta \lambda_j)^{2k} (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0)))^2} + O\left(\frac{1}{\eta \lambda_0} \frac{\eta n^{7/2}}{\sqrt{m} \lambda_0 \delta^2}\right),
\end{aligned}$$

where (a) follows from the triangle inequality and equation 202, (b) follows from $(\mathbf{I} - \eta \mathbf{H}^*)^k$ has the eigen-decomposition $(\mathbf{I} - \eta \mathbf{H}^*)^k = \sum_{j=1}^n (1 - \eta \lambda_j)^k \mathbf{v}_j \mathbf{v}_j^\top$ and $\mathbf{y} - \mathbf{u}(0) = \sum_{j=1}^n (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0))) \mathbf{v}_j$, (c) follows from equation 203, (d) follows from $\max_{k \geq 0} \{k(1 - \eta \lambda_0/4)^{k-1}\} = O(1/(\eta \lambda_0))$. \square

B.6 PROOF OF THEOREM 2.5

B.6.1 BACKGROUND ON RADEMACHER COMPLEXITY

Before we prove Theorem 2.5 stated in Section 2.3, we introduce Rademacher Complexity and the theorem derived from it.

Define a loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$. For a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we define the population loss over true model distribution \mathcal{D} and the empirical loss over n samples $S = \{(\mathbf{x}_j, y_j)\}_{j=1}^n$ from \mathcal{D} , respectively, as

$$\begin{aligned}
\mathcal{L}_{\mathcal{D}}(f) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f(\mathbf{x}), y)] \\
\mathcal{L}_S(f) &= \sum_{j=1}^n [\ell(f(\mathbf{x}_j), y_j)].
\end{aligned} \tag{103}$$

Then, Rademacher complexity of a function class \mathcal{F} mapping \mathbb{R}^d to \mathbb{R} is expressed as

$$\mathcal{R}_S(\mathcal{F}) := \frac{1}{n} \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \sum_{j=1}^n \epsilon_j f(\mathbf{x}_j) \right], \tag{104}$$

where $\epsilon = (\epsilon_1, \dots, \epsilon_n)^\top \in \mathbb{R}^n$ includes i.i.d. random variables $\epsilon_j \sim \text{unif}(\{1, -1\})$ for $j \in \{1, \dots, n\}$. This provides an upper bound of generalization error as the following theorem given from Mohri et al. (2018).

Theorem B.4. Suppose the α -Lipschitz loss function $\ell(\cdot, \cdot)$ is bounded in $[0, \beta]$ in the first argument. Then, with probability at least $1 - \delta$ over sample S of size n ,

$$\sup_{f \in \mathcal{F}} \{\mathcal{L}_{\mathcal{D}}(f) - \mathcal{L}_S(f)\} \leq 2\alpha \mathcal{R}_S(\mathcal{F}) + 3\beta \sqrt{\frac{\log(2/\delta)}{2n}}. \tag{105}$$

B.6.2 PROOF OF THEOREM 2.5

In this section, we now prove Theorem 2.5 stated in Section 2.3. We first show some technical lemmas.

As a result of Lemma 19, we can obtain the following upper bound of the magnitude of trained NN output, when the neural network is over-parameterized.

Lemma 11. Let the input data $\{\mathbf{x}_j\}_{j=1}^n$ and label data $\{y_j\}_{j=1}^n$ of n training samples independently follow model distribution $\mathcal{D}(\mathbf{x}, y)$. We consider the same setting as Theorem F.1, i.e., $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{1, \dots, n\}} \|\mathbf{x}_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$. Then, for input sample \mathbf{x} obtained from $\mathcal{D}(\mathbf{x}, y)$ and for every $k \geq 0$, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$|f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x})| = O\left(\frac{\kappa}{\delta}\right). \tag{106}$$

Proof of Lemma 11. The proof is similar with that for Lemma 4. Define a set $\Gamma(\mathbf{W}(0), R') := \{\check{\mathbf{W}} = [\check{\mathbf{w}}_1, \dots, \check{\mathbf{w}}_m] \in \mathbb{R}^{m \times d} \mid \max_{r \in \{m\}} \|\check{\mathbf{w}}_r - \mathbf{w}_r(0)\| \leq R'\}$. Then, for any matrix $\check{\mathbf{W}} = [\check{\mathbf{w}}_1, \dots, \check{\mathbf{w}}_m] \in \mathbb{R}^{m \times d}$ belonging to $\Gamma(\mathbf{W}(0), R')$, it follows that for any $\mathbf{x} \in \mathbb{R}^d$,

$$\begin{aligned} \mathbb{E}_{\mathbf{a}}[f_{\check{\mathbf{W}}}(\mathbf{x})^2] &= \mathbb{E}_{\mathbf{a}} \left[\frac{1}{m} \left(\sum_{r \in \{m\}} a_r \sigma(\check{\mathbf{w}}_r^\top \mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{a}} \left[\frac{1}{m} \left(\sum_{r \in \{m\}} \sigma(\check{\mathbf{w}}_r^\top \mathbf{x})^2 + \sum_{r, r' \in \{m\} \times \{m\}, r \neq r'} a_r a_{r'} \sigma(\check{\mathbf{w}}_r^\top \mathbf{x}) \sigma(\check{\mathbf{w}}_{r'}^\top \mathbf{x}) \right) \right] \\ &\stackrel{(a)}{=} \frac{1}{m} \left(\sum_{r \in \{m\}} \sigma(\check{\mathbf{w}}_r^\top \mathbf{x})^2 \right) + \frac{1}{m} \left(\sum_{r, r' \in \{m\} \times \{m\}, r \neq r'} \mathbb{E}_{\mathbf{a}}[a_r a_{r'}] \sigma(\check{\mathbf{w}}_r^\top \mathbf{x}) \sigma(\check{\mathbf{w}}_{r'}^\top \mathbf{x}) \right) \\ &= \frac{1}{m} \left(\sum_{r \in \{m\}} \sigma(\check{\mathbf{w}}_r^\top \mathbf{x})^2 \right), \end{aligned} \quad (107)$$

where (a) follows from $\check{\mathbf{w}}_r$ and $\check{\mathbf{w}}_{r'}$ are independent of the random vector \mathbf{a} (i.e., $\check{\mathbf{w}}_r$ and $\check{\mathbf{w}}_{r'}$ are only depending on $\mathbf{W}(0)$ and R' as $\check{\mathbf{W}}$ is an arbitrary matrix satisfying $\check{\mathbf{W}} \in \Gamma(\mathbf{W}(0), R')$). Thus, by using Markov's inequality, we obtain with probability at least $1 - \delta$ over the random initialization of \mathbf{a} ,

$$f_{\check{\mathbf{W}}}(\mathbf{x})^2 \leq \frac{1}{\delta m} \sum_{r=1}^m \sigma(\check{\mathbf{w}}_r^\top \mathbf{x})^2. \quad (108)$$

Using equation 108, the following inequalities hold with probability at least $1 - \Omega(\delta)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\begin{aligned} f_{\check{\mathbf{W}}}(\mathbf{x})^2 &\stackrel{(a)}{\leq} \frac{1}{\delta m} \sum_{r=1}^m \sigma((\check{\mathbf{w}}_r - \mathbf{w}_r(0) + \mathbf{w}_r(0))^\top \mathbf{x})^2 \\ &\leq \frac{1}{\delta m} \sum_{r=1}^m (\|\check{\mathbf{w}}_r - \mathbf{w}_r(0)\|^2) + \frac{1}{\delta m} \sum_{r=1}^m \left| \mathbf{w}_r(0)^\top \frac{\mathbf{x}}{\|\mathbf{x}\|} \right|^2 \\ &\stackrel{(b)}{\leq} \frac{R'^2}{\delta} + \frac{1}{\delta m} \sum_{r=1}^m \left| \mathbf{w}_r(0)^\top \frac{\mathbf{x}}{\|\mathbf{x}\|} \right|^2 \\ &\stackrel{(c)}{\leq} \frac{R'^2}{\delta} + \frac{1}{\delta m} \sum_{r=1}^m \|\mathbf{w}_r(0)\|^2 \\ &\stackrel{(d)}{\leq} \frac{R'^2}{\delta} + \frac{\kappa^2}{\delta^2}, \end{aligned} \quad (109)$$

where (a) follows from equation 108, (b) follows from the fact that $\check{\mathbf{W}}$ belongs to $\Gamma(\mathbf{W}(0), R)$ (i.e., $\|\check{\mathbf{w}}_r - \mathbf{w}_r(0)\| \leq R$), (c) follows from Cauchy-Schwarz inequality, and (d) follows from the fact that $\mathbb{E}[\sum_{r=1}^m \|\mathbf{w}_r(0)\|^2] = m\kappa^2$ and Markov's inequality (i.e., $\sum_{r=1}^m \|\mathbf{w}_r(0)\|^2 = m\kappa^2/\delta$ holds with probability at least δ).

On the other hand, it follows from Lemma 19 that $\mathbf{W}(k)$ belongs to $\Gamma(\mathbf{W}(0), R')$ with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, where $R' = \sqrt{\frac{c\kappa^2 n^2}{m\lambda_0^2 \delta}}$ for some constant c . This implies $\check{\mathbf{W}}$ in equation 109 can be replaced by $\mathbf{W}(k)$ (i.e., $f_{\check{\mathbf{W}}}(\mathbf{x})$ in equation 109 can be replaced by $f_{\mathbf{W}(k)}(\mathbf{x})$).

By using the union bound over the above statement (i.e., $f_{\check{\mathbf{W}}}(\mathbf{x})$ in equation 109 can be replaced by $f_{\mathbf{W}(k)}(\mathbf{x})$) and the inequality in equation 109 and setting $R' = \sqrt{\frac{c\kappa^2 n^2}{m\lambda_0^2 \delta}}$, it follows that with probability at least $1 - \Omega(\delta)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\begin{aligned} |f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x})|^2 &\stackrel{(c)}{=} O \left(\frac{1}{\delta} \left(\frac{\kappa n}{\sqrt{m\lambda_0} \sqrt{\delta}} \right)^2 + \frac{\kappa^2}{\delta^2} \right) \\ &\stackrel{(d)}{=} O \left(\frac{\kappa^2}{\delta^2} \right), \end{aligned} \quad (110)$$

The proof is completed by rescaling δ to a constant such that equation 110 holds with probability at least $1 - \delta$. \square

We can obtain the following lemma 21 by modifying Lemma 5.3 in Arora et al. (2019a). Note that Lemma 5.3 in Arora et al. (2019a) provides an upper bound of distance between trained NN weights and its initial ones. Lemma 21 is a result obtained by removing the terms affected by κ from this upper bound of Lemma 5.3 in Arora et al. (2019a). Therefore, this lemma is our major contribution to prove Theorem 2.5.

Lemma 12 (Modification/revision of Lemma 5.3 in Arora et al. (2019a)). *Consider the same setting as Theorem F.1, i.e., $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, set $\{\mathbf{x}_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|\mathbf{x}_j\| \leq 1$, $\eta = O\left(\frac{\lambda_0}{n^2}\right)$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for all $k \geq 0$*

$$\begin{aligned} & \bullet \text{ (Lemma 19) } \|\mathbf{w}_r(k) - \mathbf{w}_r(0)\| = O\left(\frac{\kappa n}{\sqrt{m} \lambda_0 \sqrt{\delta}}\right) \quad (:= R), \forall r \in \{1 : m\} \\ & \bullet \|\mathbf{W}(k) - \mathbf{W}(0)\| \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^*)^{-1} (\mathbf{y} - \mathbf{u}(0))} + O\left(\sqrt{\frac{n^2 \sqrt{\log(n/\delta)}}{\lambda_0^2 \sqrt{m}}} + \sqrt{\frac{n^3}{\sqrt{m} \lambda_0^3 \delta^{3/2}}} + \frac{n^4}{\sqrt{m} \lambda_0^3 \delta^2}\right) \end{aligned}$$

Proof of Lemma 21. The first part is proved by using Lemma 19. The rest is to prove the second part.

From Proposition F.1, we get

$$\mathbf{u}(k) - \mathbf{y} = -(\mathbf{I} - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) + \mathbf{e}(k), \quad (111)$$

where

$$\|\mathbf{e}(k)\| = O\left(k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} \left(\frac{\eta n^{7/2}}{\sqrt{m} \lambda_0 \delta^2}\right)\right). \quad (112)$$

We apply (207) to (135), which is

$$\text{vec}(\mathbf{W}(k+1)) = \text{vec}(\mathbf{W}(k)) - \eta \mathbf{Z}(k)(\mathbf{u}(k) - \mathbf{y}), \quad (113)$$

and for $k \in \{0, \dots, K-1\}$ we obtain

$$\begin{aligned} & \text{vec}(\mathbf{W}(k)) - \text{vec}(\mathbf{W}(0)) \\ &= -\sum_{k=0}^{K-1} \eta \mathbf{Z}(k)(\mathbf{u}(k) - \mathbf{y}) \\ &= \sum_{k=0}^{K-1} \eta \mathbf{Z}(k)((\mathbf{I} - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) - \mathbf{e}(k)) \\ &= \sum_{k=0}^{K-1} \eta \mathbf{Z}(0)(\mathbf{I} - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) + \sum_{k=0}^{K-1} \eta (\mathbf{Z}(k) - \mathbf{Z}(0))(\mathbf{I} - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) - \sum_{k=0}^{K-1} \eta \mathbf{Z}(k) \mathbf{e}(k). \end{aligned} \quad (114)$$

Then, we bound the first term of (210) as

$$\begin{aligned} & \left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(0)(\mathbf{I} - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) \right\|^2 \\ &= \|\mathbf{Z}(0) \mathbf{T}(\mathbf{y} - \mathbf{u}(0))\|^2 \\ &= (\mathbf{y} - \mathbf{u}(0))^\top \mathbf{T} \mathbf{H}(0) \mathbf{T} (\mathbf{y} - \mathbf{u}(0)) \\ &\leq (\mathbf{y} - \mathbf{u}(0))^\top \mathbf{T} \mathbf{H}^* \mathbf{T} (\mathbf{y} - \mathbf{u}(0)) + \|\mathbf{H}^* - \mathbf{H}(0)\|_2 \|\mathbf{y} - \mathbf{u}(0)\|_2^2 \|\mathbf{T}\|_2^2 \\ &\stackrel{(a)}{=} (\mathbf{y} - \mathbf{u}(0))^\top \mathbf{T} \mathbf{H}^* \mathbf{T} (\mathbf{y} - \mathbf{u}(0)) + O\left(\frac{n \sqrt{\log(n/\delta)}}{\sqrt{m}}\right) \left(\frac{n}{\delta}\right) \left(\frac{1}{\lambda_0}\right)^2 \\ &= (\mathbf{y} - \mathbf{u}(0))^\top \mathbf{T} \mathbf{H}^* \mathbf{T} (\mathbf{y} - \mathbf{u}(0)) + O\left(\frac{n^2 \sqrt{\log(n/\delta)}}{\delta \lambda_0^2 \sqrt{m}}\right), \end{aligned}$$

where $\mathbf{T} := \sum_{i=1}^n \sum_{k=0}^{K-1} \eta (1 - \eta \lambda_i)^k \mathbf{v}_i \mathbf{v}_i^\top = \sum_{i=1}^n \frac{1 - (1 - \eta \lambda_i)^K}{\lambda_i} \mathbf{v}_i \mathbf{v}_i^\top$ ($\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$ are orthonormal eigenvectors of \mathbf{H}^* and $\lambda_1, \dots, \lambda_n$ are the corresponding eigenvalues) and (a) follows from Lemma 10 and

