

---

# Subquadratic Overparameterization for Shallow Neural Networks

---

Chaehwan Song<sup>1\*</sup>Ali Ramezani-Kebrya<sup>1\*</sup>Thomas Pethick<sup>1</sup>Armin Eftekhari<sup>2†</sup>Volkan Cevher<sup>1</sup><sup>1</sup>Laboratory for Information and Inference Systems (LIONS), EPFL    <sup>2</sup>Umea University

ali.ramezani@epfl.ch

## Abstract

Overparameterization refers to the important phenomenon where the width of a neural network is chosen such that learning algorithms can provably attain zero loss in nonconvex training. The existing theory establishes such global convergence using various initialization strategies, training modifications, and width scalings. In particular, the state-of-the-art results require the width to scale quadratically with the number of training data under standard initialization strategies used in practice for best generalization performance. In contrast, the most recent results obtain linear scaling either with requiring initializations that lead to the “lazy-training”, or training only a single layer. In this work, we provide an analytical framework that allows us to adopt standard initialization strategies, possibly avoid lazy training, and train all layers simultaneously in basic shallow neural networks while attaining a desirable subquadratic scaling on the network width. We achieve the desiderata via Polyak-Lojasiewicz condition, smoothness, and standard assumptions on data, and use tools from random matrix theory.

## 1 Introduction

Training a neural network involves solving a nonconvex optimization problem, which, in theory, might trap first-order methods such as gradient descent to fall in bad local minima or saddle points. However, empirical evidence suggests that first-order methods with random initialization can consistently find a global minimum, even with randomized labels [46]. Demystifying this observation is of central interest to deep learning.

Recently, a line of research [48, 4, 11, 30, 41, 12, 39] suggests that such an empirical success can possibly be explained by the *overparameterization* of neural networks, whose number of parameters exceeds the number of training data  $n$ . In particular, gradient descent converges linearly fast to a global optimum in a number of problems with models that have wide hidden layers [48, 12, 41].

Despite of these remarkable results, the natural key question “*How much should we overparameterize a neural network?*” remains open even for the toy example of two-layer neural networks. On one hand, it is widely accepted that, for two-layer neural networks, the number of parameters should grow linearly with  $n$  (e.g., [22, 39]). On the other hand, theoretical results either require much more parameters, or they are established under restrictive settings. Specifically,

---

\*Equal contributions.

†This work was done while Armin Eftekhari was at EPFL.

**Table 1:** Scaling with the number of training data in the overparameterization regime. QL=quadratic loss, CLL=convex and Lipschitz loss, SD=separable data.

Depth	Algorithm	Setting	Activation	Scaling	Reference
2	GD on layer 1	QL	ReLU	$\tilde{\Omega}(n^2)$	Oymak and Soltanolkotabi [39]
$L$	GD on layer $L$	CLL	ReLU	$\tilde{\Omega}(n)$	Kawaguchi and Huang [22]
2	GD	SD	ReLU	$\tilde{\Omega}(n^2)$	Song and Yang [41]
2	GD	SD and QL	ReLU	$\tilde{\Omega}(n^6)$	Du et al. [12]
$L$	GD	SD and QL	ReLU	$\tilde{\Omega}(n^8 L^{12})$	Zou and Gu [47]
2	GD	QL	Smooth	$\tilde{\Omega}(n^{\frac{3}{2}})$	<b>This paper</b>

- Kawaguchi and Huang [22] has proven the ideal  $\tilde{\Omega}(n)$  scaling for deep neural networks. However, they apply gradient descent only to the last layer, which is not the case in practical scenarios.
- A similar issue exists in [41, 39], where the authors have shown that  $\tilde{\Omega}(n^2)$  parameters suffice for two-layer neural networks, but only the first layers are trained. Furthermore, even with infinite width, Oymak and Soltanolkotabi [39] cannot guarantee zero training error with probability approaching to one.

The goal of this paper is to close the gap between theory and practice, without resorting to unrealistic assumptions such as those discussed above. We sharpen the results of Oymak and Soltanolkotabi [39] by proving that, with proper random initialization of each layer, training error approaches to zero with high probability, exponentially fast in the width of the network. In addition, we show that only  $\tilde{\Omega}(n^{\frac{3}{2}})$  parameters suffice such that gradient descent converges to a global minimum with linear rate, which improves upon the state-of-the-art by a factor of  $\tilde{O}(n^{\frac{1}{2}})$ . We summarize the bounds on the number of parameters in terms of  $n$  in Table 1.

While our analysis on gradient descent focuses on training error, it has been observed that overparameterization can lead to poor *generalization*. In particular, [7, 45, 15] have observed the phenomenon of *lazy training*. Chizat et al. [7] has explained lazy training as a model behaves similar to its linearization around the initialization. It is known that an overparameterized neural network is likely to be trapped in the lazy regime since the parameters will hardly vary over the course of training with gradient descent [12, 30, 48]. The same phenomenon has been observed for infinitely wide neural networks [20]. In this paper, we provide theoretical guidance to possibly avoid lazy training through proper initialization. Experimental results confirm that lazy training might be avoided with our theoretically inspired initialization so that the issues reported in [7] do not apply.

## 1.1 Summary of contributions

- We first focus on a general minimization problem assuming that the loss function satisfies Polyak-Łojasiewicz (PL) condition. We find sufficient conditions in terms of initialization for the convergence of gradient flow and gradient descent to a global minimum.
- We then focus on the special problem of training a two-layer neural network with quadratic loss and smooth activation, and show that  $\tilde{\Omega}(n^{\frac{3}{2}})$  parameters are sufficient for gradient descent to converge to a global minimum with linear rate and probability approaching to one. We achieve *linear scaling* for the width when the number of input features is in  $\tilde{\Omega}(\sqrt{n})$ .
- We theoretically guide how to initialize the parameters of a neural network in the overparameterized regime of interest while possibly avoiding lazy training.

## 1.2 Further related work

In terms of techniques, our paper is closely related to [38, 39]. Similar to our Theorem 3, Oymak and Soltanolkotabi [39, Theorem 2.1] showed that gradient descent converges with linear rate when the Jacobian of the nonlinear mapping has smooth deviations, and the number of parameters grows quadratically with  $n$ . However, Oymak and Soltanolkotabi [39] assumed that gradient descent updates only the first layer. In this paper, we consider the case where gradient descent updates both layers simultaneously, and show that it suffices to have  $\tilde{\Omega}(n^{\frac{3}{2}})$  parameters with a linear rate of convergence.

ReLU is an important instance of activation functions that does not satisfy the smoothness assumption. A line of research aims to relax this assumption by instead assuming the data is separable. For shallow neural networks, Du et al. [12] proved that gradient descent finds a global minimum if the width of the network scales  $\tilde{\Omega}(n^6)$  assuming that no two data points are parallel. In a similar setting, Song and Yang [41] established convergence to a global minimum with the sufficient width of  $\tilde{\Omega}(n^2)$ . As a result, in the absence of the smoothness assumption, these papers require substantially more number of parameters to guarantee convergence to a global minimum.

The theoretical bounds for deep neural networks are even worse. For instance, Allen-Zhu et al. [1] required the total number of parameters of  $\Omega(n^{24}L^{12})$  where  $L$  is the number of layers. Zou and Gu [47] improved the scaling to  $\Omega(n^8L^{12})$ . In our setting, *i.e.*,  $L = 2$ , these bounds become vacuous in most interesting regimes. Further, in [22], the authors showed that  $\tilde{\Omega}(n)$  parameters is enough to achieve global convergence under the assumption that gradient descent updates only the last layer, which essentially reduces the problem to a simple least-squares regression.

Recently, Ji and Telgarsky [21], Chen et al. [5] showed that a polylogarithmic width suffices to achieve convergence for shallow and deep neural networks in an ergodic sense. We note that this is a weaker notion of convergence compared to the one we consider.

Li et al. [29] showed that gradient descent along with early stopping are robust to label noise on a constant fraction of labels in an overparameterized network. However, only the first layer is optimized in [29]. For possibly overparameterized and linear networks, Eftekhari [14] showed that gradient flow can successfully avoid lazy training assuming that the network has a layer with a single neuron. We note that our analysis does not require those restrictions.

Under an assumption similar to PL condition, Zou et al. [48] studied the problem of binary classification for a deep network with ReLU activation, which is a different problem compared to ours. In [42], the authors proved that gradient descent with overparameterization achieves zero-approximation when the underlying function that generates the labels has low-rank approximation. Their scaling requires perfect information about the target function, which is not the case in our paper. Under a variant of Xavier initialization, Daniely [9] found near optimal scaling for a binary classification problem trained by stochastic gradient descent. We note that the setting considered in our paper is more challenging than binary classification. Our results establish a new state-of-the-art on the required number of parameters in a nonrestrictive setting when both layers are trained at the same time. Recently, Nguyen and Mondelli [35] obtained subquadratic scaling for a deep neural network with pyramidal structure under an initialization that leads to lazy training. Our results do not have such restrictions.

Mean-field analysis was used to approximate a target distribution of parameters of a neural network by the empirical distributions [34, 33]. However, these results do not provide useful bounds on the scaling in terms of  $n$ , which is our focus in this paper.

Liu et al. [32] established global convergence when the function to minimize satisfies a variant of PL condition (local PL condition) assuming the map is Lipschitz continuous, which is not the case in our paper. Liu et al. [31] characterized the constancy of the neural tangent kernel via scaling properties of the norm of the Hessian matrix of the network. In this work, we focus on obtaining a sufficient number of parameters for gradient descent to converge to a global minimum with linear rate.

**Notation.** We use  $\|\cdot\|$  to represent the Euclidean norm of a vector and Frobenius norm of a matrix. We use  $\nabla$  to denote the Jacobian of a vector-valued and gradient of a scalar-valued function and  $\nabla\Phi(a)\{b\}$  to represent the directional derivative of  $\Phi$  along  $b$ . We use  $\odot$  and  $\otimes$  to denote the Hadamard (entry-wise) product and Kronecker product, respectively. For  $A \in \mathbb{R}^{m \times n}$  and  $t \in \mathbb{Z}_+$ , we denote  $A^{*t} \in \mathbb{R}^{m^t \times n}$  with its  $a$ -th column defined as  $\text{vec}(x_a \otimes \cdots \otimes x_a) \in \mathbb{R}^{m^t}$ . We use lower-case bold font to denote vectors. Sets and scalars are represented by calligraphic and standard fonts, respectively. We use  $[n]$  to denote  $\{1, \dots, n\}$  for an integer  $n$ . We use  $\tilde{O}$  and  $\tilde{\Omega}$  to hide logarithmic factors and use  $\lesssim$  to ignore terms up to constant and logarithmic factors.

## 2 Problem, definitions, and assumptions

In this section, we set up a general compositional optimization problem. Then we focus on the special case of shallow neural networks in Section 5.

Let  $\mathbf{w} \in \mathbb{R}^d$  denote a parameter vector where  $d$  denotes the number of parameters. In a neural network,  $\mathbf{w}$  consists of weights and biases of all layers. We consider the minimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} h(\mathbf{w}) \quad (1)$$

where  $h : \mathbb{R}^d \rightarrow \mathbb{R}_+$  is the composition of a loss function  $f : \mathbb{R}^{\bar{d}} \rightarrow \mathbb{R}_+$  and a nonlinear and nonconvex function  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^{\bar{d}}$ :

$$h(\mathbf{w}) = f(\Phi(\mathbf{w})) = f(\mathbf{z}) \quad (2)$$

where  $\mathbf{z} = \Phi(\mathbf{w})$ .

Before providing the details, let us highlight the simple idea behind the argument (see also [38]). Let  $\mathbf{w}_0$  and  $\bar{\mathbf{w}}$  denote the initial point and limit point when the gradient descent algorithm is run with some learning rate, respectively. The precise formulation of gradient descent is provided in Section 4. Let  $\nabla\Phi^*(\bar{\mathbf{w}}) : \mathbb{R}^{\bar{d}} \rightarrow \mathbb{R}^d$  denote the adjoint operator of  $\nabla\Phi(\bar{\mathbf{w}})$ . Since  $\bar{\mathbf{w}}$  is a first-order stationary point of  $h$ , we have

$$0 = \nabla h(\bar{\mathbf{w}}) = \nabla\Phi^*(\bar{\mathbf{w}}) \{\nabla f(\bar{\mathbf{z}})\}$$

where  $\bar{\mathbf{z}} = \Phi(\bar{\mathbf{w}})$ . Suppose that  $\nabla\Phi^*(\bar{\mathbf{w}})$  is a nonsingular operator. Then  $\nabla f(\bar{\mathbf{z}}) = 0$ . If  $\bar{\mathbf{z}}$  is a global minimizer of  $f$ , then  $\bar{\mathbf{w}}$  is a global minimizer of  $h$ . To prove global convergence, it suffices to show that  $\nabla\Phi^*$  is nonsingular within a neighborhood of the initialization  $\mathbf{w}_0$ , and that points reached by gradient descent remain within this neighborhood. We will prove that both statements hold with high probability for shallow neural networks.

We first define two notions that are useful to state a key lemma for our main results:

**Definition 1** (Near-isometry). *A linear mapping  $T : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$  is  $(\mu, \nu)$ -near-isometry if there exist  $0 < \mu \leq \nu$  such that*

$$\mu \leq \sigma_{\min}(T) \leq \sigma_{\max}(T) \leq \nu. \quad (3)$$

**Definition 2** (Smoothness). *Let  $\beta_\psi > 0$ . A function  $\psi : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$  is  $\beta_\psi$ -smooth, if for all  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{d_1}$ , we have*

$$\sigma_{\max}(\nabla\psi(\mathbf{u}) - \nabla\psi(\mathbf{v})) \leq \beta_\psi \|\mathbf{u} - \mathbf{v}\|. \quad (4)$$

The following lemma shows that a smooth function, which is near-isometry at initialization, remains near-isometry for all nearby points of the initialization.

**Lemma 1.** *Suppose that  $\Phi$  is  $\beta_\Phi$ -smooth and  $\nabla\Phi^*(\mathbf{w}_0)$  is  $(\mu_\Phi, \nu_\Phi)$ -near-isometry. Then, for all  $\mathbf{w} \in \text{ball}(\mathbf{w}_0, \rho_\Phi)$ , we have*

$$\frac{\mu_\Phi}{2} \leq \sigma_{\min}(\nabla\Phi^*(\mathbf{w})) \leq \sigma_{\max}(\nabla\Phi^*(\mathbf{w})) \leq \frac{3\nu_\Phi}{2} \quad (5)$$

where

$$\rho_\Phi = \frac{\mu_\Phi}{2\beta_\Phi}. \quad (6)$$

Intuitively, if  $\nabla\Phi^*(\mathbf{w}_0)$  is a  $(\mu_\Phi, \nu_\Phi)$ -near-isometry, then one would expect  $\nabla\Phi^*$  to remain near-isometry for all nearby points.

**Definition 3** (PL condition [3]). *A function  $\psi : \mathbb{R}^{d_1} \rightarrow \mathbb{R}$  satisfies the PL condition if there exists  $\alpha_\psi > 0$  such that, for all  $\mathbf{u} \in \mathbb{R}^{d_1}$ , we have*

$$\psi(\mathbf{u}) \leq \frac{\|\nabla\psi(\mathbf{u})\|^2}{2\alpha_\psi}. \quad (7)$$

We note that strongly convex functions satisfy a minor variation of the PL condition in (7).

