049

050

051

052

053

054

000

Convergence Analysis of Natural Gradient Descent for Over-parameterized Physics-Informed Neural Networks

Anonymous Authors¹

Abstract

In the context of over-parameterization, there is a line of work demonstrating that randomly initialized (stochastic) gradient descent (GD) converges to a globally optimal solution at a linear convergence rate for the quadratic loss function. However, the learning rate of GD for training two-layer neural networks exhibits poor dependence on the sample size and the Gram matrix, leading to a slow training process. In this paper, we show that for training two-layer ReLU³ Physics-Informed Neural Networks (PINNs), the learning rate can be improved from $\mathcal{O}(\lambda_0)$ to $\mathcal{O}(1/\|\boldsymbol{H}^{\infty}\|_2)$, implying that GD actually enjoys a faster convergence rate. Despite such improvements, the convergence rate is still tied to the least eigenvalue of the Gram matrix, leading to slow convergence. We then develop the positive definiteness of Gram matrices with general smooth activation functions and provide the convergence analysis of natural gradient descent (NGD) in training two-layer PINNs, demonstrating that the learning rate can be $\mathcal{O}(1)$ and at this rate, the convergence rate is independent of the Gram matrix. In particular, for smooth activation functions, the convergence rate of NGD is quadratic.

1. Introduction

In recent years, neural networks have achieved remarkable breakthroughs in the fields of image recognition (He et al., 2016), natural language processing (Devlin et al., 2018), reinforcement learning (Silver et al., 2016), and so on. Moreover, due to the flexibility and scalability of neural networks, researchers are paying much attention in exploring new methods involving neural networks for handling problems in scientific computing. One long-standing and essential problem in this area is solving partial differential equations (PDEs) numerically. Classical numerical methods, such as finite difference, finite volume and finite elements methods, suffer from the curse of dimensionality when solving high-dimensional PDEs. Due to this drawback, various methods involving neural networks have been proposed for solving different type PDEs (Müller & Zeinhofer, 2023; Raissi et al., 2019; Yu et al., 2018; Zang et al., 2020; Siegel et al., 2023). Among them, the most representative approach is Physics-Informed Neural Networks (PINNs) (Raissi et al., 2019). In the framework of PINNs, one incorporate PDE constraints into the loss function and train the neural network with it. With the use of automatic differentiation, the neural network can be efficiently trained by first-order or second-order methods.

In the applications of neural networks, one inevitable issue is the selection of the optimization methods. First-order methods, such as gradient descent (GD) and stochastic gradient descent (SGD), are widely used in optimizing neural networks as they only calculate the gradient, making them computationally efficient. In addition to first-order methods, there has been significant interest in utilizing second-order optimization methods to accelerate training. These methods have proven to be applicable not only to regression problems, as demonstrated in Martens & Grosse (2015), but also to problems related to PDEs, as shown in Müller & Zeinhofer (2023); Raissi et al. (2019).

As for the convergence aspect of the optimization methods, it has been shown that gradient descent algorithm can even achieve zero training loss under the setting of over-parameterization, which refers to a situation where a model has more parameters than necessary to fit the data (Du et al., 2018; 2019; Allen-Zhu et al., 2019a;b; Arora et al., 2019; Li & Liang, 2018; Zou et al., 2020; Cao & Gu, 2019). These works are based on the idea of neural tangent kernel (NTK)(Jacot et al., 2018), which shows that training multi-layer fully-connected neural networks via gradient descent is equivalent to performing a certain kernel method as the width of every layer goes to infinity. As for the finite width neural networks, with more refined analysis, it can be shown that the parameters are closed to the initializations throughout the entire training process when

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

the width is large enough. This directly leads to the linear convergence for GD. Despite these attractive convergence results, the learning rate depends on the sample size and the 058 Gram matrix, so it needs to be sufficiently small to guaran-059 tee convergence in practice. However, doing so results in 060 a slow training process. In contrast to first-order methods, 061 the second-order method NGD has been shown to enjoy 062 fast convergence for the L^2 regression problems, as demon-063 strated in Zhang et al. (2019); Cai et al. (2019). However, 064 the convergence of NGD in the context of training PINNs 065 is still an open problem. In this paper, we demonstrate that 066 when training PINNs, NGD indeed enjoys a faster conver-067 gence rate. 068

1.1. Contributions

069

070

The main contributions of our work can are summarized as follows:

- For the PINNs, we simultaneously improve both the learning rate η of gradient descent and the requirement for the width m. The improvements rely on a new recursion formula for gradient descent, which is similar to that for regression problems. Specifically, we can improve the learning rate $\eta = \mathcal{O}(\lambda_0)$ required in Gao et al. (2023) to $\eta = \mathcal{O}(1/||\mathbf{H}^{\infty}||_2)$ and the requirement for the width m, i.e. $m = \tilde{\Omega}\left(\frac{(n_1+n_2)^2}{\lambda_0^4\delta^3}\right)$, can be improved to $m = \tilde{\Omega}\left(\frac{1}{\lambda_0^4}(\log(\frac{n_1+n_2}{\delta}))\right)$, where $\tilde{\Omega}$ indicates that some terms involving $\log(m)$ are omitted.
- We present a framework for demonstrating the positive definiteness of Gram matrices for a variety of commonly used smooth activation functions, including the logistic function, softplus function, hyperbolic tangent function, and others. This conclusion is not only applicable to the PDE we have considered but can also be naturally extended to other forms of PDEs.
- 093 • We provide the convergence results for natural gra-094 dient descent (NGD) in training over-parameterized 095 two-layer PINNs with ReLU³ activation functions and 096 smooth activation functions. Due to the distinct op-097 timization dynamics of NGD compared to GD, the 098 learning rate can be $\mathcal{O}(1)$. Consequently, the conver-099 gence rate is independent of n and λ_0 , leading to faster 100 convergence. Moreover, when the activation function is smooth, NGD can achieve a quadratic convergence rate.

105 1.2. Related Works

104

First-order methods. There are mainly two approaches to studying the optimization of neural networks and understanding why first-order methods can find a global minimum.

One approach is to analyze the optimization landscape, as demonstrated in Jin et al. (2017); Ge et al. (2015). It has been shown that gradient descent can find a global minimum in polynomial time if the optimization landscape possesses certain favorable geometric properties. However, some unrealistic assumptions in these works make it challenging to generalize the findings to practical neural networks. Another approach to understand the optimization of neural networks is by analyzing the optimization dynamics of firstorder methods. For the two-layer ReLU neural networks, as shown in Du et al. (2018), randomly initialized gradient descent converges to a globally optimal solution at a linear rate, provided that the width m is sufficiently large and no two inputs are parallel. Later, these results were extended to deep fully-connected feedforward neural networks and ResNet with smooth activation functions (Du et al., 2019). Results for both shallow and deep neural networks depend on the stability of the Gram matrices throughout the training process, which is crucial for convergence to the global minimum. In addition to regression and classification problems, Gao et al. (2023) demonstrated the convergence of the gradient descent for two-layer PINNs through a similar analysis of optimization dynamics. However, both Du et al. (2018) and Gao et al. (2023) require a sufficiently small learning rate and a large enough network width to achieve convergence. In this work, we conduct a refined analysis of gradient descent for PINNs, resulting in milder requirements for the learning rate and network width.

Second-order methods. Although second-order methods possess better convergence rate, they are rarely used in training deep neural networks due to the prohibitive computational cost. As a variant of the Gauss-Newton method, natural gradient descent (NGD) is more efficient in practice. Meanwhile, as shown in Zhang et al. (2019) and Cai et al. (2019), NGD also enjoys faster convergence rate for the L^2 regression problems compared to gradient descent. Müller & Zeinhofer (2023) proposed energy natural gradient descent for PINNs and deep Ritz method, demonstrating experimentally that this method yields solutions that are more accurate than those obtained through GD, Adam or BFGS. After observing the ill-conditioned loss landscape of PINNs, Rathore et al. (2024) introduced a novel second-order optimizer, Nys-NewtonCG (NNCG), showing that NNCG can significantly improve the solution returned by Adam+L-BFGS. Moreover, under the assumption that the PL*-condition holds, Rathore et al. (2024) demonstrated that the convergence rate of their algorithm is independent of the condition number, which is similar with our result. However, although the PL*condition holds for over-parameterized neural networks in the context of regression problems (Liu et al., 2022), it remains unclear whether this condition holds for PINNs. In this paper, we provide the convergence analysis for NGD in training two-layer PINNs with ReLU³ activation functions or smooth activation functions, showing that it indeedconverges at a faster rate.

1.3. Notations

112 113

127

128 129

130

131

132

133

134

135

136

137 138

139 140

141

142 143

144 145

147

155

156 157

114 We denote $[n] = \{1, 2, \dots, n\}$ for $n \in \mathbb{N}$. Given a set 115 S, we denote the uniform distribution on S by $Unif\{S\}$. 116 We use $I{E}$ to denote the indicator function of the event 117 E. For two positive functions $f_1(n)$ and $f_2(n)$, we use 118 $f_1(n) = \mathcal{O}(f_2(n)), f_2(n) = \Omega(f_1(n)) \text{ or } f_1(n) \leq f_2(n)$ 119 to represent $f_1(n) \leq C f_2(n)$, where C is a universal con-120 stant C. A universal constant means a constant independent 121 of any variables. Throughout the paper, we use boldface to 122 denote vectors. Given $x_1, \dots, x_d \in \mathbb{R}$, we use (x_1, \dots, x_d) 123 or $[x_1, \cdots, x_d]$ to denote a row vector with *i*-th component 124 x_i for $i \in [d]$ and then $(x_1, \cdots, x_d)^T \in \mathbb{R}^d$ is a column 125 vector. 126

1.4. Organization of this Paper

In Section 2, we provide the problem setup for training two-layer PINNs. We then present the improved convergence results of gradient descent for PINNs in Section 3. In Section 4, we analyze the convergence of natural gradient descent in training two-layer PINNs with ReLU³ activation functions and smooth activation functions. We conclude in Section 5, and the detailed proofs are provided in the Appendix for readability and brevity.

2. Problem Setup

In this section, we consider the same setup as Gao et al. (2023), focusing on the PDE with the following form.

$$\begin{cases} \frac{\partial u}{\partial x_0}(\boldsymbol{x}) - \sum_{i=1}^d \frac{\partial^2 u}{\partial x_i^2}(\boldsymbol{x}) = f(\boldsymbol{x}), \, \boldsymbol{x} \in (0, T) \times \Omega, \\ u(\boldsymbol{x}) = g(\boldsymbol{x}), \, \boldsymbol{x} \in \{0\} \times \Omega \cup [0, T] \times \partial \Omega, \end{cases}$$
(1)

148 where $\Omega \subset \mathbb{R}^d$ is an open and bounded domain, $x = (x_0, x_1, \cdots, x_d)^T \in \mathbb{R}^{d+1}$ and $x_0 \in [0, T]$ is the time 150 variable. In the following, we assume that $||x||_2 \leq 1$ for 151 $x \in [0, T] \times \overline{\Omega}$ and f, g are bounded continuous functions.

Moreover, we consider a two-layer neural network of the
following form.

$$\phi(\boldsymbol{x}; \boldsymbol{w}, \boldsymbol{a}) = \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \sigma(\boldsymbol{w}_r^T \tilde{\boldsymbol{x}}), \quad (2)$$

158 where $\boldsymbol{w} = (\boldsymbol{w}_1^T, \cdots, \boldsymbol{w}_m^T)^T \in \mathbb{R}^{m(d+2)}, \boldsymbol{a} = (a_1, \cdots, a_m)^T \in \mathbb{R}^m$ and for $r \in [m], \boldsymbol{w}_r \in \mathbb{R}^{d+2}$ 160 is the weight vector of the first layer, a_r is the output 161 weight and $\sigma(\cdot)$ is the ReLU³ activation function. Here, 162 $\tilde{\boldsymbol{x}} = (\boldsymbol{x}^T, 1)^T \in \mathbb{R}^{d+2}$ is the augmented vector from \boldsymbol{x} and 164 in the following, we write \boldsymbol{x} for $\tilde{\boldsymbol{x}}$ for brevity. In the framework of PINNs, given training samples $\{x_p\}_{p=1}^{n_1}$ and $\{y_j\}_{j=1}^{n_2}$ that are from interior and boundary respectively, we aim to minimize the following empirical loss function.

$$L(\boldsymbol{w}, \boldsymbol{a}) :=$$

$$\sum_{p=1}^{n_1} \frac{1}{2n_1} \left(\frac{\partial \phi}{\partial x_0}(\boldsymbol{x}_p; \boldsymbol{w}, \boldsymbol{a}) - \sum_{i=1}^d \frac{\partial^2 \phi}{\partial x_i^2}(\boldsymbol{x}_p; \boldsymbol{w}, \boldsymbol{a}) - f(\boldsymbol{x}_p) \right)^2$$

$$+ \sum_{j=1}^{n_2} \frac{1}{2n_2} \left(\phi(\boldsymbol{y}_j; \boldsymbol{w}, \boldsymbol{a}) - g(\boldsymbol{y}_j) \right)^2.$$
(3)

Similar to that for the L^2 regression problems, we initialize the first layer vector $\boldsymbol{w}_r(0) \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$, output weight $a_r \sim Unif(\{-1,1\})$ for $r \in [m]$ and fix the output weights. Then the gradient descent updates the hidden weights by the following formulations:

$$\boldsymbol{w}_r(k+1) = \boldsymbol{w}_r(k) - \eta \frac{\partial L(\boldsymbol{w}(k), \boldsymbol{a})}{\partial \boldsymbol{w}_r}$$
(4)

for all $r \in [m]$ and $k \in \mathbb{N}$, where $\eta > 0$ is the learning rate. For brevity, we write L(w) for L(w, a).

To simplify the notations, for the residuals of interior and boundary, we denote them by $s_p(w)$ and $h_j(w)$ respectively, i.e.,

$$s_p(\boldsymbol{w}) = \frac{1}{\sqrt{n_1}} \left(\frac{\partial \phi}{\partial x_0}(\boldsymbol{x}_p; \boldsymbol{w}) - \sum_{i=1}^d \frac{\partial^2 \phi}{\partial x_i^2}(\boldsymbol{x}_p; \boldsymbol{w}) - f(\boldsymbol{x}_p) \right)$$
(5)

and

$$h_j(\boldsymbol{w}) = \frac{1}{\sqrt{n_2}} (\phi(\boldsymbol{y}_j; \boldsymbol{w}) - g(\boldsymbol{y}_j)).$$
(6)

Then the empirical loss function can be written as

$$L(\boldsymbol{w}) = \frac{1}{2} \left(\|\boldsymbol{s}(\boldsymbol{w})\|_{2}^{2} + \|\boldsymbol{h}(\boldsymbol{w})\|_{2}^{2} \right),$$
(7)

where

$$\boldsymbol{s}(\boldsymbol{w}) = (s_1(\boldsymbol{w}), \cdots, s_{n_1}(\boldsymbol{w}))^T \in \mathbb{R}^{n_1}$$
 (8)

and

$$\boldsymbol{h}(\boldsymbol{w}) = (h_1(\boldsymbol{w}), \cdots, h_{n_2}(\boldsymbol{w}))^T \in \mathbb{R}^{n_2}.$$
 (9)

At this time, we have

$$\frac{\partial L(\boldsymbol{w})}{\partial \boldsymbol{w}_r} = \sum_{p=1}^{n_1} s_p(\boldsymbol{w}) \frac{\partial s_p(\boldsymbol{w})}{\partial \boldsymbol{w}_r} + \sum_{j=1}^{n_2} h_j(\boldsymbol{w}) \frac{\partial h_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r}$$
(10)

and the Gram matrix H(w) is defined as $H(w) = D^T D$, where

$$\boldsymbol{D} := \left(\frac{\partial s_1(\boldsymbol{w})}{\partial \boldsymbol{w}}, \cdots, \frac{\partial s_{n_1}(\boldsymbol{w})}{\partial \boldsymbol{w}}, \frac{\partial h_1(\boldsymbol{w})}{\partial \boldsymbol{w}}, \cdots, \frac{\partial h_{n_2}(\boldsymbol{w})}{\partial \boldsymbol{w}}\right)$$
(11)

165 3. Improved Results of GD for Two-Layer 166 PINNs 167

To simplify the analysis, we make the following assumptionson the training data.

170Assumption 3.1. For $p \in [n_1]$ and $j \in [n_2]$, $||\boldsymbol{x}_p||_2 \leq \sqrt{2}$,171 $\sqrt{2}$, $||\boldsymbol{y}_j||_2 \leq \sqrt{2}$, where all inputs have been augmented.172Assumption 3.2. No two samples in $\{\boldsymbol{x}_p\}_{p=1}^{n_1} \cup \{\boldsymbol{y}_j\}_{j=1}^{n_2}$ 173are parallel.

175 Under Assumption 3.2, Lemma 3.3 in Gao et al. (2023) 176 implies that the Gram matrix $H^{\infty} := \mathbb{E}_{w \sim \mathcal{N}(\mathbf{0}, I)}[H(w)]$ 177 is strictly positive definite and we let $\lambda_0 = \lambda_{min}(\mathbf{H}^{\infty})$. 178 Similar to the case of the regression problem in Du et al. 179 (2018), H^{∞} plays an important role in the optimization 180 process. Specifically, under over-parameterization and ran-181 dom initialization, we have two facts that (1) at initialization 182 $\|\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}\|_2 = \mathcal{O}(1/\sqrt{m})$ and (2) for any iteration 183 $k \in \mathbb{N}, \|\boldsymbol{H}(k) - \boldsymbol{H}(0)\|_2 = \mathcal{O}(1/\sqrt{m}).$ The following 184 two lemmas can be used to verify these two facts, which are 185 crucial in the convergence analysis. 186

1800 1877 **Lemma 3.3.** If $m = \Omega\left(\frac{d^4}{\lambda_0^2}\log\left(\frac{n_1+n_2}{\delta}\right)\right)$, we have that 1888 with probability at least $1 - \delta$, $\|\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}\|_2 \leq \frac{\lambda_0}{4}$ and 1899 $\lambda_{min}(\boldsymbol{H}(0)) \geq \frac{3}{4}\lambda_0$.