Lemma 15. Then,

$$\begin{aligned}
& \left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(0) (I - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) \right\| \\
& \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top T \mathbf{H}^* T (\mathbf{y} - \mathbf{u}(0))} + O\left(\frac{n^2 \sqrt{\log(n/\delta)}}{\lambda_0^2 \sqrt{m}}\right) \\
& \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top T \mathbf{H}^* T (\mathbf{y} - \mathbf{u}(0))} + O\left(\sqrt{\frac{n^2 \sqrt{\log(n/\delta)}}{\lambda_0^2 \sqrt{m}}}\right). \tag{115}
\end{aligned}$$

By using

$$T \mathbf{H}^* T = \sum_{i=1}^n \left(\frac{1 - (1 - \eta \lambda_i)^K}{\lambda_i} \right)^2 \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \preceq \sum_{i=1}^n \frac{1}{\lambda_i} \mathbf{v}_i \mathbf{v}_i^\top = (\mathbf{H}^*)^{-1}, \tag{116}$$

we get

$$\left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(0) (I - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) \right\| \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (H^\infty)^{-1} (\mathbf{y} - \mathbf{u}(0))} + O\left(\sqrt{\frac{n^2 \sqrt{\log(n/\delta)}}{\lambda_0^2 \sqrt{m}}}\right). \tag{117}$$

We bound the second term of (210) as

$$\begin{aligned}
& \left\| \sum_{k=0}^{K-1} \eta (\mathbf{Z}(k) - \mathbf{Z}(0)) (1 - \eta \mathbf{H}^*)^k (\mathbf{y} - \mathbf{u}(0)) \right\| \\
& \leq \sum_{k=0}^{K-1} \eta \|\mathbf{Z}(k) - \mathbf{Z}(0)\|_2 \left\| (I - \eta \mathbf{H}^*)^k \right\|_2 \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \leq \sum_{k=0}^{K-1} \eta \|\mathbf{Z}(k) - \mathbf{Z}(0)\|_2 \|\mathbf{I} - \eta \mathbf{H}^*\|_2^k \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \leq \eta \|\mathbf{Z}(K) - \mathbf{Z}(0)\|_2 \sum_{k=0}^{K-1} (1 - \eta \lambda_0)^k \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \stackrel{(a)}{=} O\left(\sqrt{\frac{n^2}{\sqrt{m} \lambda_0 \delta^{3/2}}}\right) \sum_{k=0}^{K-1} \eta (1 - \eta \lambda_0)^k \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \stackrel{(b)}{=} O\left(\sqrt{\frac{n^3}{\sqrt{m} \lambda_0^3 \delta^{3/2}}}\right), \tag{118}
\end{aligned}$$

where (a) follows from Lemma 20 and (b) follows from $\sum_{k=0}^{K-1} \eta (1 - \eta \lambda_0)^k = \frac{1 - (1 - \eta \lambda_0)^K}{\lambda_0} \leq \frac{1}{\lambda_0}$.

By using (208) and the fact that $\|\mathbf{Z}(k)\| \leq \sqrt{n}$, we bound the third term of (210) as

$$\begin{aligned}
& \left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(k) e(k) \right\| = O\left(\sum_{k=0}^{K-1} \eta \sqrt{n} \cdot \left[k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} \left(\frac{\eta n^{7/2}}{\sqrt{m} \lambda_0 \delta^2}\right)\right]\right) \\
& = O\left(\left(\sum_{k=0}^{K-1} k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1}\right) \left(\frac{\eta^2 n^4}{\sqrt{m} \lambda_0 \delta^2}\right)\right) \\
& \stackrel{(a)}{=} O\left(\frac{n^4}{\sqrt{m} \lambda_0^3 \delta^2}\right), \tag{119}
\end{aligned}$$

where (a) follows from the fact that $\sum_{k=0}^{K-1} k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} \leq \sum_{k=0}^{\infty} k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} = \left(\frac{4}{\eta \lambda_0}\right)^2$.

Therefore, by applying (213), (214), and (215) to equation 210,

$$\begin{aligned}
& \|\mathbf{W}(k) - \mathbf{W}(0)\| \\
& \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^*)^{-1} (\mathbf{y} - \mathbf{u}(0))} + O\left(\sqrt{\frac{n^2 \sqrt{\log(n/\delta)}}{\lambda_0^2 \sqrt{m}}} + \sqrt{\frac{n^3}{\sqrt{m} \lambda_0^3 \delta^{3/2}}} + \frac{n^4}{\sqrt{m} \lambda_0^3 \delta^2}\right)
\end{aligned}$$

□

We introduce the following lemma (i.e., Lemma 22) proved by Arora et al. (2019a). This lemma shows Rademacher complexity can be upper bounded by a term depending on the distance between the trained weight and its initial value.

Lemma 13 (Lemma 5.4 in Arora et al. (2019a)). *Given $R > 0$, we assume that the input data $\{\mathbf{x}_j\}_{j=1}^n$ is given as $\|\mathbf{x}_j\| \leq 1$ for $j \in \{1, \dots, n\}$. Consider the following function class in equation 216 with $\mathbf{W}(0) =: [\mathbf{w}_1(0), \dots, \mathbf{w}_m(0)]$*

$$\mathcal{F}_{R,B}^{\mathbf{W}(0),\mathbf{a}} := \{f_{\hat{\mathbf{W}},\mathbf{a}} : \hat{\mathbf{W}} = [\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_m], \|\hat{\mathbf{w}}_r - \mathbf{w}_r(0)\| \leq R (\forall r \in \{1, \dots, m\}), \|\hat{\mathbf{W}} - \mathbf{W}(0)\| \leq B\}. \quad (120)$$

Then it follows that with a probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, for every $B > 0$, the empirical Rademacher complexity $\mathcal{R}_S(\mathcal{F}_{R,B}^{\mathbf{W}(0),\mathbf{a}})$ based on the function class in equation 216 is bounded as

$$\begin{aligned} \mathcal{R}_S(\mathcal{F}_{R,B}^{\mathbf{W}(0),\mathbf{a}}) &:= \frac{1}{n} \mathbb{E}_{\epsilon \in \{\pm 1\}^n} \left[\sup_{f \in \mathcal{F}_{R,B}^{\mathbf{W}(0),\mathbf{a}}} \sum_{j=1}^n \epsilon_j f(\mathbf{x}_j) \right] \\ &\leq \frac{B}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{2}{\delta}}{m} \right)^{1/4} \right) + \frac{2R^2 \sqrt{m}}{\kappa} + R \sqrt{2 \log \frac{2}{\delta}}. \end{aligned}$$

Now, by using Lemmas 19, 11, and 22, we prove Theorem 2.5, which is given as Theorem F.4.

Theorem B.5. [Theorem 2.5, modification/revision of Theorem 5.1 in Arora et al. (2019a)] *Suppose that all conditions except $\lambda_0 > 0$ in Theorem 2.1 hold and we fix a failure probability $\delta \in (0, 1)$. Suppose further that $\kappa = \Theta(1)$ and $m = \tilde{\Omega}(\text{poly}(n, \lambda_0^{-1}, \delta^{-1}))$. Suppose also that $\lambda_0 > 0$ holds with probability at least $1 - \delta/3$ for n i.i.d. training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ from true model distribution \mathcal{D} . Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ and the training samples, it follows that for any $k \geq \Omega(\frac{1}{\eta \lambda_0} \log \frac{n}{\delta})$,*

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \frac{1}{2} |y - f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x})|^2 = O\left(\sqrt{\frac{2(\mathbf{y} - \mathbf{u}(0))^\top \mathbf{H}^{*-1}(\mathbf{y} - \mathbf{u}(0))}{n}}\right) + O\left(\sqrt{\frac{\log \frac{n}{\lambda_0 \delta}}{n}}\right). \quad (121)$$

Proof of Theorem F.4. We consider a loss function $\ell(a, b) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ as $\ell(a, b) = (a - b)^2/2$. We assume that this loss function $\ell(a, b)$ is α -Lipschitz in the first argument, this function is bounded in $[0, \beta]$, and α and β follow $O(1)$. We will prove that this assumption holds at the end of the proof.

Using the loss function and equation 204, we can define the population loss over true model distribution \mathcal{D} and the empirical loss over n samples \mathcal{S} , respectively, as

$$\begin{aligned} \mathcal{L}_{\mathcal{D}}(f) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(f(\mathbf{x}), y)] = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \frac{1}{2} (f(\mathbf{x}) - y)^2 \\ \mathcal{L}_{\mathcal{S}}(f) &= \sum_{j=1}^n [\ell(f(\mathbf{x}_j), y_j)] = \sum_{j=1}^n \left[\frac{1}{2} (f(\mathbf{x}_j) - y_j)^2 \right], \end{aligned} \quad (122)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes a trained neural network to be specified (i.e., $f_{\mathbf{W}(k)}$).

We assume that $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$ holds with probability at least $1 - \delta/3$. We also assume the following conditions hold: $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $\eta = O(\frac{\lambda_0}{n^2})$, and $m = \Omega(\frac{n^6}{\lambda_0^4 \delta^3})$.

With probability at least $1 - \delta/6$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, the followings hold simultaneously:

- Optimization succeeds: By using Theorem F.1 with equation 220, and the fact that $\|\mathbf{y} - \mathbf{u}(0)\| = O(\frac{\sqrt{n}}{\sqrt{\delta}})$ (which is obtained from Lemma 15 and the assumption $\kappa = \Theta(1)$), if $k = \Omega(\frac{1}{\eta \lambda_0} \log \frac{n}{\delta})$, it follows that

$$L(\mathbf{W}(k)) \leq (1 - \frac{\eta \lambda_0}{2})^k O\left(\frac{n}{\delta}\right) \leq \frac{1}{2}. \quad (123)$$

Then,

$$\begin{aligned} \mathcal{L}_{\mathcal{S}}(f_{\mathbf{W}(k)}) &:= \frac{1}{2n} \sum_{q=1}^n |f_{\mathbf{W}(k)}(\mathbf{x}_q) - y_q|^2 \\ &= \frac{1}{n} L(\mathbf{W}(k)) \\ &\stackrel{(a)}{=} O\left(\frac{1}{n}\right), \end{aligned} \quad (124)$$

where (a) follows from equation 221.

- From Lemma 21, we get $\|\mathbf{w}_r(k) - \mathbf{w}_r(0)\| = R (\forall r \in \{1 : m\})$ where $R = O\left(\frac{n}{\sqrt{m}\lambda_0\sqrt{\delta}}\right)$, and $\|\mathbf{W}(k) - \mathbf{W}(0)\| \leq B$ where $B = \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^*)^{-1} (\mathbf{y} - \mathbf{u}(0))} + O\left(\sqrt{\frac{n^2 \sqrt{\log(n/\delta)}}{\lambda_0^2 \sqrt{m}}} + \sqrt{\frac{n^3}{\sqrt{m}\lambda_0^3 \delta^{3/2}}} + \frac{n^4}{\sqrt{m}\lambda_0^3 \delta^2}\right)$. Note that $B \leq O\left(\sqrt{\frac{n}{\lambda_0}}\right)$.
- Let $B_j = j$ ($j = 1, 2, \dots$). For all i , the function class $\mathcal{F}_{R, B_j}^{\mathbf{W}(0), \mathbf{a}}$ has Rademacher complexity, which is upper bounded by

$$\mathcal{R}_S(\mathcal{F}_{R, B_j}^{\mathbf{W}(0), \mathbf{a}}) \leq \frac{B_j}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{20}{\delta}}{m}\right)^{1/4}\right) + \frac{2R^2 \sqrt{m}}{\kappa} + R \sqrt{2 \log \frac{20}{\delta}}. \quad (125)$$

Let j^* be the smallest integer such that $B \leq B_{j^*}$. Then we have $B_{j^*} \leq B + 1$ and $j^* \leq O\left(\sqrt{\frac{n}{\lambda_0}}\right)$. Note that $f_{\mathbf{W}(0), \mathbf{a}} \in \mathcal{F}_{R, B_{j^*}}^{\mathbf{W}(0), \mathbf{a}}$. And we get

$$\begin{aligned} \mathcal{R}_S(\mathcal{F}_{R, B_{j^*}}^{\mathbf{W}(0), \mathbf{a}}) &\leq \frac{B + 1}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{20}{\delta}}{m}\right)^{1/4}\right) + O(R^2 \sqrt{m}) + R \sqrt{2 \log \frac{20}{\delta}} \\ &= \frac{A}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{20}{\delta}}{m}\right)^{1/4}\right) + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\frac{\text{poly}(n, \frac{1}{\lambda_0}, \frac{1}{\delta})}{m^{1/4}}\right) + O(R^2 \sqrt{m}) + R \sqrt{2 \log \frac{20}{\delta}} \\ &= \frac{A}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{20}{\delta}}{m}\right)^{1/4}\right) + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\frac{\text{poly}(n, \frac{1}{\lambda_0}, \frac{1}{\delta})}{m^{1/4}}\right) \\ &= \frac{A}{\sqrt{2n}} + O\left(\frac{\sqrt{\sqrt{n}\lambda_0^{-1}\sqrt{n}}}{\sqrt{2n}} \left(\frac{\log \frac{20}{\delta}}{m}\right)^{1/4}\right) + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\frac{\text{poly}(n, \frac{1}{\lambda_0}, \frac{1}{\delta})}{m^{1/4}}\right) \\ &= \frac{A}{\sqrt{2n}} + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\frac{\text{poly}(n, \frac{1}{\lambda_0}, \frac{1}{\delta})}{m^{1/4}}\right), \end{aligned} \quad (126)$$

where $A = \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^*)^{-1} (\mathbf{y} - \mathbf{u}(0))}$. From the result of Rademacher complexity (Theorem F.3) and the union bound over a finite set $\{1 : B_{j^*}\}$, with probability at least $1 - \delta/6$, the following inequality holds for all $j \in \{1, 2, \dots, j^*\}$.

$$\sup_{f \in \mathcal{F}_{R, B_j}^{\mathbf{W}(0), \mathbf{a}}} \{\mathcal{L}_{\mathcal{D}}(f) - \mathcal{L}_S(f)\} \leq 2\alpha \mathcal{R}_S(\mathcal{F}_{R, B_j}^{\mathbf{W}(0), \mathbf{a}}) + O\left(\beta \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right) \quad (127)$$

By using the union bound jointly to consider equation 222, equation 126, and equation 227, we obtain the fact that with probability at least $1 - 5\delta/6$, the followings are satisfied at the same time.

$$\begin{aligned} \mathcal{L}_S(f_{\mathbf{W}(k), \mathbf{a}}) &= O\left(\frac{1}{n}\right) \\ f_{\mathbf{W}(k), \mathbf{a}} &\in \mathcal{F}_{R, B_{j^*}}^{\mathbf{W}(0), \mathbf{a}} \\ \mathcal{R}_S(\mathcal{F}_{R, B_{j^*}}^{\mathbf{W}(0), \mathbf{a}}) &= \sqrt{\frac{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^*)^{-1} (\mathbf{y} - \mathbf{u}(0))}{2n}} + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\frac{\text{poly}(n, \frac{1}{\lambda_0}, \frac{1}{\delta})}{m^{1/4}}\right) \\ \sup_{f \in \mathcal{F}_R^{\mathbf{W}(0), \mathbf{a}}} \{\mathcal{L}_{\mathcal{D}}(f) - \mathcal{L}_S(f)\} &\leq 2\alpha \mathcal{R}_S(\mathcal{F}_R^{\mathbf{W}(0), \mathbf{a}}) + O\left(\beta \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right) \end{aligned} \quad (128)$$

If the assumption that α and β follow $O(1)$ holds with probability at least $1 - \delta/6$, by using the union bound, it follows that with probability at least $1 - \delta$,

$$\begin{aligned}
& \mathcal{L}_{\mathcal{D}}(f_{\mathbf{W}^{(k), \mathbf{a}}}) \\
&= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \frac{1}{2} \left| y - f_{\mathbf{W}^{(k), \mathbf{a}}}(\mathbf{x}) \right|^2 \\
&\stackrel{(a)}{=} O\left(\frac{1}{n}\right) + O\left(\alpha \mathcal{R}_S(\mathcal{F}_{R, B_{j^*}}^{\mathbf{W}^{(0), \mathbf{a}}})\right) + O\left(\beta \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right) \\
&\stackrel{(b)}{=} O\left(\frac{1}{n}\right) + O\left(\mathcal{R}_S(\mathcal{F}_{R, B_{j^*}}^{\mathbf{W}^{(0), \mathbf{a}}})\right) + O\left(\sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right) \\
&\stackrel{(c)}{=} O\left(\sqrt{\frac{(\mathbf{y} - \mathbf{u}(0))^{\top} (\mathbf{H}^*)^{-1} (\mathbf{y} - \mathbf{u}(0))}{2n}} + \frac{\text{poly}(n, \frac{1}{\lambda_0}, \frac{1}{\delta})}{m^{1/4}} + \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right), \tag{129}
\end{aligned}$$

where (a) follow from equation 228, (b) follows from the assumption that α and β follow $O(1)$ for n , and (c) follow from equation 228.