In our analysis, we will assume that  $\Phi$  and  $f$  satisfy the following properties:

**Assumption 1** (Basic assumptions for  $\Phi, f$ ).

- $\Phi$  is twice-differentiable and  $\beta_\Phi$ -smooth.
- $f$  is twice-differentiable, satisfies the PL condition with  $\alpha_f$ , and  $\min f(\mathbf{z}) = 0$ .

Despite  $f$  satisfies the PL condition, the nonconvex  $\Phi$  can render  $h$  nonconvex, and hence difficult to minimize in theory. However, we show that fast convergence of gradient descent to a global minimum can be established with appropriate initialization.

The intuition behind these assumptions is that to achieve nonsingularity of  $\nabla\Phi^*$ , we approximate  $\nabla\Phi^*(\mathbf{w}_0)$  at initialization and bound  $\nabla\Phi^*(\mathbf{w}_0) - \nabla\Phi^*(\mathbf{w}_i)$  at iteration  $i$  using the fact that  $\|\mathbf{w}_0 - \mathbf{w}_i\|$  is sufficiently small by the overparameterization. In the special case of shallow neural networks, we expect a similar argument applies even when the activation function is ReLU. Adapting our analysis for such extensions is an interesting area of future work.

### 3 Gradient flow

In this section, we consider gradient flow, which can be viewed as the limit of gradient descent for infinitesimally small learning rates. Inspired by the analysis of gradient flow, we provide an upper bound on the length of the trajectory traversed by gradient descent iterates and then find a sufficient condition in terms of initialization to establish its convergence to a global minimum. We focus on gradient descent in Section 4.

Let  $t \geq 0$  and consider the gradient flow, which is initialized at  $\mathbf{w}_0 \in \mathbb{R}^d$  and traverses the curve  $\gamma: \mathbb{R}_+ \rightarrow \mathbb{R}^d$ , given by

$$\dot{\gamma}(t) = \frac{d\gamma(t)}{dt} = -\nabla h(\gamma(t)) \quad (8)$$

where  $\gamma(0) = \mathbf{w}_0$ .

We now calculate the length of the curve  $\gamma$ . Suppose that  $\nabla\Phi^*(\mathbf{w}_0)$  is  $(\mu_\Phi, \nu_\Phi)$ -near-isometry. Using Lemma 1, in the following lemma, we control the length inside of ball  $(\mathbf{w}_0, \rho_\Phi)$ . See Appendix B for the proof.

**Lemma 2.** *Let  $t \geq 0$  and let  $\ell(t)$  denote the length of the curve  $\gamma$  in (8), restricted to the interval  $[0, t]$ . Let  $t_\Phi \in (0, \infty]$  be the smallest value such that  $\gamma(t_\Phi) \notin \text{ball}(\mathbf{w}_0, \rho_\Phi)$ . Suppose  $\nabla\Phi^*(\mathbf{w}_0)$  is  $(\mu_\Phi, \nu_\Phi)$ -near-isometry. Then, for all  $t \leq t_\Phi$ , we have*

$$\ell(t) = O\left(\frac{\nu_\Phi \sqrt{h(\mathbf{w}_0)}}{\mu_\Phi^2 \sqrt{\alpha_f}}\right).$$

Lemma 2 implies that if the objective value at initialization,  $h(\mathbf{w}_0)$ , is sufficiently small, then we can localize gradient flows to a region around  $\mathbf{w}_0$ . Combining with Lemma 1, we show that the limit point of gradient flow is a global minimum. This theorem is formally stated below.

**Theorem 1** (Gradient flow). *Let  $\mathbf{w}_0 \in \mathbb{R}^d$ . Suppose that  $\Phi$  and  $f$  satisfy Assumption 1 and  $\nabla\Phi^*(\mathbf{w}_0)$  is  $(\mu_\Phi, \nu_\Phi)$ -near-isometry. If  $\mathbf{w}_0$  satisfies*

$$h(\mathbf{w}_0) = O\left(\frac{\alpha_f \mu_\Phi^6}{\beta_\Phi^2 \nu_\Phi^2}\right), \quad (9)$$

*then the gradient flow  $\gamma$  in (8) converges to a global minimum.*

*Proof of Theorem 1.* Proper initialization in (9) ensures  $\ell(t_\Phi) < \rho_\Phi$ , which implies that

$$\|\gamma(t_\Phi) - \mathbf{w}_0\| = \|\gamma(t_\Phi) - \gamma(0)\| < \rho_\Phi. \quad (10)$$

Therefore,  $\gamma(t) \in \text{ball}(\mathbf{w}_0, \rho_\Phi)$  for all  $t \geq 0$ , and the length of  $\gamma$  is upper bounded by  $\rho_\Phi$  using Lemma 2. Hence, the gradient flow  $\gamma$  converges, i.e., the limit point  $\bar{\mathbf{w}} \in \mathbb{R}^d$  exists and satisfies

$$\|\bar{\mathbf{w}} - \mathbf{w}_0\| \leq \rho_\Phi. \quad (11)$$

Combining (5) and (11), we have

$$\frac{\mu_\Phi}{2} \leq \sigma_{\min}(\nabla\Phi^*(\bar{\mathbf{w}})) \leq \sigma_{\max}(\nabla\Phi^*(\bar{\mathbf{w}})) \leq \frac{3\nu_\Phi}{2}.$$

In particular, we note that  $\nabla\Phi^*(\bar{\mathbf{w}})$  is nonsingular. So we have  $\nabla f(\bar{\mathbf{z}}) = 0$ . Since  $f$  satisfies the PL condition in (7),  $\bar{\mathbf{z}}$  is a global minimizer of  $f$ , and  $\bar{\mathbf{w}}$  is a global minimizer of  $h$  in (1).  $\square$

## 4 Gradient descent

We now view gradient descent as the discretization of gradient flow, and show that a similar argument as in Section 3 holds for gradient descent.

Let  $\eta > 0$  denote the learning rate and let  $i \geq 0$ . The gradient descent update rule is given by

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \eta \nabla h(\mathbf{w}_i). \quad (12)$$

To study gradient descent, in addition to the previous assumptions on  $\Phi$  and  $f$  for the case of gradient flow described in Theorem 1, we also assume that  $f$  is smooth, *i.e.*, there exists  $\beta_f \geq 0$  such that, for all  $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^{\bar{d}}$ , we have

$$f(\mathbf{z}) - f(\mathbf{z}') \leq \langle \mathbf{z} - \mathbf{z}', \nabla f(\mathbf{z}') \rangle + \frac{\beta_f}{2} \|\mathbf{z} - \mathbf{z}'\|^2.$$

Smoothness of  $f$  allows safe discretization of gradient flow without deviating too much from its trajectory. The following result is the analogue of Theorem 1 for gradient descent; see Appendix C for the proof.

**Theorem 2** (Gradient descent). *Let  $\mathbf{w}_0 \in \mathbb{R}^d$ . Suppose that  $\Phi$  and  $f$  satisfy Assumption 1,  $f$  is  $\beta_f$ -smooth, and  $\nabla \Phi^*(\mathbf{w}_0)$  is  $(\mu_\Phi, \nu_\Phi)$ -near-isometry. Suppose that gradient descent is executed with sufficiently small learning rate*

$$\eta = O\left(\frac{1}{\beta_\Phi \|\nabla f(\Phi(\mathbf{w}_0))\| + \beta_f \mu_\Phi^2 + \beta_f \nu_\Phi^2}\right), \quad (13)$$

and  $\mathbf{w}_0$  satisfies (9).

Then the sequence of iterates  $\{\mathbf{w}_i\}_{i \geq 0}$  converges to a global minimum of  $h$  exponentially fast.

In addition, the rate of convergence is given by

$$h(\mathbf{w}_i) \leq (1 - C\eta\alpha_f\mu_\Phi^2)^i h(\mathbf{w}_0) \quad (14)$$

where  $C$  is a universal constant.

To prove Theorem 2, we first compute the length of the trajectory traversed by gradient descent iterates. We then use the smoothness of  $f$  and follow the descent inequality to lower bound  $f(\mathbf{z}_i) - f(\mathbf{z}_{i+1})$ . Finally, we compute the local Lipschitz constant of  $f$ .

**Remark 1.** *The idea of initializing a nonconvex problem close to a global minimum has a long history in nonconvex optimization, particularly in matrix factorization; see [6] and references therein. The observation that the length of the learning trajectory is short in the overparameterization regime has a precedent in [12, 38]. From an algorithmic perspective, the idea of linearizing  $\Phi$  when minimizing  $h = f \circ \Phi$  is studied in nonlinear regression and the Gauss-Newton method [36].*

In order to apply Theorem 2, the key step is to verify that  $h(\mathbf{w}_0)$  satisfies (9). In Section 5, we focus on the special case of shallow neural networks and improve the state of the art.

## 5 Shallow neural networks

In this section, we consider the problem of training shallow neural networks with gradient descent. Our strategy is to cast this problem as a special case of problem (1) and then apply Theorem 2 to establish global convergence. We start with the formal problem statement.

### 5.1 Setup, assumptions, and initialization

Consider a shallow neural network with  $d_0$  inputs, one hidden layer that consists of  $d_1$  hidden nodes, and  $d_2$  outputs. This shallow network is specified by the map

$$\begin{aligned} \mathbb{R}^{d_0} &\mapsto \mathbb{R}^{d_2} \\ \mathbf{x} &\mapsto V \cdot \phi(W\mathbf{x}), \end{aligned} \quad (15)$$



where  $W \in \mathbb{R}^{d_1 \times d_0}$ ,  $V \in \mathbb{R}^{d_2 \times d_1}$ , and  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is an activation function, which is applied entry-wise. Let  $\mathbf{x}_i \in \mathbb{R}^{d_0}$  and  $y_i \in \mathbb{R}^{d_2}$  denote the  $i$ -th training data and label, respectively, for  $i \in [n]$ . By concatenating the training data and their labels, we form the matrices  $X \in \mathbb{R}^{d_0 \times n}$  and  $Y \in \mathbb{R}^{d_2 \times n}$ . Let denote  $\Theta = (W, V) \in \mathbb{R}^{d_1 \times d_0} \times \mathbb{R}^{d_2 \times d_1}$  and  $Z = \Phi(\Theta) = V \cdot \phi(WX) \in \mathbb{R}^{d_2 \times n}$ . The fitting problem can be cast as (1) where

$$h(\Theta) = f(\Phi(\Theta)) = \|V\phi(WX) - Y\|^2. \quad (16)$$

**Remark 2.** We assume that the activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is twice-differentiable. Despite this assumption excludes the popular ReLU, it is still possible to apply our results to smooth approximations of ReLU such as the softplus or Gaussian error Linear Units (GeLU) [18, 35]. We note that softplus [13] or GeLU [10] often achieve similar or superior performance compared to the ReLU [8, 16, 24, 23, 44].

**Definition 4** (Hermite norm [37]). Let  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ . The Hermite norm of  $\phi$  is given by  $\|\phi\|_{\mathcal{H}} = \sqrt{\sum_{i=0}^{\infty} c_i^2}$  where  $c_i$  denotes the  $i$ -th Hermite coefficients of  $\phi$  given by:

$$c_i = \langle \phi, q_i \rangle_{\mathcal{H}} = \frac{1}{\sqrt{2\pi}} \int \phi(x) q_i(x) \exp\left(-\frac{x^2}{2}\right) dx$$

and  $q_i : \mathbb{R} \rightarrow \mathbb{R}$  is the  $i$ -th Hermite polynomial for  $i \geq 0$ .

In this section, we assume that  $\phi$ ,  $f$ , and data satisfy the following properties:

**Assumption 2** (Assumptions for shallow neural networks).

- $\phi$  is twice-differentiable,  $\phi(0) = 0$ ,  $\sup_a |\dot{\phi}(a)| = \dot{\phi}_{\max} < \infty$ ,  $\sup_a |\ddot{\phi}(a)| = \ddot{\phi}_{\max} < \infty$ , and  $\|\phi\|_{\mathcal{H}} < \infty$ . The loss function  $f$  is quadratic (16).
- $\|\mathbf{x}_i\| = 1$ ,  $\|Y\| \leq 1$ , and  $\sigma_{\max}(V_k) = O\left(\frac{\dot{\phi}_{\max}}{\ddot{\phi}_{\max}}\right)$  for  $i \in [n]$  and  $k \geq 0$ .

The assumption on  $\phi$  hold for GeLU, sigmoid, and tanh. The assumption  $\phi(0) = 0$  is to simplify the derivations and we suspect that it can be removed at the expense of more complicated expressions. The bounded Hermite norm is a mild assumption, which is used to obtain an upper bound on  $\sigma_{\max}(\phi(W^0 X))$  in terms of the Hermite coefficients of  $\phi$ . See Appendix E.1 for details. The assumption on the data is fairly mild and standard in the overparameterization literature as we can always normalize the data [30, 21]. Similar boundedness assumptions to the last assumption are commonly used in nonconvex optimization to guarantee convergence [25]. Moreover, such a bound naturally holds by applying a projection step to the gradient descent update rule, which we plan to adopt as a future work.

**Initialization.** We first consider the initialization scheme:

$$W_0 \sim \mathcal{N}(0, \omega_1^2), \quad V_0 \sim \mathcal{N}(0, \omega_2^2). \quad (17)$$

In Section 6, we study the implications of our initialization and show how to possibly avoid lazy training by varying  $(\omega_1, \omega_2)$ .

## 5.2 Main results for shallow neural networks

For shallow networks as described above, we verify in Appendix D that the key conditions in Lemma 1 hold with high probability. Combining with Theorem 2, we establish the global convergence guarantees. The proof in Appendix E uses standard tools from random matrix theory to control the random variables involved with initialization. We first estimate variables  $\mu_{\Phi}, \nu_{\Phi}$  defined in Definition 1 and  $\beta_{\Phi}$  in (4) for the neural network described in Section 5.1.

**Lemma 3** (Estimation of  $\mu_{\Phi}, \nu_{\Phi}, \beta_{\Phi}$ ). Suppose that a shallow neural network, which is constructed in Section 5.1, satisfies Assumption 2. Then we have

$$\begin{aligned} \mu_{\Phi} &= \sigma_{\min}(\phi(W_0 X)), \\ \nu_{\Phi} &= \dot{\phi}_{\max} \sigma_{\max}(X) \sigma_{\max}(V_0) + \sigma_{\max}(\phi(W_0 X)), \\ \beta_{\Phi} &= \sqrt{2} \sigma_{\max}(X) \left( \dot{\phi}_{\max} + \ddot{\phi}_{\max} \chi_{\max} \right) \end{aligned} \quad (18)$$

where  $\chi_{\max} = \sup_V \sigma_{\max}(V)$ .