190 Remark 3.4. Under the premise of deriving the same con-191 clusion as Lemma 3.3, Lemma 3.5 in Gao et al. (2023) 192 requires that $m = \tilde{\Omega}\left(\frac{(n_1+n_2)^4}{(n_1n_2)^2\lambda_0^2}\left(\log(\frac{1}{\delta})\right)^7\right)$, where some 193 194 terms involving log(m) are omitted. In contrast, on one hand, our conclusion is independent of n_1 and n_2 , and on 196 the other hand, our conclusion exhibits a clear dependence 197 on d. Moreover, the method in Gao et al. (2023) involves truncating the Gaussian distribution and then applying Ho-199 effding's inequality, which is quite complicated. In contrast, we utilize the concentration inequality for sub-Weibull ran-200 201 dom variables, which serves as a simple framework for this 202 class of problems.

203 **Lemma 3.5.** Let $R \in (0, 1]$, if $w_1(0), \dots, w_m(0)$ are i.i.d. 204 generated from $\mathcal{N}(\mathbf{0}, \mathbf{I})$, then with probability at least $1 - \delta - n_1 e^{-mR}$, the following holds. For any set of weight 206 vectors $w_1, \dots, w_m \in \mathbb{R}^{d+1}$ that satisfy for any $r \in [m]$, 207 $\|w_r - w_r(0)\|_2 < R$, then

$$\|\boldsymbol{H}(\boldsymbol{w}) - \boldsymbol{H}(0)\|_F < CM^2R,\tag{12}$$

211 where $M = 2(d+2)\log(2m(d+2)/\delta)$ and C is a universal 212 constant.

209

210

213 Remark 3.6. Lemma 3.6 in Gao et al. (2023) shows that 214 when $\|\boldsymbol{w}_r - \boldsymbol{w}_r(0)\|_2 \leq R = \tilde{\mathcal{O}}\left(\frac{\lambda_0\delta}{(n_1+n_2)(\log m)^3}\right)$ holds 216 for all r in[m], then $\|\boldsymbol{H}(\boldsymbol{w}) - \boldsymbol{H}(0)\|_2 \leq \frac{\lambda_0}{4}$. In contrast, 217 Lemma 3.5 only requires $R = \mathcal{O}\left(\frac{\lambda_0}{d^2(\log(m/\delta)^2}\right)$ to reach 218 same result. For the L^2 regression problem, as shown in Du et al. (2018), the convergence of gradient descent requires that the learning rate $\eta = O(\lambda_0/n^2)$, where *n* is the sample size of the regression problem. It is evident that this requirement on the learning rate is difficult to satisfy in practical scenarios, since λ_0 is unknown and n^2 is too large . For PINNs, Gao et al. (2023) follows the methodology of Du et al. (2018), thus inheriting similarly stringent requirements on the learning rate. Indeed, such stringent requirement stems from an inadequate decomposition method for the residual. Specifically, in Gao et al. (2023), the decomposition for the residual in the (k + 1)-th iteration is same as the one in Du et al. (2018), i.e.,

$$\begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} = \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} + \left[\begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} - \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} \right],$$
(13)

which leads to the requirements that $\eta = O(\lambda_0)$ and $m = Poly(n_1, n_2, 1/\delta)$. Thus, it requires a new approach to achieve the improvements for η and m. In fact, we can derive the following recursion formula.

Lemma 3.7. For all $k \in \mathbb{N}$, we have

$$\begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} = (\boldsymbol{I} - \eta \boldsymbol{H}(k)) \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} + \boldsymbol{I}_1(k), \quad (14)$$

where

$$I_1(k) = (I_1^1(k), \cdots, I_1^{n_1+n_2}(k))^T \in \mathbb{R}^{n_1+n_2}$$

and for $p \in [n_1]$,
$$I_1^p(k) = s_p(k+1) - s_p(k) - \left\langle \frac{\partial s_p(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k) \right\rangle,$$
(15)

for
$$i \in [n_2]$$
.

$$I_1^{n_1+j}(k) = h_j(k+1) - h_j(k) - \left\langle \frac{\partial h_j(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k) \right\rangle$$
(16)

In the recursion formula (14), $I_1(k)$ serves as a residual term. From the proof, we can see that $||I_1(k)||_2 = O(1/\sqrt{m})$ and thus, as *m* becomes large enough, only the term $I - \eta H(k)$ is significant. This observation is the reason for the requirement of η . With these facts in mind, we arrive at our main result.

Theorem 3.8. Under Assumption 3.1 and Assumption 3.2, if we set the number of hidden nodes

$$m = \Omega\left(\frac{d^8}{\lambda_0^4}\log^6\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right)$$

and the learning rate $\eta = \mathcal{O}\left(\frac{1}{\|\mathbf{H}^{\infty}\|_{2}}\right)$, then with probability at least $1 - \delta$ over the random initialization, the gradient descent algorithm satisfies

$$\left\| \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} \right\|_{2}^{2} \le \left(1 - \frac{\eta \lambda_{0}}{2} \right)^{k} \left\| \begin{pmatrix} \boldsymbol{s}(0) \\ \boldsymbol{h}(0) \end{pmatrix} \right\|_{2}^{2}$$
(17)

220 for all $k \in \mathbb{N}$.

221 Remark 3.9. It may be confusing that Gao et al. (2023) has 222 used the same method in Du et al. (2018), yet it only requires 223 $\eta = \mathcal{O}(\lambda_0)$. Actually, it is because that the loss function of 224 PINN has been normalized. If we let $n_1 = n_2 = n$ and \widetilde{H}^{∞} 225 be the Gram matrix induced by unnormalized loss function 226 of PINN, then $\lambda_{min}(\boldsymbol{H}^{\infty}) = \lambda_{min}(\widetilde{\boldsymbol{H}}^{\infty})/n$, leading to 227 the convergence rate similar to that of regression problem. 228 At this point, due to the normalization of loss function, 229 $\|\boldsymbol{H}^{\infty}\|_{2}$ can be bounded by the trace of \boldsymbol{H}^{∞} , which is 230 an explicit constant that is independent of the sample size 231 n_1, n_2 . Therefore, in practice, we can set the learning rate 232 to satisfy the theoretical convergence requirement, bridging 233 the gap between theory and practice. 234

Remark 3.10. Regarding the impact of dimensionality on 235 the convergence of GD for PINNs, Theorem 3.8 consists 236 of two parts: one is explicit, namely d^8 , and the other is 237 implicit, specifically λ_0^4 , whose relationship with dimension-238 ality remains unclear. The explicit impact arises from the 239 form of the PDE. For instance, the PDE (1) we consider 240 contains $\mathcal{O}(d)$ terms. In concurrent work, hoon Song et al. 241 (2024) investigated the impact of dimensionality and the or-242 der of PDEs on convergence under the setting of continuous 243 gradient flow. Specifically, for PDEs of order 2, the form 244 they considered includes $\mathcal{O}(d^2)$ terms, and the explicit de-245 pendence on dimensionality is d^{28} . On the one hand, hoon 246 Song et al. (2024) only addressed the continuous gradient 247 flow case, while the discrete case requires more refined anal-248 ysis as stated in Du et al. (2018). On the other hand, our 249 results can naturally extend to the PDEs they considered, 250 where the explicit dependence on dimensionality is d^{16} , 251 which is better than the result in hoon Song et al. (2024). 252 Investigating the lower bounds of the smallest eigenvalue of 253 the NTK for PINNs, similar to what has been done for deep 254 ReLU neural networks in Nguyen et al. (2021), represents a 255 promising direction for future research. 256

Similar to Du et al. (2018) and Gao et al. (2023), we prove
Theorem 3.8 by induction. Our induction hypothesis is the
following convergence rate of the empirical loss and upper
bounds for the weights.
Condition 1. At the t th iteration, we have that for each

Condition 1. At the t-th iteration, we have that for each $r \in [m], \|\boldsymbol{w}_r(t)\|_2 \leq B$ and

262

263

264

265

266

267

268

269

270

271

272

273

274

$$L(t) \le \left(1 - \frac{\eta \lambda_0}{2}\right)^t L(0), \tag{18}$$

where $B = \sqrt{2(d+2)\log\left(\frac{2m(d+2)}{\delta}\right)} + 1$ and L(k) is an abbreviation of $L(\boldsymbol{w}(k))$.

From the update formula of gradient descent, we can directly derive the following corollary, which indicates that under over-parameterization, the weights are closed to their initializations. **Corollary 3.11.** If Condition 1 holds for $t = 0, \dots, k$, then we have for every $r \in [m]$,

$$\|\boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(0)\|_{2} \leq \frac{CB^{2}\sqrt{L(0)}}{\sqrt{m\lambda_{0}}},$$
 (19)

where C is a universal constant.

Proof Sketch: Assume that Condition 1 holds for $t = 0, \dots, k$, it suffices to demonstrate that Condition 1 also holds for t = k + 1.

From the recursion formula (14), we have that

$$\begin{aligned} \left\| \begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} \right\|_{2}^{2} \\ &= \left\| (\boldsymbol{I} - \eta \boldsymbol{H}(k)) \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} + \boldsymbol{I}_{1}(k) \right\|_{2}^{2} \\ &\leq \|\boldsymbol{I} - \eta \boldsymbol{H}(k)\|_{2}^{2} \left\| \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} \right\|_{2}^{2} + \|\boldsymbol{I}_{1}(k)\|_{2}^{2} \\ &+ 2 \|\boldsymbol{I} - \eta \boldsymbol{H}(k)\|_{2} \left\| \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} \right\|_{2} \| \boldsymbol{I}_{1}(k)\|_{2} , \end{aligned}$$
(20)

where the inequality follows from the Cauchy's inequality.

Combining Corollary 3.11 with Lemma 3.5, we can deduce that when *m* is large enough, we have $\|\boldsymbol{H}(k) - \boldsymbol{H}(0)\|_2 \leq \lambda_0/4$. Thus, $\lambda_{min}(\boldsymbol{H}(k)) \geq \lambda_0/2$ and $\boldsymbol{I} - \eta \boldsymbol{H}(k)$ is positive definite when $\eta = \mathcal{O}(1/\|\boldsymbol{H}^\infty\|_2)$. On the other hand, with Corollary 3.11, we can derive that $\|\boldsymbol{I}_1(k)\|_2 = \mathcal{O}(\eta\sqrt{L}(k)/\sqrt{m})$. Plugging these results into (20), we have

$$\begin{split} \left\| \begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} \right\|_{2}^{2} \\ &= \left(\left(1 - \frac{\eta \lambda_{0}}{2} \right)^{2} + \mathcal{O}\left(\frac{\eta^{2}}{m}\right) + \mathcal{O}\left(\frac{\eta}{\sqrt{m}}\right) \right) \left\| \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} \right\|_{2}^{2} \\ &\leq \left(1 - \frac{\eta \lambda_{0}}{2} \right) \left\| \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} \right\|_{2}^{2}, \end{split}$$

$$(21)$$

where the last inequality holds when m is large enough.

4. Convergence of NGD for Two-Layer PINNs

Although we have improved the learning rate of gradient descent for PINNs, it may still be necessary to set the learning rates to be sufficiently small for some complex PDEs. Because, although for all PDEs, $Trace(\mathbf{H}^{\infty})$ is an explicit constant, it depends on the form of the PDE, and for complex PDEs, it may be quite large. Moreover, the convergence rate $1 - \frac{\eta \lambda_0}{2}$ also depends on λ_0 , which depends on the sample size and may be extremely small. Zhang et al. (2019) and Cai et al. (2019) have provided the convergence results for natural gradient descent (NGD) in training overparameterized two-layer neural networks for L^2 regression

275 problems. They showed that the maximal learning rate can 276 be O(1) and the convergence rate is independent of λ_0 , 277 which result in a faster convergence rate. However, their 278 methods cannot generalize directly to PINNs. In the section, 279 we conduct the convergence analysis of NGD for PINNs 280 and demonstrate that it results in a faster convergence rate 281 for PINNs compared to gradient descent.

In this section, we consider the same setup as described in
Section 2. Specifically, we focus on the PDE of the form
given in (1) and follow the same initialization as described
in Section 2. During the training process, we fix the output
weight *a* and update the hidden weights via NGD. The optimization objective is the empirical loss function presented
in (7), which is defined as follows:

$$L(\boldsymbol{w}) = \frac{1}{2} \left(\|\boldsymbol{s}(\boldsymbol{w})\|_{2}^{2} + \|\boldsymbol{h}(\boldsymbol{w})\|_{2}^{2} \right), \qquad (22)$$

where s(w) and h(w) are defined in (8) and (9), respectively.

The NGD gives the following update rule:

$$\boldsymbol{w}(k+1) = \boldsymbol{w}(k) - \eta \boldsymbol{J}(k)^T \left(\boldsymbol{J}(k) \boldsymbol{J}(k)^T \right)^{-1} \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix},$$
(23)

where

290

291

292

293

294

295

296 297

299

300

301

303

304

306

307

308

309

312

313

314

318

319

324

327

329

$$\boldsymbol{J}(k) = \left(\boldsymbol{J}_1(k)^T, \cdots, \boldsymbol{J}_{n_1+n_2}(k)^T\right)^T \in \mathbb{R}^{(n_1+n_2) \times m(d+2)}$$

is the Jacobian matrix for the whole dataset and $\eta > 0$ is the learning rate. Specifically, for $p \in [n_1]$,

$$\boldsymbol{J}_{p}(k) = \left[\left(\frac{\partial s_{p}(k)}{\partial \boldsymbol{w}_{1}} \right)^{T}, \cdots, \left(\frac{\partial s_{p}(k)}{\partial \boldsymbol{w}_{m}} \right)^{T} \right] \in \mathbb{R}^{1 \times m(d+2)}$$
(24)

and for $j \in [n_2]$,

$$\boldsymbol{J}_{n_1+j}(k) = \left[\left(\frac{\partial h_j(k)}{\partial \boldsymbol{w}_1} \right)^T, \cdots, \left(\frac{\partial h_j(k)}{\partial \boldsymbol{w}_m} \right)^T \right] \in \mathbb{R}^{1 \times m(d+2)}$$
(25)

Remark 4.1. Zhang et al. (2019) and Cai et al. (2019) have independently and concurrently established the convergence of NGD in the context of regression problems. The difference lies in the fact that Zhang et al. (2019) focused on ReLU activation functions, whereas Cai et al. (2019) considered smooth activation functions and consistently set the learning rate to 1. Here, following Zhang et al. (2019), we refer to this approach as NGD. In Cai et al. (2019), the authors derived this method based on NTK kernel regression and termed it the Gram-Gauss-Newton (GGN) method.

 $\frac{325}{326}$ For the activation function of the two-layer neural network

$$\phi(\boldsymbol{x}; \boldsymbol{w}, \boldsymbol{a}) = \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \sigma(\boldsymbol{w}_r^T \boldsymbol{x}), \quad (26)$$

we consider settings where $\sigma(\cdot)$ is either the ReLU³ activation function or a smooth activation function satisfying the following assumption.

Assumption 4.2. There exists a constant c > 0 such that $\sup_{z \in \mathbb{R}} |\sigma^{(3)}(z)| \le c$ and for any $z, z' \in \mathbb{R}$,

$$|\sigma^{(k)}(z) - \sigma^{(k)}(z')| \le c|z - z'|, \tag{27}$$

where $k \in \{0, 1, 2, 3\}$. Moreover, $\sigma(\cdot)$ is analytic and is not a polynomial function.

Lemma 4.3. If no two samples in $\{x_p\}_{p=1}^{n_1} \cup \{y_j\}_{j=1}^{n_2}$ are parallel, then the Gram matrix \mathbf{H}^{∞} is strictly positive definite for activation functions that satisfy Assumption 4.2, i.e., $\lambda_0 := \lambda_{min}(\mathbf{H}^{\infty}) > 0.$

Remark 4.4. Assumption 4.2 holds for various commonly used activation function, including logistic function $\sigma(z) = 1/(1 + e^{-z})$, softplus function $\sigma(z) = \log(1 + e^{z})$, hyperbolic tangent function $\sigma(z) = (e^{z} - e^{-z})/(e^{z} + e^{-z})$ and others. Compared to the ReLU³ activation function, these smooth activation functions are more popular for PINNs because solving PDEs typically requires high-order derivatives.

Unlike the approach for gradient descent, Zhang et al. (2019) focus on the change of the Jacobian matrix for NGD rather than the Gram matrix. More precisely, they demonstrate that J(w) is stable with respect to w, where J(w) is the Jacobian matrix with weight vector $\boldsymbol{w} = (\boldsymbol{w}_1^T, \cdots, \boldsymbol{w}_m^T)^T$. Roughly speaking, they show that when $\|\boldsymbol{w} - \boldsymbol{w}(0)\|_2$ is small, then $\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2$ is also proportionately small. However, this approach is not applicable to PINNs, because the loss function involves derivatives. Roughly speaking, the stability considered in Zhang et al. (2019) is more global in nature, whereas ours is local. Since the subsequent conclusions require the boundedness of local weights, we do not use this stability. Moreover, from Theorem 1 in Zhang et al. (2019), we can see that this stability imposes additional con-2) straints on the learning rate. Therefore, we instead focus on the stability of J(w) with respect to each individual weight vector w_r , which provides a more targeted approach.