Therefore, as we assume that $m = \tilde{\Omega}(\text{poly}(n, \lambda_0^{-1}, \delta^{-1}))$, we get equation 217 from equation 229.

Now we will prove the assumption that α and β follow $O(1)$. From Lemma 11, with probability at least $1 - \delta/6$, $|f_{\mathbf{W}^{(k), \mathbf{a}}}(\mathbf{x})|$ in equation 229 follows $O(\frac{\kappa}{\delta}) = O(1)$ for every $k \geq 0$ and $\mathbf{x} \sim \mathcal{D}$. On the other hand, y follows $O(1)$ for n . This is because y is independent of n , as y is i.i.d. sample of the model $\mathcal{D}(\mathbf{x}, y)$. These imply that $|y - f_{\mathbf{W}^{(k), \mathbf{a}}}(\mathbf{x})|$ in equation 229 follows $O(1)$ for every $k \geq 0$ and $(\mathbf{x}, y) \sim \mathcal{D}$. Therefore, α and β follow $O(1)$. \square

B.7 PROOF OF COROLLARIES IN SECTION 2.2

In this section, we prove Corollaries 2.2, 2.3, and 2.4, which are given sequentially as follows.

Proof of Corollary 2.2. Theorem 2.3 suggests that Theorem 2.1 does not hold, if $\kappa = O(n^\alpha)$ holds for some constant $\alpha < 0$ and $\lambda_0 = O(n) > 0$ holds. It is because condition $m = \Omega(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}) = \Omega(n^{6-4\gamma-2\alpha})$ in Theorem 2.1 implies condition $m = \Omega(n^{3-2\gamma})$ in Theorem 2.3 if $\kappa = O(n^\alpha)$ holds for some constant $\alpha < 0$ and $\lambda_0 = O(n) > 0$ holds. \square

Proof of Corollary 2.3. In order for the error term κ/δ in equation 6 in Corollary 2.1 to converges to zero as n increases, the condition $\kappa = o(\delta) = o(1)$ for n should be satisfied. That is, κ should follow $o(1)$ for n in order for Corollary 2.1 to guarantee that the training error converges to zero. However, Corollary 2.1 does not hold under this condition of κ if $\lambda_0 = O(n) > 0$ holds. This is because Corollary 2.2 implies that Corollary 2.1 does not hold if $\kappa = o(1)$ for n holds and $\lambda_0 = O(n) > 0$ holds, as Corollary 2.1 is derived from Theorem 2.1. Therefore, if $\lambda_0 = O(n) > 0$ holds, there exists no instance of κ satisfying both zero convergence of training error and correctness. \square

Proof of Corollary 2.4. In order for the error term $\frac{\sqrt{n}\kappa}{\lambda_0 \delta}$ in equation 7 in Corollary 2.2 to converge to zero as n increases, the condition $\kappa = o(\frac{\lambda_0 \delta}{\sqrt{n}}) = o(\frac{\lambda_0}{\sqrt{n}})$ should hold. However, Corollary 2.2 implies that Corollary 2.2 does not hold if $\kappa = o(1)$ for n holds and $\lambda_0 = O(n) > 0$ holds, as Corollary 2.2 is derived from Theorem 2.1. As the condition $\kappa = o(\lambda_0/\sqrt{n})$ implies $\kappa = o(1)$ for n if $\lambda_0 = O(\sqrt{n})$, Corollary 2.2 fails to guarantee both of zero generalization error and correctness if $\lambda_0 = O(\sqrt{n})$ holds. \square

B.8 EXPERIMENT DETAILS: FIGURE 1

All experiments are executed under the Tensorflow environment with NVIDIA Titan RTX GPU. Note that λ_0 denotes the minimum eigenvalue of \mathbf{H}^* , where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$ is set of n input training samples. To plot Figure 1, we randomly select the same number of training samples for each class from 10 categories. We repeat this task 500 times; we showed the average value of λ_0 as the blue line, and shaded around this line by using the minimum/maximum values as borders.

C ADDITIONAL NOTATIONS AND SETUP FOR PROOF OF THEOREM 3.1

The key idea for proof of Theorem 3.1 is similar to the proof of Theorem 2.5, but we provide the full derivations for completeness of the paper.

C.1 ADDITIONAL NOTATIONS AND SETUP

For notational simplicity, we let $d = d_h$ and $l = d_o$.

Additional notations. We are given sets $x = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ and $y = [y_1, \dots, y_n] \in \mathbb{R}^{l \times n}$ of n input and label training samples, $\{x_j\}_{j=1}^n$ and $\{y_j\}_{j=1}^n$, respectively. For $i \in \{l\}$, we denote by $g_i \in \mathbb{R}^n$ each row vector of training label matrix $y \in \mathbb{R}^{l \times n}$ such that $y = [y_1, \dots, y_n] = [g_1, \dots, g_l]^\top \in \mathbb{R}^{l \times n}$. The spectral norm is denoted by $\|\cdot\|_2$. Thus, (i, j) -th entry of H^* is given by

$$[H^*]_{ij} := \mathbb{E}_{w \sim \mathcal{N}(0, I_d)} [x_i^\top x_j \mathbb{1}\{w^\top x_i \geq 0, w^\top x_j \geq 0\}].$$

Setup for subnetwork. As we define in Section 3, we consider the neural network $F(W) := [f_W(x_1), \dots, f_W(x_n)] \in \mathbb{R}^{l \times n}$, which has training parameter $W = [w_1, \dots, w_m] \in \mathbb{R}^{d \times m}$ and as input training set $x \in \mathbb{R}^{d \times n}$, and denote by \bar{m} an integer satisfying $m = \bar{m}l$. Define submatrix $W_i = [w_{\bar{m} \cdot (i-1)+1}, \dots, w_{\bar{m} \cdot i}] \in \mathbb{R}^{\bar{m} \times d}$ of $W = [W_1^\top, \dots, W_l^\top]^\top$ for $i \in \{l\}$.

For $i \in \{l\}$, we define function $f_{W_i}(x) : \mathbb{R}^d \mapsto \mathbb{R}$ as follows.

$$f_{W_i}(x) := \frac{1}{\sqrt{\bar{m}}} \sum_{r \in \{\bar{m} \cdot (i-1)+1 : \bar{m} \cdot i\}} a_i[r] \sigma(w_r^\top x) \quad (130)$$

Note that the weight matrix in the first layer in the neural network $f_{W_i}(x)$ includes only submatrix W_i of W and the weight matrix in the second layer in the neural network $f_{W_i}(x)$ includes only vector a_i of $a = [a_1, \dots, a_l]$. For this reason, we call $f_{W_i}(x)$ the i th subnetwork of original/two-layer NN $f_W(x) : \mathbb{R}^d \mapsto \mathbb{R}^l$ given in equation 3, as $f_W(x) = [f_{W_1}(x), \dots, f_{W_l}(x)]^\top$ holds by the definitions of $f_W(x)$ and $f_{W_i}(x)$ in equation 3 and equation 130. Thus, it follows that

$$F(W) = [(f_{W_1}(x_1), \dots, f_{W_1}(x_n))^\top, \dots, (f_{W_l}(x_1), \dots, f_{W_l}(x_n))^\top]^\top \in \mathbb{R}^{l \times n}. \quad (131)$$

For $i \in \{l\}$, we define vector $u_i \in \mathbb{R}^n$ of the i th network outputs for n training samples as

$$u_i := (f_{W_i}(x_1), \dots, f_{W_i}(x_n))^\top. \quad (132)$$

GD optimization in whole network. We recall weight matrix $W(k)$ at the k th step of gradient descent as

$$W(k+1) := W(k) - \eta \frac{\partial L(W)}{\partial W} \Big|_{W=W(k)}. \quad (133)$$

Furthermore, we define $Z(k) := \frac{\partial u(k)}{\partial \text{vec}(W(k))} \in \mathbb{R}^{md \times n}$. Thus, $Z(k)$ is derived as

$$Z(k) = \frac{1}{\sqrt{m}} \begin{pmatrix} \mathbb{I}_{1,1}(k)a_1x_1 & \dots & \mathbb{I}_{1,n}(k)a_1x_n \\ \vdots & & \vdots \\ \mathbb{I}_{m,1}(k)a_mx_1 & \dots & \mathbb{I}_{m,n}(k)a_mx_n \end{pmatrix} \in \mathbb{R}^{md \times n}, \quad (134)$$

where $\mathbb{I}_{p,q}(k) := \mathbb{1}\{x_q^\top w_p(k) \geq 0\}$.

Then, equation 138 can be expressed as

$$\text{vec}(W(k+1)) = \text{vec}(W(k)) - \eta Z(k)(u(k) - y). \quad (135)$$

GD optimization in subnetwork. By using equation 132, we define the following loss

$$L_i(W_i) := (u_i - g_i)^\top (u_i - g_i)/2. \quad (136)$$

Define weight matrix $W_i(k)$ at the k th step of GD to minimize equation 136 as

$$W_i(k+1) := W_i(k) - \eta \frac{\partial L_i(W_i)}{\partial W_i} \Big|_{W_i=W_i(k)}. \quad (137)$$

Define $\bar{w}_j^i(k)$ as the j th column vector of $W_i(k)$ such that $[\bar{w}_1^i(k), \dots, \bar{w}_{\bar{m}}^i(k)] = W_i(k)$. We define $u_i(k) \in \mathbb{R}^n$ as the vector $u_i \in \mathbb{R}^n$ obtained at the k -th step of GD in equation 137. That is,

$$u_i(k) := (f_{W_i(k)}(x_1), \dots, f_{W_i(k)}(x_n))^\top. \quad (138)$$

Furthermore, we define $Z_i := \frac{\partial u_i}{\partial \text{vec}(\mathbf{W}_i)} \in \mathbb{R}^{\bar{m}d \times n}$ and $Z_i(k) := \frac{\partial u_i(k)}{\partial \text{vec}(\mathbf{W}_i(k))} \in \mathbb{R}^{\bar{m}d \times n}$. Thus, $Z_i(k)$ is derived as

$$Z_i(k) = \frac{1}{\sqrt{\bar{m}}} \begin{pmatrix} \mathbb{I}_{1,1}^i(k) a_i[1] x_1 & \dots & \mathbb{I}_{1,n}^i(k) a_i[1] x_n \\ \dots & \dots & \dots \\ \mathbb{I}_{\bar{m},1}^i(k) a_i[\bar{m}] x_1 & \dots & \mathbb{I}_{\bar{m},n}^i(k) a_i[\bar{m}] x_n \end{pmatrix} \in \mathbb{R}^{\bar{m}d \times n}, \quad (139)$$

where $\mathbb{I}_{p,q}^i(k) := \mathbb{1}\{x_q^\top \bar{w}_p^i(k) \geq 0\}$. Then, equation 137 can be expressed as

$$\text{vec}(\mathbf{W}_i(k+1)) = \text{vec}(\mathbf{W}_i(k)) - \eta Z_i(k)(u_i(k) - g_i). \quad (140)$$

GD optimization in whole network. As we show in Section ??, we define weight matrix $\mathbf{W}(k)$ at the k th step of GD to minimize $L(\mathbf{W}) := \sum_{i=1}^l L_i(\mathbf{W}_i) = \frac{1}{2} \|\mathbf{y} - F(\mathbf{W})\|^2$ given in equation 15 as

$$\mathbf{W}(k+1) := \mathbf{W}(k) - \eta \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}}|_{\mathbf{W}=\mathbf{W}(k)}, \quad (141)$$

where $L_i(\mathbf{W}_i)$ is given in equation 136.

C.2 INITIALIZATION

We initialize parameter $\mathbf{W}(0) = [\mathbf{W}_1(0)^\top, \dots, \mathbf{W}_l(0)^\top]^\top \in \mathbb{R}^{d \times m}$ in equation 137 such that each element of $\mathbf{W}(0)$ is i.i.d. sample following $\mathcal{N}(0, \kappa^2)$. We initialize parameter $\mathbf{a} = [a_1, \dots, a_l]$ in equation 130 as $a_i[r] \sim \text{unif}(\{-1, 1\})$ for $i \in \{l\}$ and $r \in \{\bar{m} \cdot i - \bar{m} + 1 : \bar{m} \cdot i\}$, otherwise $a_i[r] = 0$.

D EQUIVALENCE BETWEEN ORIGINAL NETWORK AND SUBNETWORKS

We prove in Theorem D.1 that the original/trained NN $f_{\mathbf{W}(k)}(x)$ is equivalent to the set of l trained subnetworks $[f_{\mathbf{W}_1(k)}(x), \dots, f_{\mathbf{W}_l(k)}(x)]^\top$ for any k th step of GD.

Theorem D.1. *If $\mathbf{W}(0) = [\mathbf{W}_1(0)^\top, \dots, \mathbf{W}_l(0)^\top]^\top$, then $\mathbf{W}(k) = [\mathbf{W}_1(k), \dots, \mathbf{W}_l(k)]$ for all $k \geq 0$. That is, $f_{\mathbf{W}(k)}(x) = [f_{\mathbf{W}_1(k)}(x), \dots, f_{\mathbf{W}_l(k)}(x)]^\top$.*

Proof of Theorem D.1. It suffices to show that $\mathbf{W}(k) = [\mathbf{W}_1(k), \dots, \mathbf{W}_l(k)]$ for all k . It follows that

$$\begin{aligned} \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} &= \left[\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}_1}, \dots, \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}_d} \right] \\ &\stackrel{(a)}{=} \left[\sum_{i=1}^l \frac{\partial L_i(\mathbf{W}_i)}{\partial \mathbf{W}_1}, \dots, \sum_{i=1}^l \frac{\partial L_i(\mathbf{W}_i)}{\partial \mathbf{W}_d} \right] \\ &\stackrel{(b)}{=} \left[\frac{\partial L_1(\mathbf{W}_1)}{\partial \mathbf{W}_1}, \dots, \frac{\partial L_d(\mathbf{W}_d)}{\partial \mathbf{W}_d} \right], \end{aligned}$$

where (a) follows from the loss definition (136) and (15), and (b) follows from the fact that if $i \neq j$,

$$\frac{\partial L_i(\mathbf{W}_i)}{\partial \mathbf{W}_j} = 0,$$

given that $a_i[r] = 0$ for $i \in \{d\}$ and $r \notin \{\bar{m} \cdot i - \bar{m} + 1 : \bar{m} \cdot i\}$. Thus, for all k ,

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}}|_{\mathbf{W}=\mathbf{W}(k)} = \left[\frac{\partial L_1(\mathbf{W}_1)}{\partial \mathbf{W}_1}|_{\mathbf{W}_1=\mathbf{W}_1(k)}, \dots, \frac{\partial L_d(\mathbf{W}_d)}{\partial \mathbf{W}_d}|_{\mathbf{W}_d=\mathbf{W}_d(k)} \right]$$

so that $\mathbf{W}(k) = [\mathbf{W}_1(k), \dots, \mathbf{W}_l(k)]$. \square

E ASSUMPTION

All assumptions used in this paper were presented as follows. The assumptions used were specified in each theory as not all the assumptions were used in each theory.