**Remark 3.** The terms  $\sigma_{\min}(\phi(W_0X))$  and  $\sigma_{\max}(\phi(W_0X))$  in (18) play a critical role in our analysis. In [12, 41], strict positivity of the eigenvalues of Gram matrix is the primary tool to show the convergence. Oymak and Soltanolkotabi [39] also followed a similar argument using the neural network covariance matrix. The underlying intuition seems similar to Lemma 3. However, the resulting bounds are different since gradient descent updates  $(W, V)$  simultaneously in our problem setup, which is more realistic.

By combining Lemma 3 and the results on global convergence of gradient descent in Section 4, we establish global convergence for shallow neural network.

**Theorem 3** (Shallow network with gradient descent). *Consider the shallow network described in Section 5.1 that satisfies Assumption 2 and  $\tau^{r_1}|\phi(a)| \leq |\phi(\tau a)| \leq \tau^{r_2}|\phi(a)|$  for all  $a$ ,  $0 < \tau < 1$ , and some constants  $r_1, r_2$ .<sup>3</sup> Suppose that  $\Theta_0$  is randomly initialized as in (17) with  $\omega_1$  and  $\omega_2$ , which satisfy*

$$\omega_1\omega_2 \lesssim \frac{1}{\sqrt{d_0d_1}}, \quad (19)$$

and suppose that the hidden layer width  $d_1$  satisfies

$$d_1 = \tilde{\Omega} \left( \xi(\mathcal{C}_\delta, t, \phi, \{c_i\}_{i \geq 0}) \frac{\sigma_{\max}(X)^2 \sqrt{n}}{\sigma_{\min}(X^{*t})^3} \right) \quad (20)$$

where  $\mathcal{C}_\delta$  is a set of constants,  $\xi$  is a term independent to  $d_0, n$ ,  $t$  is a constant such that  $n \simeq d_0^t$ , and  $X^{*t} \in \mathbb{R}^{d_0^t \times n}$  is derived from Khatri-Rao product with its  $a$ -th column defined as  $\text{vec}(x_a \otimes \dots \otimes x_a) \in \mathbb{R}^{d_0^t}$ . Then gradient descent converges to a global minimum exponentially fast with probability at least  $1 - \psi(\phi, \xi, d_0, d_1, d_2, X)$ .<sup>4</sup> See Appendix E.6 for the exact expressions of  $\xi$  and  $\psi$ .

**Remark 4.** Theorem 3 shows that, with sufficient degree of overparameterization, gradient descent finds a global minimum, except with an arbitrary small probability. Note that we need two conditions for Theorem 3 to hold, both of which are related to the overparameterization of the network. The condition (19) is for the concentration of random matrices, to make  $\psi$  arbitrary small, and (20) is for the locality of gradient descent.

### 5.3 Order analysis

We first decompose the random matrix  $\phi(X^\top W_0^\top)\phi(W_0X)$  into independent random matrices. We then apply concentration inequalities to establish an upper bound on  $\sigma_{\max}(\phi(W_0X))$  and a lower bound on  $\sigma_{\min}(\phi(W_0X))$  through the Hermite decomposition of  $\phi(W_0X)$  and note that with high probability,

$$\sqrt{\frac{c_t^2}{t!}} d_1 \sigma_{\min}(X^{*t}) \lesssim \sigma_{\min}(\phi(W_0X)) \lesssim \sigma_{\max}(\phi(W_0X)) \lesssim \sqrt{c_0^2 d n}.$$

We also find an upper bound on  $h(\Theta_0)$  at initialization. Substituting  $\nu_\Phi, \mu_\Phi, \beta_\Phi$  into (9), we obtain the sufficient condition in (20). We note that  $\xi(\mathcal{C}_\delta, t, \phi, \{c_i\}_{i \geq 0})$  can be viewed as a constant w.r.t.  $d_0, d_1$ , and  $n$ . For  $t = 1$ , it requires  $n \simeq d_0$ , which is not a common setting in practice. For  $t \geq 2$ , we suppose that  $n \simeq d_0^t$ , which is the case in practice and estimate  $\sigma_{\max}(X) \simeq \sqrt{\frac{n}{d_0}}$  and  $\sigma_{\min}(X^{*t}) \simeq \sqrt{\frac{n}{d_0^t}} \simeq 1$  along the lines of [39, Section 2.1]. Substituting  $\sigma_{\max}(X)$  and  $\sigma_{\min}(X^{*t})$  into (20), we have

$$d_1 \gtrsim \frac{n^{\frac{3}{2}}}{d_0}. \quad (21)$$

Therefore, the overall overparameterization degree becomes  $d_0d_1 \simeq \tilde{\Omega}(n^{\frac{3}{2}})$ , which is sufficient for gradient descent to find a global minimum at a linear rate except with an arbitrary small probability. We note that an optimal linear scaling for the width  $d_1 \simeq \tilde{O}(n)$  is sufficient when the number of input features is sufficiently large  $d_0 \simeq \tilde{\Omega}(\sqrt{n})$ , which improves upon the results of [39] by a factor of  $\tilde{O}(n^{\frac{1}{2}})$ . Furthermore, unlike [39], we adopt standard initialization strategies in Theorem 3.

<sup>3</sup>The last assumption holds for popular activation functions such as sigmoid, tanh, and ELU, and can be relaxed if  $\omega_1 = 1$  in (17).

<sup>4</sup> $\psi$  can be arbitrary small.



## 6 Lazy training and experimental evaluation

Following the theoretically motivated initialization in Theorem 3, we set  $\omega_1\omega_2 \simeq \frac{1}{\sqrt{d_0d_1}}$ . This gives rise to a broad family of initialization schemes as one varies the ratio  $\omega_2/\omega_1$ . Interestingly, we note that popular initialization schemes such as LeCun [27] and He initialization [17] belong to this family. The purpose of this section is to empirically investigate the impact the choice of this ratio has on generalization of shallow networks.

To this end, we will look at the generalization error of varying initializations in the more practical setting of stochastic gradient descent (SGD). Specifically, we fix the product of the weight initialization  $\omega_1\omega_2$  and then proceed by varying  $\omega_2$ . To ensure that perfect generalization is possible, we adopt the teacher-student setup, where, for the teacher network, we train a two-layer fully connected neural network, on MNIST [26] until SGD reaches zero training error. The student networks are trained for 300 epochs to ensure convergence. The results are shown in Figure 1. We use mean-square loss and a smooth activation function (GeLU [18]) for the student network to match the problem setup as closely as possible.

In Figure 1, we observe that while SGD achieves zero training error for every  $\omega_2$ , as suggested by Theorem 3 applicable in the full batch setting, the generalization ability increases as the ratio  $\omega_2/\omega_1$  grows. It is also interesting to observe that the popular He initialization scheme corresponds to a rather balanced ratio that lies at the boundary of the well-performing region. In our experiments, we used He initialization to fix the value  $\omega_1\omega_2$ . This tendency suggests that a wide family of initialization schemes could generalize well as long as the ratio  $\omega_2/\omega_1$  is not too small.

**Comment on lazy training.** It is important to address the so called lazy regime when generalization is of concern. Let

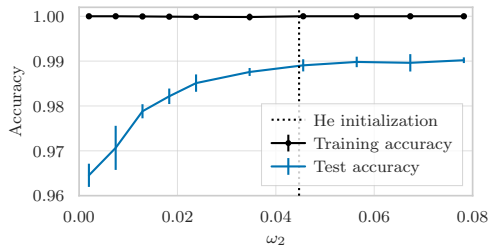
$$\tilde{h}(\Theta) := h(\Theta_0) + \langle \nabla h(\Theta_0), \Theta - \Theta_0 \rangle$$

be the linearized function of  $h$  around  $\Theta_0$  and let  $\Theta_i$  and  $\tilde{\Theta}_i$  denote the iterates of gradient descent at time  $i$ . The lazy training regime refers to the case where the training trajectory stays close to this linearization, i.e.  $\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\| \simeq 0$  for all  $i$  [7]. Such a linearization occurs in infinitely wide neural networks [7], which have been shown to generalize well in some settings [2, 28]. However, in our case of subquadratic (finite) width, the lazy regime might lead to poor generalization. To gain insight on when we cannot avoid it with certainty, let us make a simple rewriting of our network assuming  $\phi$  is homogeneous:

$$\Phi(\Theta) = \alpha V \phi(WX)$$

with  $V_0 \sim \mathcal{N}(0, 1)$  and  $W_0 \sim \mathcal{N}(0, 1)$  where the standard deviations are pulled out as a scaling factor  $\alpha = \omega_1\omega_2 \simeq 1/\sqrt{d_0d_1}$ . This seems to fit into an example in [7, Appendix A.2] suggesting lazy training as  $d_1 \rightarrow \infty$ . However, their results require an odd activation function and infinite width, while our activation function is required not to be odd (see the proof in Appendix E) and our results are under subquadratic (finite) width. Instead, to study lazy training, we explicitly compute an upper bound on  $\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\|$  in Appendix F following [7, Theorem 2.3].

It turns out that the upper bound becomes  $\infty$  when  $\omega_1 \ll \omega_2$ , and it becomes zero when  $\omega_1 \gg \omega_2$ . Our analysis suggests that shallow neural networks can avoid lazy training provided that  $\omega_2/\omega_1 \rightarrow \infty$ . This analysis is corroborated by the empirical results showing that the generalization capability improves as  $\omega_2$  grows in Figure 1. On the other hand, if  $\omega_2/\omega_1 \rightarrow 0$ , then lazy training is bound to happen asymptotically. For details, see Appendix F. Finally, we note that we have not theoretically claimed that our initialization is guaranteed to be non-lazy, since doing so would require establishing a lower bound on  $\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\|$ , which is an interesting problem for future work. Instead, our discussion above only provides a necessary condition for non-lazy training, and a sufficient condition for lazy training.



**Figure 1:** Training and test error on MNIST for different  $\omega_2$ . Error bars indicates the 95% confidence interval computed over 5 independent runs. The setup details are provided in Appendix G.

## 7 Conclusions and future work

In this paper, we prove the linear convergence of first-order methods on subquadratically overparameterized two-layer neural networks with smooth activation functions. Our theoretical analysis is compatible with standard initialization strategies, which can potentially avoid lazy training. We train both layers simultaneously and achieve a desirable subquadratic scaling on the width of the network. In particular, we note that a linear scaling for the width  $d_1 \simeq \tilde{O}(n)$  is sufficient when the number of input features is sufficiently large  $d_0 \simeq \tilde{\Omega}(\sqrt{n})$ . We use tools from random matrix theory under standard assumptions on data and leverage on the assumption that the loss satisfies Polyak-Łojasiewicz condition. We carefully find an explicit upper bound and lower bound on singular values of the outputs of the first layer at initialization with high probability under general initialization.

It is natural to ask whether we can attain similar degree of overparameterization with nonsmooth activation functions such as ReLU. We plan to adapt our analysis for such extensions as a future work. While our analysis provides a necessary condition for avoiding lazy training, it is interesting to develop sufficient conditions in the future. In particular, developing lower bounds on  $\|h(\Theta_i) - \tilde{h}(\hat{\Theta}_i)\|$  will be a key to fully characterize lazy training.

Finally, as a theoretical work, we do not anticipate any potential negative societal impacts of our paper. However, the long-term impacts of our work may depend on how machine learning algorithms are used in society.

## Acknowledgments and Disclosure of Funding

The authors would like to thank Fabian Latorre, Fanghui Liu, and Paul Rolland for helpful discussions.

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement n° 725594 - time-data). This project was sponsored by the Department of the Navy, Office of Naval Research (ONR) under a grant number N62909-17-1-2111. This work was supported by Hasler Foundation Program: Cyber Human Systems (project number 16066). Research was sponsored by the Army Research Office and was accomplished under Grant Number W911NF-19-1-0404.

## References

- [1] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning (ICML)*, 2019.
- [2] Sanjeev Arora, Simon S. Du, Zhiyuan Li, Ruslan Salakhutdinov, Ruosong Wang, and Dingli Yu. Harnessing the power of infinitely wide deep nets on small-data tasks. In *International Conference on Learning Representations (ICLR)*, 2020.
- [3] Jérôme Bolte, Trong Phong Nguyen, Juan Peypouquet, and Bruce W Suter. From error bounds to the complexity of first-order descent methods for convex functions. *Mathematical Programming*, 165:471–507, 2017.
- [4] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a ConvNet with Gaussian inputs. In *International Conference on Machine Learning (ICML)*, 2017.
- [5] Zixiang Chen, Yuan Cao, Difan Zou, and Quanquan Gu. How much over-parameterization is sufficient to learn deep ReLU networks? In *International Conference on Learning Representations (ICLR)*, 2021.
- [6] Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing (TSP)*, 67:5239–5269, 2019.
- [7] Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in neural information processing systems (NeurIPS)*, 2019.

- [8] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (ELUs). In *International Conference on Learning Representations (ICLR)*, 2016.
- [9] Amit Daniely. Neural networks learning and memorization with (almost) no over-parameterization. In *Advances in neural information processing systems (NeurIPS)*, 2020.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805v2*, 2019.
- [11] Simon S. Du and Jason D. Lee. On the power of over-parametrization in neural networks with quadratic activation. In *International Conference on Machine Learning (ICML)*, 2018.
- [12] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations (ICLR)*, 2019.
- [13] Charles Dugas, Yoshua Bengio, François Bélisle, Claude Nadeau, and René Garcia. Incorporating second-order functional knowledge for better option pricing. In *Advances in neural information processing systems (NeurIPS)*, 2000.
- [14] Armin Eftekhari. Training linear neural networks: Non-local convergence and complexity results. In *International Conference on Machine Learning (ICML)*, 2020.
- [15] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of lazy training of two-layers neural networks. In *Advances in neural information processing systems (NeurIPS)*, 2019.
- [16] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of Wasserstein GANs. In *Advances in neural information processing systems (NeurIPS)*, 2017.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [18] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (GELU). *arXiv preprint arXiv:1606.08415v4*, 2020.
- [19] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of American Statistical Association*, 58:13–30, 1963.
- [20] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems (NeurIPS)*, 2018.
- [21] Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow ReLU networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- [22] Kenji Kawaguchi and Jiaoyang Huang. Gradient descent finds global minima for generalizable deep neural networks of practical sizes. In *Annual Allerton Conference on Communication, Control, and Computing*, 2019.
- [23] Youngjin Kim, Minjung Kim, and Gunhee Kim. Memorization precedes generation: Learning unsupervised GANs with memory networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [24] Abhishek Kumar, Prasanna Sattigeri, and Tom Fletcher. Semi-supervised learning with GANs: Manifold invariance with improved inference. In *Advances in neural information processing systems (NeurIPS)*, 2017.
- [25] Fabian Latorre, Armin Eftekhari, and Volkan Cevher. Fast and provable ADMM for learning with generative priors. In *Advances in neural information processing systems (NeurIPS)*, 2019.