Lemma 4.5. Let $R \in (0,1]$, if $w_1(0), \dots, w_m(0)$ are i.i.d. generated $\mathcal{N}(\mathbf{0}, \mathbf{I})$, then with probability at least $1 - P(\delta, m, R)$ the following holds. For any set of weight vectors $w_1, \dots, w_m \in \mathbb{R}^{d+2}$ that satisfy for any $r \in [m]$, $\|w_r - w_r(0)\|_2 < R$, then

(1) when $\sigma(\cdot)$ is the ReLU³ activation function, we have that

$$\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2 \le CM\sqrt{R},\tag{28}$$

where C is a universal constant, $M = 2(d+2)\log(2m(d+2)/\delta)$ and

$$P(\delta, m, R) = \delta + n_1 e^{-mR}; \qquad (29)$$

(2) when $\sigma(\cdot)$ satisfies Assumption 4.2, we have that

$$\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2 \le CdR \tag{30}$$

for $m \geq \log^2(1/\delta)$, where C is a universal constant and $P(\delta, m, R) = \delta.$

332

344 345

347

348

349

351

352

354

356

357

361

363

369

371

373

374

375

377

379

381

As shown in Lemma 4.5, for $ReLU^3$ activation function, 333 $\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2 = \mathcal{O}(\sqrt{R})$, whereas for smooth activa-334 tion function, $\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2 = \mathcal{O}(R)$. Since R is is 335 sufficiently small, $\mathcal{O}(R)$ is more favorable than $\mathcal{O}(\sqrt{R})$. In 336 fact, the difference of (28) and (30) arises from the continu-337 ity of $\sigma^{'''}(\cdot)$. 338

339 Remark 4.6. For the regression problems, it is shown in 340 Zhang et al. (2019) that when $\sigma(\cdot)$ is the ReLU activation 341 function, then with probability at least $1 - \delta$, for all weight 342 vectors \boldsymbol{w} that satisfy $\|\boldsymbol{w} - \boldsymbol{w}(0)\|_2 \leq R'$, the following 343 holds.

$$\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2 \lesssim \frac{(R')^{1/3}}{\delta^{1/3}m^{1/6}}.$$

Setting $R = R' / \sqrt{m}$ in Lemma 4.5, then $\|\boldsymbol{w} - \boldsymbol{w}(0)\|_2 \leq$ R' and (28) becomes

$$\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2 \lesssim \frac{\log(\frac{1}{\delta})(R')^{1/2}}{m^{1/4}}.$$

Since $R' = \mathcal{O}(\|\boldsymbol{y} - \boldsymbol{u}(0)\|_2/\sqrt{\lambda_0})$ for regression problems, our method results in a less favorable dependence on R' and more favorable dependence on m and δ . This can improve $m = Poly(1/\delta)$ to $m = Poly(\log(1/\delta))$ for the regression problems.

358 More importantly, the stability considered in Zhang et al. 359 (2019) results in that the learning rate must satisfy that $\eta \leq \frac{1-C}{(1+C)^2}$, where $0 \leq C < 1/2$ is a constant appearing 360 in the stability of Jacobian matrix. This requirement for the learning rate may be difficult to satisfy, as C is unknown.

With the stability of Jacobian matrix, we can derive the 364 following convergence results.

Theorem 4.7. Let L(k) = L(w(k)), then the following 367 conclusions hold.

(1) When $\sigma(\cdot)$ is the ReLU³ activation function, under Assumption 3.2, we set

$$m = \Omega\left(\frac{1}{(1-\eta)^2}\frac{d^8}{\lambda_0^4}\log^6\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right)$$

and $\eta \in (0, 1)$, then with probability at least $1 - \delta$ over the random initialization for all $k \in \mathbb{N}$

$$L(k) \le (1 - \eta)^k L(0).$$
 (31)

(2) When $\sigma(\cdot)$ satisfies Assumption 4.2, under Assumption 3.2, we set

$$m = \Omega\left(\frac{1}{1-\eta}\frac{d^6}{\lambda_0^3}\log^2\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right)$$

$$m = \Omega\left(\frac{1}{1-\eta}\frac{d^6}{\lambda_0^3}\log^2\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right)$$

and $\eta \in (0, 1)$, then with probability at least $1 - \delta$ over the random initialization for all $k \in \mathbb{N}$

$$L(k) \le (1 - \eta)^k L(0).$$
 (32)

In Theorem 4.7, the requirements of m with ReLU³ and smooth activation functions exhibit different dependencies on λ_0 and d. The discrepancy is primarily due to the distinct formulations presented in (28) and (30) of Lemma 4.5.

Remark 4.8. We first compare our results with those of NGD for L^2 regression problems. Given that the convergence results are the same, our focus shifts to examining the necessary conditions for the width m. As demonstrated in Zhang et al. (2019) and Cai et al. (2019), it is required that $m = \Omega\left(\frac{n^4}{\lambda_0^4 \delta^3}\right)$ for ReLU activation function and $m = \Omega\left(\max\left\{\frac{n^4}{\lambda_0^4}, \frac{n^2 d \log(n/\delta)}{\lambda_0^2}\right\}\right)$ for smooth activation function. Clearly, our result has a worse dependence on d, which is inevitable due to the involvement of derivatives in the loss function. Moreover, our requirement for m appears to be almost independent of n, primarily because our loss function has been normalized. With smooth activation functions, in addition to the dependence on d, Theorem 4.7 (2) only requires that $m = \hat{\Omega}(\lambda_0^{-3})$. However, Cai et al. (2019) demands a more stringent condition, requiring that $m = \Omega(\lambda_0^{-4}).$

Continuing our analysis, we contrast our results with those of GD for PINNs. Roughly speaking, Gao et al. (2023) has shown that when $\sigma(\cdot)$ is the ReLU³ activation function, $m = \widetilde{\Omega}\left(\frac{(n_1+n_2)^2}{\lambda_0^4\delta^3}\right)$ and $\eta = \mathcal{O}(\lambda_0)$, then the convergence result (17) holds. It is evident that our result, i.e, Theorem 4.7 (1), has a milder dependence on n_1, n_2 and δ . Furthermore, the learning rate and convergence rate are independent of λ_0 , resulting in faster convergence.

Comparing with our results in Section 3, the requirement for m in Theorem 4.7 (1) is the same as in Theorem 3.8, when we make η less close to 1. On the other hand, since $\eta = \mathcal{O}(1)$ and the convergence rate only depends on η , NGD can lead to faster convergence than GD.

Note that as η approaches 1, the width m tends to infinity, thus, the convergence results in Theorem 4.7 become vacuous. In fact, when $\eta = 1$, NGD can enjoy a second-order convergence rate even though m is finite, provided that $\sigma(\cdot)$ satisfies Assumption 4.2.

Corollary 4.9. Under Assumption 3.2 and Assumption 4.2, set $\eta = 1$ and

$$m = \Omega\left(\frac{d^6}{\lambda_0^3}\log^2\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right),$$

then with probability at least $1 - \delta$, we have

$$\left\| \begin{pmatrix} \boldsymbol{s}(t+1) \\ \boldsymbol{h}(t+1) \end{pmatrix} \right\|_{2} \leq \frac{CB^{4}}{\sqrt{m\lambda_{0}^{3}}} \left\| \begin{pmatrix} \boldsymbol{s}(t) \\ \boldsymbol{h}(t) \end{pmatrix} \right\|_{2}^{2}$$

for all $t \in \mathbb{N}$, where C is a universal constant and $B = \sqrt{2(d+2)\log(2m(d+2)/\delta)} + 1$.

387

388

389

390

395

396 397

398 399

400

401

402

406

407

408

409

413

414

415

416 417

418 419

420

421

422

423

424

425

426

Remark 4.10. Cai et al. (2019) has demonstrated the secondorder convergence for regression problems with smooth activation functions. Specifically, it is shown in Cai et al. (2019) that

$$\| \boldsymbol{y} - \boldsymbol{u}(t+1) \|_2 \lesssim rac{n^{3/2}}{\sqrt{m}\lambda_0^2} \| \boldsymbol{y} - \boldsymbol{u}(t) \|_2^2$$

Actually, when applying our method used in Corollary 4.9, we can get a more satisfactory result as follows.

$$\| \boldsymbol{y} - \boldsymbol{u}(t+1) \|_2 \lesssim rac{n^{3/2}}{\sqrt{m\lambda_0^3}} \| \boldsymbol{y} - \boldsymbol{u}(t) \|_2^2.$$

Instead of inducing on the convergence rate of the empirical loss function, as shown in Condition 1, we perform induction on the movements of the hidden weights as follows.

403 the movements of the indden weights as follows. 404 *Condition* 2. At the *t*-th iteration, we have $||w_r(t)||_2 \le B$ 405 and

$$\|\boldsymbol{w}_{r}(t) - \boldsymbol{w}_{r}(0)\|_{2} \le \frac{CB^{2}\sqrt{L(0)}}{\sqrt{m\lambda_{0}}} := R^{\prime}$$

for all $r \in [m]$, where C is a universal constant and B =

$$\begin{array}{l} 410\\ 411\\ 412 \end{array} \quad \sqrt{2(d+2)\log\left(\frac{2m(d+2)}{\delta}\right) + 1}. \end{array}$$

With Condition 2, we can directly derive the following convergence rate of the empirical loss function.

Corollary 4.11. If Condition 2 holds for $t = 0, \dots, k$ and $R' \leq R$ and $R'' \leq \sqrt{1 - \eta} \sqrt{\lambda_0}$, then

$$L(t) \le (1 - \eta)^t L(0),$$

holds for $t = 0, \dots, k$, where R is the constant in Lemma 4.5 and $R'' = CM\sqrt{R}$ is in (28) when σ is the ReLU³ activation function, R'' = CdR is in (30) when σ satisfies Assumption 4.2.

Proof Sketch: First, let $u(t) = \begin{pmatrix} s(t) \\ h(t) \end{pmatrix}$, then from the updating formula of NGD (23), we have

where the second equality is from the fundamental theorem of calculus and $\boldsymbol{w}(s) = s\boldsymbol{w}(t+1) + (1-s)\boldsymbol{w}(t) = \boldsymbol{w}(t) - s\eta \boldsymbol{J}(t)^T \boldsymbol{H}(t)^{-1} \boldsymbol{u}(t).$

In the proof, we assume that Condition 2 holds for $t = 0, \dots, k$. Then from Corollary 4.11, to prove Theorem 4.7, it suffices to demonstrate that this condition also holds for t = k + 1. Here, we primarily explain the process from Condition 2 to Corollary 4.11, while other content is placed in the appendix.

Note that $\frac{\partial u(w(t))}{\partial w} = J(t)$, thus $I_1(t) = \eta u(t)$. Plugging this into (33) yields that

$$u(t+1) = (1-\eta)u(t) + I_2(t).$$
 (34)

From equation (34), we can see the difference between NGD and GD. Recall that the iteration formula for GD is

$$\boldsymbol{u}(t+1) = (1 - \eta \boldsymbol{H}(t))\boldsymbol{u}(t) + \boldsymbol{I}_1(t).$$

Precisely because of this, the convergence rate of GD is inevitably influenced by λ_0 , whereas that of NGD is not.

From the stability of the Jacobian matrix, we can deduce that $\|\mathbf{I}_2(t)\|_2 = \mathcal{O}(\eta \|\mathbf{u}(t)\|_2 / \sqrt{m})$. Plugging this into (34) yields that

$$\begin{aligned} \|\boldsymbol{u}(t+1)\|_{2}^{2} \\ \leq \|(1-\eta)\boldsymbol{u}(t)\|_{2}^{2} + \|\boldsymbol{I}_{2}(t)\|_{2}^{2} + 2(1-\eta)\|\boldsymbol{u}(t)\|_{2}\|\boldsymbol{I}_{2}(t)\|_{2} \\ = \left((1-\eta)^{2} + \mathcal{O}\left(\frac{\eta^{2}}{m}\right) + 2(1-\eta)\mathcal{O}\left(\frac{\eta}{\sqrt{m}}\right)\right)\|\boldsymbol{u}(t)\|_{2}^{2} \\ \leq (1-\eta)\|\boldsymbol{u}(t)\|_{2}^{2}, \end{aligned}$$
(35)

where the last inequality holds if m is large enough.

5. Conclusion and Discussion

In this paper, we have improved the conditions required for the convergence of gradient descent for PINNs, showing that gradient descent actually achieves a better convergence rate. Furthermore, we demonstrate that natural gradient descent can find the global optima of two-layer PINNs with ReLU³ or smooth activation functions for a class of second-order linear PDEs. Compared to gradient descent, natural gradient descent exhibits a faster convergence rate and its maximal learning rate is $\mathcal{O}(1)$. However, natural gradient descent is quite expensive in terms of computation and memory in training neural networks. As a result, several cost-effective variants have been proposed, such as K-FAC (Martens & Grosse, 2015) and mini-batch natural gradient descent. It would be interesting to investigate the convergence of these methods for PINNs. Additionally, extending the convergence analysis to deep neural networks and studying the generalization bounds of trained PINNs are important directions for future research.

440 Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

441

442

443

444

445

446

447

448

449

450

451

452

- Allen-Zhu, Z., Li, Y., and Liang, Y. Learning and generalization in overparameterized neural networks, going beyond two layers. *Advances in neural information processing systems*, 32, 2019a.
- Allen-Zhu, Z., Li, Y., and Song, Z. A convergence theory for
 deep learning via over-parameterization. In *International conference on machine learning*, pp. 242–252. PMLR,
 2019b.
- Arora, S., Du, S., Hu, W., Li, Z., and Wang, R. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pp. 322–332. PMLR, 2019.
- Cai, T., Gao, R., Hou, J., Chen, S., Wang, D., He, D., Zhang,
 Z., and Wang, L. Gram-gauss-newton method: Learning
 overparameterized neural networks for regression problems. arXiv preprint arXiv:1905.11675, 2019.
- 468
 469
 470
 470
 471
 471
 471
 472
 473
 474
 474
 474
 475
 475
 476
 476
 477
 477
 478
 478
 479
 479
 470
 470
 470
 470
 470
 470
 471
 470
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
 471
- 472
 473 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert:
 474 Pre-training of deep bidirectional transformers for lan475 guage understanding. *arXiv preprint arXiv:1810.04805*,
 476 2018.
- Du, S., Lee, J., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. In *International conference on machine learning*, pp. 1675–1685. PMLR, 2019.
- 482 Du, S. S., Zhai, X., Poczos, B., and Singh, A. Gradient
 483 descent provably optimizes over-parameterized neural
 484 networks. *arXiv preprint arXiv:1810.02054*, 2018.
- Gao, Y., Gu, Y., and Ng, M. Gradient descent finds the
 global optima of two-layer physics-informed neural networks. In *International Conference on Machine Learning*,
 pp. 10676–10707. PMLR, 2023.
- Ge, R., Huang, F., Jin, C., and Yuan, Y. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on learning theory*, pp. 797–842. PMLR, 2015.

- Giné, E. and Nickl, R. *Mathematical foundations of infinitedimensional statistical models*. Cambridge university press, 2021.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- hoon Song, C., Park, Y., and Kang, M. How does pde order affect the convergence of pinns? In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. How to escape saddle points efficiently. In *International conference on machine learning*, pp. 1724–1732. PMLR, 2017.
- Kuchibhotla, A. K. and Chakrabortty, A. Moving beyond sub-gaussianity in high-dimensional statistics: Applications in covariance estimation and linear regression. *Information and Inference: A Journal of the IMA*, 11(4): 1389–1456, 2022.
- Li, Y. and Liang, Y. Learning overparameterized neural networks via stochastic gradient descent on structured data. *Advances in neural information processing systems*, 31, 2018.
- Liu, C., Zhu, L., and Belkin, M. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:85–116, 2022.
- Martens, J. and Grosse, R. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417. PMLR, 2015.
- Müller, J. and Zeinhofer, M. Achieving high accuracy with pinns via energy natural gradient descent. In *International Conference on Machine Learning*, pp. 25471– 25485. PMLR, 2023.
- Nguyen, Q., Mondelli, M., and Montufar, G. F. Tight bounds on the smallest eigenvalue of the neural tangent kernel for deep relu networks. In *International Conference on Machine Learning*, pp. 8119–8129. PMLR, 2021.
- Raissi, M., Perdikaris, P., and Karniadakis, G. E. Physicsinformed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear

- partial differential equations. Journal of Computational physics, 378:686-707, 2019.
- Rathore, P., Lei, W., Frangella, Z., Lu, L., and Udell, M. Challenges in training pinns: A loss landscape perspec-tive. In Forty-first International Conference on Machine Learning, 2024.
- Siegel, J. W., Hong, Q., Jin, X., Hao, W., and Xu, J. Greedy training algorithms for neural networks and applications to pdes. Journal of Computational Physics, 484:112084, 2023.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. nature, 529(7587):484-489, 2016.
- Yu, B. et al. The deep ritz method: a deep learning-based nu-merical algorithm for solving variational problems. Com-munications in Mathematics and Statistics, 6(1):1-12, 2018.
- Zang, Y., Bao, G., Ye, X., and Zhou, H. Weak adversarial networks for high-dimensional partial differential equa-tions. Journal of Computational Physics, 411:109409, 2020.
- Zhang, G., Martens, J., and Grosse, R. B. Fast convergence of natural gradient descent for over-parameterized neural networks. Advances in Neural Information Processing Systems, 32, 2019.
- Zou, D., Cao, Y., Zhou, D., and Gu, Q. Gradient descent op-timizes over-parameterized deep relu networks. Machine learning, 109:467-492, 2020.