Assumption 1. *There exist R and φ such that with probability at least $1 - \delta$, the following equation 142 is satisfied for each $i \in \{1 : l\}$,*

$$\frac{1}{n\bar{m}} \sum_{p \in \{1:n\}} \sum_{r \in \mathcal{M}_i} \mathbb{1}\{|w_r(0)^\top x_p| \leq R\} \leq \varphi, \quad (142)$$

where $\bar{m} = m/l$ and $\mathcal{M}_i = \{\bar{m} \cdot i - \bar{m} + 1 : \bar{m} \cdot i\}$.

Assumption 2. There exists φ' such that with probability at least $1 - \delta$, the following equation 143 is satisfied for each $i \in \{1 : l\}$,

$$\frac{1}{n\bar{m}} \sum_{p \in \{1:n\}} \sum_{r \in \mathcal{M}_i} \mathbb{1}\{|w_r(0)^\top x_p| \leq R'\} \leq \varphi', \quad (143)$$

where $\bar{m} = m/l$, $\mathcal{M}_i = \{\bar{m} \cdot i - \bar{m} + 1 : \bar{m} \cdot i\}$, and $R' = O(\frac{\kappa n}{\sqrt{m\lambda_0}})$.

Assumption 3. With probability at least $1 - \delta$, the following equation 144 is satisfied for each $i \in \{1 : l\}$,

$$\frac{1}{n\bar{m}} \sum_{p \in \{1:n\}} \sum_{r \in \mathcal{M}_i} |w_r(0)^\top x_p|^2 \leq \kappa^2, \quad (144)$$

where $\bar{m} = m/l$, $\mathcal{M}_i = \{\bar{m} \cdot i - \bar{m} + 1 : \bar{m} \cdot i\}$, and κ is a constant invariant of n or m (i.e., $\kappa = \Theta(1)$).

Assumption 4.

$$\bar{m} = \Omega\left(\frac{n^4}{\min(\lambda_0^2, \lambda_0^4)}\right)$$

Assumption 5.

$$\varphi' = O\left(\frac{\min(\lambda_0^2, \lambda_0^3)}{n^4}\right)$$

Assumption 6. Let the input data $\{x_j\}_{j=1}^n$ and label data $\{y_j\}_{j=1}^n$ of n training samples independently follow model distribution $\mathcal{D}(x, y)$. Then, for input sample x obtained from $\mathcal{D}(x, y)$ and for every $k \geq 0$, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, the following equation 145 is satisfied for each $i \in \{1 : l\}$,

$$|f_{\mathbf{W}(k)}(x)[i]| = O(1). \quad (145)$$

F PROOF OF THEOREMS 2.4 AND 2.5 WHEN THE NETWORK HAS A SCALAR OUTPUT AS $l = 1$

We simplify $\mathbf{H}_1(0)$ as $\mathbf{H}(0)$.

F.1 MODIFICATION OF THEOREM 4.1 IN DU ET AL. (2018)

In order to prove Theorem 2.5 stated in Section 2.3, we first prove Theorem F.1 in this section, since Theorem 2.5 is proved by using the result of Theorem F.1. Theorem F.1 is the result of extending the condition for κ to $\kappa = \Theta(1)$ from $\kappa = 1$, which is given in Theorem 4.1 in Du et al. (2018). Therefore, most of the proof processes for Theorem F.1 (and its technical lemmas) are already proved in Du et al. (2018); we provide them in this section for completeness.

To prove Theorem F.1, we first introduce some technical lemmas.

The following lemma provides an upper bound of the magnitude of the initial NN output.

Lemma 14. Suppose that set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$. Then, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\|\mathbf{u}(0)\|^2 = O\left(\frac{n\kappa^2}{\delta}\right). \quad (146)$$

Proof of Lemma 14. It follows that

$$\begin{aligned} \mathbb{E}_{\mathbf{a}}[\|\mathbf{u}(0)\|^2] &= \mathbb{E}_{\mathbf{a}}[\| (f_{\mathbf{W}(0)}(x_1), \dots, f_{\mathbf{W}(0)}(x_n)) \|^2] \\ &= \mathbb{E}_{\mathbf{a}}\left[\sum_{j=1}^n |f_{\mathbf{W}(0)}(x_j)|^2\right] \\ &= \mathbb{E}_{\mathbf{a}}\left[\frac{1}{m} \sum_{j=1}^n \left| \sum_{r \in \{1:m\}} a_r \sigma(w_r^\top x_j) \right|^2\right] \\ &= \frac{1}{m} \sum_{j=1}^n \sum_{r \in \{1:m\}} \left| \sigma(w_r^\top x_j) \right|^2. \end{aligned} \quad (147)$$

Furthermore, it follows that with probability at least $1 - \delta$,

$$\begin{aligned} \frac{1}{m} \sum_{j=1}^n \sum_{r \in \{1:m\}} \left| \sigma(w_r^\top x_j) \right|^2 &\leq \frac{1}{m} \sum_{j=1}^n \sum_{r \in \{1:m\}} \left| w_r^\top x_j \right|^2 \\ &\stackrel{(a)}{\leq} n\kappa^2, \end{aligned} \quad (148)$$

where (a) follows from the assumption equation 144.

By combining equation 147 and equation 148, we obtain the fact that with probability at least $1 - \delta$,

$$\mathbb{E}_{\mathbf{a}} [\|\mathbf{u}(0)\|^2] = O(n\kappa^2).$$

Therefore, by using Markov's inequality, $\|\mathbf{u}(0)\|^2 = O(n\kappa^2/\delta)$ is satisfied with probability at least $1 - \Omega(\delta)$. The proof is completed by rewriting $\Omega(\delta)$ as δ . \square

Then, by using Lemma 14, we can also obtain an upper bound of gap between the initial NN output and label as Lemma 15.

Lemma 15. Suppose that set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$. If $\|\mathbf{y}\| = O(\sqrt{n})$ is satisfied, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\|\mathbf{y} - \mathbf{u}(0)\|^2 = O\left(\frac{\max(\kappa^2, 1)n}{\delta}\right).$$

Proof of Lemma 15. It follows from Lemma 14 that with probability at least $1 - \delta$,

$$2 \max(\|\mathbf{y}\|^2, \|\mathbf{u}(0)\|^2) = O\left(\max\left(1, \frac{\kappa^2}{\delta}\right)n\right) = O\left(\frac{\max(\kappa^2, 1)n}{\delta}\right). \quad (149)$$

Then, the proof is completed by applying the following inequality to equation 149.

$$\|\mathbf{y} - \mathbf{u}(0)\|^2 \leq \|\mathbf{y}\|^2 + \|\mathbf{u}(0)\|^2 \leq 2 \max(\|\mathbf{y}\|^2, \|\mathbf{u}(0)\|^2)$$

\square

The following lemma (i.e., Lemma 16) gives an upper bound of the gap between each trained weight vector and its initialization, when the training loss is reduced by the GD optimization. This lemma is the result of extending the condition for κ to arbitrary $\kappa > 0$ from $\kappa = 1$, which is given in Corollary 4.1 in Du et al. (2018).

Lemma 16 (Variant of Corollary 4.1 in Du et al. (2018)). We are given arbitrary $\kappa > 0$. Suppose that set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$. If the following condition holds for $k' \in \{0, 1, \dots, k-1\}$,

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq \left(1 - \frac{\eta\lambda_0}{2}\right)^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2, \quad (150)$$

then for every $r \in \{m\}$,

$$\|w_r(k) - w_r(0)\| \leq \frac{4\sqrt{n} \|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m}\lambda_0} \quad (= R'), \quad (151)$$

where $w_j(k)$ is the column of $\mathbf{W}(k) =: [w_1(k), \dots, w_m(k)]$ at the k -th step of GD.

Proof of Lemma 16. Since

$$\frac{\partial L(\mathbf{W})}{\partial w_r} = \frac{1}{\sqrt{m}} \sum_{q=1}^n (u_q - y_q) a_r x_q \mathbf{1}(w_r^\top x_q \geq 0),$$

we get

$$\left\| \frac{\partial L(\mathbf{W}(k'))}{\partial w_r(k')} \right\| \leq \frac{\sqrt{n}}{\sqrt{m}} \|\mathbf{y} - \mathbf{u}(k')\|.$$

Thus, we have

$$\begin{aligned}
\|w_r(k) - w_r(0)\| &\leq \eta \sum_{k'=0}^{k-1} \left\| \frac{\partial L(\mathbf{W}(k'))}{\partial w_r(k')} \right\| \\
&\leq \eta \sum_{k'=0}^{k-1} \frac{\sqrt{n}}{\sqrt{m}} \|\mathbf{y} - \mathbf{u}(k')\| \\
&\stackrel{(a)}{\leq} \frac{\sqrt{n}}{\sqrt{m}} \sum_{k'=0}^{k-1} \eta \left(1 - \frac{\eta \lambda_0}{2}\right)^{k'/2} \|\mathbf{y} - \mathbf{u}(0)\| \\
&\leq \frac{\sqrt{n}}{\sqrt{m}} \sum_{k'=0}^{k-1} \eta \left(1 - \frac{\eta \lambda_0}{4}\right)^{k'} \|\mathbf{y} - \mathbf{u}(0)\| \\
&\leq \frac{4\sqrt{n} \|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m} \lambda_0},
\end{aligned}$$

where (a) follows from equation 150. \square

The following lemma (i.e. Lemma 17) is a direct extension of Lemma 3.2 in Du et al. (2018) with respect to κ ; we further specify κ in Lemma 17 as Du et al. (2018) assume that $\kappa = 1$. This result provides that the induced Gram matrix \mathbf{H} is lower bounded by λ_0 and remains near from the Gram matrix $\mathbf{H}(0)$.

Lemma 17 (Variant of Lemma 3.2 in Du et al. (2018)). *Define matrix $\mathbf{H}(k) \in \mathbb{R}^{n \times n}$ such that p, q -th entry of $\mathbf{H}(k)$ is given by*

$$H_{pq}(k) := \frac{1}{m} x_p^\top x_q \sum_{r=1}^m [\mathbb{1}\{w_r(k)^\top x_p \geq 0, w_r(k)^\top x_q \geq 0\}], \quad (152)$$

where $w_j(k)$ is the j th column vector of $\mathbf{W}(k)$ such that $[w_1(k), \dots, w_m(k)] = \mathbf{W}(k)$. Suppose that the assumptions 1 and 3 hold. Then, with probability at least $1 - \delta$, the following holds. For any set of weight vectors $w_1, \dots, w_m \in \mathbb{R}^d$ that satisfies $\|w_r(0) - w_r\| \leq R$ for any $r \in \{m\}$, a positive constant R , then the matrix $\mathbf{H} \in \mathbb{R}^{n \times n}$ whose p, q -th entry is defined by

$$H_{pq} := \frac{1}{m} x_p^\top x_q \sum_{r=1}^m [\mathbb{1}\{w_r^\top x_p \geq 0, w_r^\top x_q \geq 0\}] \quad (153)$$

satisfies $\|\mathbf{H} - \mathbf{H}(0)\|_2 < 2n^2\varphi$ and $\lambda_{\min}(\mathbf{H}) \geq \lambda_{\min}(\mathbf{H}(0)) - 2n^2\varphi$, where $\mathbf{H}(0)$ is defined in equation 152 and $\lambda_{\min}(\mathbf{H})$ is the smallest eigenvalue of \mathbf{H} .

Proof of Lemma 17. The following event is defined as

$$\mathcal{E}_{qr} := \{\exists w : \|w - w_r(0)\| \leq R, \mathbb{1}\{x_q^\top w_r(0) \geq 0\} \neq \mathbb{1}\{x_q^\top w \geq 0\}\}. \quad (154)$$

This event happens if and only if $|w_r(0)^\top x_q| \leq R$.

On the other hand, for any $(p, q) \in \{n\}^2$, it follows that

$$\begin{aligned}
|H_{pq}(0) - H_{pq}| &= \frac{1}{m} |x_p^\top x_q \sum_{r=1}^m (\mathbb{1}\{w_r(0)^\top x_p \geq 0, w_r(0)^\top x_q \geq 0\} - \mathbb{1}\{w_r^\top x_p \geq 0, w_r^\top x_q \geq 0\})| \\
&\leq \frac{1}{m} \sum_{r=1}^m \mathbb{1}\{\mathcal{E}_{pr} \cup \mathcal{E}_{qr}\}.
\end{aligned} \quad (155)$$

By summing equation 155 over (p, q) ,

$$\sum_{pq} |H_{pq}(0) - H_{pq}| \leq \frac{1}{m} \sum_{pq} \sum_{r=1}^m \mathbb{1}\{\mathcal{E}_{pr} \cup \mathcal{E}_{qr}\} \leq \frac{2n}{m} \sum_p \sum_{r=1}^m \mathbb{1}\{\mathcal{E}_{pr}\} = \frac{2n}{m} \sum_p \sum_{r=1}^m \mathbb{1}\{|w_r(0)^\top x_p| \leq R\}. \quad (156)$$

By applying the assumption equation 142, with probability at least $1 - \delta$, we get

$$\frac{2n}{m} \sum_p \sum_{r=1}^m \mathbb{1}\{|w_r(0)^\top x_p| \leq R\} \leq 2n^2\varphi. \quad (157)$$

Then, it follows from applying equation 157 to equation 156 that

$$\|\mathbf{H} - \mathbf{H}(0)\|_2 \leq \sum_{pq} |H_{pq}(0) - H_{pq}| \leq 2n^2\varphi. \quad (158)$$

Finally, we can obtain a lower bound of the smallest eigenvalue of $\mathbf{H}(\lambda_{\min}(\mathbf{H}))$ by plugging in equation 158 as follows

$$\lambda_{\min}(\mathbf{H}) \geq \lambda_{\min}(\mathbf{H}(0)) - \|\mathbf{H} - \mathbf{H}(0)\|_2 \geq \lambda_{\min}(\mathbf{H}(0)) - 2n^2\varphi. \quad (159)$$

□

The following lemma (i.e. Lemma 18) is a direct extension of Lemma 4.1 in Du et al. (2018) with respect to κ ; we further specify κ in Lemma 18 as Du et al. (2018) assume $\kappa = 1$. We include the proof of Lemma 18 for completeness.

Lemma 18 (Variant of Lemma 4.1 in Du et al. (2018)). *Let $S_q := \{r \in \{m\} : \mathbb{1}\{\mathcal{E}_{qr}\} = 0\}$ and $(S_q)^\perp := \{m\} \setminus S_q$, where \mathcal{E}_{qr} is defined in equation 154. Then, from the assumption equation 142, with probability at least $1 - \delta$ over the random initialization of $\mathbf{W}(0)$, we have $\sum_{q=1}^n |(S_q)^\perp| \leq nm\varphi$.*

Proof of Lemma 18. From the assumption equation 142,

$$\sum_q |(S_q)^\perp| = \sum_q \sum_{r=1}^m \mathbb{1}\{\mathcal{E}_{qr}\} \leq nm\varphi. \quad (160)$$

□

By using Lemmas 17 and 18, we prove the following theorem (i.e. Theorem F.1). Note that Theorem F.1 is a direct extension of Theorem 4.1 in Du et al. (2018) with respect to κ (from $\kappa = 1$ to $\kappa = \Theta(1)$).