- [26] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86:2278–2324, 1998.
- [27] Yann A. LeCun, Léon Bottou, Genevieve B. Orr, and Klaus-Robert Müller. *Efficient BackProp*. In *Neural networks: Tricks of the Trade*. Springer, 2012.
- [28] Jaehoon Lee, Samuel S. Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. In *Advances in neural information processing systems (NeurIPS)*, 2020.
- [29] Mingchen Li, Mahdi Soltanolkotabi, and Samet Oymak. Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.
- [30] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in neural information processing systems (NeurIPS)*, 2018.
- [31] Chaoyue Liu, Libin Zhu, and Mikhail Belkin. On the linearity of large non-linear models: when and why the tangent kernel is constant. In *Advances in neural information processing systems (NeurIPS)*, 2020.
- [32] Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. In *Advances in neural information processing systems (NeurIPS)*, 2020.
- [33] Yiping Lu, Chao Ma, Yulong Lu, Jianfeng Lu, and Lexing Ying. A mean field analysis of deep resnet and beyond: Towards provably optimization via overparameterization from depth. In *International Conference on Machine Learning (ICML)*, 2020.
- [34] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, 2019.
- [35] Quynh Nguyen and Marco Mondelli. Global convergence of deep networks with one wide layer followed by pyramidal topology. In *Advances in neural information processing systems (NeurIPS)*, 2020.
- [36] J. Nocedal and S. Wright. *Numerical Optimization*. Springer New York, 2006.
- [37] Frank W. J. Olver, Daniel W. Lozier, Ronald F. Boisvert, and Charles W. Clark. *NIST Handbook of Mathematical Functions Paperback and CD-ROM*. Cambridge University Press, 2010.
- [38] Samet Oymak and Mahdi Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *International Conference on Machine Learning (ICML)*, 2019.
- [39] Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 1:84–105, 2020.
- [40] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [41] Zhao Song and Xin Yang. Quadratic suffices for over-parametrization via matrix Chernoff bound. *arXiv preprint arXiv:1906.03593v2*, 2019.
- [42] Lili Su and Pengkun Yang. On learning over-parameterized neural networks: A functional approximation perspective. In *Advances in neural information processing systems (NeurIPS)*, 2019.

- [43] Roman Vershynin. *Introduction to the Non-asymptotic Analysis of Random Matrices*. Cambridge University Press, 2012.
- [44] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853v2*, 2020.
- [45] Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for understanding neural networks. In *Advances in neural information processing systems (NeurIPS)*, 2019.
- [46] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations (ICLR)*, 2017.
- [47] Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. In *Advances in neural information processing systems (NeurIPS)*, 2019.
- [48] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep ReLU networks. *arXiv preprint arXiv:1811.08888v3*, 2018.

## A Proof of Lemma 1

Intuitively, if  $\nabla\Phi^*(\mathbf{w}_0)$  is a  $(\mu_\Phi, \nu_\Phi)$ -near-isometry, then one would expect  $\nabla\Phi^*$  to remain near-isometry for all nearby points. Formally, let  $A, B \in R^{m \times n}$  and let singular values of a matrix are ordered such that  $\sigma_i(A) \geq \sigma_j(A)$  and  $\sigma_i(B) \geq \sigma_j(B)$  for  $1 \leq i \leq j \leq \min\{m, n\}$ . Using Weyl's inequality and for  $i + j - 1 \leq \min\{m, n\}$ , we have:

$$\sigma_{i+j-1}(A + B) \leq \sigma_i(A) + \sigma_j(B). \quad (22)$$

More formally, suppose that  $\mathbf{w} \in \mathbb{R}^d$  satisfies

$$\|\mathbf{w} - \mathbf{w}_0\| \leq \frac{\mu_\Phi}{2\beta_\Phi} = \rho_\Phi. \quad (23)$$

If  $\nabla\Phi^*(\mathbf{w}_0)$  is  $(\mu_\Phi, \nu_\Phi)$ -isometry in the sense of Definition 1, then applying Weyl's inequality (22) along with using smoothness and (23), we have

$$\begin{aligned} \sigma_{\min}(\nabla\Phi^*(\mathbf{w})) &\geq \sigma_{\min}(\nabla\Phi^*(\mathbf{w}_0)) - \sigma_{\max}(\nabla\Phi^*(\mathbf{w}) - \nabla\Phi^*(\mathbf{w}_0)) \\ &\geq \mu_\Phi - \beta_\Phi \|\mathbf{w} - \mathbf{w}_0\| \\ &\geq \frac{\mu_\Phi}{2}. \end{aligned}$$

Using a similar argument, we establish an upper bound  $\sigma_{\max}(\nabla\Phi^*(\mathbf{w}))$ :

$$\sigma_{\max}(\nabla\Phi^*(\mathbf{w})) \leq \sigma_{\max}(\nabla\Phi^*(\mathbf{w}_0)) + \sigma_{\max}(\nabla\Phi^*(\mathbf{w}) - \nabla\Phi^*(\mathbf{w}_0)) \leq \nu_\Phi + \frac{\mu_\Phi}{2} \leq \frac{3\nu_\Phi}{2}.$$

## B Proof of Lemma 2

Let  $t \geq 0$  and denote

$$\zeta(t) = \Phi(\gamma(t)) \quad (24)$$

so we have

$$h(\gamma(t)) = f(\Phi(\gamma(t))) = f(\zeta(t)). \quad (25)$$

Taking the first-order derivative w.r.t.  $t$ , we have

$$\begin{aligned} \dot{\zeta}(t) &= \nabla\Phi(\gamma(t)) \{\dot{\gamma}(t)\} \\ &= -\nabla\Phi(\gamma(t)) \{\nabla h(\gamma(t))\}. \end{aligned} \quad (26)$$

Note that we have

$$\begin{aligned}
\frac{dh(\gamma(t))}{dt} &= \nabla h(\gamma(t)) \{\dot{\gamma}(t)\} \\
&= -\nabla h(\gamma(t)) \{\nabla h(\gamma(t))\} \\
&= -\|\nabla h(\gamma(t))\|^2.
\end{aligned} \tag{27}$$

Length of the segment of the curve  $\gamma_K$  restricted to the interval  $[0, t]$  is given by

$$\begin{aligned}
\ell(t) &= \int_0^t \|\dot{\gamma}(\tau)\| d\tau \\
&= \int_0^t \|\nabla h(\gamma(\tau))\| d\tau \\
&\leq \int_0^t \sigma_{\max}(\nabla \Phi^*(\gamma(\tau))) \cdot \|\nabla f(\zeta(\tau))\| d\tau \\
&\lesssim \nu_{\Phi} \int_0^t \|\nabla f(\zeta(\tau))\| d\tau.
\end{aligned} \tag{28}$$

To control the norm in the last line of (28), we note that

$$\begin{aligned}
-\frac{d\sqrt{f(\zeta(\tau)) - f(\zeta(t))}}{d\tau} &= -\frac{\frac{df(\zeta(\tau))}{d\tau}}{2\sqrt{f(\zeta(\tau)) - f(\zeta(t))}} \\
&= -\frac{\langle \nabla f(\zeta(\tau)), \dot{\zeta}(\tau) \rangle}{2\sqrt{f(\zeta(\tau)) - f(\zeta(t))}} \\
&= \frac{\langle \nabla f(\zeta(\tau)), \nabla \Phi(\gamma(\tau)) \{\nabla h(\gamma(\tau))\} \rangle}{2\sqrt{f(\zeta(\tau)) - f(\zeta(t))}} \\
&= \frac{\|\nabla h(\gamma(\tau))\|^2}{2\sqrt{f(\zeta(\tau)) - f(\zeta(t))}} \\
&\geq \frac{\sigma_{\min}^2(\nabla \Phi^*(\gamma(\tau))) \cdot \|\nabla f(\zeta(\tau))\|^2}{2\sqrt{f(\zeta(\tau)) - f(\zeta(t))}} \\
&\gtrsim \frac{\mu_{\Phi}^2 \cdot \|\nabla f(\zeta(\tau))\|^2}{\sqrt{f(\zeta(\tau)) - f(\zeta(t))}} \\
&\gtrsim \frac{\sqrt{\alpha_f} \mu_{\Phi}^2 \cdot \|\nabla f(\zeta(\tau))\|^2}{\|\nabla f(\zeta(\tau))\|} \\
&= \sqrt{\alpha_f} \mu_{\Phi}^2 \cdot \|\nabla f(\zeta(\tau))\|,
\end{aligned} \tag{29}$$

provided that the denominators are nonzero. Substituting (29) into (28), the desired length is bounded by

$$\begin{aligned}
\ell(t) &\lesssim \nu_{\Phi} \int_0^t \|\nabla f(\zeta(\tau))\| d\tau \\
&\lesssim -\frac{\nu_{\Phi}}{\mu_{\Phi}^2 \sqrt{\alpha_f}} \int_0^t \frac{d\sqrt{f(\zeta(\tau)) - f(\zeta(t))}}{d\tau} d\tau \\
&= \frac{\nu_{\Phi}}{\mu_{\Phi}^2 \sqrt{\alpha_f}} \left( \sqrt{f(\zeta(0))} - \sqrt{f(\zeta(t))} \right) \\
&\leq \frac{\nu_{\Phi} \sqrt{f(\zeta(0))}}{\mu_{\Phi}^2 \sqrt{\alpha_f}} \\
&= \frac{\nu_{\Phi} \sqrt{h(\gamma(0))}}{\mu_{\Phi}^2 \sqrt{\alpha_f}}
\end{aligned}$$



$$= \frac{\nu_\Phi \sqrt{h(\mathbf{w}_0)}}{\mu_\Phi^2 \sqrt{\alpha_f}},$$

which completes the proof of Lemma 2.

## C Proof of Theorem 2

The proof is along the lines of Theorem 1. We first compute the length of the trajectory traversed by gradient descent iterates. Formally, let  $I$  denote the first iteration such that  $\mathbf{w}_I \notin \text{ball}(\mathbf{w}_0, \rho_\Phi)$ . The length of the trajectory traced by  $\{\mathbf{w}_i\}_{i=0}^I$  is upper bounded by

$$\begin{aligned} \ell(I) &:= \sum_{i=0}^{I-1} \|\mathbf{w}_{i+1} - \mathbf{w}_i\| \\ &= \eta \sum_{i=0}^{I-1} \|\nabla h(\mathbf{w}_i)\| \\ &\lesssim \eta \nu_\Phi \sum_{i=0}^{I-1} \|\nabla f(\mathbf{z}_i)\|. \end{aligned} \tag{30}$$

This following lemma is useful for our proof.

**Lemma 4.** *Suppose  $\mathbf{u}, \mathbf{v} \in \text{ball}(\mathbf{w}_0, \rho_\Phi)$ . Then we have  $\|\Phi(\mathbf{u}) - \Phi(\mathbf{v})\| \leq \frac{3\nu_\Phi}{2} \|\mathbf{u} - \mathbf{v}\|$ .*

*Proof.* Using Lemma 1, we establish a bound on  $\|\Phi(\mathbf{u}) - \Phi(\mathbf{v})\|$ :

$$\begin{aligned} \|\Phi(\mathbf{u}) - \Phi(\mathbf{v})\| &= \left\| \int_0^1 \nabla \Phi(\mathbf{v} + t(\mathbf{u} - \mathbf{v})) (\mathbf{u} - \mathbf{v}) dt \right\| \\ &\leq \int_0^1 \|\nabla \Phi(\mathbf{v} + t(\mathbf{u} - \mathbf{v})) (\mathbf{u} - \mathbf{v})\| dt \\ &\leq \frac{3\nu_\Phi}{2} \|\mathbf{u} - \mathbf{v}\|. \end{aligned}$$

□

Let  $i \leq I - 2$ . To control the upper bound in (30), we use the smoothness of  $f$  and Lemma 4 to obtain a standard “descent inequality” as:

$$\begin{aligned} f(\mathbf{z}_i) - f(\mathbf{z}_{i+1}) &\geq \langle \mathbf{z}_i - \mathbf{z}_{i+1}, \nabla f(\mathbf{z}_i) \rangle - \frac{\beta_f}{2} \|\mathbf{z}_{i+1} - \mathbf{z}_i\|^2 \\ &= \langle \Phi(\mathbf{w}_i) - \Phi(\mathbf{w}_{i+1}), \nabla f(\mathbf{z}_i) \rangle - \frac{\beta_f}{2} \|\Phi(\mathbf{w}_{i+1}) - \Phi(\mathbf{w}_i)\|^2 \\ &= \langle \nabla \Phi(\mathbf{w}_i) \{\mathbf{w}_i - \mathbf{w}_{i+1}\}, \nabla f(\mathbf{z}_i) \rangle - \frac{\beta_f}{2} \|\Phi(\mathbf{w}_{i+1}) - \Phi(\mathbf{w}_i)\|^2 \\ &\quad - \langle \Phi(\mathbf{w}_{i+1}) - \Phi(\mathbf{w}_i) - \nabla \Phi(\mathbf{w}_i) \{\mathbf{w}_{i+1} - \mathbf{w}_i\}, \nabla f(\mathbf{z}_i) \rangle \\ &\geq \langle \nabla \Phi(\mathbf{w}_i) \{\mathbf{w}_i - \mathbf{w}_{i+1}\}, \nabla f(\mathbf{z}_i) \rangle - \frac{\beta_f}{2} \|\Phi(\mathbf{w}_{i+1}) - \Phi(\mathbf{w}_i)\|^2 \\ &\quad - \frac{\beta_\Phi}{2} \|\mathbf{w}_{i+1} - \mathbf{w}_i\|^2 \|\nabla f(\mathbf{z}_i)\| \\ &\geq \langle \nabla \Phi(\mathbf{w}_i) \{\mathbf{w}_i - \mathbf{w}_{i+1}\}, \nabla f(\mathbf{z}_i) \rangle - \frac{1}{2} \|\mathbf{w}_{i+1} - \mathbf{w}_i\|^2 \left( \beta_\Phi \|\nabla f(\mathbf{z}_i)\| + \frac{9\beta_f \nu_\Phi^2}{4} \right) \\ &= \eta \langle \nabla \Phi(\mathbf{w}_i) \{\nabla h(\mathbf{w}_i)\}, \nabla f(\mathbf{z}_i) \rangle - \frac{\eta^2}{2} \|\nabla h(\mathbf{w}_i)\|^2 \left( \beta_\Phi \|\nabla f(\mathbf{z}_i)\| + \frac{9\beta_f \nu_\Phi^2}{4} \right) \\ &= \eta \|\nabla h(\mathbf{w}_i)\|^2 - \frac{\eta^2}{2} \|\nabla h(\mathbf{w}_i)\|^2 \left( \beta_\Phi \|\nabla f(\mathbf{z}_i)\| + \frac{9\beta_f \nu_\Phi^2}{4} \right) \end{aligned}$$

$$\begin{aligned}
&= \eta \|\nabla h(\mathbf{w}_i)\|^2 \left( 1 - \frac{\eta \beta_\Phi \|\nabla f(\mathbf{z}_i)\|}{2} - \frac{9\eta \beta_f \nu_\Phi^2}{8} \right) \\
&\gtrsim \eta \mu_\Phi^2 \|\nabla f(\mathbf{z}_i)\|^2 \quad (\text{chain rule and Lemma 1})
\end{aligned}$$

where the fourth inequality holds since  $\|\Phi(\mathbf{a}) - \Phi(\mathbf{b}) - \nabla\Phi(\mathbf{b})(\mathbf{a} - \mathbf{b})\| \leq \frac{\beta_\Phi}{2} \|\mathbf{b} - \mathbf{a}\|^2$  for  $\beta_\Phi$ -smooth  $\Phi$ , and the last line holds provided that  $\eta$  satisfies:

$$\eta \lesssim \frac{1}{\beta_\Phi \max_i \|\nabla f(\mathbf{z}_i)\| + \beta_f \nu_\Phi^2}. \quad (31)$$