Appendix

Before the proofs, we first define the event

$$A_{ir} := \{ \exists \boldsymbol{w} : \| \boldsymbol{w} - \boldsymbol{w}_r(0) \|_2 \le R, I\{ \boldsymbol{w}^T \boldsymbol{x}_i \ge 0 \} \ne I\{ \boldsymbol{w}_r(0)^T \boldsymbol{x}_i \ge 0 \} \}$$
(36)

for all $i \in [n]$.

Note that the event happens if and only if $|w_r(0)^T x_i| < ||x_i||_2 R$, thus by the anti-concentration inequality of Gaussian distribution, we have

$$P(A_{ir}) = P_{z \sim \mathcal{N}(0, \|\boldsymbol{x}_i\|_2^2)}(|z| < R) = P_{z \sim \mathcal{N}(0, 1)}(|z| < R) \le \frac{2R}{\sqrt{2\pi}}.$$
(37)

Let $S_i = \{r \in [m] : I\{A_{ir}\} = 0\}$ and $S_i^{\perp} = [m] \backslash S_i$.

Then, we need to recall that

$$\frac{\partial s_p(\boldsymbol{w})}{\partial \boldsymbol{w}_r} = \frac{a_r}{\sqrt{mn_1}} \left[\sigma^{\prime\prime}(\boldsymbol{w}_r^T \boldsymbol{x}_p) w_{r0} \boldsymbol{x}_p + \sigma^{\prime}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \begin{pmatrix} 1\\ \mathbf{0}_{d+1} \end{pmatrix} - \sigma^{\prime\prime\prime}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \|\boldsymbol{w}_{r1}\|_2^2 \boldsymbol{x}_p - 2\sigma^{\prime\prime}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \begin{pmatrix} 0\\ \boldsymbol{w}_{r1} \end{pmatrix} \right]$$
(38)

and

$$\frac{\partial h_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r} = \frac{a_r}{\sqrt{mn_2}} \sigma'(\boldsymbol{w}_r^T \boldsymbol{y}_j) \boldsymbol{y}_j.$$
(39)

A. Proof of Section 3

A.1. Proof of Lemma 3.3

Proof. In the following, we aim to bound $\|\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}\|_{F}$, as $\|\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}\|_{2} \leq \|\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}\|_{F}$. Note that the entries of $\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}$ have three forms as follows.

$$\sum_{r=1}^{m} \left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial s_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}(0)} \left[\sum_{r=1}^{m} \left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial s_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right], \tag{40}$$

$$\sum_{r=1}^{m} \left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}(0)} \left[\sum_{r=1}^{m} \left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right]$$
(41)

583 and 584

$$\sum_{r=1}^{m} \left\langle \frac{\partial h_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}} \left[\sum_{r=1}^{m} \left\langle \frac{\partial h_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right].$$
(42)

587 For the first form (40), to simplify the analysis, we let

$$\boldsymbol{Z}_{r}(i) = \sigma^{''}(\boldsymbol{w}_{r}(0)^{T}\boldsymbol{x}_{i})w_{r0}(0)\boldsymbol{x}_{i} + \sigma^{'}(\boldsymbol{w}_{r}(0)^{T}\boldsymbol{x}_{i})\begin{pmatrix}1\\\boldsymbol{0}_{d+1}\end{pmatrix}$$
$$-\sigma^{'''}(\boldsymbol{w}_{r}(0)^{T}\boldsymbol{x}_{p})\|\boldsymbol{w}_{r1}(0)\|_{2}^{2}\boldsymbol{x}_{p} - 2\sigma^{''}(\boldsymbol{w}_{r}(0)^{T}\boldsymbol{x}_{i})\begin{pmatrix}0\\\boldsymbol{w}_{r1}(0)\end{pmatrix}$$

594 and

$$X_r(ij) = \langle \mathbf{Z}_r(i), \mathbf{Z}_r(j) \rangle_{\mathcal{F}}$$

then

$$\sum_{r=1}^{m} \left\langle \frac{\partial s_p(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial s_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}} \left[\sum_{r=1}^{m} \left\langle \frac{\partial s_p(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial s_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right] = \frac{1}{n_1 m} \sum_{r=1}^{m} \left[X_r(ij) - \mathbb{E} X_r(ij) \right].$$

601 Note that $|X_r(ij)| \lesssim 1 + \|\boldsymbol{w}_r(0)\|_2^4$, thus

603
604
$$\|X_r(ij)\|_{\psi_{\frac{1}{2}}} \lesssim 1 + \|\|\boldsymbol{w}_r(0)\|_2^4\|_{\psi_{\frac{1}{2}}} \lesssim 1 + \|\|\boldsymbol{w}_r(0)\|_2^2\|_{\psi_1}^2 \lesssim d^2.$$

605 Here, for more details on the Orlicz norm, see the remarks after Lemma C.1.

For the centered random variable, the property of $\psi_{\frac{1}{2}}$ quasi-norm implies that

$$\|X_r(ij) - \mathbb{E}[X_r(ij)]\|_{\psi_{\frac{1}{2}}} \lesssim \|X_r(ij)\|_{\psi_{\frac{1}{2}}} + \|\mathbb{E}[X_r(ij)]\|_{\psi_{\frac{1}{2}}} \lesssim d^2.$$

Therefore, applying Lemma C.1 yields that with probability at least $1 - \delta$,

$$\left|\sum_{r=1}^{m} \frac{1}{m} \left[X_r(ij) - \mathbb{E}X_r(ij) \right] \right| \lesssim \frac{d^2}{\sqrt{m}} \sqrt{\log\left(\frac{1}{\delta}\right)} + \frac{d^2}{m} \left(\log\left(\frac{1}{\delta}\right) \right)^2,$$

which directly yields that

$$\left|\sum_{r=1}^{m} \left\langle \frac{\partial s_p(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial s_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}(0)} \left[\sum_{r=1}^{m} \left\langle \frac{\partial s_p(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial s_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right] \right| \lesssim \frac{d^2}{n_1 \sqrt{m}} \sqrt{\log\left(\frac{1}{\delta}\right)} + \frac{d^2}{n_1 m} \left(\log\left(\frac{1}{\delta}\right)\right)^2$$
(43)

Similarly, for the second form (41) and third form (42), we can deduce that

$$\left|\left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}(0)} \left[\left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right] \right\|_{\psi_{\frac{1}{2}}} \lesssim \frac{d^2}{\sqrt{n_1 n_2} m_1^2}$$

and

$$\left\|\left\langle\frac{\partial h_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}\right\rangle - \mathbb{E}_{\boldsymbol{w}(0)}\left[\left\langle\frac{\partial h_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}\right\rangle\right]\right\|_{\psi_{\frac{1}{2}}} \lesssim \frac{d^2}{n_2 m}$$

Thus applying Lemma C.1 yields that with probability at least $1 - \delta$,

$$\left|\sum_{r=1}^{m} \left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}(0)} \left[\sum_{r=1}^{m} \left\langle \frac{\partial s_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right] \right| \lesssim \frac{d^2}{\sqrt{n_1 n_2} \sqrt{m}} \sqrt{\log\left(\frac{1}{\delta}\right)} + \frac{d^2}{\sqrt{n_1 n_2} m} \log\left(\frac{1}{\delta}\right)$$
(44)

and with probability at least $1 - \delta$,

$$\left|\sum_{r=1}^{m} \left\langle \frac{\partial h_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle - \mathbb{E}_{\boldsymbol{w}(0)} \left[\sum_{r=1}^{m} \left\langle \frac{\partial h_i(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r}, \frac{\partial h_j(\boldsymbol{w}(0))}{\partial \boldsymbol{w}_r} \right\rangle \right] \right| \lesssim \frac{d^2}{n_2 \sqrt{m}} \sqrt{\log\left(\frac{1}{\delta}\right)} + \frac{d^2}{n_2 m} \log\left(\frac{1}{\delta}\right).$$
(45)

Combining (43), (44) and (45), we can deduce that with probability at least $1 - \delta$,

$$\begin{split} \|\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}\|_{2}^{2} \\ &\leq \|\boldsymbol{H}(0) - \boldsymbol{H}^{\infty}\|_{F}^{2} \\ &\lesssim \frac{d^{4}}{m} \log\left(\frac{n_{1} + n_{2}}{\delta}\right) + \frac{d^{4}}{m^{2}} \left(\log\left(\frac{n_{1} + n_{2}}{\delta}\right)\right)^{4} \\ &\lesssim \frac{d^{4}}{m} \log\left(\frac{n_{1} + n_{2}}{\delta}\right). \end{split}$$

Thus when $\sqrt{\frac{d^4}{m}\log\left(\frac{n_1+n_2}{\delta}\right)} \lesssim \frac{\lambda_0}{4}$, i.e.,

$$m = \Omega\left(\frac{d^4}{\lambda_0^2}\log\left(\frac{n_1+n_2}{\delta}\right)\right),$$

657 we have $\lambda_{min}(\boldsymbol{H}(0)) \geq \frac{3}{4}\lambda_0$.

A.2. Proof of Lemma 3.5

Proof. We first reformulate the term $\frac{\partial s_p(k)}{\partial w_r}$ in (38) as follows.

$$\frac{\partial s_p(\boldsymbol{w})}{\partial \boldsymbol{w}_r} = \frac{a_r}{\sqrt{mn_1}} \left[\sigma^{\prime\prime}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \begin{pmatrix} w_{r0} \boldsymbol{x}_{p0} \\ w_{r0} \boldsymbol{x}_{p1} - 2\boldsymbol{w}_{r1} \end{pmatrix} + \sigma^{\prime}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \begin{pmatrix} 1 \\ \mathbf{0}_{d+1} \end{pmatrix} - \sigma^{\prime\prime\prime}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \|\boldsymbol{w}_{r1}\|_2^2 \boldsymbol{x}_p \right]$$

It suffices to bound $\|\boldsymbol{H}(\boldsymbol{w}) - \boldsymbol{H}(0)\|_F$, which can in turn allows us to bound each entry of $\boldsymbol{H}(\boldsymbol{w}) - \boldsymbol{H}(0)$. For $i \in [n_1]$ and $j \in [n_1]$, we have that

$$\begin{split} H_{ij}(\boldsymbol{w}) &= \sum_{r=1}^{m} \left\langle \frac{\partial s_i(\boldsymbol{w})}{\partial \boldsymbol{w}_r}, \frac{\partial s_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r} \right\rangle \\ &= \frac{1}{n_1 m} \sum_{r=1}^{m} \left\langle \sigma^{''}(\boldsymbol{w}_r^T \boldsymbol{x}_i) \begin{pmatrix} w_{r0} x_{i0} \\ w_{r0} \boldsymbol{x}_{i1} - 2 \boldsymbol{w}_{r1} \end{pmatrix} + \sigma^{'}(\boldsymbol{w}_r^T \boldsymbol{x}_i) \begin{pmatrix} 1 \\ \mathbf{0}_{d+1} \end{pmatrix} - \sigma^{'''}(\boldsymbol{w}_r^T \boldsymbol{x}_i) \|\boldsymbol{w}_{r1}\|_2^2 \boldsymbol{x}_i, \\ &\sigma^{''}(\boldsymbol{w}_r^T \boldsymbol{x}_j) \begin{pmatrix} w_{r0} x_{j0} \\ w_{r0} \boldsymbol{x}_{j1} - 2 \boldsymbol{w}_{r1} \end{pmatrix} + \sigma^{'}(\boldsymbol{w}_r^T \boldsymbol{x}_j) \begin{pmatrix} 1 \\ \mathbf{0}_{d+1} \end{pmatrix} - \sigma^{'''}(\boldsymbol{w}_r^T \boldsymbol{x}_j) \|\boldsymbol{w}_{r1}\|_2^2 \boldsymbol{x}_j \rangle \end{split}$$

After expanding the inner product term, we can find that although it has nine terms, it only consists of six classes. For simplicity, we use the following six symbols to represent the corresponding classes.

$$\sigma^{''}\sigma^{''},\sigma^{''}\sigma^{'},\sigma^{'}\sigma^{'},\sigma^{'''}\sigma^{''},\sigma^{'''}\sigma^{''},\sigma^{'''}\sigma^{'''}.$$

For instance, $\sigma'' \sigma'$ represents

$$\left\langle \sigma^{''}(\boldsymbol{w}_{r}^{T}\boldsymbol{x}_{i}) \begin{pmatrix} w_{r0}\boldsymbol{x}_{i0} \\ w_{r0}\boldsymbol{x}_{i1} - 2\boldsymbol{w}_{r1} \end{pmatrix}, \sigma^{'}(\boldsymbol{w}_{r}^{T}\boldsymbol{x}_{j}) \begin{pmatrix} 1 \\ \boldsymbol{0}_{d+1} \end{pmatrix} \right\rangle, \left\langle \sigma^{'}(\boldsymbol{w}_{r}^{T}\boldsymbol{x}_{i}) \begin{pmatrix} 1 \\ \boldsymbol{0}_{d+1} \end{pmatrix}, \sigma^{''}(\boldsymbol{w}_{r}^{T}\boldsymbol{x}_{j}) \begin{pmatrix} w_{r0}\boldsymbol{x}_{j0} \\ w_{r0}\boldsymbol{x}_{j1} - 2\boldsymbol{w}_{r1} \end{pmatrix} \right\rangle.$$

In fact, when bounding the corresponding terms for $H_{ij}(w) - H_{ij}(0)$, the first four classes can be grouped into one category. They are of the form $f_1(w)f_2(w)f_3(w)f_4(w)$, where for each $i \ (1 \le i \le 4)$, $f_i(w)$ is Lipschitz continuous with respect to $\|\cdot\|_2$ and $\|f_i(w)\| \le \|w\|_2$ (Note that $\sigma'(\cdot) = (\sigma''(\cdot))^2$). On the other hand, when $\|w_1 - w_2\|_2 \le R \le 1$, we can deduce that $\|f_i(w)\| = \|f_i(w)\|_2 \|f_$

$$|f_1(\boldsymbol{w}_1)f_2(\boldsymbol{w}_1)f_3(\boldsymbol{w}_1)f_4(\boldsymbol{w}_1) - f_1(\boldsymbol{w}_2)f_2(\boldsymbol{w}_2)f_3(\boldsymbol{w}_2)f_4(\boldsymbol{w}_2)| \lesssim R(\|\boldsymbol{w}_1\|_2^3 + 1).$$

Thus, for the terms in $H_{ij}(\boldsymbol{w}) - H_{ij}(0)$ that belong to the first four classes, we can deduce that they are less than $CR(\|\boldsymbol{w}_r(0)\|_2^3 + 1)$, where C is a universal constant.

For the classes $\sigma''' \sigma''$ and $\sigma''' \sigma'$, they are both involving σ''' that is not Lipschitz continuous. To make it precise, we write the class $\sigma''' \sigma''$ explicitly as follows.

$$\sigma^{''}(oldsymbol{w}_r^Toldsymbol{x}_i)\sigma^{'''}(oldsymbol{w}_r^Toldsymbol{x}_j)\|oldsymbol{w}_{r1}\|_2^2 igg(egin{array}{c} w_{r0}x_{i0}\ w_{r0}oldsymbol{x}_{i1}-2oldsymbol{w}_{r1} igg)^Toldsymbol{x}_j.$$

Note that when $\|\boldsymbol{w}_r - \boldsymbol{w}_r(0)\|_2 < R$, we have that

$$|\sigma^{'''}(\boldsymbol{w}_r^T \boldsymbol{x}_j) - \sigma^{'''}(\boldsymbol{w}_r(0)^T \boldsymbol{x}_j)| = |I\{\boldsymbol{w}_r^T \boldsymbol{x}_j \ge 0\} - I\{\boldsymbol{w}_r(0)^T \boldsymbol{x}_j \ge 0\}| \le I\{A_{jr}\},$$

where the event A_{jr} has been defined in (36).

Thus, we can deduce that for the terms in $H_{ij}(w) - H_{ij}(0)$ that belong to the classes $\sigma''' \sigma''$ and $\sigma''' \sigma'$, they are less than

$$C\left[(I\{A_{ir}\}+I\{A_{jr}\})(\|\boldsymbol{w}_{r}(0)\|_{2}^{3}+1)+R(\|\boldsymbol{w}_{r}(0)\|_{2}^{3}+1)\right],$$

where C is a universal constant.