Theorem F.1. (Modification of Theorem 4.1 in Du et al. (2018)) *Suppose that the assumptions 1 and 3 hold with $\varphi = O(\frac{\lambda_0}{n^{\frac{1}{2}}})$, $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega(\frac{n^2\kappa^2}{\lambda_0^2 R^2 \delta})$, set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^{\frac{1}{2}}})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$. The DNN parameter $\mathbf{W}(k)$ is optimized via the gradient descent with the step size $\eta = O(\frac{\lambda_0}{n^{\frac{1}{2}}})$. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for $k \in \{0, 1, 2, \dots\}$,*

$$\|\mathbf{y} - \mathbf{u}(k)\|^2 \leq (1 - \frac{\eta\lambda_0}{2})^k \|\mathbf{y} - \mathbf{u}(0)\|^2. \quad (161)$$

Proof of Theorem F.1. This proof is based on that of Theorem 4.1 in Du et al. (2018). To do this, we use the induction hypothesis. We assume that $k = 0$. Then, equation 162 holds for $k' \in \{0, 1, \dots, k\} = \{0\}$.

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq (1 - \frac{\eta\lambda_0}{2})^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2 \quad (162)$$

Next, we assume that k is an integer satisfying $k > 0$. We assume that for $k' \in \{0, 1, \dots, k\}$, it holds

$$\|\mathbf{y} - \mathbf{u}(k')\|^2 \leq (1 - \frac{\eta\lambda_0}{2})^{k'} \|\mathbf{y} - \mathbf{u}(0)\|^2. \quad (163)$$

The gradient descent of training loss $L(\mathbf{W})$ with respect to the parameter w_r for $r \in \{m\}$ can be derived as

$$\frac{\partial L(\mathbf{W})}{\partial w_r} = \frac{1}{\sqrt{m}} \sum_{q=1}^n (u_q - y_q) a_r x_q \mathbb{1}\{w_r^\top x_q \geq 0\}. \quad (164)$$

We define the event

$$\mathcal{E}_{qr} := \{\exists \mathbf{w} : \|\mathbf{w} - w_r(0)\| \leq R', \mathbb{1}\{x_q^\top w_r(0) \geq 0\} \neq \mathbb{1}\{x_q^\top \mathbf{w} \geq 0\}\}.$$

And we define $S_q := \{r \in \{m\} : \mathbb{1}\{\mathcal{E}_{qr}\} = 0\}$ and $(S_q)^\perp := \{m\} \setminus S_q$. Then,

$$\begin{aligned} \mathbf{u}_q(k+1) - \mathbf{u}_q(k) &= \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r (\sigma(w_r(k+1)^\top x_q) - \sigma(w_r(k)^\top x_q)) \\ &= \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \left(\sigma \left((w_r(k) - \eta \frac{\partial L(\mathbf{W}(k))}{\partial w_r(k)})^\top x_q \right) - \sigma(w_r(k)^\top x_q) \right) \\ &=: I_1^q + I_2^q, \end{aligned}$$

where

$$I_1^q := \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r \left(\sigma \left((w_r(k) - \eta \frac{\partial L(\mathbf{W}(k))}{\partial w_r(k)})^\top x_q \right) - \sigma(w_r(k)^\top x_q) \right)$$

$$I_2^q := \frac{1}{\sqrt{m}} \sum_{r \in (S_q)^\perp} a_r \left(\sigma \left((w_r(k) - \eta \frac{\partial L(\mathbf{W}(k))}{\partial w_r(k)})^\top x_q \right) - \sigma(w_r(k)^\top x_q) \right).$$

Then, it follows that for some positive constant C , with probability at least $1 - \Omega(\delta)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,

$$\begin{aligned} |I_2^q| &\leq \frac{\eta}{\sqrt{m}} \sum_{r \in (S_q)^\perp} \left| \frac{\partial L(\mathbf{W}(k))}{\partial w_r(k)}^\top x_q \right| \\ &\leq \frac{\eta |S_q|}{\sqrt{m}} \max_{r \in \{m\}} \left\| \frac{\partial L(\mathbf{W}(k))}{\partial w_r(k)} \right\| \\ &\stackrel{(a)}{\leq} \frac{\eta \sqrt{n} |S_q| \|\mathbf{y} - \mathbf{u}(k)\|}{m}, \end{aligned} \quad (165)$$

where (a) follows from equation 164.

To analyze I_1^q , by Lemma 16 and the assumption equation 163, we obtain that $\|w_r(k+1) - w_r(0)\| \leq R'$ and $\|w_r(k) - w_r(0)\| \leq R'$ for all $r \in S_q$. Note that $R' < R$, which is equivalent to

$$R' := \frac{4\sqrt{n} \|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m}\lambda_0} < R \quad \Rightarrow \quad m = \Omega\left(\frac{n \|\mathbf{y} - \mathbf{u}(0)\|^2}{\lambda_0^2 R^2}\right).$$

Note that from Lemma 15 and the assumption $\kappa = \Theta(1)$, $\|\mathbf{y} - \mathbf{u}(0)\|^2 = O(\kappa^2 n / \delta)$ with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$. Thus, it follows that

$$m = \Omega\left(\frac{n \|\mathbf{y} - \mathbf{u}(0)\|^2}{\lambda_0^2 R^2}\right) = \Omega\left(\frac{n^2 \kappa^2}{\lambda_0^2 R^2 \delta}\right). \quad (166)$$

Since $R' < R$, for $r \in S_q$,

$$\mathbb{1}\{w_r(k+1)^\top x_q \geq 0\} = \mathbb{1}\{w_r(k)^\top x_q \geq 0\}.$$

Thus,

$$\begin{aligned} I_1^q &= \frac{\eta}{m} \sum_{p=1}^n x_q^\top x_p (u_p(k) - y_p(k)) \sum_{r \in S_q} \mathbb{1}\{w_r(k+1)^\top x_q \geq 0, w_r(k+1)^\top x_p \geq 0\} \\ &= -\eta \sum_{p=1}^n (u_p(k) - y_p(k)) (H_{qp}(k) - H_{qp}^\perp(k)), \end{aligned}$$

where

$$\begin{aligned} H_{ij}(k) &:= \frac{1}{m} x_i^\top x_j \sum_{r=1}^m [\mathbb{1}\{w_r(k)^\top x_i \geq 0, w_r(k)^\top x_j \geq 0\}] \\ H_{ij}^\perp(k) &:= \frac{1}{m} x_i^\top x_j \sum_{r \in (S_q)^\perp} [\mathbb{1}\{w_r(k)^\top x_i \geq 0, w_r(k)^\top x_j \geq 0\}]. \end{aligned} \quad (167)$$

By using Lemma 18, it follows that with probability at least $1 - \delta$ over the random initialization of $\mathbf{W}(0)$,

$$\left\| \mathbf{H}^\perp(k) \right\|_2 \leq \sum_{(q,p)=(1,1)}^{(n,n)} |\mathbf{H}_{qp}^\perp(k)| \leq \frac{n \sum_{q=1}^n |(S_q)^\perp|}{m} \leq n^2 \varphi. \quad (168)$$

By using equation 164, we get

$$\|\mathbf{u}(k+1) - \mathbf{u}(k)\|_2^2 \leq \eta^2 \sum_{q=1}^n \left(\sum_{r=1}^m \left\| \frac{\partial L(\mathbf{W}(k))}{\partial w_r(k)} \right\|_2 \right)^2 \leq \eta^2 n^2 \|\mathbf{y} - \mathbf{u}(k)\|^2. \quad (169)$$

Now, using Lemma 17, the fact that $R' < R$, and the assumption that $\varphi = O(\frac{\lambda_0}{n^2}) < \frac{\lambda_0}{4n^2}$, we get

$$\lambda_{\min}(\mathbf{H}(k)) \geq \lambda_0 - 2n^2 \varphi > \frac{\lambda_0}{2}, \quad (170)$$

where $\lambda_{\min}(\mathbf{H}(k))$ is the smallest eigenvalue of $\mathbf{H}(k)$. Then, by using union bound, the following inequalities hold with probability at least $1 - \Omega(\delta)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$.

$$\begin{aligned}
\|\mathbf{y} - \mathbf{u}(k+1)\|_2^2 &= \|\mathbf{y} - \mathbf{u}(k)\|_2^2 - 2(\mathbf{y} - \mathbf{u}(k))^\top (\mathbf{u}(k+1) - \mathbf{u}(k)) + \|\mathbf{u}(k+1) - \mathbf{u}(k)\|_2^2 \\
&= \|\mathbf{y} - \mathbf{u}(k)\|_2^2 - 2\eta(\mathbf{y} - \mathbf{u}(k))^\top \mathbf{H}(k)(\mathbf{y} - \mathbf{u}(k)) \\
&\quad + 2\eta(\mathbf{y} - \mathbf{u}(k))^\top \mathbf{H}^\perp(k)(\mathbf{y} - \mathbf{u}(k)) - 2(\mathbf{y} - \mathbf{u}(k))^\top \mathbf{I}_2 + \|\mathbf{u}(k+1) - \mathbf{u}(k)\|_2^2 \\
&\stackrel{(a)}{\leq} (1 - \eta\lambda_0 + 2\eta n^2\varphi + 2\eta n\sqrt{n}\varphi + \eta^2 n^2) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \\
&\stackrel{(b)}{\leq} (1 - \eta\lambda_0 + \frac{1}{5}\lambda_0\eta + O(\frac{\lambda_0\eta}{\sqrt{n}}) + \eta^2 n^2) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \\
&\stackrel{(c)}{\leq} (1 - \eta\lambda_0 + \frac{1}{5}\lambda_0\eta + O(\frac{\lambda_0\eta}{\sqrt{n}}) + \frac{1}{5}\lambda_0\eta) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \\
&\leq (1 - \frac{\eta\lambda_0}{2}) \|\mathbf{y} - \mathbf{u}(k)\|_2^2, \tag{171}
\end{aligned}$$

where $\mathbf{I}_2 := (\mathbf{I}_2^1, \dots, \mathbf{I}_2^n)^\top$, (a) follows from equation 170, equation 165 (Applying Lemma 18, $\sum_q \frac{\eta\sqrt{n}|(S_q)^\perp| \|\mathbf{y} - \mathbf{u}(k)\|}{m} \leq \eta n\sqrt{n}\varphi \|\mathbf{y} - \mathbf{u}(k)\|$), equation 168, and equation 169, (b) follows from the assumption that $\varphi = O(\frac{\lambda_0}{n^2}) < \frac{\lambda_0}{10n^2}$, and (c) follows from the definition of step size $\eta = O(\frac{\lambda_0}{n^2})$ (i.e., η can be set less than $\lambda_0/(5n^2)$). We can rescale δ to a constant such that the following condition equation 172 holds with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$.

$$\|\mathbf{y} - \mathbf{u}(k+1)\|_2^2 \stackrel{(b)}{\leq} (1 - \frac{\eta\lambda_0}{2}) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \tag{172}$$

Therefore, by using the induction hypothesis with equation 172, with probability at least $1 - \delta$, it follows that for $k \in \{0, 1, 2, \dots\}$,

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq (1 - \frac{\eta\lambda_0}{2})^k \|\mathbf{y} - \mathbf{u}(0)\|_2^2. \tag{173}$$

□

F.2 PROOF OF THEOREM 2.4 WHEN THE NETWORK HAS A SCALAR OUTPUT AS $l = 1$

In this section, we prove Theorem 2.4. We first show some technical lemmas.

The following lemma (i.e., Lemma 19) gives an upper bound of the gap between each trained weight vector and its initialization. This is the result of fixing κ in Lemma C.1 in Arora et al. (2019a) as $\kappa = \Theta(1)$.

Lemma 19 (Specific case of Lemma C.1 in Arora et al. (2019a) and Corollary of Lemma 16). *Under same setting as Theorem F.1, i.e., the assumptions 1 and 3 hold with $\varphi = O(\frac{\lambda_0}{n^2})$, $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega(\frac{n^2\kappa^2}{\lambda_0^2 R^2 \delta})$, set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, it follows that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$,*

$$\|w_r(k) - w_r(0)\| \leq \frac{4\sqrt{n} \|\mathbf{y} - \mathbf{u}(0)\|}{\sqrt{m}\lambda_0} = O\left(\frac{\kappa n}{\sqrt{m}\lambda_0\sqrt{\delta}}\right) \quad (:= R') \tag{174}$$

Proof of Lemma 19. The condition (150) is satisfied if the conditions in Theorem F.1 hold. Then, the proof is completed by combining Lemma 16 and the fact that $\|\mathbf{y} - \mathbf{u}(0)\| = O(\frac{\kappa\sqrt{n}}{\sqrt{\delta}})$ holds with probability at least $1 - \delta$ (which is obtained from Lemma 15 and the assumption $\kappa = \Theta(1)$). □

The following lemma (i.e., Lemma 20) is the result of fixing κ in Lemma C.2 in Arora et al. (2019a) as $\kappa = \Theta(1)$. Therefore, we omit the proof of Lemma 20 as Lemma 20 is a specific case of Lemma C.2 in Arora et al. (2019a).

Lemma 20 (Modification of Lemma C.2 in Arora et al. (2019a)). *Let the assumption 2 hold. Under same setting as Theorem F.1, i.e., the assumptions 1 and 3 hold with $\varphi = O(\frac{\lambda_0}{n^2})$, $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega(\frac{n^2\kappa^2}{\lambda_0^2 R^2 \delta})$, set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, it follows that with probability at least $1 - \delta$ over the random*

initialization, for all $k \geq 0$ we have

$$\begin{aligned}\|\mathbf{H}(k) - \mathbf{H}(0)\| &= O\left(\frac{n^2\delta}{m} + n^2\varphi'\right), \\ \|\mathbf{Z}(k) - \mathbf{Z}(0)\| &= O\left(\sqrt{\frac{n\delta}{m}} + n\varphi'\right).\end{aligned}\tag{175}$$

Proof of Lemma 20. Define $R' = O(\frac{\kappa n}{\sqrt{m}\lambda_0\sqrt{\delta}})$. From Lemma 19, we obtain that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, we have $\|w_r(k) - w_r(0)\| \leq R'$ for all $r \in \{1 : m\}$ and $k \geq 0$. On the other hand, the following event is defined as

$$\mathcal{A}_{rq} := \{|w_r(0)^\top x_q| \leq R'\}, \quad q \in \{n\}, r \in \{m\}.\tag{176}$$

Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for all r, p

$$\mathbb{I}\{\mathbb{I}_{r,p}(k) \neq \mathbb{I}_{r,p}(0)\} = \mathbb{I}\{\mathcal{A}_{rp}\}.$$

It is because the following event \mathcal{E}_{qr} happens if and only if $|w_r(0)^\top x_q| \leq R'$.

$$\mathcal{E}_{qr} := \{\exists \mathbf{w} : \|\mathbf{w} - w_r(0)\| \leq R', \mathbb{I}\{x_q^\top w_r(0) \geq 0\} \neq \mathbb{I}\{x_q^\top \mathbf{w} \geq 0\}\}$$

Then, it follows that for any $p, q \in \{n\}$,

$$\begin{aligned}|H_{pq}(k) - H_{pq}(0)| &= \left| \frac{x_p^\top x_q}{m} \sum_{r=1}^m (\mathbb{I}_{r,p}(k)\mathbb{I}_{r,q}(k) - \mathbb{I}_{r,p}(0)\mathbb{I}_{r,q}(0)) \right| \\ &\leq \frac{1}{m} \sum_{r=1}^m \left(\mathbb{I}\{\mathbb{I}_{r,p}(k) \neq \mathbb{I}_{r,p}(0)\} + \mathbb{I}\{\mathbb{I}_{r,q}(k) \neq \mathbb{I}_{r,q}(0)\} \right) \\ &\leq \frac{1}{m} \sum_{r=1}^m \left(\mathbb{I}\{\mathcal{A}_{rp}\} + \mathbb{I}\{\mathcal{A}_{rq}\} + 2\mathbb{I}\{\|w_r(k) - w_r(0)\| > R'\} \right).\end{aligned}\tag{177}$$

This event happens if and only if $|w_r(0)^\top x_q| \leq R'$. On the other hand, for any $(p, q) \in \{n\}^2$, it follows that with probability at least $1 - \Omega(\delta)$,

$$\begin{aligned}\|\mathbf{H}(k) - \mathbf{H}(0)\| &= \sum_{pq} |H_{pq}(k) - H_{pq}(0)| \\ &\leq \frac{1}{m} \sum_{pq} \sum_{r=1}^m \left(\mathbb{I}\{\mathcal{A}_{rp}\} + \mathbb{I}\{\mathcal{A}_{rq}\} + 2\mathbb{I}\{\|w_r(k) - w_r(0)\| > R'\} \right) \\ &\stackrel{(a)}{\leq} \frac{2n^2\delta}{m} + \frac{2n}{m} \sum_{r=1}^m \sum_p \mathbb{I}\{\mathcal{A}_{rp}\} \\ &\stackrel{(b)}{\leq} \frac{2n^2\delta}{m} + 2n^2\varphi',\end{aligned}\tag{178}$$

where (a) follows from Lemma 19 (with probability at least $1 - \delta$, $\|w_r(k) - w_r(0)\| \leq R'$ for all $r \in \{m\}$ and all $k \geq 0$) and (b) follows from the assumption *equation 143*.