We now use the bound above to find an upper bound on  $\sqrt{f(\mathbf{z}_i) - f(\mathbf{z}_{I-1})} - \sqrt{f(\mathbf{z}_{i+1}) - f(\mathbf{z}_{I-1})}$ :

$$\begin{aligned}
\sqrt{f(\mathbf{z}_i) - f(\mathbf{z}_{I-1})} - \sqrt{f(\mathbf{z}_{i+1}) - f(\mathbf{z}_{I-1})} &= \frac{f(\mathbf{z}_i) - f(\mathbf{z}_{i+1})}{\sqrt{f(\mathbf{z}_i) - f(\mathbf{z}_{I-1})} + \sqrt{f(\mathbf{z}_{i+1}) - f(\mathbf{z}_{I-1})}} \\
&\gtrsim \frac{\eta \mu_\Phi^2 \|\nabla f(\mathbf{z}_i)\|^2}{\sqrt{f(\mathbf{z}_i) - f(\mathbf{z}_{I-1})} + \sqrt{f(\mathbf{z}_{i+1}) - f(\mathbf{z}_{I-1})}} \\
&\geq \frac{\eta \mu_\Phi^2 \|\nabla f(\mathbf{z}_i)\|^2}{2\sqrt{f(\mathbf{z}_i) - f(\mathbf{z}_{I-1})}} \\
&\geq \frac{\eta \sqrt{\alpha_f} \mu_\Phi^2 \|\nabla f(\mathbf{z}_i)\|^2}{\sqrt{2} \|\nabla f(\mathbf{z}_i)\|} \\
&= \frac{\eta \sqrt{\alpha_f} \mu_\Phi^2}{\sqrt{2}} \|\nabla f(\mathbf{z}_i)\|.
\end{aligned} \quad (32)$$

Substituting (32) into (30), we have

$$\begin{aligned}
\ell(I) &\lesssim \eta \nu_\Phi \sum_{i=0}^{I-1} \|\nabla f(\mathbf{z}_i)\| \\
&\lesssim \frac{\nu_\Phi}{\sqrt{\alpha_f} \mu_\Phi^2} \sum_{i=0}^{I-2} \left( \sqrt{f(\mathbf{z}_i) - f(\mathbf{z}_{I-1})} - \sqrt{f(\mathbf{z}_{i+1}) - f(\mathbf{z}_{I-1})} \right) + \eta \nu_\Phi \|\nabla f(\mathbf{z}_{I-1})\| \\
&\lesssim \frac{\nu_\Phi}{\sqrt{\alpha_f} \mu_\Phi^2} \sqrt{f(\mathbf{z}_0) - f(\mathbf{z}_{I-1})} + \eta \nu_\Phi \|\nabla f(\mathbf{z}_{I-1})\| \\
&\leq \frac{\nu_\Phi \sqrt{f(\mathbf{z}_0)}}{\sqrt{\alpha_f} \mu_\Phi^2} + \eta \nu_\Phi \|\nabla f(\mathbf{z}_{I-1})\|.
\end{aligned} \quad (33)$$

Note that

$$f(\mathbf{z}_0) = h(\mathbf{w}_0) \lesssim \frac{\alpha_f \mu_\Phi^6}{\beta_\Phi^2 \nu_\Phi^2}$$

and scaling down the learning rate sufficiently to control the second term in the upper bound ensure that

$$\ell(I) \leq \frac{\rho_\Phi}{2} = \frac{\mu_\Phi}{4\beta_\Phi}.$$

Hence, the gradient descent iterates satisfy:

$$\{\mathbf{w}_i\}_{i \geq 0} \in \text{ball}(\mathbf{w}_0, \rho_\Phi),$$

which implies that the limit  $\bar{\mathbf{w}}$  exists and is globally optimal. In the following, we simplify the expression for  $\eta$  in (31). Since the iterates of gradient flow remain within a ball of radius  $\rho_\Phi$ , we can compute the local Lipschitz constant of  $f$  as

$$\begin{aligned}
\max_i \|\nabla f(\mathbf{z}_i)\| &\leq \|\nabla f(\mathbf{z}_0)\| + \max_i \|\nabla f(\mathbf{z}_i) - \nabla f(\mathbf{z}_0)\| \\
&\leq \|\nabla f(\mathbf{z}_0)\| + \beta_f \max_i \|\mathbf{z}_i - \mathbf{z}_0\| \\
&= \|\nabla f(\mathbf{z}_0)\| + \beta_f \max_i \|\Phi(\mathbf{w}_i) - \Phi(\mathbf{w}_0)\| \\
&= \|\nabla f(\mathbf{z}_0)\| + \frac{3\beta_f\nu_\Phi}{2} \max_i \|\mathbf{w}_i - \mathbf{w}_0\| \\
&\leq \|\nabla f(\mathbf{z}_0)\| + \frac{3\beta_f\nu_\Phi}{2} \cdot \rho_\Phi \\
&= \|\nabla f(\mathbf{z}_0)\| + \frac{3\beta_f\mu_\Phi\nu_\Phi}{4\beta_\Phi}.
\end{aligned} \tag{34}$$

Substituting (34) into (31), an upper bound on  $\eta$  is given by

$$\eta \lesssim \frac{1}{\beta_\Phi \|\nabla f(\mathbf{z}_0)\| + \beta_f \mu_\Phi \nu_\Phi + \beta_f \nu_\Phi^2} \leq \frac{1}{\beta_\Phi \|\nabla f(\mathbf{z}_0)\| + \beta_f \mu_\Phi^2 + \beta_f \nu_\Phi^2} \tag{35}$$

where the last inequality holds since  $\mu_\Phi \leq \nu_\Phi$ .

Finally, using (7), we prove the linear convergence to the limit point  $\bar{\mathbf{w}}$ :

$$\begin{aligned}
h(\mathbf{w}_{i+1}) &= h(\mathbf{w}_{i+1}) - h(\mathbf{w}_i) + h(\mathbf{w}_i) \\
&= f(\mathbf{z}_{i+1}) - f(\mathbf{z}_i) + h(\mathbf{w}_i) \\
&\leq -C\eta\mu_\Phi^2 \|\nabla f(\mathbf{z}_i)\|^2 + h(\mathbf{w}_i) \\
&\leq (1 - C\eta\alpha_f\mu_\Phi^2)h(\mathbf{w}_i)
\end{aligned} \tag{36}$$

where  $C$  is a universal constant. This completes the proof of Theorem 2.

## D Proof of Lemma 3

We first obtain the expression for adjoint operator  $\nabla\Phi^*(\Theta) : \mathbb{R}^{d_2 \times n} \rightarrow \mathbb{R}^{d_1 \times d_0} \times \mathbb{R}^{d_2 \times d_1}$ . Let  $\Delta_W \in \mathbb{R}^{d_1 \times d_0}$ ,  $\Delta_V \in \mathbb{R}^{d_2 \times d_1}$ , and  $\Delta \in \mathbb{R}^{d_2 \times n}$ . We expand  $\Phi$  as follow:

$$\begin{aligned}
\Phi(W + \Delta_W, V) &\approx \Phi(W, V) + \nabla_W \Phi(\Delta_W), \\
\Phi(W, V + \Delta_V) &\approx \Phi(W, V) + \nabla_V \Phi(\Delta_V)
\end{aligned} \tag{37}$$

where

$$\nabla_W \Phi(\Delta_W) = V \left( \dot{\phi}(WX) \odot \Delta_W X \right), \quad \nabla_V \Phi(\Delta_V) = \Delta_V \phi(WX),$$

$\odot$  stands for the Hadamard (entry-wise) product, and  $\dot{\phi}(WX)$  is the derivative of  $\phi$  calculated at each entry of the matrix  $WX$ . The operator  $\nabla\Phi(\Theta)$  is given by  $(\Delta_W, \Delta_V) \rightarrow \nabla_W \Phi(\Delta_W) + \nabla_V \Phi(\Delta_V)$ .

Using the cyclic property of the trace operator and trace  $((A \odot B)C) = \text{trace}((A \odot C^\top)B^\top)$ , we have

$$\begin{aligned}
\langle \Delta, \nabla_W \Phi(\Delta_W) \rangle &= \left\langle \left( \dot{\phi}(WX) \odot V^\top \Delta \right) X^\top, \Delta_W \right\rangle, \\
\langle \Delta, \nabla_V \Phi(\Delta_V) \rangle &= \langle \Delta_V, \Delta \phi(X^\top W^\top) \rangle.
\end{aligned} \tag{38}$$

Substituting (38), the adjoint operator is given by

$$\nabla\Phi^*(\Theta) : \Delta \rightarrow \left( \left( \dot{\phi}(WX) \odot V^\top \Delta \right) X^\top, \Delta \phi(X^\top W^\top) \right). \tag{39}$$

Suppose that there exist  $\dot{\phi}_{\max}, \ddot{\phi}_{\max} < \infty$  such that

$$\sup_a |\dot{\phi}(a)| \leq \dot{\phi}_{\max}, \quad \sup_a |\ddot{\phi}(a)| \leq \ddot{\phi}_{\max}. \tag{40}$$

**Lemma 5.** Let  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times k}$ . Then, we have

$$\sigma_{\min}(A)\|B\| \leq \|AB\| \leq \sigma_{\max}(A)\|B\|.$$

Using Lemma 5 and triangular inequality, we note that

$$\begin{aligned} \|\nabla\Phi^*(\Theta, \Delta)\| &\leq \left\| \left( \dot{\phi}(WX) \odot (V^\top \Delta) \right) X^\top \right\| + \|\Delta\phi(X^\top W^\top)\| \\ &\leq \dot{\phi}_{\max}\sigma_{\max}(X)\sigma_{\max}(V)\|\Delta\| + \sigma_{\max}(\phi(WX))\|\Delta\|. \end{aligned} \quad (41)$$

Similarly, we have this lower bound:

$$\|\nabla\Phi^*(\Theta, \Delta)\| \geq \sigma_{\min}(\phi(WX))\|\Delta\|. \quad (42)$$

Substituting  $\Theta_0 = (W_0, V_0)$  into (41) and (42),  $\mu_\Phi$  and  $\nu_\Phi$  are given by:

$$\begin{aligned} \sigma_{\max}(\nabla\Phi^*(\Theta_0)) &\leq \dot{\phi}_{\max}\sigma_{\max}(X)\sigma_{\max}(V_0) + \sigma_{\max}(\phi(W_0X)) =: \nu_\Phi, \\ \sigma_{\min}(\nabla\Phi^*(\Theta_0)) &\geq \sigma_{\min}(\phi(W_0X)) =: \mu_\Phi. \end{aligned} \quad (43)$$

In the following, we find the smoothness parameter  $\beta_\Phi$  in (4). Let  $\Theta, \hat{\Theta} \in \mathbb{R}^{d_1 \times d_0} \times \mathbb{R}^{d_2 \times d_1}$ . We note that  $\|\nabla\Phi(\Theta, \Delta) - \nabla\Phi(\hat{\Theta}, \Delta)\| \leq U_1 + U_2$  where

$$\begin{aligned} U_1 &= \|V(\dot{\phi}(W^\top X) \odot (\Delta_W^\top X)) - \hat{V}(\dot{\phi}(\hat{W}^\top X) \odot (\Delta_W^\top X))\| \\ U_2 &= \|\Delta_V\phi(W^\top X) - \Delta_V\phi(\hat{W}^\top X)\|. \end{aligned} \quad (44)$$

Let us denote

$$\sigma_{\max}(\hat{V}) \leq \chi_{\max}. \quad (45)$$

An upper bound on  $U_1$  in (44) is given by:

$$\begin{aligned} U_1 &\leq \|(V - \hat{V})(\dot{\phi}(W^\top X) \odot (\Delta_W^\top X))\| + \|\hat{V}(\dot{\phi}(W^\top X) \odot (\Delta_W^\top X)) - \hat{V}(\dot{\phi}(\hat{W}^\top X) \odot (\Delta_W^\top X))\| \\ &\leq \dot{\phi}_{\max}\sigma_{\max}(X)\|V - \hat{V}\|\|\Delta_W\| + \sigma_{\max}(X)\sigma_{\max}(\hat{V})\|\dot{\phi}(W^\top X) - \dot{\phi}(\hat{W}^\top X)\|_\infty\|\Delta_W\| \\ &\leq \dot{\phi}_{\max}\sigma_{\max}(X)\|V - \hat{V}\|\|\Delta_W\| + \ddot{\phi}_{\max}\sigma_{\max}(X)\|X\|_\infty\sigma_{\max}(\hat{V})\|W - \hat{W}\|\|\Delta_W\| \\ &\leq \dot{\phi}_{\max}\sigma_{\max}(X)\|V - \hat{V}\|\|\Delta_W\| + \ddot{\phi}_{\max}\chi_{\max}\sigma_{\max}(X)\|W - \hat{W}\|\|\Delta_W\|. \end{aligned}$$

An upper bound on  $U_2$  in (44) is given by:

$$U_2 \leq \dot{\phi}_{\max}\sigma_{\max}(X)\|W - \hat{W}\|\|\Delta_V\|.$$

Substituting the upper bounds on  $U_1$  and  $U_2$ , an upper bound on  $\sigma_{\max}(\nabla\Phi(\Theta) - \nabla\Phi(\hat{\Theta}))$  is given by

$$\begin{aligned} \sigma_{\max}(\nabla\Phi(\Theta) - \nabla\Phi(\hat{\Theta})) &\leq \sigma_{\max}(X) \left( \dot{\phi}_{\max} + \ddot{\phi}_{\max}\chi_{\max} \right) \|W - \hat{W}\| + \sigma_{\max}(X)\dot{\phi}_{\max}\|V - \hat{V}\| \\ &\leq \sqrt{2}\sigma_{\max}(X) \left( \dot{\phi}_{\max} + \ddot{\phi}_{\max}\chi_{\max} \right) \|\Theta - \hat{\Theta}\| \end{aligned}$$

where the last inequality holds since

$$\|W - \hat{W}\| + \|V - \hat{V}\| \leq \sqrt{2}\sqrt{\|W - \hat{W}\|^2 + \|V - \hat{V}\|^2}.$$

Finally,  $\beta_\Phi$  in (4) is given by

$$\beta_\Phi = \sqrt{2}\sigma_{\max}(X) \left( \dot{\phi}_{\max} + \ddot{\phi}_{\max}\chi_{\max} \right). \quad (46)$$

## E Proof of Theorem 3

This is our setup:  $\min_{\Theta \in \mathbb{R}^{d_1 \times d_0} \times \mathbb{R}^{d_2 \times d_1}} h(\Theta)$  where

$$h(\Theta) = \|V\phi(WX) - Y\|^2.$$

Note that  $\alpha_f = \beta_f = 2$ .

Suppose that there exists  $\chi_{\max} < \infty$  such that, for all  $i \geq 0$ , we have

$$\sigma_{\max}(V_i) \leq \chi_{\max}.$$

The details of  $\chi_{\max}$  later will be provided in Section E.6.