Similarly, for the last class $\sigma''' \sigma'''$ that are of the form

$$\sigma^{\prime\prime\prime}(oldsymbol{w}_r^Toldsymbol{x}_i)\sigma^{\prime\prime\prime}(oldsymbol{w}_r^Toldsymbol{x}_j)\|oldsymbol{w}_{r1}\|_2^4oldsymbol{x}_i^Toldsymbol{x}_j)$$

we can deduce that

$$\begin{aligned} & |\sigma^{'''}(\boldsymbol{w}_{r}^{T}\boldsymbol{x}_{i})\sigma^{'''}(\boldsymbol{w}_{r}^{T}\boldsymbol{x}_{j})\|\boldsymbol{w}_{r1}\|_{2}^{4}\boldsymbol{x}_{i}^{T}\boldsymbol{x}_{j} - \sigma^{'''}(\boldsymbol{w}_{r}(0)^{T}\boldsymbol{x}_{i})\sigma^{'''}(\boldsymbol{w}_{r}(0)^{T}\boldsymbol{x}_{j})\|\boldsymbol{w}_{r1}(0)\|_{2}^{4}\boldsymbol{x}_{i}^{T}\boldsymbol{x}_{j}|\\ & \lesssim I\{A_{ir}\lor A_{jr}\}\|\boldsymbol{w}_{r}(0)\|_{2}^{4} + R(\|\boldsymbol{w}_{r}(0)\|_{2}^{3} + 1).\end{aligned}$$

Combining the upper bounds for the terms in the six classes, we have that

j

$$|H_{ij}(\boldsymbol{w}) - H_{ij}(0)| \lesssim \frac{1}{n_1} \left[\frac{1}{m} \left(R \sum_{r=1}^m \|\boldsymbol{w}_r(0)\|_2^3 \right) + \frac{1}{m} \sum_{r=1}^m (I\{A_{ir}\} + I\{A_{jr}\}) (\|\boldsymbol{w}_r(0)\|_2^4 + \|\boldsymbol{w}_r(0)\|_2^3 + 1) + R \right]$$

$$\lesssim \frac{1}{n_1} \left[\frac{1}{m} \left(R \sum_{r=1}^m \|\boldsymbol{w}_r(0)\|_2^4 \right) + \frac{1}{m} \sum_{r=1}^m (I\{A_{ir}\} + I\{A_{jr}\}) (\|\boldsymbol{w}_r(0)\|_2^4 + 1) + R \right],$$
(46)

where the last inequality follows from that $\|\boldsymbol{w}_r(0)\|_2^3 \lesssim \|\boldsymbol{w}_r(0)\|_2^4 + 1$ due to Young's inequality for products.

729
730 Now, we focus on the term
$$\frac{1}{m} \sum_{r=1}^{m} I\{A_{ir}\} \| \boldsymbol{w}_r(0) \|_2^4$$
.
731

Since

$$P\left(|w_{ri}(0)|^2 \ge 2\log\left(\frac{2}{\delta}\right)\right) \le \delta$$

and then

$$P\left(\|\boldsymbol{w}_r(0)\|_2^2 \ge 2(d+2)\log\left(\frac{2(d+2)}{\delta}\right)\right) \le \delta$$

This implies that

$$P\left(\exists r \in [m], \|\boldsymbol{w}_r(0)\|_2^2 \ge 2(d+2)\log\left(\frac{2m(d+2)}{\delta}\right)\right) \le \delta.$$
(47)

Applying Bernstein's inequality for the first term yields that with probability at least $1 - e^{-mR}$,

$$\frac{1}{m}\sum_{r=1}^{m}I\{A_{ir}\}\leq 4R.$$

Moreover, from (47), we have that with probability at least $1 - \delta$, $I\{||\boldsymbol{w}_r(0)||_2^2 > M\} = 0$ holds for all $r \in [m]$. Thus from (46), with probability at least $1 - \delta - n_1 e^{-mR}$, we have that for any $i \in [n_1]$ and $j \in [n_1]$,

$$|H_{ij}(\boldsymbol{w}) - H_{ij}(0)| \lesssim rac{1}{n_1} \left[RM^2 + RM^2 + R
ight]$$

 $\lesssim rac{1}{n_1} M^2 R.$

For $i \in [n_1], j \in [n_1 + 2, n_2]$ and $i \in [n_1 + 1, n_2], j \in [n_2]$, from the form of $\frac{\partial h_j(w)}{\partial w_r}$, i.e.,

$$\frac{\partial h_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r} = \frac{a_r}{\sqrt{n_2 m}} \sigma'(\boldsymbol{w}_r^T \boldsymbol{y}_j) \boldsymbol{y}_j$$

770 we can obtain similar results for the terms $\left\langle \frac{\partial s_i}{\partial \boldsymbol{w}}, \frac{\partial h_j}{\partial \boldsymbol{w}} \right\rangle$ and $\left\langle \frac{\partial h_i}{\partial \boldsymbol{w}}, \frac{\partial h_j}{\partial \boldsymbol{w}} \right\rangle$.

With all results above, we have that with probability at least $1 - \delta - n_1 e^{-mR}$,

$$\|\boldsymbol{H}(\boldsymbol{w}) - \boldsymbol{H}(0)\|_F \lesssim M^2 R.$$

A.3. Proof of Lemma 3.7

779 Proof. First, we have

$$s_p(k+1) - s_p(k) = \left[s_p(k+1) - s_p(k) - \left\langle \frac{\partial s_p(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k) \right\rangle \right] + \left\langle \frac{\partial s_p(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k) \right\rangle$$

$$:= I_1^p(k) + I_2^p(k).$$
(48)

For the second term $I_2^p(k)$, from the updating rule of gradient descent, we have that

$$I_{2}^{p}(k) = \left\langle \frac{\partial s_{p}(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k) \right\rangle$$

$$= \left\langle \frac{\partial s_{p}(k)}{\partial \boldsymbol{w}}, -\eta \frac{\partial L(k)}{\partial \boldsymbol{w}} \right\rangle$$

$$= -\sum_{r=1}^{m} \eta \left\langle \frac{\partial s_{p}(k)}{\partial \boldsymbol{w}_{r}}, \frac{\partial L(k)}{\partial \boldsymbol{w}_{r}} \right\rangle$$

$$= -\sum_{r=1}^{m} \eta \left\langle \frac{\partial s_{p}(k)}{\partial \boldsymbol{w}_{r}}, \sum_{t=1}^{n_{1}} s_{t}(k) \frac{\partial s_{t}(k)}{\partial \boldsymbol{w}_{r}} + \sum_{j=1}^{n_{2}} h_{j}(k) \frac{\partial h_{j}(k)}{\partial \boldsymbol{w}_{r}} \right\rangle$$

$$= -\eta \left[\sum_{t=1}^{n_{1}} \left\langle \frac{\partial s_{p}(k)}{\partial \boldsymbol{w}_{r}}, \frac{\partial s_{t}(k)}{\partial \boldsymbol{w}_{r}} \right\rangle s_{t}(k) + \sum_{j=1}^{n_{2}} \left\langle \frac{\partial s_{p}(k)}{\partial \boldsymbol{w}_{r}}, \frac{\partial h_{j}(k)}{\partial \boldsymbol{w}_{r}} \right\rangle h_{j}(k) \right]$$

$$= -\eta [\boldsymbol{H}(k)]_{p} \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix},$$
(49)

802 where $[\boldsymbol{H}(k)]_p$ denotes the *p*-row of $\boldsymbol{H}(k)$.

Similarly, for h(k), we have

$$h_{j}(k+1) - h_{j}(k) = \left[h_{j}(k+1) - h_{j}(k) - \left\langle\frac{\partial h_{j}(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k)\right\rangle\right] + \left\langle\frac{\partial h_{j}(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k)\right\rangle$$
$$:= I_{1}^{n_{1}+j}(k) + I_{2}^{n_{1}+j}(k)$$
(50)

and

$$I_2^{n_1+j}(k) = -\eta [\boldsymbol{H}(k)]_{n_1+j} \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix}.$$
(51)

Combining (48), (49), (50) and (51) yields that

$$\begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} - \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} = \boldsymbol{I}_1(k) + \boldsymbol{I}_2(k)$$
$$= \boldsymbol{I}_1(k) - \eta \boldsymbol{H}(k) \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix}.$$

A simple transformation directly leads to

$$\begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} = (\boldsymbol{I} - \eta \boldsymbol{H}(k)) \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} + \boldsymbol{I}_1(k).$$

A.4. Proof of Theorem 3.8

 Proof. For the sake of completeness in the proof, we restate Condition 1 and Corollary 3.11 from the main text, and label them as Condition 3 and Corollary A.1, respectively.

Condition 3. At the t-th iteration, we have that for each $r \in [m]$, $\|\boldsymbol{w}_r(t)\|_2 \leq B$ and

$$L(t) \le \left(1 - \frac{\eta \lambda_0}{2}\right)^t L(0), \tag{52}$$

where $B = \sqrt{2(d+2)\log\left(\frac{2m(d+2)}{\delta}\right)} + 1$ and L(k) is an abbreviation of $L(\boldsymbol{w}(k))$.

From (47), we know that with probability at least $1 - \delta$, $\|\boldsymbol{w}_r(0)\|_2 \leq \sqrt{2(d+2)\log\left(\frac{2m(d+2)}{\delta}\right)}$ holds for all $r \in [m]$. Thus, if we can prove that $\boldsymbol{w}_r(t)$ is closed enough to $\boldsymbol{w}_r(0)$, then $\|\boldsymbol{w}_r(t)\|_2 \leq B$ holds.

Corollary A.1 (Lemma 4.1 in (Gao et al., 2023)). If Condition 3 holds for $t = 0, \dots, k$, then we have for every $r \in [m]$,

$$\|\boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(0)\|_{2} \le \frac{CB^{2}\sqrt{L(0)}}{\sqrt{m\lambda_{0}}} := R^{'},$$
(53)

where C is a universal constant.

Corollary A.1 implies that when m is large enough, we have $\|\boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(0)\|_2 \le 1$ and then $\|\boldsymbol{w}_r(k+1)\|_2 \le B$. Thus, in induction, we only need to prove that (52) also holds for t = k + 1, which relies on the recursion formula (14).

Recall that the recursion formula is

$$\begin{pmatrix} \boldsymbol{s}(k+1) \\ \boldsymbol{h}(k+1) \end{pmatrix} = (\boldsymbol{I} - \eta \boldsymbol{H}(k)) \begin{pmatrix} \boldsymbol{s}(k) \\ \boldsymbol{h}(k) \end{pmatrix} + \boldsymbol{I}_1(k) .$$

From Corollary A.1 and Lemma 3.5, taking $CM^2R < \frac{\lambda_0}{4}$ in (12) and $R' \leq R$ in (53) yields that $\lambda_{min}(\boldsymbol{H}(k)) \geq \lambda_{min}(\boldsymbol{H}(0)) - \frac{\lambda_0}{4} \geq \frac{\lambda_0}{2}$ and

$$\|\boldsymbol{H}(k)\|_{2} \leq \|\boldsymbol{H}(0)\|_{2} + \frac{\lambda_{0}}{4} \leq \|\boldsymbol{H}^{\infty}\|_{2} + \frac{\lambda_{0}}{2} \leq \frac{3}{2}\|\boldsymbol{H}^{\infty}\|_{2}.$$

Therefore, if we take $\eta \leq \frac{2}{3} \frac{1}{\|\boldsymbol{H}^{\infty}\|_2}$, then $\boldsymbol{I} - \eta \boldsymbol{H}(k)$ is positive definite and $\|\boldsymbol{I} - \eta \boldsymbol{H}(k)\|_2 \leq 1 - \frac{\eta \lambda_0}{2}$.

Combining these facts with the recursion formula, we have that

 $\left\| \left(s(k+1) \right) \right\|^2$

$$\| \left(\boldsymbol{h}(k+1) \right) \|_{2} = \| \left(\boldsymbol{I} - \eta \boldsymbol{H}(k) \right) \left(\frac{\boldsymbol{s}(k)}{\boldsymbol{h}(k)} \right) \|_{2}^{2} + \| \boldsymbol{I}_{1}(k) \|_{2}^{2} + 2 \left\langle \left(\boldsymbol{I} - \eta \boldsymbol{H}(k) \right) \left(\frac{\boldsymbol{s}(k)}{\boldsymbol{h}(k)} \right), \boldsymbol{I}_{1}(k) \right\rangle$$

$$\leq \left(1 - \frac{\eta \lambda_{0}}{2} \right)^{2} \| \left(\frac{\boldsymbol{s}(k)}{\boldsymbol{h}(k)} \right) \|_{2}^{2} + \| \boldsymbol{I}_{1}(k) \|_{2}^{2} + 2 \left(1 - \frac{\eta \lambda_{0}}{2} \right) \| \left(\frac{\boldsymbol{s}(k)}{\boldsymbol{h}(k)} \right) \|_{2} \| \boldsymbol{I}_{1}(k) \|_{2}.$$
(54)

Thus, it remains only to bound $\|I_1(k)\|_2$.

872 For $I_1(k)$, recall that $I_1(k) = (I_1^1(k), \dots, I_1^{n_1}(k), I_1^{n_1+1}(k), \dots, I_1^{n_1+n_2}(k))^T \in \mathbb{R}^{n_1+n_2}$ and for $p \in [n_1]$,

$$I_1^p(k) = s_p(k+1) - s_p(k) - \left\langle \frac{\partial s_p(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k) \right\rangle,$$

for $j \in [n_2]$,

$$I_1^{n_1+j}(k) = h_j(k+1) - h_j(k) - \left\langle \frac{\partial h_j(k)}{\partial \boldsymbol{w}}, \boldsymbol{w}(k+1) - \boldsymbol{w}(k) \right\rangle.$$

880 Recall that

$$s_p(k) = \frac{1}{\sqrt{n_1}} \left(\frac{1}{\sqrt{m}} \left(\sum_{r=1}^m a_r \sigma'(\boldsymbol{w}_r(k)^T x_p) w_{r0}(k) - a_r \sigma''(\boldsymbol{w}_r(k)^T x_p) \|\boldsymbol{w}_{r1}(k)\|_2^2 \right) - f(x_p) \right)$$

885 and

$$\begin{aligned} \frac{\partial s_p(k)}{\partial \boldsymbol{w}_r} &= \frac{a_r}{\sqrt{n_1 m}} \left[\boldsymbol{\sigma}^{\prime\prime}(\boldsymbol{w}_r(k)^T \boldsymbol{x}_p) \boldsymbol{w}_{r0}(k) \boldsymbol{x}_p + \boldsymbol{\sigma}^{\prime}(\boldsymbol{w}_r(k)^T \boldsymbol{x}_p) \begin{pmatrix} 1\\ \boldsymbol{0}_{d+2} \end{pmatrix} - \boldsymbol{\sigma}^{\prime\prime\prime}(\boldsymbol{w}_r(k)^T \boldsymbol{x}_p) \|\boldsymbol{w}_{r1}(k)\|_2^2 \boldsymbol{x}_p \\ &- 2\boldsymbol{\sigma}^{\prime\prime}(\boldsymbol{w}_r(k)^T \boldsymbol{x}_p) \begin{pmatrix} 0\\ \boldsymbol{w}_{r1}(k) \end{pmatrix} \right]. \end{aligned}$$

Define $\chi_{pr}^1(k) := \sigma'(\boldsymbol{w}_r(k)^T \boldsymbol{x}_p) w_{r0}(k)$ and $\chi_{pr}^2(k) := \sigma''(\boldsymbol{w}_r(k)^T \boldsymbol{x}_p) \|\boldsymbol{w}_{r1}(k)\|_2^2$, i.e., $\chi_{pr}^1(k)$ and $\chi_{pr}^2(k)$ are related to the operators $\frac{\partial u}{\partial t}$ and Δu respectively.

Then define

$$\hat{\chi}_{pr}^{1}(k) = \chi_{pr}^{1}(k+1) - \chi_{pr}^{1}(k) - \left\langle \frac{\partial \chi_{pr}^{1}(k)}{\partial \boldsymbol{w}_{r}}, \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k) \right\rangle$$

899 and

$$\hat{\chi}_{pr}^2(k) = \chi_{pr}^2(k+1) - \chi_{pr}^2(k) - \left\langle \frac{\partial \chi_{pr}^2(k)}{\partial \boldsymbol{w}_r}, \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k) \right\rangle$$

At this time, we have

$$I_1^p(k) = \frac{1}{\sqrt{n_1 m}} \sum_{r=1}^m a_r \left[\hat{\chi}_{pr}^1(k) - \hat{\chi}_{pr}^2(k) \right].$$

The purpose of defining $\hat{\chi}_{pr}^1(k)$ and $\hat{\chi}_{pr}^1(k)$ in this way is to enable us to handle the terms related to the operators $\frac{\partial u}{\partial t}$ and Δu separately.

909 We first recall some definitions. For $p \in [n_1]$,

$$A_{p,r} = \{ \exists \boldsymbol{w} : \| \boldsymbol{w} - \boldsymbol{w}_r(0) \|_2 \le R, I\{ \boldsymbol{w}^T \boldsymbol{x}_p \ge 0 \} \neq I\{ \boldsymbol{w}_r(0)^T \boldsymbol{x}_p \ge 0 \} \}$$

and $S_p = \{r \in [m] : I\{A_{p,r} = 0\}\}, S_p^{\perp} = [n_1] \backslash S_p$.