Similarly, it follows that with probability at least $1 - \Omega(\delta)$,

$$\begin{aligned}\|\mathbf{Z}(k) - \mathbf{Z}(0)\|^2 &= \frac{1}{m} \sum_q \sum_r |\mathbb{I}_{r,q}(k) - \mathbb{I}_{r,q}(0)|^2 \\ &= \frac{1}{m} \sum_q \sum_r \mathbb{I}\{\mathbb{I}_{r,q}(k) \neq \mathbb{I}_{r,q}(0)\} \\ &\leq \frac{1}{m} \sum_q \sum_{r=1}^m \left(\mathbb{I}\{\mathcal{A}_{rq}\} + \mathbb{I}\{\|w_r(k) - w_r(0)\| > R'\} \right) \\ &\stackrel{(a)}{\leq} \frac{n\delta}{m} + n\varphi',\end{aligned}\tag{180}$$

where (a) follows from the assumption *equation 143*.

Therefore, the proof is completed by rescaling $\Omega(\delta)$ to δ . \square

Then, by using the above lemmas, we prove the following proposition. This result is a revision of Theorem 4.1 in Arora et al. (2019a) by removing a κ -affected value (i.e., $(1 - \eta\lambda_0)^k \frac{\sqrt{n\kappa}}{\delta}$) in the original bound given as in (33) in Arora et al. (2019a). Therefore, this proposition is our major contribution to prove Theorem 2.4.

Proposition F.1 (Modification/revision of Theorem 4.1 in Arora et al. (2019a)). *Let the assumption 2 hold. Under same setting as Theorem F.1, i.e., the assumptions 1 and 3 hold with $\varphi = O(\frac{\lambda_0}{n^2})$, $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega(\frac{n^2\kappa^2}{\lambda_0^2 R^2 \delta})$, set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, it follows with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ that for all $k \in \{0, 1, \dots\}$,*

$$\mathbf{u}(k) - \mathbf{y} = -(I - \eta \mathbf{H}(0))^k \mathbf{y} + \mathbf{e}(k), \quad (181)$$

where

$$\|\mathbf{e}(k)\| = O\left(k(1 - \frac{\eta\lambda_0}{4})^{k-1} \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi}{\sqrt{\delta}} + \frac{n^2\eta\varphi}{\sqrt{\delta}}\right)\right). \quad (182)$$

Proof of Proposition F.1. We define $u_q(k) := f_{\mathbf{W}(k)}(x_q)$ is the q th entry of $\mathbf{u}(k) := (f_{\mathbf{W}(k)}(x_1), \dots, f_{\mathbf{W}(k)}(x_n))^\top$. Then,

$$u_q(k+1) - u_q(k) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r [\sigma(w_r(k+1)^\top x_q) - \sigma(w_r(k)^\top x_q)]. \quad (183)$$

Define $R' = O(\frac{\kappa n}{\sqrt{m\lambda_0}\sqrt{\delta}})$. From Lemma 19, we obtain that with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, we have $\|w_r(k) - w_r(0)\| \leq R'$ for all $r \in \{1 : m\}$ and $k \geq 0$. On the other hand, the following event is defined as

$$\mathcal{A}_{rq} := \{|w_r(0)^\top x_q| \leq R'\}, \quad q \in \{n\}, r \in \{m\}. \quad (184)$$

Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for all r, p

$$\mathbb{I}\{\mathbb{I}_{r,p}(k) \neq \mathbb{I}_{r,p}(0)\} = \mathbb{I}\{\mathcal{A}_{rp}\}.$$

It is because the following event \mathcal{E}_{qr} happens if and only if $|w_r(0)^\top x_q| \leq R'$.

$$\mathcal{E}_{qr} := \{\exists \mathbf{w} : \|\mathbf{w} - w_r(0)\| \leq R', \mathbb{I}\{x_q^\top w_r(0) \geq 0\} \neq \mathbb{I}\{x_q^\top \mathbf{w} \geq 0\}\}$$

Let $S_q := \{r \in \{1 : m\} : \mathbb{I}\{\mathcal{A}_{qr}\} = 0\}$ and $S_q^\perp := \{1 : m\} \setminus S_q$. Then from the assumption equation 143, with probability at least $1 - \delta$ over the random initialization of $\mathbf{W}(0)$, we have

$$\sum_{q=1}^n |(S_q)^\perp| \leq nm\varphi'. \quad (185)$$

From (183), we get

$$\begin{aligned} u_q(k+1) - u_q(k) &= \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r [\sigma(w_r(k+1)^\top x_q) - \sigma(w_r(k)^\top x_q)] \\ &\quad + \frac{1}{\sqrt{m}} \sum_{r \in S_q^\perp} a_r [\sigma(w_r(k+1)^\top x_q) - \sigma(w_r(k)^\top x_q)] \end{aligned} \quad (186)$$

We denote the second term as $\dot{e}_q(k)$

$$\begin{aligned} |\dot{e}_q(k)| &= \left| \frac{1}{\sqrt{m}} \sum_{r \in S_q^\perp} a_r [\sigma(w_r(k+1)^\top x_q) - \sigma(w_r(k)^\top x_q)] \right| \\ &\leq \frac{1}{\sqrt{m}} \sum_{r \in S_q^\perp} \|w_r(k+1) - w_r(k)\| \\ &= \frac{\eta}{\sqrt{m}} \sum_{r \in S_q^\perp} \left\| \frac{\partial L_i(\mathbf{W}(k))}{\partial w_r(k)} \right\| \\ &\leq \frac{\eta}{m} \sum_{r \in S_q^\perp} \sum_{j=1}^n |y_j - u_j(k)| \\ &\leq \frac{\eta\sqrt{n}|S_q^\perp|}{m} \|\mathbf{y} - \mathbf{u}(k)\|. \end{aligned} \quad (187)$$

For the first term in (186),

$$\begin{aligned}
& \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r [\sigma(w_r(k+1)^\top x_q) - \sigma(w_r(k)^\top x_q)] \\
&= \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r \mathbb{1}\{w_r(k)^\top x_q \geq 0\} (w_r(k+1) - w_r(k))^\top x_q \\
&= \frac{1}{\sqrt{m}} \sum_{r \in S_q} a_r \mathbb{1}\{w_r(k)^\top x_q \geq 0\} \left(-\frac{\eta}{\sqrt{m}} \sum_{p=1}^n (u_p(k) - y_p) a_r x_p \mathbb{1}\{w_r(k)^\top x_p \geq 0\} \right)^\top x_q \\
&= -\frac{\eta}{m} \sum_{p=1}^n (u_p(k) - y_p) x_p^\top x_q \sum_{r \in S_q} \mathbb{1}\{w_r(k)^\top x_p \geq 0\} \mathbb{1}\{w_r(k)^\top x_q \geq 0\} + \bar{\epsilon}_q(k) \\
&= -\eta \sum_{p=1}^n (u_p(k) - y_p) H_{qp}(k) + \bar{\epsilon}_q(k), \tag{188}
\end{aligned}$$

where

$$\bar{\epsilon}_q(k) = \frac{\eta}{m} \sum_{p=1}^n (u_p(k) - y_p) x_p^\top x_q \sum_{r \in S_q^\perp} \mathbb{1}\{w_r(k)^\top x_p \geq 0\} \mathbb{1}\{w_r(k)^\top x_q \geq 0\}. \tag{189}$$

Then,

$$|\bar{\epsilon}_q(k)| \leq \frac{\eta \sqrt{n} |S_q^\perp|}{m} \|\mathbf{y} - \mathbf{u}(k)\|. \tag{190}$$

Combining (186), (187), (188) and (190),

$$u_q(k+1) - u_q(k) = -\eta \sum_{p=1}^n (u_p(k) - y_p) H_{qp}(k) + \dot{\epsilon}_q(k) + \bar{\epsilon}_q(k) \tag{191}$$

which gives

$$\mathbf{u}(k+1) - \mathbf{u}(k) = -\eta \mathbf{H}(k)(\mathbf{u}(k) - \mathbf{y}) + \boldsymbol{\epsilon}(k), \tag{192}$$

where $\boldsymbol{\epsilon}(k) = \dot{\boldsymbol{\epsilon}}(k) + \bar{\boldsymbol{\epsilon}}(k)$. Note that by using equation 185,

$$\|\boldsymbol{\epsilon}(k)\| \leq \|\boldsymbol{\epsilon}(k)\|_1 = \sum_{q=1}^n |\dot{\epsilon}_q(k) + \bar{\epsilon}_q(k)| \leq \sum_{q=1}^n \frac{2\eta \sqrt{n} |S_q^\perp|}{m} \|\mathbf{y} - \mathbf{u}(k)\| = O\left(\eta n \sqrt{n} \varphi'\right) \|\mathbf{y} - \mathbf{u}(k)\|. \tag{193}$$

We rewrite (192) as

$$\mathbf{u}(k+1) - \mathbf{u}(k) = -\eta \mathbf{H}(0)(\mathbf{u}(k) - \mathbf{y}) + \boldsymbol{\zeta}(k), \tag{194}$$

where $\boldsymbol{\zeta}(k) = \eta(\mathbf{H}(0) - \mathbf{H}(k))(\mathbf{u}(k) - \mathbf{y}) + \boldsymbol{\epsilon}(k)$. Then, we get

$$\mathbf{u}(k) - \mathbf{y} = -(1 - \eta \mathbf{H}(0))^k (\mathbf{y} - \mathbf{u}(0)) + \sum_{t=0}^{k-1} (I - \eta \mathbf{H}(0))^t \boldsymbol{\zeta}(k-1-t). \tag{195}$$

From (193) and Lemmas 20, we bound $\boldsymbol{\zeta}(k)$ as

$$\begin{aligned}
\|\boldsymbol{\zeta}(k)\| &\leq \eta \|\mathbf{H}(0) - \mathbf{H}(k)\|_2 \|\mathbf{y} - \mathbf{u}(k)\| + O\left(\eta n \sqrt{n} \varphi'\right) \|\mathbf{y} - \mathbf{u}(k)\| \\
&= O\left(\frac{2n^2 \delta}{m} + 2n^2 \varphi' + \eta n \sqrt{n} \varphi'\right) \|\mathbf{y} - \mathbf{u}(k)\|. \tag{196}
\end{aligned}$$

Then,

$$\begin{aligned}
& \left\| \sum_{t=0}^{k-1} (I - \eta \mathbf{H}^{(0)})^t \zeta(k-1-t) \right\| \\
& \leq \sum_{t=0}^{k-1} \|I - \eta \mathbf{H}(0)\|^t \|\zeta(k-1-t)\| \\
& \leq \sum_{t=0}^{k-1} \|I - \eta \mathbf{H}(0)\|^t \|\zeta(k-1-t)\| \\
& \leq \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t \|\zeta(k-1-t)\| \\
& \stackrel{(a)}{\leq} \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t O\left(\frac{2n^2 \delta}{m} + 2n^2 \varphi' + \eta n \sqrt{n} \varphi'\right) \|\mathbf{y} - \mathbf{u}(k-1-t)\| \\
& \stackrel{(b)}{\leq} \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t O\left(\frac{2n^2 \delta}{m} + 2n^2 \varphi' + \eta n \sqrt{n} \varphi'\right) \left(1 - \frac{\eta \lambda_0}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\| \\
& \stackrel{(c)}{\leq} \sum_{t=0}^{k-1} (1 - \eta \lambda_0)^t O\left(\frac{2n^2 \delta}{m} + 2n^2 \varphi' + \eta n \sqrt{n} \varphi'\right) \left(1 - \frac{\eta \lambda_0}{4}\right)^k O\left(\frac{\sqrt{n}}{\sqrt{\delta}}\right) \\
& \leq k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} O\left(\frac{2n^{5/2} \sqrt{\delta}}{m} + \frac{2n^{5/2} \varphi'}{\sqrt{\delta}} + \frac{n^2 \eta \varphi'}{\sqrt{\delta}}\right), \tag{197}
\end{aligned}$$

where (a) follows from (196), (b) follows from Theorem F.1 that

$$\|\mathbf{y} - \mathbf{u}(k)\| \leq \left(1 - \frac{\eta \lambda_0}{2}\right)^{k/2} \|\mathbf{y} - \mathbf{u}(0)\| \leq \left(1 - \frac{\eta \lambda_0}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|, \tag{198}$$

and (c) follows from Lemma 15 and the assumption $\kappa = \Theta(1)$ (i.e., $\|\mathbf{y} - \mathbf{u}(0)\| = O(\frac{\sqrt{n}}{\sqrt{\delta}})$).

By applying (197) to (195), under same setting as Theorem F.1, it follows that for $k \geq 0$, with probability at least $1 - \delta$,

$$\mathbf{u}(k) - \mathbf{y} = -(I - \eta \mathbf{H}(0))^k (\mathbf{y} - \mathbf{u}(0)) + \mathbf{e}(k), \tag{199}$$

where

$$\|\mathbf{e}(k)\| = O\left(k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} \left(\frac{2n^{5/2} \sqrt{\delta}}{m} + \frac{2n^{5/2} \varphi'}{\sqrt{\delta}} + \frac{n^2 \eta \varphi'}{\sqrt{\delta}}\right)\right). \tag{200}$$

□

As a simple corollary of Proposition F.1, now we can prove Theorem 2.4 as follows.

Theorem F.2. [Theorem 2.4, modification/revision of Theorem 4.1 in Arora et al. (2019a)] Let the assumption 2 hold. Under same setting as Theorem F.1, i.e., the assumptions 1 and 3 hold with $\varphi = O(\frac{\lambda_0}{n^2})$, $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\frac{n^2 \kappa^2}{\lambda_0^2 R^2 \delta}\right)$, set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$, it follows with probability at least $1 - \delta$ for $\delta \in (0, 1)$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ that for all $k \geq 0$,

$$\frac{1}{n} \|\mathbf{y} - \mathbf{u}(k)\|^2 = \frac{1}{n} \sum_{i=1}^n (1 - \eta \lambda_i)^{2k} \left(v_i^\top (\mathbf{y} - \mathbf{u}(0))\right)^2 + O\left(\left[\frac{n^4 \sqrt{\delta}}{m \lambda_0^2} + \frac{n^4 \varphi'}{\lambda_0^2 \sqrt{\delta}} + \frac{\varphi' n^{3/2}}{\lambda_0 \sqrt{\delta}}\right]^2\right), \tag{201}$$

where $v_1, \dots, v_n \in \mathbb{R}^n$ are orthonormal eigenvectors of $\mathbf{H}(0)$ and $\lambda_1, \dots, \lambda_n$ are the corresponding eigenvalues.