In Lemma 3, we have shown that

$$\begin{aligned} \mu_{\Phi} &= \sigma_{\min}(\phi(W_0X)), \\ \nu_{\Phi} &= \dot{\phi}_{\max} \sigma_{\max}(X) \sigma_{\max}(V_0) + \sigma_{\max}(\phi(W_0X)), \\ \beta_{\Phi} &= \sqrt{2} \sigma_{\max}(X) \left( \dot{\phi}_{\max} + \ddot{\phi}_{\max} \chi_{\max} \right). \end{aligned}$$

In order to apply Theorem 2, we now establish high-probability bounds on random quantities  $\mu_{\Phi}$ ,  $\nu_{\Phi}$ , and  $h(\Theta_0)$  given the initialization in (17).

### E.1 Estimating $\mu_{\Phi}, \nu_{\Phi}$

We now estimate the random quantities  $\mu_{\Phi}, \nu_{\Phi}$  in our neural network setting. The key quantities to estimate are  $\sigma_{\min}(\phi(W_0X))$  and  $\sigma_{\max}(\phi(W_0X))$ . To that end, we consider Hermite decomposition of the activation function  $\phi$ .

We start with the basic definition of Hermite polynomial and its properties. Let  $i \geq 0$  and let  $q_i : \mathbb{R} \rightarrow \mathbb{R}$  denote the  $i$ -th Hermite polynomial. Note that  $q_i$ 's form an orthogonal basis for the Hilbert space of functions.:

$$\mathcal{H} = \left\{ u : \mathbb{R} \rightarrow \mathbb{R} \mid \int u^2(x) \exp\left(-\frac{x^2}{2}\right) dx < \infty \right\},$$

which is equipped with the inner product

$$\langle u, v \rangle_{\mathcal{H}} = \frac{1}{\sqrt{2\pi}} \int u(x)v(x) \exp\left(-\frac{x^2}{2}\right) dx$$

for  $u, v \in \mathcal{H}$ . We consider probabilist's convention of Hermite polynomial. Specifically, for  $i, j \geq 0$ , we have

$$\langle q_i, q_j \rangle_{\mathcal{H}} = \begin{cases} i! & i = j, \\ 0 & i \neq j. \end{cases} \quad (47)$$

Using the above orthogonal basis to decompose  $\phi(W_0X)$ , we have

$$\phi(W_0X) = \sum_{i=0}^{\infty} \frac{c_i}{i!} \cdot q_i(W_0X) \quad (48)$$

where  $c_i = \langle \phi, q_i \rangle_{\mathcal{H}}$  and each matrix  $q_i(W_0X) \in \mathbb{R}^{d_1 \times n}$  is formed by applying  $q_i$  entry-wise to the matrix  $W_0X$ . Let us denote

$$M_0 := \phi(X^{\top} W_0^{\top}) \phi(W_0X).$$

Let  $0 < \tau < 1$ . Suppose there are constants  $r_1, r_2$  such that  $\tau^{r_1} |\phi(a)| \leq |\phi(\tau a)| \leq \tau^{r_2} |\phi(a)|$  for all  $a$ . In the following, we first obtain  $\mathbb{E}[\tilde{M}_0] = \mathbb{E}[\phi(X^{\top} \tilde{W}_0^{\top}) \phi(\tilde{W}_0X)]$  with  $\tilde{W}_0 \sim \mathcal{N}(0, 1)$  and then obtain a lower bound on  $\sigma_{\min}(\mathbb{E}[\tilde{M}_0])$  and an upper bound on  $\sigma_{\min}(\mathbb{E}[M_0])$  by scaling the variance.

Applying Hermite decomposition (48) and taking expectation, we have

$$\begin{aligned}\mathbb{E}[\tilde{M}_0] &= \mathbb{E}\left[\phi(X^\top \tilde{W}_0^\top)\phi(\tilde{W}_0 X)\right] \\ &= \sum_{i,j=0}^{\infty} \frac{c_i c_j}{i!j!} \mathbb{E}[q_i(X^\top \tilde{W}_0^\top)q_j(\tilde{W}_0 X)]\end{aligned}\quad (49)$$

where the expectation is w.r.t. the random matrix  $\tilde{W}_0$ . Let  $\mathbf{x}_a \in \mathbb{R}^{d_0}$  denote the  $a$ -th column of the training data  $X$ . Each summand in (49) is an  $n \times n$  matrix where

$$\left[\mathbb{E}[q_i(X^\top \tilde{W}_0^\top)q_j(\tilde{W}_0 X)]\right]_{a,b} = \sum_{c=1}^{d_1} \mathbb{E}\left[q_i(\mathbf{x}_a^\top \tilde{W}_{0,c,\rightarrow})q_j(\tilde{W}_{0,c,\rightarrow}^\top \mathbf{x}_b)\right], \quad (50)$$

where  $\tilde{W}_{0,c,\rightarrow}$  is the  $c$ -th row of  $\tilde{W}_0$  for  $a, b \in [n]$ .

In summand on the RHS of (50), we note that there is a linear combination of  $\tilde{W}_0$ 's elements inside of each Hermite polynomial.

We use the properties of Hermite polynomials [37][§18.18.11]:

$$\frac{(a_1^2 + \dots + a_r^2)^{\frac{i}{2}}}{i!} \tilde{q}_i\left(\frac{a_1 x_1 + \dots + a_r x_r}{(a_1^2 + \dots + a_r^2)^{\frac{1}{2}}}\right) = \sum_{s_1 + \dots + s_r = i} \frac{a_1^{s_1} \dots a_r^{s_r}}{s_1! \dots s_r!} \tilde{q}_{s_1}(x_1) \dots \tilde{q}_{s_r}(x_r) \quad (51)$$

where  $\tilde{q}_i$ 's form an orthogonal basis, equipped with the inner product  $\langle u, v \rangle_{\tilde{\mathcal{H}}} = \frac{1}{\sqrt{\pi}} \int u(x)v(x) \exp(-x^2) dx$ . This basis follows the physicist's convention of Hermite polynomial.

Since  $\tilde{q}_i$  and  $q_i$  are rescalings of the other, we can replace  $q_i$ 's into (51). Note that we have  $\|\mathbf{x}_a\|_2 = 1$  for all  $a \in [n]$ . Then we have

$$q_i(\mathbf{x}_a^\top \tilde{W}_{0,c,\rightarrow}) = i! \sum_{s_1 + \dots + s_{d_0} = i} \frac{x_{a,1}^{s_1} \dots x_{a,d_0}^{s_{d_0}}}{s_1! \dots s_{d_0}!} q_{s_1}(\tilde{W}_{0,c,1}) \dots q_{s_{d_0}}(\tilde{W}_{0,c,d_0}) \quad (52)$$

where  $x_{a,k}$  and  $\tilde{W}_{0,c,k}$  are  $k$ -th entry of  $\mathbf{x}_a$  and  $\tilde{W}_{0,c,\rightarrow}$  for  $k \in [d_0]$ . Using the expansion in (52), we expand (50) as follows:

$$\begin{aligned}\zeta_{i,j}(a,b) &= i!j! \sum_{s_1 + \dots + s_{d_0} = i} \sum_{s'_1 + \dots + s'_{d_0} = j} \frac{x_{a,1}^{s_1} \dots x_{a,d_0}^{s_{d_0}}}{s_1! \dots s_{d_0}!} \cdot \frac{x_{b,1}^{s'_1} \dots x_{b,d_0}^{s'_{d_0}}}{s'_1! \dots s'_{d_0}!} \rho_{\mathbf{s},\mathbf{s}'}(\tilde{W}_{0,c,\rightarrow}) \\ &= \begin{cases} (i!)^2 \sum_{s_1 + \dots + s_{d_0} = i} \frac{(x_{a,1} x_{b,1})^{s_1} \dots (x_{a,d_0} x_{b,d_0})^{s_{d_0}}}{s_1! \dots s_{d_0}!} & i = j, \\ 0 & i \neq j \end{cases} \\ &= \begin{cases} i! \sum_{s_1 + \dots + s_{d_0} = i} \binom{i}{s_1, \dots, s_{d_0}} (x_{a,1} x_{b,1})^{s_1} \dots (x_{a,d_0} x_{b,d_0})^{s_{d_0}} & i = j, \\ 0 & i \neq j \end{cases}\end{aligned}\quad (53)$$

where  $\zeta_{i,j}(a,b) = \mathbb{E}\left[q_i(\mathbf{x}_a^\top \tilde{W}_{0,c,\rightarrow})q_j(\tilde{W}_{0,c,\rightarrow}^\top \mathbf{x}_b)\right]$ ,

$$\rho_{\mathbf{s},\mathbf{s}'}(\tilde{W}_{0,c,\rightarrow}) = \mathbb{E}\left[q_{s_1}(\tilde{W}_{0,c,1}) \dots q_{s_{d_0}}(\tilde{W}_{0,c,d_0}) \cdot q_{s'_1}(\tilde{W}_{0,c,1}) \dots q_{s'_{d_0}}(\tilde{W}_{0,c,d_0})\right],$$

$\mathbf{s} = [s_1, \dots, s_{d_0}]$ , and  $\mathbf{s}' = [s'_1, \dots, s'_{d_0}]$ .

To simplify the expression in (53), we define  $X^{*i} \in \mathbb{R}^{d_0^i \times n}$  where the  $a$ -th column is given by

$$X_a^{*i} = \text{vec}(\mathbf{x}_a \otimes \dots \otimes \mathbf{x}_a) \in \mathbb{R}^{d_0^i},$$

which is also called Khatri-Rao product. For  $i = 0$ , we use the convention that  $X^{*0} = \mathbf{1}\mathbf{1}^\top \in \mathbb{R}^{n \times n}$ .

We can rewrite (53) as follows:

$$\zeta_{i,j}(a,b) = \begin{cases} i! \langle X_a^{*i}, X_b^{*i} \rangle & i = j \\ 0 & i \neq j. \end{cases} \quad (54)$$



Substituting (54) back into (50), we find that

$$\begin{aligned} \left[ \mathbb{E}[q_i(X^\top \tilde{W}_0^\top) q_j(\tilde{W}_0 X)] \right]_{a,b} &= \sum_{c=1}^{d_1} \mathbb{E} \left[ q_i(\mathbf{x}_a^\top \tilde{W}_{0,c,\rightarrow}) q_j(\tilde{W}_{0,c,\rightarrow}^\top \mathbf{x}_b) \right] \\ &= \begin{cases} d_1 i! \langle X_a^{*i}, X_b^{*i} \rangle & i = j \\ 0 & i \neq j. \end{cases} \end{aligned} \quad (55)$$

Substituting (55) into (49), we have

$$\mathbb{E}[\tilde{M}_0] = d_1 \left( c_0^2 \mathbf{1}\mathbf{1}^\top + c_1^2 X^\top X + \sum_{i=2}^{\infty} \frac{c_i^2}{i!} (X^{*i})^\top X^{*i} \right). \quad (56)$$

We now establish an upper bound on  $\sigma_{\max} \left( \sum_{i=2}^{\infty} \frac{c_i^2}{i!} (X^{*i})^\top X^{*i} \right)$ :

$$\begin{aligned} \sigma_{\max} \left( \sum_{i=2}^{\infty} \frac{c_i^2}{i!} (X^{*i})^\top X^{*i} \right) &\leq \sum_{i=2}^{\infty} \frac{c_i^2}{i!} \sigma_{\max}((X^{*i})^\top X^{*i}) \\ &\leq c_\infty^2 \sigma_{\max}^2(X) \end{aligned} \quad (57)$$

where  $c_\infty$  is given by

$$c_\infty^2 = \sum_{i=2}^{\infty} \frac{c_i^2}{i!},$$

which is finite provided that  $\|\phi\|_{\mathcal{H}}$  is bounded.

Using (57), we now establish an upper bound on  $\sigma_{\max}(\mathbb{E}[\tilde{M}_0])$ :

$$\sigma_{\max}(\mathbb{E}[\tilde{M}_0]) \lesssim d_1 (nc_0^2 + (c_1^2 + c_\infty^2) \sigma_{\max}^2(X)).$$

Moreover, suppose there exists some  $t$  such that  $\sigma_{\min}(X^{*t}) > 0$ . This requires to have  $d_0^t \geq n$ . Putting together the lower bound on  $\sigma_{\min}(\mathbb{E}[\tilde{M}_0])$  and the upper bound on  $\sigma_{\min}(\mathbb{E}[\tilde{M}_0])$ , noting  $W_0 = \omega_1 \tilde{W}_0$  and applying  $\tau^{r_1} \phi(a) \leq \phi(\tau a) \leq \tau^{r_2} \phi(a)$ , we have

$$\omega_1^{2r_1} d_1 \frac{c_t^2}{t!} \sigma_{\min}^2(X^{*t}) \lesssim \sigma_{\min}(\mathbb{E}[M_0]) \leq \sigma_{\max}(\mathbb{E}[M_0]) \lesssim \omega_1^{2r_2} d_1 (nc_0^2 + (c_1^2 + c_\infty^2) \sigma_{\max}^2(X)). \quad (58)$$

## E.2 Concentration of the random matrix $M_0$

To see how well the random matrix  $M_0$  concentrates about its expectation, note that

$$\begin{aligned} M_0 &= \phi(X^\top W_0^\top) \phi(W_0 X) \\ &= \sum_{i=1}^{d_1} \phi(X^\top W_{0,i,\rightarrow}^\top) \phi(W_{0,i,\rightarrow} X) \\ &= \sum_{i=1}^{d_1} A_i \end{aligned} \quad (59)$$

where  $\{A_i\}_{i=1}^{d_1} \subset \mathbb{R}^{n \times n}$  are independent random matrices.

Consider the event  $\mathcal{E}_1$  that

$$\max_{i \in [d_1]} \|W_{0,i,\rightarrow}\|_2 \lesssim k_1 \omega_1 \sqrt{d_0 \log d_1}, \quad \max_{i \in [d_1]} \|V_{0,i,\downarrow}\|_2 \lesssim k_2 \omega_2 \sqrt{d_2 \log d_1} \quad (60)$$

where  $V_{0,i,\downarrow}$  is the  $i$ -th column of  $V_0$ . Note that  $W_{0,i,\rightarrow} \in \mathbb{R}^{d_0}$  and  $V_{0,i,\downarrow} \in \mathbb{R}^{d_2}$  are random zero-mean Gaussian vectors whose entries' variances are  $\omega_1^2$  and  $\omega_2^2$ , respectively. Therefore, with an

application of the scalar Bernstein inequality [43, Proposition 5.16], followed by the union bound, we observe that the event  $\mathcal{E}_1$  happens except with a probability of at most

$$p_1 := d_1^{-Ck_1d_0} + d_1^{-Ck_2d_2}, \quad (61)$$

for a universal constant  $C$  with sufficiently large  $k_1, k_2$ .