For $\hat{\chi}_{pr}^1(k)$, from its definition, we have that

$$\begin{aligned} \hat{\chi}_{pr}^{1}(k) &= \sigma'(\boldsymbol{w}_{r}(k+1)^{T}\boldsymbol{x}_{p})w_{r0}(k+1) - \sigma'(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})w_{r0}(k) \\ &- \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p} \rangle \sigma''(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})w_{r0}(k) - (w_{r0}(k+1) - w_{r0}(k))\sigma'(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p}) \\ &= (\sigma'(\boldsymbol{w}_{r}(k+1)^{T}\boldsymbol{x}_{p}) - \sigma'(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p}))w_{r0}(k+1) - \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p} \rangle \sigma''(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})w_{r0}(k). \end{aligned}$$

926 From the mean value theorem, we can deduce that there exists $\zeta(k) \in \mathbb{R}$ such that

$$\sigma^{'}(\boldsymbol{w}_{r}(k+1)^{T}\boldsymbol{x}_{p}) - \sigma^{'}(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p}) = \sigma^{''}(\zeta(k))\langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p}\rangle$$

930 and

$$egin{aligned} |\sigma^{''}(\zeta(k)) - \sigma^{''}(oldsymbol{w}_r(k)^Toldsymbol{x}_p)| &\leq |\zeta(k) - oldsymbol{w}_r(k)^Toldsymbol{x}_p| \ &\leq \sqrt{2} \|oldsymbol{w}_r(k+1) - oldsymbol{w}_r(k)\|_2. \end{aligned}$$

935 Then, for $\hat{\chi}_{pr}^1(k)$, we can rewrite it as follows.

$$\begin{split} \hat{\chi}_{pr}^{1}(k) &= \sigma^{''}(\zeta(k)) \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p} \rangle w_{r0}(k+1) - \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p} \rangle \sigma^{''}(\boldsymbol{w}_{r}(k)^{T} \boldsymbol{x}_{p}) w_{r0}(k) \\ &= \left[\left(\sigma^{''}(\zeta(k)) - \sigma^{''}(\boldsymbol{w}_{r}(k)^{T} \boldsymbol{x}_{p}) \right) \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p} \rangle w_{r0}(k+1) \right] \\ &+ \left[\langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p} \rangle \sigma^{''}(\boldsymbol{w}_{r}(k)^{T} \boldsymbol{x}_{p}) (w_{r0}(k+1) - w_{r0}(k)) \right]. \end{split}$$

This implies that

$$|\hat{\chi}_{pr}^{1}(k)| \lesssim B \| \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k) \|_{2}^{2}.$$

For $\hat{\chi}_{pr}^2(k)$, we write it as follows explicitly.

$$\hat{\chi}_{pr}^{2}(k) = \sigma^{''}(\boldsymbol{w}_{r}(k+1)^{T}\boldsymbol{x}_{p})\|\boldsymbol{w}_{r1}(k+1)\|_{2}^{2} - \sigma^{''}(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})\|\boldsymbol{w}_{r1}(k)\|_{2}^{2} - \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{x}_{p} \rangle \sigma^{'''}(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})\|\boldsymbol{w}_{r1}(k)\|_{2}^{2} - 2\langle \boldsymbol{w}_{r1}(k+1) - \boldsymbol{w}_{r1}(k), \boldsymbol{w}_{r1}(k) \rangle \sigma^{''}(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p}).$$
(55)

953 Note that for the term $\sigma''(\boldsymbol{w}_r(k)^T \boldsymbol{w}_p) \| \boldsymbol{w}_{r1}(k) \|_2^2$, we can rewrite it as follows.

954
$$\sigma''(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})\|\boldsymbol{w}_{r1}(k)\|_{2}^{2}$$
955
$$=\sigma''(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})\|\boldsymbol{w}_{r1}(k) - \boldsymbol{w}_{r1}(k+1) + \boldsymbol{w}_{r1}(k+1)\|_{2}^{2}$$
957
$$=\sigma''(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})[\|\boldsymbol{w}_{r1}(k) - \boldsymbol{w}_{r1}(k+1)\|_{2}^{2} + \|\boldsymbol{w}_{r1}(k+1)\|_{2}^{2} - 2\langle \boldsymbol{w}_{r1}(k+1) - \boldsymbol{w}_{r1}(k), \boldsymbol{w}_{r1}(k+1)\rangle],$$
958
$$\sigma''(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{x}_{p})[\|\boldsymbol{w}_{r1}(k) - \boldsymbol{w}_{r1}(k+1)\|_{2}^{2} + \|\boldsymbol{w}_{r1}(k+1)\|_{2}^{2} - 2\langle \boldsymbol{w}_{r1}(k+1) - \boldsymbol{w}_{r1}(k), \boldsymbol{w}_{r1}(k+1)\rangle],$$

959 where the first term $\sigma''(\boldsymbol{w}_r(k)^T \boldsymbol{x}_p) \| \boldsymbol{w}_{r1}(k) - \boldsymbol{w}_{r1}(k+1) \|_2^2 = \mathcal{O}(B \| \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k) \|_2^2).$

Plugging (56) into (55) yields that

$$\begin{aligned}
\begin{aligned}
962 \\
963 \\
964 \\
964 \\
964 \\
964 \\
964 \\
964 \\
965 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
966 \\
967 \\
968 \\
969 \\
967 \\
968 \\
969 \\
970 \\
968 \\
969 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\
970 \\$$

975 Thus, we only need to consider the term

$$\sigma^{\prime\prime}(\boldsymbol{w}_r(k+1)^T\boldsymbol{x}_p) - \sigma^{\prime\prime}(\boldsymbol{w}_r(k)^T\boldsymbol{x}_p) - \langle \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k), \boldsymbol{x}_p \rangle \sigma^{\prime\prime\prime}(\boldsymbol{w}_r(k)^T\boldsymbol{x}_p).$$

979 For $r \in S_p$, since $\|\boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(0)\|_2 \leq R$, $\|\boldsymbol{w}_r(k) - \boldsymbol{w}_r(0)\|_2 \leq R$, we have that $I\{\boldsymbol{w}_r(k+1)^T\boldsymbol{x}_p \geq 0\} = I\{\boldsymbol{w}_r(k)^T\boldsymbol{x}_p \geq 0\}$, which yields that

990 001	For $r \in S_p^{\perp}$, the Lipschitz continuity of $\sigma^{''}$ implies that	
992 993	$\sigma^{\prime\prime}(\boldsymbol{w}_r(k+1)^T\boldsymbol{x}_p) - \sigma^{\prime\prime}(\boldsymbol{w}_r(k)^T\boldsymbol{x}_p) - \langle \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k), \boldsymbol{x}_p \rangle \sigma^{\prime\prime\prime}(\boldsymbol{w}_r(k)^T\boldsymbol{x}_p) = \mathcal{O}(\ \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k)\ _2).$	(59)
994 995	Combining (57), (58) and (59), we can deduce that for $r \in S_p$,	
996	$ \hat{\chi}_{nr}^2(k) \lesssim B \ oldsymbol{w}_r(k+1) - oldsymbol{w}_r(k)\ _2^2$	
997		
999	and for $r \in S_p^+$, $ ^2(l) \leq D _{(l+1)} = (l) ^2 + D^2 _{(l+1)} = (l) _{l}$	
1000	$ \chi_{pr}^{-}(k) \gtrsim B \ \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k) \ _{2}^{2} + B^{2} \ \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k) \ _{2}^{2}.$	
1001	With the estimations for $\hat{\chi}_{pr}^1(k)$ and $\hat{\chi}_{pr}^2(k)$, we have	
1003	m m	
1004	$ I_1^p(k) \le \frac{1}{\sqrt{m}} \sum_{k=1}^{m} (\hat{\chi}_{nr}^1(k) + \hat{\chi}_{nr}^2(k))$	
1005	$\sqrt{n_1 m} \sum_{r=1}^{n_1 m} \sum_{$	(60)
1007	$\leq \frac{1}{2} \sum_{k=1}^{m} B \ \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k) \ _{2}^{2} + \frac{1}{2} \sum_{k=1}^{m} B^{2} \ \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k) \ _{2}$	(00)
1008	$\sim \sqrt{n_1 m} \sum_{r=1}^{2} - \ u_r(u+1) - u_r(v) \ _2 + \sqrt{n_1 m} \sum_{r\in S_p^{\perp}} - \ u_r(u+1) - u_r(v) \ _2$	
1009		
1011	For $j \in [n_2]$, we consider $I_1^{n_1+j}(k)$, which can be written as follows.	
1012	$\langle 2L(1) \rangle$	
1013	$I_1^{n_1+j}(k) = h_j(k+1) - h_j(k) - \left\langle \boldsymbol{w}(k+1) - \boldsymbol{w}(k), \frac{\partial h_j(k)}{\partial \boldsymbol{w}} \right\rangle$	
1014		
1016	$= \sum \frac{u_r}{\sqrt{n_2m}} \left \sigma(\boldsymbol{w}_r(k+1)^T \boldsymbol{y}_j) - \sigma(\boldsymbol{w}_r(k)^T \boldsymbol{y}_j) - \langle \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k), \boldsymbol{y}_j \rangle \sigma'(\boldsymbol{w}_r(k)^T \boldsymbol{y}_j) \right .$	
1017	$r=1$ V r_2 r_2 r_2	
1018	From the mean value theorem, we have that there exists $\zeta(k) \in \mathbb{R}$ such that	
1020	· · · · · · · · · · · · · · · · · · ·	
1021	$\sigma(\boldsymbol{w}_r(k+1)^T\boldsymbol{y}_j) - \sigma(\boldsymbol{w}_r(k)^T\boldsymbol{y}_j) = \sigma^{'}(\zeta(k)) \langle \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k), \boldsymbol{y}_j \rangle$	
1022	and	
1024		
1025	$ \sigma^{'}(\zeta(k))-\sigma^{'}(oldsymbol{w}_{r}(k)^{T}oldsymbol{y}_{j}) \leq 2B \zeta(k)-oldsymbol{w}_{r}(k)^{T}oldsymbol{y}_{j} $	
1026	$\leq 2\sqrt{2}B\ oldsymbol{w}_r(k+1)-oldsymbol{w}_r(k)\ _2.$	
1028	Thus	
1029		
1030	$ \sigma(\boldsymbol{w}_r(k+1)^T\boldsymbol{y}_j) - \sigma(\boldsymbol{w}_r(k)^T\boldsymbol{y}_j) - \langle \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k), \boldsymbol{y}_j \rangle \sigma^{'}(\boldsymbol{w}_r(k)^T\boldsymbol{y}_j) $	
1031	$= \sigma^{'}(\zeta(k)) \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{y}_{j} \rangle - \sigma(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{y}_{j}) - \langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{y}_{j} \rangle \sigma^{'}(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{y}_{j}) $	
1033	$= (\sigma^{'}(\zeta(k)) - \sigma^{'}(\boldsymbol{w}_{r}(k)^{T}\boldsymbol{y}_{j}))\langle \boldsymbol{w}_{r}(k+1) - \boldsymbol{w}_{r}(k), \boldsymbol{y}_{j}\rangle $	
1034	$\lesssim B \ oldsymbol{w}_r(k+1) - oldsymbol{w}_r(k) \ _2.$	
1035	Therefore for $i \in [n]$	
1037	Therefore, for $j \in [n_2]$,	
1038	$ I_1^{n_1+j}(k) \lesssim rac{D}{\sqrt{m_2m}} \sum \ \boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k)\ _2^2.$	(61)
1039	$v \cdots z \cdots r = 1$	
1041	From the updating rule of gradient descent, we can deduce that for every $r \in [m]$,	
1042	$\ a_{m}(k+1) - a_{m}(k) \ = \ \partial L(k) \ \geq \eta B^{2} \sqrt{L(k)}$	(60)
1043 1044	$\ \boldsymbol{w}_r(\kappa+1) - \boldsymbol{w}_r(\kappa)\ _2 = \left\ -\eta \overline{\partial \boldsymbol{w}_r}\right\ _2 \gtrsim \overline{\sqrt{m}} \sqrt{L(\kappa)}.$	(02)

1045 Plugging (62) into (61) and (60), we can deduce that 1046

$$|I_1^p(k)| \lesssim \frac{B}{\sqrt{n_1 m}} \sum_{r=1}^m \|\boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k)\|_2^2 + \frac{B^2}{\sqrt{n_1 m}} \sum_{r \in S_p^\perp} \|\boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k)\|_2$$
$$\lesssim \frac{B}{\sqrt{m_1 m}} \sum_{r=1}^m \frac{\eta^2 B^4}{2} L(k) + \frac{B^2}{\sqrt{m_1 m}} \sum_{r \in S_p^\perp} \frac{\eta B^2}{\sqrt{L(k)}} \sqrt{L(k)}$$

$$\approx \frac{1}{\sqrt{n_1 m}} \sum_{r=1}^{\infty} \frac{1}{m} L(k) + \frac{1}{\sqrt{n_1 m}} \sum_{r \in S_p^{\perp}} \sqrt{m} \sqrt{L(k)}$$

$$= \frac{\eta^2 B^5 L(k)}{\sqrt{n_1 m}} + \frac{\eta B^4 \sqrt{L(k)}}{\sqrt{n_1}} \frac{1}{m} \sum_{r=1}^{m} I\{r \in S_p^{\perp}\}$$

$$\leq \frac{\eta^2 B^5 \sqrt{L(0)} \sqrt{L(k)}}{\sqrt{n_1 m}} + \frac{\eta B^4 \sqrt{L(k)}}{\sqrt{n_1}} \frac{1}{m} \sum_{r=1}^{m} I\{r \in S_p^{\perp}\}$$
(63)

058 and

$$P(A_{p,r}) \le \frac{2R}{\sqrt{2\pi}}, \ S_p = \{r \in [m] : I\{A_{p,r}\} = 0\}$$

 $|I_1^{n_1+j}(k)| \lesssim \frac{B}{\sqrt{n_2m}} \sum_{r=1}^m \|\boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(k)\|_2^2$

 $\lesssim \frac{B}{\sqrt{n_2 m}} \sum_{r=1}^m \frac{\eta^2 B^4}{m} L(k)$

 $\leq \frac{\eta^2 B^5 \sqrt{L(0)} \sqrt{L(k)}}{\sqrt{n_2 m}}.$

¹⁰⁷⁰ Thus, from Bernstein's inequality, we have that with probability at least $1 - e^{-mR}$,

$$\frac{1}{m}\sum_{r=1}^{m}I\{r\in S_{p}^{\perp}\} = \frac{1}{m}\sum_{r=1}^{m}I\{A_{pr}\} \lesssim 4R$$

Then the inequality holds for all $p \in [n_1]$ with probability at least $1 - n_1 e^{-mR}$. Plugging this into (63), we can conclude that for every $p \in [n_1]$

$$|I_1^p(k)| \lesssim \frac{\eta^2 B^5 \sqrt{L(0)} \sqrt{L(k)}}{\sqrt{n_1 m}} + \frac{\eta B^4 \sqrt{L(k)}}{\sqrt{n_1}} R.$$
(65)

(64)

Combining (64) and (65), we have that

$$\|\mathbf{I}_{1}(k)\|_{2} = \sqrt{\sum_{p=1}^{n_{1}} |I_{1}^{p}(k)|^{2} + \sum_{j=1}^{n_{2}} |I_{1}^{n_{1}+j}(k)|^{2}}$$
$$\lesssim \frac{\eta^{2} B^{5} \sqrt{L(0)} \sqrt{L(k)}}{\sqrt{m}} + \eta B^{4} \sqrt{L(k)} R.$$

1087 Plugging this into (54) yields that

 $\|\langle \mathbf{s}(k+1)\rangle\|^2$

$$\begin{split} \left\| \begin{pmatrix} \mathbf{s}_{(k+1)} \\ \mathbf{h}_{(k+1)} \end{pmatrix} \right\|_{2} \\ &\leq \left(1 - \frac{\eta \lambda_{0}}{2} \right)^{2} \left\| \begin{pmatrix} \mathbf{s}_{(k)} \\ \mathbf{h}_{(k)} \end{pmatrix} \right\|_{2}^{2} + \|\mathbf{I}_{1}(k)\|_{2}^{2} + 2\left(1 - \frac{\eta \lambda_{0}}{2} \right) \left\| \begin{pmatrix} \mathbf{s}_{(k)} \\ \mathbf{h}_{(k)} \end{pmatrix} \right\|_{2} \|\mathbf{I}_{1}(k)\|_{2} \\ &\leq \left[\left(1 - \frac{\eta \lambda_{0}}{2} \right)^{2} + C^{2} \left(\frac{\eta^{2} B^{5} \sqrt{L(0)}}{\sqrt{m}} + \eta B^{4} R \right)^{2} + 2C \left(\frac{\eta^{2} B^{5} \sqrt{L(0)}}{\sqrt{m}} + \eta B^{4} R \right) \right] \left\| \begin{pmatrix} \mathbf{s}_{(k)} \\ \mathbf{h}_{(k)} \end{pmatrix} \right\|_{2}^{2} \\ &\leq \left(1 - \frac{\eta \lambda_{0}}{2} \right) \left\| \begin{pmatrix} \mathbf{s}_{(k)} \\ \mathbf{h}_{(k)} \end{pmatrix} \right\|_{2}^{2}, \end{split}$$

where C is a universal constant and the last inequality requires that

$$\frac{\eta^2 B^5 \sqrt{L(0)}}{\sqrt{m}} \lesssim \eta \lambda_0, \ \eta B^4 R \lesssim \eta \lambda_0$$

Recall that we also require $CM^2R < \frac{\lambda_0}{4}$ for R in (12) and

$$\boldsymbol{R'} = \frac{CB^2\sqrt{L(0)}}{\sqrt{m}\lambda_0} < \boldsymbol{R}$$

for R' in (53) to make sure $\|H(k) - H(0)\|_2 \le \frac{\lambda_0}{4}$.