Proof of Theorem F.2. Our proof is based on that for Theorem 4.1 in Arora et al. (2019a).

From Proposition F.1, it follows that for all $k \in \{0, 1, \dots\}$,

$$\mathbf{u}(k) - \mathbf{y} = -(I - \eta \mathbf{H}^{(0)})^k \mathbf{y} + \mathbf{e}(k), \tag{202}$$

where

$$\|e(k)\| = O\left(k(1 - \frac{\eta\lambda_0}{4})^{k-1} \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi'}{\sqrt{\delta}} + \frac{n^2\eta\varphi'}{\sqrt{\delta}}\right)\right). \quad (203)$$

Therefore, we get

$$\begin{aligned} \|\mathbf{u}(k) - \mathbf{y}\|^2 &\stackrel{(a)}{\leq} \left\|(\mathbf{I} - \eta\mathbf{H}^{(0)})^k(\mathbf{u}(0) - \mathbf{y})\right\|^2 + \|e(k)\|^2 \\ &\stackrel{(b)}{=} \sum_{j=1}^n (1 - \eta\lambda_j)^{2k} (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0)))^2 + \|e(k)\|^2 \\ &\stackrel{(c)}{=} \sum_{j=1}^n (1 - \eta\lambda_j)^{2k} (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0)))^2 + O\left(\left[k(1 - \frac{\eta\lambda_0}{4})^{k-1} \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi'}{\sqrt{\delta}} + \frac{n^2\eta\varphi'}{\sqrt{\delta}}\right)\right]^2\right) \\ &\stackrel{(d)}{=} \sum_{j=1}^n (1 - \eta\lambda_j)^{2k} (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0)))^2 + O\left(\left[\frac{1}{\eta\lambda_0} \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi'}{\sqrt{\delta}} + \frac{n^2\eta\varphi'}{\sqrt{\delta}}\right)\right]^2\right) \\ &\stackrel{(e)}{=} \sum_{j=1}^n (1 - \eta\lambda_j)^{2k} (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0)))^2 + O\left(\left[\frac{n^{9/2}\sqrt{\delta}}{m\lambda_0^2} + \frac{n^{9/2}\varphi'}{\lambda_0^2\sqrt{\delta}} + \frac{\varphi'n^2}{\lambda_0\sqrt{\delta}}\right]^2\right), \end{aligned}$$

where (a) follows from the triangle inequality and equation 202, (b) follows from $(\mathbf{I} - \eta\mathbf{H}^{(0)})^k$ has the eigen-decomposition $(\mathbf{I} - \eta\mathbf{H}^{(0)})^k = \sum_{j=1}^n (1 - \eta\lambda_j)^k \mathbf{v}_j \mathbf{v}_j^\top$ and $\mathbf{y} - \mathbf{u}(0) = \sum_{j=1}^n (\mathbf{v}_j^\top (\mathbf{y} - \mathbf{u}(0))) \mathbf{v}_j$, (c) follows from equation 203, (d) follows from $\max_{k \geq 0} \{k(1 - \eta\lambda_0/4)^{k-1}\} = O(1/(\eta\lambda_0))$, and (e) follows from $\eta = O(\lambda_0/n^2)$. \square

F.3 PROOF OF THEOREM 2.5 WHEN THE NETWORK HAS A SCALAR OUTPUT AS $l = 1$

F.3.1 BACKGROUND ON RADEMACHER COMPLEXITY

Before we prove Theorem 2.5 stated in Section 2.3, we introduce Rademacher Complexity and the theorem derived from it.

Define a loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$. For a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we define the population loss over true model distribution \mathcal{D} and the empirical loss over n samples $\mathcal{S} = \{(x_j, y_j)\}_{j=1}^n$ from \mathcal{D} , respectively, as

$$\begin{aligned} \mathcal{L}_{\mathcal{D}}(f) &= \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(f(x), y)] \\ \mathcal{L}_{\mathcal{S}}(f) &= \sum_{j=1}^n [\ell(f(x_j), y_j)]. \end{aligned} \quad (204)$$

Then, Rademacher complexity of a function class \mathcal{F} mapping \mathbb{R}^d to \mathbb{R} is expressed as

$$\mathcal{R}_{\mathcal{S}}(\mathcal{F}) := \frac{1}{n} \mathbb{E}_{\epsilon} \left[\sup_{f \in \mathcal{F}} \sum_{j=1}^n \epsilon_j f(x_j) \right], \quad (205)$$

where $\epsilon = (\epsilon_1, \dots, \epsilon_n)^\top \in \mathbb{R}^n$ includes i.i.d. random variables $\epsilon_j \sim \text{unif}(\{1, -1\})$ for $j \in \{1, \dots, n\}$. This provides an upper bound of generalization error as the following theorem given from Mohri et al. (2018).

Theorem F.3. Suppose the α -Lipschitz loss function $\ell(\cdot, \cdot)$ is bounded in $[0, \beta]$ in the first argument. Then, with probability at least $1 - \delta$ over sample \mathcal{S} of size n ,

$$\sup_{f \in \mathcal{F}} \{\mathcal{L}_{\mathcal{D}}(f) - \mathcal{L}_{\mathcal{S}}(f)\} \leq 2\alpha \mathcal{R}_{\mathcal{S}}(\mathcal{F}) + 3\beta \sqrt{\frac{\log(2/\delta)}{2n}}. \quad (206)$$

F.3.2 PROOF OF THEOREM 2.5

In this section, we now prove Theorem 2.5 stated in Section 2.3. We first show some technical lemmas.

As a result of Lemma 19, we can obtain the following upper bound of the magnitude of trained NN output, when the neural network is over-parameterized.

We can obtain the following lemma 21 by modifying Lemma 5.3 in Arora et al. (2019a). Note that Lemma 5.3 in Arora et al. (2019a) provides an upper bound of distance between trained NN weights and its initial ones. Lemma 21 is a result obtained by removing the terms affected by κ from this upper bound of Lemma 5.3 in Arora et al. (2019a). Therefore, this lemma is our major contribution to prove Theorem 2.5.

Lemma 21 (Modification/revision of Lemma 5.3 in Arora et al. (2019a)). *Let the assumption 2 hold. Consider the same setting as Theorem F.1, i.e., the assumptions 1 and 3 hold with $\varphi = O(\frac{\lambda_0}{n^2})$, $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $m = \Omega\left(\max\left(\frac{n^6}{\lambda_0^4 \kappa^2 \delta^3}, \frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)\right)$, set $\{x_j\}_{j=1}^n$ of n training input samples is bounded as $\max_{j \in \{n\}} \|x_j\| \leq 1$, $\eta = O(\frac{\lambda_0}{n^2})$, and $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$. Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, it follows that for all $k \geq 0$*

- (Lemma 19) $\|w_r(k) - w_r(0)\| = O\left(\frac{\kappa n}{\sqrt{m} \lambda_0 \sqrt{\delta}}\right) \quad (:= R), \forall r \in \{1 : m\}$
- $\|\mathbf{W}(k) - \mathbf{W}(0)\| \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^{(0)})^{-1} (\mathbf{y} - \mathbf{u}(0))} + O\left(\sqrt{\frac{n^2 \delta}{m \lambda_0^2} + \frac{n^2 \varphi'}{\lambda_0^2}} + \frac{1}{\eta \lambda_0^2} \left(\frac{2n^3 \sqrt{\delta}}{m} + \frac{2n^3 \varphi'}{\sqrt{\delta}} + \frac{n^{5/2} \eta \varphi'}{\sqrt{\delta}}\right)\right)$

Proof of Lemma 21. The first part is proved by using Lemma 19. The rest is to prove the second part.

From Proposition F.1, we get

$$\mathbf{u}(k) - \mathbf{y} = -(\mathbf{I} - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) + \mathbf{e}(k), \quad (207)$$

where

$$\|\mathbf{e}(k)\| = O\left(k \left(1 - \frac{\eta \lambda_0}{4}\right)^{k-1} \left(\frac{2n^{5/2} \sqrt{\delta}}{m} + \frac{2n^{5/2} \varphi'}{\sqrt{\delta}} + \frac{n^2 \eta \varphi'}{\sqrt{\delta}}\right)\right). \quad (208)$$

We apply (207) to (135), which is

$$\text{vec}(\mathbf{W}(k+1)) = \text{vec}(\mathbf{W}(k)) - \eta \mathbf{Z}(k)(\mathbf{u}(k) - \mathbf{y}), \quad (209)$$

and for $k \in \{0, \dots, K-1\}$ we obtain

$$\begin{aligned} & \text{vec}(\mathbf{W}(k)) - \text{vec}(\mathbf{W}(0)) \\ &= -\sum_{k=0}^{K-1} \eta \mathbf{Z}(k)(\mathbf{u}(k) - \mathbf{y}) \\ &= \sum_{k=0}^{K-1} \eta \mathbf{Z}(k)((\mathbf{I} - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) - \mathbf{e}(k)) \\ &= \sum_{k=0}^{K-1} \eta \mathbf{Z}(0)(\mathbf{I} - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) + \sum_{k=0}^{K-1} \eta (\mathbf{Z}(k) - \mathbf{Z}(0))(\mathbf{I} - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) - \sum_{k=0}^{K-1} \eta \mathbf{Z}(k) \mathbf{e}(k). \end{aligned} \quad (210)$$

Then we bound the first term of (210) as

$$\left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(0)(\mathbf{I} - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) \right\|^2 = (\mathbf{y} - \mathbf{u}(0))^\top \mathbf{T} \mathbf{H}^{(0)} \mathbf{T} (\mathbf{y} - \mathbf{u}(0)),$$

where $\mathbf{T} := \sum_{i=1}^n \sum_{k=0}^{K-1} \eta (1 - \eta \lambda_i)^k \mathbf{v}_i \mathbf{v}_i^\top = \sum_{i=1}^n \frac{1 - (1 - \eta \lambda_i)^K}{\lambda_i} \mathbf{v}_i \mathbf{v}_i^\top$ ($v_1, \dots, v_n \in \mathbb{R}^n$ are orthonormal eigenvectors of \mathbf{H}^* and $\lambda_1, \dots, \lambda_n$ are the corresponding eigenvalues). Then,

$$\left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(0)(\mathbf{I} - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) \right\| = \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top \mathbf{T} \mathbf{H}^{(0)} \mathbf{T} (\mathbf{y} - \mathbf{u}(0))} \quad (211)$$

By using

$$\mathbf{T} \mathbf{H}^{(0)} \mathbf{T} = \sum_{i=1}^n \left(\frac{1 - (1 - \eta \lambda_i)^K}{\lambda_i}\right)^2 \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \preceq \sum_{i=1}^n \frac{1}{\lambda_i} \mathbf{v}_i \mathbf{v}_i^\top = (\mathbf{H}^{(0)})^{-1}, \quad (212)$$

we get

$$\left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(0)(\mathbf{I} - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) \right\| \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^*)^{-1} (\mathbf{y} - \mathbf{u}(0))}. \quad (213)$$

We bound the second term of (210) as

$$\begin{aligned}
& \left\| \sum_{k=0}^{K-1} \eta (\mathbf{Z}(k) - \mathbf{Z}(0)) (1 - \eta \mathbf{H}^{(0)})^k (\mathbf{y} - \mathbf{u}(0)) \right\| \\
& \leq \sum_{k=0}^{K-1} \eta \|\mathbf{Z}(k) - \mathbf{Z}(0)\|_2 \left\| (I - \eta \mathbf{H}^{(0)})^k \right\|_2 \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \leq \sum_{k=0}^{K-1} \eta \|\mathbf{Z}(k) - \mathbf{Z}(0)\|_2 \left\| I - \eta \mathbf{H}^{(0)} \right\|_2^k \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \leq \eta \|\mathbf{Z}(K) - \mathbf{Z}(0)\|_2 \sum_{k=0}^{K-1} (1 - \eta \lambda_0)^k \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \stackrel{(a)}{=} O \left(\sqrt{\frac{n\delta}{m} + n\varphi'} \right) \sum_{k=0}^{K-1} \eta (1 - \eta \lambda_0)^k \|\mathbf{y} - \mathbf{u}(0)\|_2 \\
& \stackrel{(b)}{=} O \left(\sqrt{\frac{n^2\delta}{m\lambda_0^2} + \frac{n^2\varphi'}{\lambda_0^2}} \right), \tag{214}
\end{aligned}$$

where (a) follows from Lemma 20 and (b) follows from $\sum_{k=0}^{K-1} \eta (1 - \eta \lambda_0)^k = \frac{1 - (1 - \eta \lambda_0)^K}{\lambda_0} \leq \frac{1}{\lambda_0}$.

By using (208) and the fact that $\|\mathbf{Z}(k)\| \leq \sqrt{n}$, we bound the third term of (210) as

$$\begin{aligned}
& \left\| \sum_{k=0}^{K-1} \eta \mathbf{Z}(k) e(k) \right\| = O \left(\sum_{k=0}^{K-1} \eta \sqrt{n} \cdot \left[k(1 - \frac{\eta \lambda_0}{4})^{k-1} \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi'}{\sqrt{\delta}} + \frac{n^2\eta\varphi'}{\sqrt{\delta}} \right) \right] \right) \\
& = O \left(\left(\sum_{k=0}^{K-1} k(1 - \frac{\eta \lambda_0}{4})^{k-1} \right) \left(\eta \sqrt{n} \cdot \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi'}{\sqrt{\delta}} + \frac{n^2\eta\varphi'}{\sqrt{\delta}} \right) \right) \right) \\
& \stackrel{(a)}{=} O \left(\frac{1}{\eta \lambda_0^2} \left(\frac{2n^3\sqrt{\delta}}{m} + \frac{2n^3\varphi'}{\sqrt{\delta}} + \frac{n^{5/2}\eta\varphi'}{\sqrt{\delta}} \right) \right), \tag{215}
\end{aligned}$$

where (a) follows from the fact that $\sum_{k=0}^{K-1} k(1 - \frac{\eta \lambda_0}{4})^{k-1} \leq \sum_{k=0}^{\infty} k(1 - \frac{\eta \lambda_0}{4})^{k-1} = (\frac{4}{\eta \lambda_0})^2$.

Therefore, by applying (213), (214), and (215) to equation 210,

$$\begin{aligned}
\|\mathbf{W}(k) - \mathbf{W}(0)\| & \leq \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^{(0)})^{-1} (\mathbf{y} - \mathbf{u}(0))} \\
& + O \left(\sqrt{\frac{n^2\delta}{m\lambda_0^2} + \frac{n^2\varphi'}{\lambda_0^2}} + \frac{1}{\eta \lambda_0^2} \left(\frac{2n^3\sqrt{\delta}}{m} + \frac{2n^3\varphi'}{\sqrt{\delta}} + \frac{n^{5/2}\eta\varphi'}{\sqrt{\delta}} \right) \right)
\end{aligned}$$

□

We introduce the following lemma (i.e., Lemma 22) proved by Arora et al. (2019a). This lemma shows Rademacher complexity can be upper bounded by a term depending on the distance between the trained weight and its initial value.