Let  $i \in [d_1]$ . Conditioned on the event  $\mathcal{E}_1$ , an upper bound on  $\|\phi(X^\top W_{0,i,\rightarrow})\|_2$  is given by:

$$\|\phi(X^\top W_{0,i,\rightarrow})\|_2 \lesssim \dot{\phi}_{\max} \sigma_{\max}(X) k_1 \omega_1 \sqrt{d_0 \log d_1}. \quad (62)$$

Moreover, we have

$$\begin{aligned} \sigma_{\max}(A_i) &= \|\phi(X^\top W_{0,i,\rightarrow})\|_2^2 \\ &= \|\phi(X^\top W_{0,i,\rightarrow}) - \phi(0)\|_2^2 \\ &\lesssim \dot{\phi}_{\max}^2 \sigma_{\max}^2(X) k_1^2 \omega_1^2 d_0 \log d_1. \end{aligned} \quad (63)$$

We now focus on the concentration of  $\sigma_{\min}(M_0)$  and  $\sigma_{\max}(M_0)$ . We use a concentration property, which provides the tail bound of  $\tilde{f}(W) = \phi(X^\top W^\top) \phi(WX)$  with multivariate Gaussian input  $W$ . In the following lemma, we show that  $\tilde{f}$  is a Lipschitz function, and its Lipschitz constant explains how  $\tilde{f}(W)$  concentrates around its mean.

**Lemma 6.** *Let  $\tilde{f}(W) = \phi(X^\top W^\top) \phi(WX)$ . Suppose  $W$  satisfies (60). Then  $\tilde{f}$  is  $\kappa$ -Lipschitz function with constant  $\kappa = 4\dot{\phi}_{\max}^2 \sigma_{\max}^2(X) k_1 \omega_1 \sqrt{d_0 \log d_1}$ . So we have*

$$\|\tilde{f}(W) - \tilde{f}(W')\| < 4\dot{\phi}_{\max}^2 \sigma_{\max}^2(X) k_1 \omega_1 \sqrt{d_0 \log d_1} \cdot \|W - W'\|.$$

*Proof.* Note that  $\tilde{f}(W_0) = M_0$  and  $\tilde{f}$  can be represented as

$$\tilde{f}(X) = \sum_{i=1}^{d_1} f_i(W_{i,\rightarrow})$$

where  $f_i$  is given by  $f_i(W_{i,\rightarrow}) = \phi(X^\top W_{i,\rightarrow}^\top) \phi(W_{i,\rightarrow} X)$ . We prove that each  $f_i$  is  $\kappa$ -Lipschitz, which implies that  $\tilde{f}$  is also  $\kappa$ -Lipschitz.

We note that  $f_i$ 's can be expressed as a composition of three functions:

$$f_i(\mathbf{v}) = (g_1 \circ g_2 \circ g_3)(\mathbf{v})$$

where  $g_1, g_2$ , and  $g_3$  are given by

$$g_1(\mathbf{v}) = \mathbf{v}\mathbf{v}^\top, \quad g_2(\mathbf{v}) = \phi(\mathbf{v}), \quad g_3(\mathbf{v}) = \mathbf{v}X. \quad (64)$$

It is clear that  $g_2$  is  $\dot{\phi}_{\max}$ -Lipschitz, and  $g_3$  is  $\sigma_{\max}(X)$ -Lipschitz from their definitions. Lipschitz constant of  $g_1$  comes from the domain bound as follows:

$$\begin{aligned} \|g_1(\mathbf{v} + \delta\mathbf{v}) - g_1(\mathbf{v})\| &= \|\delta\mathbf{v}\mathbf{v}^\top + \mathbf{v}\delta\mathbf{v}^\top + \delta\mathbf{v}\delta\mathbf{v}^\top\| \\ &\leq 2\|\delta\mathbf{v}\mathbf{v}^\top\| + \|\delta\mathbf{v}\delta\mathbf{v}^\top\| \\ &\leq (2\|\mathbf{v}\| + \|\delta\mathbf{v}\|) \cdot \|\delta\mathbf{v}\|. \end{aligned} \quad (65)$$

A bound on  $(2\|\mathbf{v}\| + \|\delta\mathbf{v}\|)$  is obtained in (62). Then  $g_1$  is  $\kappa_1$ -Lipschitz function with  $\kappa_1 = 4\dot{\phi}_{\max} \sigma_{\max}(X) k_1 \omega_1 \sqrt{d_0 \log d_1}$ . Therefore, all  $g_1, g_2$  and  $g_3$  are Lipschitz function, so their composition  $f_i$  is also Lipschitz function with constant  $\kappa = 4\dot{\phi}_{\max}^2 \sigma_{\max}^2(X) k_1 \omega_1 \sqrt{d_0 \log d_1}$ , which completes the proof.  $\square$

**Lemma 7.** *Let  $\mathbf{z} \in \mathbb{R}^d$  denote a Gaussian random vector. Then we have  $\Pr\{\|\mathbf{z} - \mathbb{E}[\mathbf{z}]\| > t \mid \mathcal{E}_2\} \lesssim \exp(-t^2)$  where  $\mathcal{E}_2$  is the event that  $\|\mathbf{z}\|$  is bounded.*

We can focus on the tail distribution of  $M_0 = \tilde{f}(W_0)$ . Using Lemmas 6 and 7, we have

$$\Pr\{\|M_0 - \mathbb{E}[M_0]\| > t \mid \mathcal{E}_1\} \lesssim \exp(-k_3^2) \quad (66)$$

where  $t = k_3 4\phi_{\max}^2 \sigma_{\max}^2(X) k_1 \omega_1 \sqrt{d_0 \log d_1}$  with some constant  $k_3$ .

Using (66), we now establish a tail bound on  $\sigma_{\min}(M_0)$ :

$$\begin{aligned} \Pr\{\sigma_{\min}(M_0) \leq (1 - \delta_1)\sigma_{\min}(\mathbb{E}[M_0]) \mid \mathcal{E}_1\} &\leq \Pr\{|\sigma_{\min}(M_0) - \sigma_{\min}(\mathbb{E}[M_0])| \geq \delta_1 \sigma_{\min}(\mathbb{E}[M_0]) \mid \mathcal{E}_1\} \\ &\leq \Pr\{\sigma_{\min}(M_0 - \mathbb{E}[M_0]) \geq \delta_1 \sigma_{\min}(\mathbb{E}[M_0]) \mid \mathcal{E}_1\} \\ &\leq \Pr\{\sigma_{\max}(M_0 - \mathbb{E}[M_0]) \geq \delta_1 \sigma_{\min}(\mathbb{E}[M_0]) \mid \mathcal{E}_1\} \\ &\leq \Pr\{\|M_0 - \mathbb{E}[M_0]\| \geq \delta_1 \sigma_{\min}(\mathbb{E}[M_0]) \mid \mathcal{E}_1\} \\ &\lesssim p_2 \end{aligned}$$

where

$$p_2 = \exp\left(-\left(\frac{\delta_1 \sigma_{\min}(\mathbb{E}[M_0])}{4\phi_{\max}^2 \sigma_{\max}^2(X) k_1 \omega_1 \sqrt{d_0 \log d_1}}\right)^2\right).$$

Similarly, we obtain

$$\Pr\{\sigma_{\max}(M_0) \geq (1 + \delta_2)\sigma_{\max}(\mathbb{E}[M_0]) \mid \mathcal{E}_1\} \lesssim p_3$$

where

$$p_3 = \exp\left(-\left(\frac{\delta_2 \sigma_{\max}(\mathbb{E}[M_0])}{4\phi_{\max}^2 \sigma_{\max}^2(X) k_1 \omega_1 \sqrt{d_0 \log d_1}}\right)^2\right).$$

Putting these bounds together with (58), we have :

$$\begin{aligned} \omega_1^{r_1} \sqrt{(1 - \delta_1) \frac{c_t^2}{t!} d_1} \sigma_{\min}(X^{*t}) &\leq \sigma_{\min}(\phi(W_0 X)) \\ \sigma_{\max}(\phi(W_0 X)) &\leq \sqrt{(1 + \delta_2) \omega_1^{r_2}} (\sqrt{(c_1^2 + c_\infty^2) d_1} \sigma_{\max}(X) + |c_0| \sqrt{d_1 n}) \end{aligned} \quad (67)$$

except with a probability of at most  $p_1 + p_2 + p_3$ .

With establishing the bounds on  $\sigma_{\min}(\phi(W_0 X))$  and  $\sigma_{\max}(\phi(W_0 X))$ , we can finally estimate  $\mu_\Phi, \nu_\Phi$  as follows:

### E.3 Lower bound on $\mu_\Phi$

A lower bound on  $\mu_\Phi$  is given by

$$\omega_1^{r_1} \sqrt{(1 - \delta_1) \frac{c_t^2}{t!} d_1} \sigma_{\min}(X^{*t}) \leq \sigma_{\min}(\phi(W_0 X)) = \mu_\Phi, \quad (68)$$

except with a probability of at most  $p_1 + p_2$ .

### E.4 Upper bound on $\nu_\Phi$

Since  $\nu_\Phi = \dot{\phi}_{\max} \sigma_{\max}(X) \sigma_{\max}(V_0) + \sigma_{\max}(\phi(W_0 X))$ , we obtain a bound on  $\sigma_{\max}(V_0)$ :

Since  $V_0$  is a Gaussian random matrix, we have

$$\sigma_{\max}(V_0) \leq \omega_2 (2\sqrt{d_1} + \sqrt{d_2}) \lesssim \omega_2 \sqrt{d_1} \quad (69)$$

except with a probability of at most  $p_4 = \exp(-C d_1)$  where  $C$  is a universal constant [43][Corollary 5.35].

Combining (69) with the upper bound on  $\sigma_{\max}(\phi(W_0 X))$ , we have

$$\begin{aligned} \nu_\Phi &= \dot{\phi}_{\max} \sigma_{\max}(X) \sigma_{\max}(V_0) + \sigma_{\max}(\phi(W_0 X)) \\ &\lesssim \omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \sqrt{d_1} + \omega_1^{r_2} \sqrt{(1 + \delta_2) (c_1^2 + c_\infty^2) d_1} \sigma_{\max}(X) + \omega_1^{r_2} |c_0| \sqrt{(1 + \delta_2) d_1 n} \end{aligned}$$

except with a probability of at most  $p_1 + p_3 + p_4$ .

### E.5 Upper bound on $h(\Theta_0)$

In this section, we bound  $h(\Theta_0)$ . Using  $\|\mathbf{a} + \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$ , we have

$$\begin{aligned} h(\Theta_0) &= \|V_0\phi(W_0X) - Y\|^2 \\ &\leq 2\|V_0\phi(W_0X)\|^2 + 2\|Y\|^2. \end{aligned} \quad (70)$$

To upper bound the random norm in (70), we first decompose  $V_0\phi(W_0X)$  into terms including  $W_{0,i,\rightarrow} \in \mathbb{R}^{d_0}$  and  $V_{0,i,\downarrow} \in \mathbb{R}^{d_2}$  as follows:

$$V_0\phi(W_0X) = \sum_{i=1}^{d_1} B_i \quad (71)$$

where  $B_i = V_{0,i,\downarrow}\phi(W_{0,i,\rightarrow}^\top X) \in \mathbb{R}^{d_2 \times n}$ 's are independent random matrices for  $i \in [d_1]$ .

Conditioned on the event  $\mathcal{E}_1$  defined in (60), we bound  $\|B_i\|$ :

$$\begin{aligned} \|B_i\| &= \|V_{0,i,\downarrow}\|_2 \|\phi(W_{0,i,\rightarrow}^\top X)\|_2 \\ &\leq \|V_{0,i,\downarrow}\|_2 \cdot \dot{\phi}_{\max} \sigma_{\max}(X) k_1 \omega_1 \sqrt{d_0 \log d_1} \\ &\leq \omega_1 \omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) k_1 k_2 \sqrt{d_0 d_2} \log d_1 \end{aligned} \quad (72)$$

for  $i \leq d_1$ .

Substituting the upper bound in 71 into 72 and applying the Hoeffding inequality [19], we have

$$\begin{aligned} \Pr\{\|V_0\phi(W_0X)\| \gtrsim u(d_0, d_1, d_2) | \mathcal{E}_1\} &= \Pr\{\|V_0\phi(W_0X) - \mathbb{E}[V_0\phi(W_0X)]\| \gtrsim u(d_0, d_1, d_2) | \mathcal{E}_1\} \\ &\leq \Pr\left\{\sum_{i=1}^{d_1} \|B_i - \mathbb{E}[B_i]\| \gtrsim u(d_0, d_1, d_2) | \mathcal{E}_1\right\} \\ &\leq p_5 \end{aligned}$$

where

$$u(d_0, d_1, d_2) = \delta_3 \omega_1 \omega_2 \dot{\phi}_{\max} k_1 k_2 \sqrt{d_0 d_1 d_2} \sigma_{\max}(X) \log d_1$$

and  $p_5 = \exp(-C\delta_3^2)$  with  $\delta_3 \geq 0$  and a universal constant  $C$ .

Therefore, under the event  $\mathcal{E}_1$ , we have

$$\begin{aligned} h(\Theta_0) &\leq 2\|V_0\phi(W_0X)\|^2 + 2\|Y\|^2 \\ &\lesssim \delta_3^2 \omega_1^2 \omega_2^2 \dot{\phi}_{\max}^2 k_1^2 k_2^2 d_0 d_1 d_2 \sigma_{\max}^2(X) \log^2 d_1 + \|Y\|^2 \end{aligned} \quad (73)$$

except with a probability of at most  $p_1 + p_5$ . It is natural to assume that  $d_2 = o(d_1)$ . We also have  $\|Y\| \leq 1$ .

Suppose that

$$\omega_1 \omega_2 \lesssim \frac{1}{\dot{\phi}_{\max} \sqrt{d_0 d_1} \log d_1}. \quad (74)$$

Substituting (74) into (73), we have

$$h(\Theta_0) \leq \delta_3^2 k_1^2 k_2^2 \sigma_{\max}^2(X) \quad (75)$$

where  $\delta_3$ ,  $k_1$ , and  $k_2$  are all constants and independent of  $d_0$ ,  $d_1$ , and  $n$ .

### E.6 Denouement

The key condition for linear rate convergence of gradient descent in (9) is

$$h(\Theta_0) \lesssim \frac{\alpha_f \mu_\Phi^6}{\beta_\Phi^2 \nu_\Phi^2}.$$

Putting everything together for the shallow neural network, with high probably, we have

$$\begin{aligned}
\alpha_f &= 2 \\
\nu_\Phi &= \omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \sqrt{d_1} + \sqrt{(1 + \delta_2) \omega_1^{2r_2} (c_1^2 + c_\infty^2) \sigma_{\max}(X) \sqrt{d_1} + |c_0| \sqrt{\omega_1^{2r_2} (1 + \delta_2) d_1 n}} \\
\mu_\Phi &= \omega_1^{r_1} \sqrt{(1 - \delta_1) \frac{c_t^2}{t!} d_1 \sigma_{\min}(X^{*t})} \\
\beta_\Phi &= \sqrt{2} \sigma_{\max}(X) \left( \dot{\phi}_{\max} + \ddot{\phi}_{\max} \chi_{\max} \right).
\end{aligned} \tag{76}$$

We note that the order of  $\sigma_{\max}(X)$  and  $\sigma_{\min}(X^{*t})$  play significant roles for the overparameterization order analysis. For  $t = 1$ , it requires  $n \simeq d_0$ , which is not a common setting in practice. In the following, we focus on  $t \geq 2$ .