Finally, with $R = O(\frac{\lambda_0}{M^2})$ and Lemma C.4 for the upper bound of L(0), m needs to satisfies that

$$m = \Omega\left(\frac{M^4 B^4 L(0)}{\lambda_0^4}\right) = \Omega\left(\frac{d^8}{\lambda_0^4} \log^6\left(\frac{md}{\delta}\right) \log\left(\frac{n_1 + n_2}{\delta}\right)\right).$$

B. Proof of Section 4

B.1. Proof of Lemma 4.3

Proof. Recall that

$$\boldsymbol{H}(\boldsymbol{w}) = \boldsymbol{D}^T \boldsymbol{D}, \quad \boldsymbol{D} = \begin{bmatrix} \frac{\partial s_1(\boldsymbol{w})}{\partial \boldsymbol{w}}, \cdots, \frac{\partial s_{n_1}(\boldsymbol{w})}{\partial \boldsymbol{w}}, \frac{\partial h_1(\boldsymbol{w})}{\partial \boldsymbol{w}}, \cdots, \frac{\partial h_{n_2}(\boldsymbol{w})}{\partial \boldsymbol{w}} \end{bmatrix},$$

and $H^{\infty} = \mathbb{E}_{w \sim \mathcal{N}(\mathbf{0}, I)} G(w).$

We denote $\varphi(\boldsymbol{x}; \boldsymbol{w}) = \sigma'(\boldsymbol{w}^T \boldsymbol{x}) w_0 - \sigma''(\boldsymbol{w}^T \boldsymbol{x}) \|\boldsymbol{w}_1\|_2^2$, where $\boldsymbol{w} = (w_0, \boldsymbol{w}_1^T)^T$, $w_0 \in \mathbb{R}, \boldsymbol{w}_1 \in \mathbb{R}^d$, then

$$rac{\partial s_p(oldsymbol{w})}{\partial oldsymbol{w}_r} = rac{1}{\sqrt{n_1}} rac{a_r}{\sqrt{m}} rac{\partial arphi(oldsymbol{x}_p;oldsymbol{w}_r)}{\partial oldsymbol{w}_r}.$$

Similarly, we denote $\psi(\boldsymbol{y}; \boldsymbol{w}) = \sigma(\boldsymbol{w}^T \boldsymbol{y})$, then

$$\frac{\partial h_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r} = \frac{1}{\sqrt{n_2}} \frac{a_r}{\sqrt{m}} \frac{\partial \psi(\boldsymbol{y}_j, \boldsymbol{w}_r)}{\partial \boldsymbol{w}_r}$$

With the notations, we can deduce that

$$\begin{array}{l} 1139\\ 1140\\ 1141\\ 1142\\ 1143\\ 1144\\ 1145\\ 1146 \end{array} \qquad H_{p,j}^{\infty} = \begin{cases} \frac{1}{n_1} \mathbb{E}_{\boldsymbol{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})} \left\langle \frac{\partial \varphi(\boldsymbol{x}_p; \boldsymbol{w})}{\partial \boldsymbol{w}}, \frac{\partial \varphi(\boldsymbol{x}_j; \boldsymbol{w})}{\partial \boldsymbol{w}} \right\rangle, & 1 \leq p \leq n_1, 1 \leq j \leq n_1, \\ \frac{1}{\sqrt{n_1 n_2}} \mathbb{E}_{\boldsymbol{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})} \left\langle \frac{\partial \varphi(\boldsymbol{x}_p; \boldsymbol{w})}{\partial \boldsymbol{w}}, \frac{\partial \psi(\boldsymbol{y}_j; \boldsymbol{w})}{\partial \boldsymbol{w}} \right\rangle, & 1 \leq p \leq n_1, n_1 + 1 \leq j \leq n_1 + n_2, \\ \frac{1}{n_2} \mathbb{E}_{\boldsymbol{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})} \left\langle \frac{\partial \psi(\boldsymbol{y}_p; \boldsymbol{w})}{\partial \boldsymbol{w}}, \frac{\partial \psi(\boldsymbol{y}_j; \boldsymbol{w})}{\partial \boldsymbol{w}} \right\rangle, & n_1 + 1 \leq p \leq n_1 + n_2, n_1 + 1 \leq j \leq n_1 + n_2, \\ \end{array}$$

where $H_{p,j}^{\infty}$ is the (p, j)-th entry of H^{∞} .

The proof of this lemma requires tools from functional analysis. Let \mathcal{H} be a Hilbert space of integrable (d+2)-dimensional vector fields on \mathbb{R}^{d+2} , i.e., $f \in \mathcal{H}$ if $\mathbb{E}_{w \sim \mathcal{N}(\mathbf{0}, I)}[\|f(w)\|_2^2] < \infty$. The inner product for any two elements f, g in \mathcal{H} is $\mathbb{E}_{\boldsymbol{w} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})}[\langle f(\boldsymbol{w}), g(\boldsymbol{w}) \rangle]$. Thus, proving \boldsymbol{H}^{∞} is strictly positive definite is equivalent to show that

$$\frac{\partial \varphi(\boldsymbol{x}_{1}; \boldsymbol{w})}{\partial \boldsymbol{w}}, \cdots, \frac{\partial \varphi(\boldsymbol{x}_{n_{1}}; \boldsymbol{w})}{\partial \boldsymbol{w}}, \frac{\partial \psi(\boldsymbol{y}_{1}; \boldsymbol{w})}{\partial \boldsymbol{w}}, \cdots, \frac{\partial \psi(\boldsymbol{y}_{n_{2}}; \boldsymbol{w})}{\partial \boldsymbol{w}} \in \mathcal{H}$$

are linearly independent. Suppose that there are $\alpha_1, \dots, \alpha_{n_1}, \beta_1, \dots, \beta_{n_2} \in \mathbb{R}$ such that

$$\alpha_1 \frac{\partial \varphi(\boldsymbol{x}_1; \boldsymbol{w})}{\partial \boldsymbol{w}} + \dots + \alpha_{n_1} \frac{\partial \varphi(\boldsymbol{x}_{n_1}; \boldsymbol{w})}{\partial \boldsymbol{w}} + \beta_1 \frac{\partial \psi(\boldsymbol{y}_1; \boldsymbol{w})}{\partial \boldsymbol{w}} + \dots + \beta_{n_2} \frac{\partial \psi(\boldsymbol{y}_{n_2}; \boldsymbol{w})}{\partial \boldsymbol{w}} = 0 \text{ in } \mathcal{H}.$$

This implies that

$$\alpha_1 \frac{\partial \varphi(\boldsymbol{x}_1; \boldsymbol{w})}{\partial \boldsymbol{w}} + \dots + \alpha_{n_1} \frac{\partial \varphi(\boldsymbol{x}_{n_1}; \boldsymbol{w})}{\partial \boldsymbol{w}} + \beta_1 \frac{\partial \psi(\boldsymbol{y}_1; \boldsymbol{w})}{\partial \boldsymbol{w}} + \dots + \beta_{n_2} \frac{\partial \psi(\boldsymbol{y}_{n_2}; \boldsymbol{w})}{\partial \boldsymbol{w}} = 0$$
(66)

holds for all $w \in \mathbb{R}^{d+1}$, as $\sigma(\cdot)$ is smooth.

We first compute the derivatives of φ and ψ . Differentiating $\psi(\boldsymbol{y}; \boldsymbol{w}) k$ times with respect to \boldsymbol{w} , we have

$$rac{\partial^k \psi(oldsymbol{y};oldsymbol{w})}{\partial oldsymbol{w}^k} = \sigma^{(k)}(oldsymbol{w}^Toldsymbol{y})oldsymbol{y}^{\otimes(k)},$$

where \otimes denotes tensor product.

For $\varphi(\boldsymbol{x}; \boldsymbol{w})$, let $\varphi_0(\boldsymbol{x}; \boldsymbol{w}) = \sigma'(\boldsymbol{w}^T \boldsymbol{x}) w_0$, $\varphi_i(\boldsymbol{x}; \boldsymbol{w}) = \sigma''(\boldsymbol{w}^T \boldsymbol{x}) w_i^2$ for $1 \le i \le d$, then

$$arphi(oldsymbol{x};oldsymbol{w}) = arphi_0(oldsymbol{x};oldsymbol{w}) - \sum_{i=1}^d arphi_i(oldsymbol{x};oldsymbol{w}).$$

Differentiating $\varphi_0(x; w)$ k times with respect to w, similar to the Leibniz rule for the k-th derivative of the product of two scalar functions, we have

$$\frac{\partial^k \varphi_0(\boldsymbol{x}; \boldsymbol{w})}{\partial \boldsymbol{w}^k} = \sigma^{(k+1)}(\boldsymbol{w}^T \boldsymbol{x}) w_0 \boldsymbol{x}^{\otimes (k)} + k \sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}) \boldsymbol{e}_0 \otimes \boldsymbol{x}^{\otimes (k-1)},$$
(67)

where $e_0 = (1, 0, \dots, 0)^T \in \mathbb{R}^{d+2}$.

Similarly, for $\varphi_i(\boldsymbol{x}; \boldsymbol{w}), 1 \leq i \leq d$, we have

$$\frac{1183}{1184} \qquad \frac{\partial^{k}\varphi_{i}(\boldsymbol{x};\boldsymbol{w})}{\partial\boldsymbol{w}^{k}} = \sigma^{(k+2)}(\boldsymbol{w}^{T}\boldsymbol{x})w_{i}^{2}\boldsymbol{x}^{\otimes(k)} + C_{k}^{1}\sigma^{(k+1)}(\boldsymbol{w}^{T}\boldsymbol{x})2w_{i}\boldsymbol{e}_{i}\otimes\boldsymbol{x}^{\otimes(k-1)} + C_{k}^{2}\sigma^{(k)}(\boldsymbol{w}^{T}\boldsymbol{x})2\boldsymbol{e}_{i}^{\otimes(2)}\otimes\boldsymbol{x}^{\otimes(k-2)}, \quad (68)$$

where $e_i \in \mathbb{R}^{d+2}$, is a vector where all other components are 0, and only the (i + 1)-th component is 1.

Combining the results in (67) and (68) for the derivatives of $\varphi_0(x; w), \dots, \varphi_d(x; w)$ yields that

Note that when no two points in $\{x_1, \dots, x_{n_1}, y_1, \dots, y_{n_2}\}$ are parallel,

$$m{x}_1^{\otimes (n_1+n_2-1)}, \cdots, m{x}_{n_1}^{\otimes (n_1+n_2-1)}, m{y}_1^{\otimes (n_1+n_2-1)}, \cdots, m{y}_{n_2}^{\otimes (n_1+n_2-1)}$$

are independent (see Lemma G.6 in (Du et al., 2018)). It motivates us to differentiate both sides in (66) (k - 1) times for w with $k = n_1 + n_2 + 1$, then we have

$$\alpha_1 \frac{\partial^k \varphi(\boldsymbol{x}_1; \boldsymbol{w})}{\partial \boldsymbol{w}^k} + \dots + \alpha_{n_1} \frac{\partial^k \varphi(\boldsymbol{x}_{n_1}; \boldsymbol{w})}{\partial \boldsymbol{w}^k} + \beta_1 \frac{\partial^k \psi(\boldsymbol{y}_1; \boldsymbol{w})}{\partial \boldsymbol{w}^k} + \dots + \beta_{n_2} \frac{\partial^k \psi(\boldsymbol{y}_{n_2}; \boldsymbol{w})}{\partial \boldsymbol{w}^k} = 0.$$
(70)

1210 we can deduce from the independence of the tensors that for any $j \in [n_2]$ and $w \in \mathbb{R}^{d+2}$,

$$\beta_j \sigma^{(n_1+n_2+1)} (\boldsymbol{w}^T \boldsymbol{y}_j) \boldsymbol{y}_j^{\otimes (2)} = 0$$

Now, we can choose a \boldsymbol{w} such that $\sigma^{(n_1+n_2+1)}(\boldsymbol{w}^T\boldsymbol{y}_j)\neq 0$, thus

$$\beta_j \sigma^{(n_1+n_2+1)}(\boldsymbol{w}^T \boldsymbol{y}_j) \|\boldsymbol{y}_j\|_2^2 = Trace\left(\beta_j \sigma^{(n_1+n_2+1)}(\boldsymbol{w}^T \boldsymbol{y}_j) \boldsymbol{y}_j^{\otimes(2)}\right) = 0,$$

which implies $\beta_j = 0$ and then this holds for all $j \in [n_2]$.

Similarly, for α_i , $i \in [n_1]$, from (67) and (68), we have

$$\alpha_i [\sigma^{(k+1)}(\boldsymbol{w}^T \boldsymbol{x}_i) w_0 \boldsymbol{x}_i^{\otimes(2)} + k \sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}_i) \boldsymbol{e}_0 \otimes \boldsymbol{x}_i - \sum_{j=1}^d \left(\sigma^{(k+2)}(\boldsymbol{w}^T \boldsymbol{x}_i) w_j^2 \boldsymbol{x}_i^{\otimes(2)} + C_k^1 \sigma^{(k+1)}(\boldsymbol{w}^T \boldsymbol{x}_i) 2 w_j \boldsymbol{e}_j \otimes \boldsymbol{x}_i + C_k^2 \sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}_i) 2 \boldsymbol{e}_j^{\otimes(2)} \right)] = 0.$$

For fixed *i*, denote $\boldsymbol{x}_i = (x_{i0}, \boldsymbol{x}_{i1}^T)^T$ with $x_{i0} \in \mathbb{R}$ and $\boldsymbol{x}_{i1} \in \mathbb{R}^{d+1}$. We consider two cases: (1) $\boldsymbol{x}_{i1} \neq 0$; (2) $\boldsymbol{x}_{i1} = 0$. Although \boldsymbol{x}_i has been augmented so that $\boldsymbol{x}_{i1} \neq 0$, we still consider Case 2 to account for scenarios where we might need to use neural networks without a bias term.

In the case(1), we can let $w_0 = 0$, thus for $k = n_1 + n_2 + 1$,

$$\begin{bmatrix} 1232 \\ 1233 \\ 1234 \\ 1235 \end{bmatrix} \alpha_i \left[k\sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}_i) \boldsymbol{e}_0 \otimes \boldsymbol{x}_i - \sum_{j=1}^d \left(\sigma^{(k+2)}(\boldsymbol{w}^T \boldsymbol{x}_i) w_j^2 \boldsymbol{x}_i^{\otimes(2)} + C_k^1 \sigma^{(k+1)}(\boldsymbol{w}^T \boldsymbol{x}_i) 2w_j \boldsymbol{e}_j \otimes \boldsymbol{x}_i + C_k^2 \sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}_i) 2\boldsymbol{e}_j^{\otimes(2)} \right) \right] = 0,$$

$$\begin{bmatrix} 1232 \\ 1235 \end{bmatrix}$$

$$which implies that the trace is 0, i.e.$$

which implies that the trace is 0, i.e.,

$$\alpha_{i} \left[k\sigma^{(k)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})x_{i,0} - \sum_{j=1}^{d} \left(\sigma^{(k+2)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})w_{j}^{2} \|\boldsymbol{x}_{i}\|_{2}^{2} + C_{k}^{1}\sigma^{(k+1)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})2w_{j}\boldsymbol{x}_{ij} + C_{k}^{2}\sigma^{(k)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})2 \right) \right] = 0$$

Rearranging it yields that

$$\alpha_i \left[k \sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}_i) \boldsymbol{x}_{i0} - \left(\sigma^{(k+2)}(\boldsymbol{w}^T \boldsymbol{x}_i) \| \boldsymbol{w}_1 \|_2^2 \| \boldsymbol{x}_i \|_2^2 + 2C_k^1 \sigma^{(k+1)}(\boldsymbol{w}^T \boldsymbol{x}_i) \boldsymbol{w}_1^T \boldsymbol{x}_{i1} + 2C_k^2 \sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}_i) \right) \right] = 0$$

Now, we can set $w_1^T x_{i1} = c$ such that $\sigma^{(k+1)}(c) \neq 0$. Since $w^T x_i = w_1^T x_{i1} = c$, we have

$$\alpha_i \left[k\sigma^{(k)}(c) \boldsymbol{x}_{i,0} - \left(\sigma^{(k+2)}(c) \| \boldsymbol{w}_1 \|_2^2 \| \boldsymbol{x}_i \|_2^2 + 2C_k^1 \sigma^{(k+1)}(c) c + 2C_k^2 \sigma^{(k)}(c) \right) \right] = 0$$

Note that the only variable in above equation is w_1 and in the hyperplane $\{w_1 : w_1^T x_{i,1} = c\}, \|w_1\|_2$ can be selected to tend infinite, thus $\alpha_i = 0$.

In the case(2), we have that

$$\begin{array}{c} 1256\\ 1257\\ 1258\\ 1258\\ 1259 \end{array} \alpha_{i} \left[\sigma^{(k+1)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})w_{0}\boldsymbol{x}_{i}^{\otimes(2)} + k\sigma^{(k)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})\boldsymbol{e}_{0}\otimes\boldsymbol{x}_{i} - \sum_{j=1}^{d} \left(\sigma^{(k+2)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})w_{j}^{2}\boldsymbol{x}_{i}^{\otimes(2)} + C_{k}^{2}\sigma^{(k)}(\boldsymbol{w}^{T}\boldsymbol{x}_{i})2\boldsymbol{e}_{j}^{\otimes(2)} \right) \right] = 0.$$

From the observation of the (t, t)-th entry of the matrix above with $t \ge 2$, we have that

$$\alpha_i(-C_k^2 \sigma^{(k)}(\boldsymbol{w}^T \boldsymbol{x}_i)2) = -2\alpha_i C_k^2 \sigma^{(k)}(w_0 x_{i0}) = 0$$

Then, taking a w_0 such that $\sigma^{(k)}(w_0 x_{i0}) \neq 0$ yields the conclusion.