Lemma 22 (Lemma 5.4 in Arora et al. (2019a)). *Given $R > 0$, we assume that the input data $\{x_j\}_{j=1}^n$ is given as $\|x_j\| \leq 1$ for $j \in \{n\}$. Consider the following function class in equation 216 with $\mathbf{W}(0) =: [w_1(0), \dots, w_m(0)]$*

$$\mathcal{F}_{R',B}^{\mathbf{W}(0),\mathbf{a}} := \{f_{\hat{\mathbf{W}},\mathbf{a}} : \hat{\mathbf{W}} = [\hat{w}_1, \dots, \hat{w}_m], \|\hat{w}_r - w_r(0)\| \leq R' (\forall r \in \{m\}), \|\hat{\mathbf{W}} - \mathbf{W}(0)\| \leq B\}. \tag{216}$$

Then it follows that with a probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, for every $B > 0$, the empirical Rademacher complexity $\mathcal{R}_S(\mathcal{F}_{R',B}^{\mathbf{W}(0),\mathbf{a}})$ based on the function class in equation 216 is bounded as

$$\begin{aligned}
\mathcal{R}_S(\mathcal{F}_{R',B}^{\mathbf{W}(0),\mathbf{a}}) & := \frac{1}{n} \mathbb{E}_{\epsilon \in \{\pm 1\}^n} \left[\sup_{f \in \mathcal{F}_{R',B}^{\mathbf{W}(0),\mathbf{a}}} \sum_{j=1}^n \epsilon_j f(x_j) \right] \\
& \leq \frac{B}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{2}{\delta}}{m} \right)^{1/4} \right) + \frac{2R'^2 \sqrt{m}}{\kappa} + R' \sqrt{2 \log \frac{2}{\delta}}.
\end{aligned}$$

Now, by using Lemmas 19 and 22, we prove Theorem 2.5, which is given as Theorem F.4.

Theorem F.4. [Theorem 2.5, modification/revision of Theorem 5.1 in Arora et al. (2019a)] Suppose that all conditions except $\lambda_0 > 0$ in Theorem F.1 hold and we fix a failure probability $\delta \in (0, 1)$. Suppose further that assumption 2, 4, 5, and 6 hold. Suppose also that $\lambda_0 > 0$ holds with probability at least $1 - \delta/3$ for n i.i.d. training samples $\{(x_i, y_i)\}_{i=1}^n$ from true model distribution \mathcal{D} . Then, with probability at least $1 - \delta$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$ and the training samples, it follows that for any $k \geq \Omega(\frac{1}{\eta\lambda_0} \log \frac{n}{\delta})$,

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \frac{1}{2} \|y - f_{\mathbf{W}(k)}(x)\|^2 = O\left(\sqrt{\frac{2(\mathbf{y} - \mathbf{u}(0))^\top \mathbf{H}^{(0)-1}(\mathbf{y} - \mathbf{u}(0))}{n}}\right) + O\left(\sqrt{\frac{\log \frac{n}{\lambda_0 \delta}}{n}}\right). \quad (217)$$

Proof of Theorem F.4. We consider a loss function $\ell(a, b) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ as $\ell(a, b) = (a - b)^2/2$. We assume that this loss function $\ell(a, b)$ is α -Lipschitz in the first argument, this function is bounded in $[0, \beta]$, and α and β follow $O(1)$. We will prove that this assumption holds at the end of the proof.

Using the loss function and equation 204, we can define the population loss over true model distribution \mathcal{D} and the empirical loss over n samples \mathcal{S} , respectively, as

$$\begin{aligned} \mathcal{L}_{\mathcal{D}}(f) &= \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f(x), y)] = \mathbb{E}_{(x,y) \sim \mathcal{D}} \frac{1}{2} (f(x) - y)^2 \\ \mathcal{L}_{\mathcal{S}}(f) &= \sum_{j=1}^n [\ell(f(x_j), y_j)] = \sum_{j=1}^n \left[\frac{1}{2} (f(x_j) - y_j)^2 \right], \end{aligned} \quad (218)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes a trained neural network to be specified (i.e., $f_{\mathbf{W}(k)}$).

We assume that $\lambda_0 = O(n^\gamma) > 0$ with a constant $\gamma \leq 1$ holds with probability at least $1 - \delta/3$. We also assume the following conditions hold: $\kappa = \Theta(1)$ for n , $\|\mathbf{y}\| = O(\sqrt{n})$, $\eta = O(\frac{\lambda_0}{n^2})$, and $m = \Omega(\frac{n^6}{\lambda_0^4 \delta^3})$.

With probability at least $1 - \delta/6$ over the random initialization of $(\mathbf{W}(0), \mathbf{a})$, the followings hold simultaneously:

- Optimization succeeds: Suppose that \hat{k} is any integer satisfying

$$\hat{k} \log(1 - \frac{\eta\lambda_0}{2}) \leq \log\left(\frac{\epsilon}{\|\mathbf{y} - \mathbf{u}(0)\|^2}\right), \quad (219)$$

where ϵ is arbitrary small constant invariant of n . As $-\eta\lambda_0(2 - \eta\lambda_0)^{-1} \leq \log(1 - \frac{\eta\lambda_0}{2})$, the following condition implies equation 219.

$$\hat{k} > ((\eta\lambda_0)/(2 - \eta\lambda_0))^{-1} \log(\|\mathbf{y} - \mathbf{u}(0)\|^2 / \epsilon) \quad (220)$$

By using Theorem F.1 with equation 220, and the fact that $\|\mathbf{y} - \mathbf{u}(0)\| = O(\frac{\sqrt{n}}{\sqrt{\delta}})$ (which is obtained from Lemma 15 and the assumption $\kappa = \Theta(1)$), if $k = \Omega(\frac{1}{\eta\lambda_0} \log \frac{n}{\delta})$, it follows that

$$L(\mathbf{W}(k)) \leq (1 - \frac{\eta\lambda_0}{2})^k O\left(\frac{n}{\delta}\right) \leq \frac{1}{2}. \quad (221)$$

Then,

$$\begin{aligned} \mathcal{L}_{\mathcal{S}}(f_{\mathbf{W}(k)}) &:= \frac{1}{2n} \sum_{q=1}^n |f_{\mathbf{W}(k)}(x_q) - y_q|^2 \\ &= \frac{1}{n} L(\mathbf{W}(k)) \\ &\stackrel{(a)}{=} O\left(\frac{1}{n}\right), \end{aligned} \quad (222)$$

where (a) follows from equation 221.

- From Lemma 21, we get $\|w_r(k) - w_r(0)\| = R'$ ($\forall r \in \{1 : m\}$) where $R' = O\left(\frac{n}{\sqrt{m\lambda_0}\sqrt{\delta}}\right)$, and $\|\mathbf{W}(k) - \mathbf{W}(0)\| \leq B$ where $B = \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^{(0)})^{-1} (\mathbf{y} - \mathbf{u}(0))} + O\left(\sqrt{\frac{n^2 \delta}{m \lambda_0^2}} + \frac{n^2 \varphi'}{\lambda_0^2} + \frac{1}{\eta \lambda_0^2} \left(\frac{2n^3 \sqrt{\delta}}{m} + \frac{2n^3 \varphi'}{\sqrt{\delta}} + \frac{n^{5/2} \eta \varphi'}{\sqrt{\delta}}\right)\right)$. Note that $B \leq O(\sqrt{\frac{n}{\lambda_0}})$.

- Let $B_j = j$ ($j = 1, 2, \dots$). For all i , the function class $\mathcal{F}_{R', B_j}^{\mathbf{W}^{(0), \mathbf{a}}}$ has Rademacher complexity, which is upper bounded by

$$\mathcal{R}_S(\mathcal{F}_{R', B_j}^{\mathbf{W}^{(0), \mathbf{a}}}) \leq \frac{B_j}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{20}{\delta}}{m} \right)^{1/4} \right) + \frac{2R'^2 \sqrt{m}}{\kappa} + R' \sqrt{2 \log \frac{20}{\delta}}. \quad (223)$$

Let j^* be the smallest integer such that $B \leq B_{j^*}$. Then we have $B_{j^*} \leq B + 1$ and $j^* \leq O(\sqrt{\frac{n}{\lambda_0}})$. Note that $f_{\mathbf{W}^{(0), \mathbf{a}}} \in \mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0), \mathbf{a}}}$. And we get

$$\begin{aligned} \mathcal{R}_S(\mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0), \mathbf{a}}}) &\leq \frac{B+1}{\sqrt{2n}} \left(1 + \left(\frac{2 \log \frac{20}{\delta}}{m} \right)^{1/4} \right) + O(R'^2 \sqrt{m}) + R' \sqrt{2 \log \frac{20}{\delta}} \\ &\leq \frac{B+1}{\sqrt{2n}} O(1) + O(R'^2 \sqrt{m}) + R' \sqrt{2 \log \frac{20}{\delta}} \\ &= O\left(\frac{B}{\sqrt{n}}\right) + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\frac{n^2}{\sqrt{m} \lambda_0^2}\right) + O\left(\frac{n}{\sqrt{m} \lambda_0}\right) \\ &= O\left(\frac{A}{\sqrt{n}}\right) + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\sqrt{\frac{n\delta}{m\lambda_0^2} + \frac{n\varphi'}{\lambda_0^2}} + \frac{1}{\eta\lambda_0^2} \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi'}{\sqrt{\delta}} + \frac{n^2\eta\varphi'}{\sqrt{\delta}} \right) + O\left(\frac{n^2 + n\lambda_0}{\sqrt{m}\lambda_0^2}\right)\right) \\ &= O\left(\frac{A}{\sqrt{n}}\right) + O\left(\frac{1}{\sqrt{n}}\right) + O\left(\sqrt{\frac{n\delta}{m\lambda_0^2} + \frac{n\varphi'}{\lambda_0^2}} + \frac{1}{\eta\lambda_0^2} \left(\frac{2n^{5/2}\sqrt{\delta}}{m} + \frac{2n^{5/2}\varphi'}{\sqrt{\delta}} + \frac{n^2\eta\varphi'}{\sqrt{\delta}} \right) + \frac{n^2 + n\lambda_0}{\sqrt{m}\lambda_0^2}\right) \\ &\stackrel{(a)}{=} O\left(\frac{A}{\sqrt{n}}\right) + O\left(\frac{1}{\sqrt{n}}\right), \end{aligned} \quad (224)$$

where $A = \sqrt{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^{(0)})^{-1} (\mathbf{y} - \mathbf{u}(0))}$ and (a) follows from the condition equation 225.

$$\begin{aligned} m &= \Omega\left(\max\left(\frac{n^2}{\lambda_0^2}, \frac{n^2}{\eta\lambda_0^2}, \frac{n^3}{\lambda_0^4}, \frac{n^3}{\lambda_0^2}\right)\right) = \Omega\left(\max\left(\frac{n^2}{\lambda_0^2}, \frac{n^4}{\lambda_0^3}, \frac{n^3}{\lambda_0^4}, \frac{n^3}{\lambda_0^2}\right)\right) = \Omega\left(\frac{n^4}{\min(\lambda_0^2, \lambda_0^4)}\right) \\ \varphi' &= O\left(\min\left(\frac{\lambda_0^2}{n^2}, \frac{\eta\lambda_0^2}{n^2}\right)\right) = O\left(\min\left(\frac{\lambda_0^2}{n^2}, \frac{\lambda_0^3}{n^4}\right)\right) = O\left(\frac{\min(\lambda_0^2, \lambda_0^3)}{n^4}\right) \\ \frac{\varphi'}{m} &= O\left(\frac{\lambda_0^2}{n^{3/2}}\right), \end{aligned} \quad (225)$$

By using $\eta = O(\frac{\lambda_0}{n^2})$, equation 225 is implied by

$$\begin{aligned} m &= \Omega\left(\frac{n^4}{\min(\lambda_0^2, \lambda_0^4)}\right), \\ \varphi' &= O\left(\frac{\min(\lambda_0^2, \lambda_0^3)}{n^4}\right), \end{aligned} \quad (226)$$

which are our assumptions.

From the result of Rademacher complexity (Theorem F.3) and the union bound over a finite set $\{1 : B_{j^*}\}$, with probability at least $1 - \delta/6$, the following inequality holds for all $j \in \{1, 2, \dots, j^*\}$.

$$\sup_{f \in \mathcal{F}_{R', B_j}^{\mathbf{W}^{(0), \mathbf{a}}}} \{\mathcal{L}_{\mathcal{D}}(f) - \mathcal{L}_S(f)\} \leq 2\alpha \mathcal{R}_S(\mathcal{F}_{R', B_j}^{\mathbf{W}^{(0), \mathbf{a}}}) + O\left(\beta \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right) \quad (227)$$

By using the union bound jointly to consider equation 222, equation 224, and equation 227, we obtain the fact that with probability at least $1 - 5\delta/6$, the followings are satisfied at the same time.

$$\begin{aligned}
\mathcal{L}_S(f_{\mathbf{W}^{(k)}, \mathbf{a}}) &= O\left(\frac{1}{n}\right) \\
f_{\mathbf{W}^{(k)}, \mathbf{a}} &\in \mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0)}, \mathbf{a}} \\
\mathcal{R}_S(\mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0)}, \mathbf{a}}) &= \sqrt{\frac{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^{(0)})^{-1} (\mathbf{y} - \mathbf{u}(0))}{2n}} + O\left(\frac{1}{\sqrt{n}}\right) \\
\sup_{f \in \mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0)}, \mathbf{a}}} \{\mathcal{L}_D(f) - \mathcal{L}_S(f)\} &\leq 2\alpha \mathcal{R}_S(\mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0)}, \mathbf{a}}) + O\left(\beta \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right)
\end{aligned} \tag{228}$$

If the assumption that α and β follow $O(1)$ holds with probability at least $1 - \delta/6$, by using the union bound, it follows that with probability at least $1 - \delta$,

$$\begin{aligned}
&\mathcal{L}_D(f_{\mathbf{W}^{(k)}, \mathbf{a}}) \\
&= \mathbb{E}_{(x, y) \sim \mathcal{D}} \frac{1}{2} |y - f_{\mathbf{W}^{(k)}}(x)|^2 \\
&\stackrel{(a)}{=} O\left(\frac{1}{n}\right) + O\left(\alpha \mathcal{R}_S(\mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0)}, \mathbf{a}})\right) + O\left(\beta \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right) \\
&\stackrel{(b)}{=} O\left(\frac{1}{n}\right) + O\left(\mathcal{R}_S(\mathcal{F}_{R', B_{j^*}}^{\mathbf{W}^{(0)}, \mathbf{a}})\right) + O\left(\sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}}\right) \\
&\stackrel{(c)}{=} O\left(\sqrt{\frac{(\mathbf{y} - \mathbf{u}(0))^\top (\mathbf{H}^{(0)})^{-1} (\mathbf{y} - \mathbf{u}(0))}{2n}}\right) + O\left(\frac{1}{\sqrt{n}}\right) + \sqrt{\frac{\log(\frac{n}{\lambda_0 \delta})}{2n}},
\end{aligned} \tag{229}$$

where (a) follow from equation 228, (b) follows from the assumption that α and β follow $O(1)$ for n , and (c) follows from equation 228.

Therefore, by using equation 145, we get equation 217 from equation 229.

Now we will prove the assumption that α and β follow $O(1)$. From our assumption equation 145, with probability at least $1 - \delta$, $|f_{\mathbf{W}^{(k)}}(x)|$ in equation 229 follows $O(1)$ for every $k \geq 0$ and $x \sim \mathcal{D}$. On the other hand, y follows $O(1)$ for n . This is because y is independent of n , as y is i.i.d. sample of the model $\mathcal{D}(x, y)$. These imply that $|y - f_{\mathbf{W}^{(k)}}(x)|$ in equation 229 follows $O(1)$ for every $k \geq 0$ and $(x, y) \sim \mathcal{D}$. Therefore, α and β follow $O(1)$. \square

G PROOF OF THEOREMS 2.4 AND 2.5 WHEN THE NETWORK HAS A VECTOR OUTPUT AS $l \geq 1$

From Theorem D.1, the original/trained NN $f_{\mathbf{W}^{(k)}}(x)$ is equivalent to the set of l trained subnetworks.

By combining this fact and each result of Theorem F.2 and Theorem F.4, Theorems 2.4 and 2.5 when $l \geq 1$ can be directly obtained respectively.