### E.7 Order analysis with $t \geq 2$

In this section, we assume  $|c_0|$  is sufficiently large such that  $|c_0| \sqrt{(1 + \delta_2) d_1 n}$  becomes the dominating term in  $\nu_\Phi$ .<sup>5</sup> Then a sufficient condition to satisfy (9) is

$$d_1^2 \gtrsim \frac{\delta_3^2 c_0^2 (1 + \delta_2) k_1^2 k_2^2 (\dot{\phi}_{\max} + \ddot{\phi}_{\max} \chi_{\max})^2 \sigma_{\max}^4(X) n t!^3}{\omega_1^{6r_1 - 2r_2} (1 - \delta_1)^3 c_t^6 \sigma_{\min}^6(X^{*t})}, \tag{77}$$

which can be written as

$$d_1 \gtrsim \sqrt{\frac{\delta_3^2 c_0^2 (1 + \delta_2) k_1^2 k_2^2 (\dot{\phi}_{\max} + \ddot{\phi}_{\max} \chi_{\max})^2 t!^3}{\omega_1^{6r_1 - 2r_2} (1 - \delta_1)^3 c_t^6}} \cdot \frac{\sqrt{n} \sigma_{\max}^2(X)}{\sigma_{\min}^3(X^{*t})}.$$

For notational simplicity, we let  $\delta_4 = \max(k_1, k_2)$  and denote  $\mathcal{C}_\delta = \{\delta_1, \delta_2, \delta_3, \delta_4\}$  and

$$\xi(\mathcal{C}_\delta, t, \phi, \{c_i\}_{i \geq 0}) = \sqrt{\frac{\delta_3^2 c_0^2 (1 + \delta_2) \delta_4^4 (\dot{\phi}_{\max} + \ddot{\phi}_{\max} \chi_{\max})^2 t!^3}{\omega_1^{6r_1 - 2r_2} (1 - \delta_1)^3 c_t^6}}. \tag{78}$$

Note that  $\xi(\mathcal{C}_\delta, t, \phi, \{c_i\}_{i \geq 0})$  can be viewed as a constant w.r.t.  $d_0$ ,  $d_1$ , and  $n$ . Then (77) can be written as:

$$d_1 = \tilde{\Omega} \left( \frac{\sqrt{n} \sigma_{\max}^2(X)}{\sigma_{\min}^3(X^{*t})} \right). \tag{79}$$

It remains to estimate  $\sigma_{\max}(X)$  and  $\sigma_{\min}(X^{*t})$  to finish the order analysis of  $d_1$ . Suppose that  $n \simeq d_0^t$ . Then, along the lines of [39][Section 2.1], we have  $\sigma_{\max}(X) \simeq \sqrt{\frac{n}{d_0}}$  and  $\sigma_{\min}(X^{*t}) \simeq \sqrt{\frac{n}{d_0^t}} \simeq 1$ .

Combining them all, we have

$$d_1 \gtrsim \xi(\mathcal{C}_\delta, t, \phi, \{c_i\}_{i \geq 0}) \frac{n^{\frac{3}{2}}}{d_0}. \tag{80}$$

Therefore, the overall overparameterization degree becomes  $d_0 d_1 \simeq \tilde{\Omega}(n^{\frac{3}{2}})$  for  $t \geq 2$ .

The exact expression of  $\psi(\phi, \xi, d_0, d_1, d_2, X)$  in Theorem 3 is given by

$$\begin{aligned}
\psi &\leq p_1 + p_2 + p_3 + p_4 + p_5 \\
&\leq d_1^{-C\delta_4 d_0} + d_1^{-C\delta_4 d_2} + e^{-\left(\frac{\delta_1 \sigma_{\min}(\mathbb{E}[M_0])}{4\phi_{\max}^2 \sigma_{\max}^2(X) \delta_4 \sqrt{d_0} \log d_1}\right)^2} + e^{-\left(\frac{\delta_2 \sigma_{\max}(\mathbb{E}[M_0])}{4\phi_{\max}^2 \sigma_{\max}^2(X) \delta_4 \sqrt{d_0} \log d_1}\right)^2} + e^{-C d_1} + e^{-C \delta_3^2}.
\end{aligned}$$

Note that  $d_1^{-C\delta_4 d_0} + d_1^{-C\delta_4 d_2} + \exp(-C d_1) + \exp(-C \delta_3^2)$  decreases exponentially, which can be sufficiently small without changing the order of  $d_1$ .

<sup>5</sup>To have a nonzero  $c_0$ , the activation function should not be an odd function.

Finally, with  $d_0 d_1 \simeq \tilde{\Omega}(n^{\frac{3}{2}})$ , the gradient descent converges to a global minimum with linear rate with probability at least  $1 - \psi$ , which can be arbitrary small.

**Order analysis without boundedness assumption on  $\sigma_{\max}(V_k)$  in Assumption 2.**

So far, we assumed  $\sigma_{\max}(V_k)$  is bounded for  $k \geq 0$ . We can relax this assumption by bounding the length of the trajectory of gradient descent as discussed in Appendix C. Recall (33):

$$\ell(I) \lesssim \frac{\nu_{\Phi} \sqrt{f(Z_0)}}{\sqrt{\alpha_f \mu_{\Phi}^2}}.$$

Using triangular inequality and substituting (33), we can obtain a bound on  $\|V_k\|$

$$\begin{aligned} \|V_k\| &\leq \|V_k - V_0\| + \|V_0\| \\ &\leq \frac{\nu_{\Phi} \sqrt{f(Z_0)}}{\sqrt{\alpha_f \mu_{\Phi}^2}} + \|V_0\| \end{aligned} \quad (81)$$

As shown in (69),  $\|V_0\| \lesssim \omega_2 \sqrt{d_1}$  with high probability over the choice of  $V_0$ . With sufficiently small  $\omega_2$ , the first term in the upper bound dominates in (81). Applying (75) and substituting (81) into (77), we have

$$\begin{aligned} d_1^3 &\gtrsim \frac{n^2 \sigma_{\max}^6(X)}{\sigma_{\min}^{10}(X^{*t})} \\ d_1 &\gtrsim \frac{n^{\frac{5}{3}}}{d_0}. \end{aligned}$$

The overall overparameterization degree becomes  $d_0 d_1 \simeq \tilde{\Omega}(n^{\frac{5}{3}})$ , which is slightly worse than the result of Theorem 3 under boundedness assumption on  $\sigma_{\max}(V_k)$ . Note that we still have a subquadratic scaling on the network width.

## F Additional discussion on lazy training in Section 6

In this section, we provide an asymptotic analysis for the term  $\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\|$  to show that there exists a regime where our initialization can avoid lazy training. Recall our setting:

$$\Phi(\Theta) = V \cdot \phi(WX)$$

where  $W \sim \mathcal{N}(0, \omega_1^2)$  and  $V \sim \mathcal{N}(0, \omega_2^2)$ . Following the theoretical guidance in (19), we set  $\omega_1 \omega_2 \simeq \frac{1}{\sqrt{d_0 d_1}}$ .

An upper bound on  $\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\|$  is given by [7, Theorem 2.3]:

$$\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\| \lesssim \frac{\text{Lip}(\nabla \Phi(\Theta))}{\text{Lip}(\Phi(\Theta))^2}. \quad (82)$$

In the following, we estimate  $\frac{\text{Lip}(\nabla \Phi(\Theta))}{\text{Lip}(\Phi(\Theta))^2}$  to find when it is not bound to be close to zero.

Substituting  $\beta_{\Phi}$  and  $\nu_{\Phi}$  expressions in (76) into the upper bound in (82) for sufficiently large  $n, c_0$ , we have

$$\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\| \lesssim \frac{\sqrt{2} \sigma_{\max}(X) (\dot{\phi}_{\max} + \ddot{\phi}_{\max} \chi_{\max})}{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \sqrt{d_1} + \omega_1^{r_2} c_0 \sqrt{(1 + \delta_2) d_1 n})^2}. \quad (83)$$

We now find an upper bound on  $\chi_{\max}$  by bounding the total length of the trajectory of gradient descent as in Appendix C where the length of the trajectory traced by gradient descent is given by (33):

$$\ell(I) \leq \frac{\nu_{\Phi} \sqrt{f(Z_0)}}{\sqrt{\alpha_f \mu_{\Phi}^2}}.$$



Using (33), (69), and (75), a bound on  $\chi_{\max}$  is given by

$$\begin{aligned}
\|V_i\|_2 &\leq \|V_i - V_0\|_F + \|V_0\|_2 \\
&\leq \frac{\nu_\Phi \sqrt{f(Z_0)}}{\sqrt{\alpha_f \mu_\Phi^2}} + \|V_0\|_2 \\
&\lesssim \frac{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) + \omega_1^{r_2} c_0 \sqrt{n}) \sigma_{\max}(X)}{\omega_1^{2r_1} \sqrt{d_1} \sigma_{\min}^2(X^{*t})} + \omega_2 \sqrt{d_1}
\end{aligned} \tag{84}$$

Therefore we have

$$\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\| \lesssim \frac{\sqrt{2} \sigma_{\max}(X) \left( \dot{\phi}_{\max} + \ddot{\phi}_{\max} \frac{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) + \omega_1^{r_2} c_0 \sqrt{n}) \sigma_{\max}(X)}{\omega_1^{2r_1} \sqrt{d_1} \sigma_{\min}^2(X^{*t})} + \omega_2 \ddot{\phi}_{\max} \sqrt{d_1} \right)}{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \sqrt{d_1} + \omega_1^{r_2} c_0 \sqrt{(1 + \delta_2) d_1 n})^2}$$

We now consider two cases: 1)  $\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \gtrsim \omega_1^{r_2} c_0 \sqrt{n}$  and 2)  $\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \lesssim \omega_1^{r_2} c_0 \sqrt{n}$ . More precisely, for the asymptomatic analysis, we consider extremal cases  $\omega_1 \gg \omega_2$  and  $\omega_1 \ll \omega_2$  and evaluate  $\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\|$  in each case:

### F.1 Regime with $\omega_2 \gg \omega_1$

In the overparameterization regime with large  $d$ , we note that  $\ddot{\phi}_{\max} \frac{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) + \omega_1^{r_2} c_0 \sqrt{n}) \sigma_{\max}(X)}{\omega_1^{2r_1} \sqrt{d_1} \sigma_{\min}^2(X^{*t})} + \omega_2 \ddot{\phi}_{\max} \sqrt{d_1} \gtrsim \dot{\phi}_{\max}$ . Then we have

$$\begin{aligned}
\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\| &\lesssim \frac{\sqrt{2} \sigma_{\max}(X) \left( \frac{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) + \omega_1^{r_2} c_0 \sqrt{n}) \sigma_{\max}(X)}{\omega_1^{2r_1} \sqrt{d_1} \sigma_{\min}^2(X^{*t})} + \omega_2 \sqrt{d_1} \right)}{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \sqrt{d_1} + \omega_1^{r_2} c_0 \sqrt{(1 + \delta_2) d_1 n})^2} \\
&\lesssim \frac{\sigma_{\max}^2(X) \left( \frac{\omega_2}{\omega_1^{2r_1} \sqrt{d_1} \sigma_{\min}^2(X^{*t})} \right)}{(\omega_2 \sigma_{\max}(X) + \omega_1^{r_2} c_0 \sqrt{n})^2 d_1} \\
&\lesssim \frac{\sigma_{\max}^2(X) \omega_2 / d_1^{\frac{3}{2}}}{\sigma_{\min}^2(X^{*t}) (\omega_1^{r_1} \omega_2 \sigma_{\max}(X) + \omega_1^{r_1+r_2} c_0 \sqrt{n})^2} \\
&\lesssim \frac{\sigma_{\max}^2(X) \omega_2 / d_1^{\frac{3}{2}}}{\left( \sigma_{\min}(X^{*t}) \sigma_{\max}(X) \frac{\omega_1^{r_1-1}}{\sqrt{d_0 d_1}} + \omega_1^{r_1+r_2} \sigma_{\min}(X^{*t}) c_0 \sqrt{n} \right)^2}.
\end{aligned}$$

We note that this upper bound above goes to  $\infty$  in the regime  $\omega_2 \gg \omega_1$ , which means that gradient descent can avoid lazy training. Note that it does not imply this training scheme is guaranteed to be non-lazy though.

### F.2 Regime with $\omega_1 \gg \omega_2$

In this regime, we have  $\ddot{\phi}_{\max} \frac{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) + \omega_1^{r_2} c_0 \sqrt{n}) \sigma_{\max}(X)}{\omega_1^{2r_1} \sqrt{d_1} \sigma_{\min}^2(X^{*t})} \lesssim \dot{\phi}_{\max} + \omega_2 \ddot{\phi}_{\max} \sqrt{d_1}$ . Then we have

$$\begin{aligned}
\|h(\Theta_i) - \tilde{h}(\tilde{\Theta}_i)\| &\lesssim \frac{\sqrt{2} \sigma_{\max}(X) (\dot{\phi}_{\max} + \omega_2 \ddot{\phi}_{\max} \sqrt{d_1})}{(\omega_2 \dot{\phi}_{\max} \sigma_{\max}(X) \sqrt{d_1} + \omega_1^{r_2} c_0 \sqrt{(1 + \delta_2) d_1 n})^2} \\
&\lesssim \frac{\sqrt{2} \sigma_{\max}(X) (\dot{\phi}_{\max} + \omega_2 \ddot{\phi}_{\max} \sqrt{d_1})}{(\omega_1^{r_2} c_0 \sqrt{d_1 n})^2}.
\end{aligned} \tag{85}$$

Note that this bound goes to 0 and lazy training is bound to happen asymptotically.

## G Implementation details of Section 6

For the experiments illustrated in Figure 1, we computed the training and test accuracy for different variants of the proposed weight initialization scheme. We considered the MNIST data set made available through the *torchvision* implementation<sup>6</sup>. We used the provided split of 60 000 training examples and 10 000 test examples which we subsequently normalized.

First, a teacher neural network was trained on this data set. The label provided by the teacher was then used to relabel both the training and test examples. For each of the weight initializations a student network was constructed and trained on the relabeled data set. The student neural network had 1 000 units in its hidden layer and used the GeLU activation function. For the loss we used the mean square error against a one-hot encoding of the true class label. We minimized this loss with stochastic gradient descent (SGD) for which there were three hyperparameter choices. As the difficulty of the data set was modest we expected a large range of these hyperparameters to work. It thus sufficed to make a reasonable guess by choosing a batch size of 128, learning rate of 0.01 and 300 epochs. The teacher neural network differed from the student network by using He initialization and cross entropy loss.

All results were implemented in PyTorch [40] and run on a Slurm cluster using a Tesla K40c GPU. We fixed  $\omega_1\omega_2 \approx 0.002259$  based on the He initialization for our particular network and varied  $\omega_2$  in the range  $[0.002, 0.1]$ . We considered 10 different initializations in this range and ran 5 experiments for each configuration of weight initialization,  $(\omega_1, \omega_2)$ . Using these independent runs we plotted the mean and standard deviation of the final training and test accuracy in Figure 1, in Section 6.

---

<sup>6</sup>This implementation uses the original MNIST source: <http://yann.lecun.com/exdb/mnist/>.