B.2. Proof of Lemma 4.5

Proof. Recall that

$$\begin{aligned} \frac{\partial s_p(\boldsymbol{w})}{\partial \boldsymbol{w}_r} &= \frac{a_r}{\sqrt{n_1 m}} \left[\sigma^{''}(\boldsymbol{w}_r^T \boldsymbol{x}_p) w_{r0} \boldsymbol{x}_p + \sigma^{'}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \begin{pmatrix} 1\\ \boldsymbol{0}_{d+1} \end{pmatrix} - \sigma^{'''}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \|\boldsymbol{w}_{r1}\|_2^2 \boldsymbol{x}_p \\ &- 2\sigma^{''}(\boldsymbol{w}_r^T \boldsymbol{x}_p) \begin{pmatrix} 0\\ \boldsymbol{w}_{r1} \end{pmatrix} \right] \end{aligned}$$

1274 and

$$\frac{\partial h_{j}(\boldsymbol{w})}{\partial \boldsymbol{w}_{r}} = \frac{a_{r}}{\sqrt{n_{2}m}}\sigma^{'}(\boldsymbol{w}_{r}^{T}\boldsymbol{y}_{j})\boldsymbol{y}_{j}$$

 $^{1277}_{1278}$ (1) When $\sigma(\cdot)$ is the ReLU 3 activation function.

1270 From the form of $\frac{\partial s_p(w)}{\partial w_r}$, we can deduce that 1280

$$\begin{aligned}
& \left\| \frac{\partial s_{p}(\boldsymbol{w})}{\partial \boldsymbol{w}_{r}} - \frac{\partial s_{p}(0)}{\partial \boldsymbol{w}_{r}} \right\|_{2} \\
& 1283 \\
& 1283 \\
& 1284 \\
& 1285 \\
& 1285 \\
& 1285 \\
& 1285 \\
& 1286 \\
& 1287 \\
& 1287 \\
& 1288 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
& 1 = 1 \\
&$$

where the second inequality follows from the fact $\|\boldsymbol{w} - \boldsymbol{w}_r(0)\|_2 < R \le 1$ and the definition of A_{pr} in (36).

1290 Similarly, we have that 1291

$$\left\|\frac{\partial h_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r} - \frac{\partial h_j(0)}{\partial \boldsymbol{w}_r}\right\|_2 \lesssim \frac{1}{\sqrt{n_2 m}} R(\|\boldsymbol{w}_r(0)\|_2 + 1).$$
(72)

1294 Combining (71) and (72), we can deduce that

$$\begin{aligned} \|\mathbf{J}(\boldsymbol{w}) - \mathbf{J}(0)\|_{F}^{2} \\ &\leq \|\mathbf{J}(\boldsymbol{w}) - \mathbf{J}(0)\|_{F}^{2} \\ &\leq \|\mathbf{J}(\boldsymbol{w}) - \mathbf{J}(0)\|_{F}^{2} \\ &= \sum_{i=1}^{n_{1}+n_{2}} \|\mathbf{J}_{i}(\boldsymbol{w}) - \mathbf{J}_{i}(0)\|_{2}^{2} \\ &= \sum_{i=1}^{n_{1}+n_{2}} \|\mathbf{J}_{i}(\boldsymbol{w}) - \mathbf{J}_{i}(0)\|_{2}^{2} \\ &= \sum_{i=1}^{m} \left(\sum_{p=1}^{n_{1}} \left\|\frac{\partial s_{p}(\boldsymbol{w})}{\partial \boldsymbol{w}_{r}} - \frac{\partial s_{p}(0)}{\partial \boldsymbol{w}_{r}}\right\|_{2}^{2} + \sum_{j=1}^{n_{2}} \left\|\frac{\partial h_{j}(\boldsymbol{w})}{\partial \boldsymbol{w}_{r}} - \frac{\partial h_{j}(0)}{\partial \boldsymbol{w}_{r}}\right\|_{2}^{2} \right) \\ &\leq \sum_{r=1}^{m} \left(\sum_{p=1}^{n_{1}} \left\|\frac{\partial s_{p}(\boldsymbol{w})}{\partial \boldsymbol{w}_{r}} - \frac{\partial s_{p}(0)}{\partial \boldsymbol{w}_{r}}\right\|_{2}^{2} + 1) + I\{A_{pr}\}(\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1)\right)^{2} + \sum_{j=1}^{n_{2}} \frac{1}{n_{2}m}(R\|\boldsymbol{w}_{r}(0)\|_{2} + R)^{2} \right) \\ &\leq \sum_{r=1}^{m} \left(\sum_{p=1}^{n_{1}} \frac{1}{n_{1}m} \left(R(\|\boldsymbol{w}_{r}(0)\|_{2} + 1) + I\{A_{pr}\}(\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1)\right)\right) \\ &\leq \frac{R^{2}}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1) + \frac{1}{n_{1}m} \sum_{p=1}^{n_{1}} \sum_{r=1}^{m} I\{A_{pr}\}(\|\boldsymbol{w}_{r}(0)\|_{2}^{4} + 1) \\ &= \frac{R^{2}}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1) \\ &+ \frac{1}{n_{1}m} \sum_{p=1}^{n_{1}} \sum_{r=1}^{m} I\{A_{pr}\} \left(\|\boldsymbol{w}_{r}(0)\|_{2}^{2} \leq M\} + \|\boldsymbol{w}_{r}(0)\|_{2}^{4} I\{\|\boldsymbol{w}_{r}(0)\|_{2}^{2} > M\} + 1\right) \\ &\leq \frac{R^{2}}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1) + \frac{M^{2}}{n_{1}m} \sum_{p=1}^{m} \sum_{r=1}^{m} I\{A_{pr}\} + \frac{1}{m} \sum_{r=1}^{m} \|\boldsymbol{w}_{r}(0)\|_{2}^{4} I\{\|\boldsymbol{w}_{r}(0)\|_{2}^{2} > M\}, \\ &\leq \frac{R^{2}}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1) + \frac{M^{2}}{n_{1}m} \sum_{p=1}^{m} \sum_{r=1}^{m} I\{A_{pr}\} + \frac{1}{m} \sum_{r=1}^{m} \|\boldsymbol{w}_{r}(0)\|_{2}^{4} I\{\|\boldsymbol{w}_{r}(0)\|_{2}^{2} > M\}, \\ &\leq \frac{R^{2}}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1) + \frac{M^{2}}{n_{1}m} \sum_{p=1}^{m} \sum_{r=1}^{m} I\{A_{pr}\} + \frac{1}{m} \sum_{r=1}^{m} \|\boldsymbol{w}_{r}(0)\|_{2}^{4} I\{\|\boldsymbol{w}_{r}(0)\|_{2}^{2} > M\}, \\ &\leq \frac{R^{2}}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1) + \frac{M^{2}}{n_{1}m} \sum_{p=1}^{m} \sum_{r=1}^{m} I\{A_{pr}\} + \frac{1}{m} \sum_{r=1}^{m} \|\boldsymbol{w}_{r}(0)\|_{2}^{4} I\{\|\boldsymbol{w}_{r}(0)\|_{2}^{2} > M\}, \\ &\leq \frac{R^{2}}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_{r}(0)\|_{2}^{2} + 1) + \frac{R^{2}}{n_{1}m} \sum_{r=1}^{m} I\{A_{pr}\} + \frac{R^{2}}{m} \sum_{r=1}^{m} \|\boldsymbol{w}_{r}(0)\|_{2}^{4} I\{\|\boldsymbol{w}_{r}(0)\|_{2}^{2} > M\}, \\ &\leq \frac{R^{2}}{m} \sum_{r=1}^{m}$$

where $M = 2(d+2)\log(2m(d+2)/\delta)$. Note that from (47), we have 1320 1321 $P\left(\exists r \in [m], \|\boldsymbol{w}_r(0)\|_2^2 \ge 2(d+2)\log\left(\frac{2m(d+2)}{\delta}\right)\right) \le \delta.$ 1322 1323 1324 On the other hand, applying Bernstein's inequality yields that with probability at least $1 - n_1 e^{-mR}$, 1325 1326 $\frac{1}{m}\sum_{r=1}^{m}I\{A_{pr}\} < 4R$ 1327 1328 1329 holds for all $p \in [n_1]$. 1330 Therefore, we have that $\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_2^2 \lesssim MR^2 + R^2 + M^2R \lesssim M^2R$ 1332 holds with probability at least $1 - \delta - n_1 e^{-mR}$. 1334 (2) Note that when σ satisfies Assumption 4.2, σ', σ'' and σ''' are all Lipschitz continuous and bounded. Thus, we can 1335 obtain that 1336 $\left\|\frac{\partial s_p(\boldsymbol{w})}{\partial \boldsymbol{w}_r} - \frac{\partial s_p(0)}{\partial \boldsymbol{w}_r}\right\|_2 \lesssim \frac{1}{\sqrt{n_1 m}} R(\|\boldsymbol{w}_r(0)\|_2^2 + \|\boldsymbol{w}_r(0)\|_2 + 1) \lesssim \frac{1}{\sqrt{n_1 m}} R(\|\boldsymbol{w}_r(0)\|_2^2 + 1),$ 1338 1339 1340 where the second inequality is from Young's inequality. 1341 Similarly, we have 1342 $\left\|\frac{\partial h_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r} - \frac{\partial h_j(0)}{\partial \boldsymbol{w}_r}\right\|_2 \lesssim \frac{1}{\sqrt{n_2 m}} R(\|\boldsymbol{w}_r(0)\|_2 + 1).$ 1343 1344 1345 Combining (73) and (74) yields that 1346 1347 $\|\boldsymbol{J}(\boldsymbol{w}) - \boldsymbol{J}(0)\|_{2}^{2}$ 1348 $\leq \sum_{r=1}^{m} \left(\sum_{r=1}^{n_1} \left\| \frac{\partial s_p(\boldsymbol{w})}{\partial \boldsymbol{w}_r} - \frac{\partial s_p(0)}{\partial \boldsymbol{w}_r} \right\|_2^2 + \sum_{i=1}^{n_2} \left\| \frac{\partial h_j(\boldsymbol{w})}{\partial \boldsymbol{w}_r} - \frac{\partial h_j(0)}{\partial \boldsymbol{w}_r} \right\|_2^2 \right)$ 1349 $\lesssim \sum_{r=1}^{m} \left(\sum_{r=1}^{n_1} \frac{1}{n_1 m} (R \| \boldsymbol{w}_r(0) \|_2^2 + R)^2 + \sum_{r=1}^{n_2} \frac{1}{n_2 m} (R \| \boldsymbol{w}_r(0) \|_2 + R)^2 \right)$ 1354 $\lesssim \frac{R^2}{m} \sum_{r=1}^{m} (\|\boldsymbol{w}_r(0)\|_2^4 + 1)$ 1358 $\lesssim R^2 \left| d^2 + \frac{d^2}{\sqrt{m}} \sqrt{\log\left(\frac{1}{\delta}\right) + \frac{d^2}{m} \left(\log\left(\frac{1}{\delta}\right)\right)^2} \right|,$ 1359 1360 1361 where the last inequality follows from the fact that $\|\|\boldsymbol{w}_r(0)\|_2^4\|_{\psi_1} \lesssim d^2$ and Lemma C.4.

(73)

(74)

1363

1369 1370

1364 **B.3.** Proof of Theorem 4.7

For the sake of completeness in the proof, we restate Condition 2 and Corollary 4.11 from the main text, and label them as Condition 4 and Corollary B.1, respectively.

Condition 4. At the t-th iteration, we have $\|\boldsymbol{w}_r(t)\|_2 \leq B$ and

$$\|\boldsymbol{w}_{r}(t) - \boldsymbol{w}_{r}(0)\|_{2} \leq \frac{CB^{2}\sqrt{L(0)}}{\sqrt{m\lambda_{0}}} := R^{\prime}$$

1372 for all $r \in [m]$, where C is a universal constant and $B = \sqrt{2(d+2)\log\left(\frac{2m(d+2)}{\delta}\right)} + 1$. 1374

1375 **Corollary B.1.** If Condition 3 holds for $t = 0, \dots, k$ and $R' \leq R$ and $R'' \lesssim \sqrt{1 - \eta} \sqrt{\lambda_0}$, then

$$L(t) \le (1-\eta)^t L(0),$$

holds for $t = 0, \dots, k$, where R is the constant in Lemma 4.5 and $R'' = CM\sqrt{R}$ in (28) when σ is the ReLU³ activation function, R'' = CdR in (30) when σ satisfies Assumption 4.2.

Thanks to Corollary B.1, it is sufficient to prove that Condition 4 also holds for t = k + 1. For readability, we defer the proof of Corollary B.1 to the end of this section. In the following, we are going to show that the Condition 4 also holds for t = k + 1, thus combining Condition 4 and Corollary B.1 leads to Theorem 4.7.

1387 1388 Proof of Theorem 4.7. Recall that we let $R'' = CM\sqrt{R}$ in (28) when σ is the ReLU³ activation function and let R'' = CdR1389 in (30) when σ satisfies Assumption 4.2.

First, we can set $R' \leq R$ and $R'' \leq \frac{\sqrt{3\lambda_0}}{6}$, since $R'' \lesssim \sqrt{1-\eta}\sqrt{\lambda_0}$. Then from Lemma 4.5 we have $\|\boldsymbol{J}(t) - \boldsymbol{J}(0)\|_2 \leq \frac{\sqrt{3\lambda_0}}{6}$, thus

$$\sigma_{min}(\boldsymbol{J}(t)) \ge \sigma_{min}(\boldsymbol{J}(0)) - \|\boldsymbol{J}(t) - \boldsymbol{J}(0)\|_2 \ge \frac{\sqrt{3\lambda_0}}{2} - \frac{\sqrt{3\lambda_0}}{6} = \frac{\sqrt{3\lambda_0}}{3}$$

and then $\lambda_{min}(\boldsymbol{H}(t)) \geq \frac{\lambda_0}{3}$ for $t = 0, \dots, k$, where $\sigma_{min}(\cdot)$ denotes the least singular value.

1398 From the updating rule of NGD, we have

$$\boldsymbol{w}_{r}(t+1) = \boldsymbol{w}_{r}(t) - \eta \left[\boldsymbol{J}(t)^{T} \right]_{r} (\boldsymbol{H}(t))^{-1} \begin{pmatrix} \boldsymbol{s}(t) \\ \boldsymbol{h}(t) \end{pmatrix}$$

1376 1377

1378

1386

1393

1394 1395

1399

1400

1401 1402

1405 1406

1409

$$\left[\boldsymbol{J}(t)^{T}\right]_{r} = \left[\frac{\partial s_{1}(t)}{\partial \boldsymbol{w}_{r}}, \cdots, \frac{\partial s_{n_{1}}(t)}{\partial \boldsymbol{w}_{r}}, \frac{\partial h_{1}(t)}{\partial \boldsymbol{w}_{r}}, \cdots, \frac{\partial h_{n_{2}}(t)}{\partial \boldsymbol{w}_{r}}\right]$$

1407 1408 Therefore, for $t = 0, \dots, k$ and any $r \in [m]$, we have

$$\begin{aligned} \| \boldsymbol{w}_{r}(t+1) - \boldsymbol{w}_{r}(t) \|_{2} \\ \| \boldsymbol{w}_{r}(t+1) - \boldsymbol{w}_{r}(t) \|_{2} \\ \leq \eta \| \left[\boldsymbol{J}(t)^{T} \right]_{r} \|_{2} \| \boldsymbol{H}(t)^{-1} \|_{2} \sqrt{L(t)} \\ \leq \frac{3\eta}{\lambda_{0}} \| \left[\boldsymbol{J}(t)^{T} \right]_{r} \|_{2} \sqrt{L(t)} \\ \\ \| \boldsymbol{H} \\$$

where the last inequality is due to Corollary B.1. 1429

Summing t from 0 to k yields that

 $\|\boldsymbol{w}_r(k+1) - \boldsymbol{w}_r(0)\|_2$ $\leq \sum_{i=1}^{k} \|\boldsymbol{w}_{r}(t+1) - \boldsymbol{w}_{r}(t)\|_{2}$ $\leq C \frac{\eta B^2}{\sqrt{m\lambda_0}} \sum_{k=0}^k (1-\eta)^{t/2} \sqrt{L(0)}$ $\leq \frac{CB^2\sqrt{L(0)}}{\sqrt{m\lambda_0}},$

where C is a universal constant.

Now, when $R' \leq 1$, we can deduce that $\|\boldsymbol{w}_r(k+1)\|_2 \leq B$, implying that Condition 4 also holds for t = k + 1. Thus, it remains only to derive the requirement for m.

Recall that we need m to satisfy that $R' = \frac{CB^2 \sqrt{L(0)}}{\sqrt{m\lambda_0}} \leq R$ and $R'' \lesssim \sqrt{1-\eta} \sqrt{\lambda_0}$.

(1) When σ is the ReLU³ activation function, in Corollary B.1, $R'' = CM\sqrt{R} \lesssim \sqrt{1-\eta}\sqrt{\lambda_0}$, implying that $R \lesssim \frac{(1-\eta)\lambda_0}{M^2}$. Then $R^{'} = \frac{CB^2\sqrt{L(0)}}{\sqrt{m}\lambda_0} \leq R$ implies that

$$m = \Omega\left(\frac{1}{(1-\eta)^2} \frac{M^4 B^4 L(0)}{\lambda_0^4}\right).$$

From Lemma C.4 for the estimation of L(0), i.e.,

$$L(0) \lesssim d^2 \log\left(\frac{n_1 + n_2}{\delta}\right),$$

we can deduce that

$$m = \Omega\left(\frac{1}{(1-\eta)^2}\frac{d^8}{\lambda_0^4}\log^6\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right)$$

(2) When σ satisfies Assumption 4.2, we have that

$$R \lesssim rac{\sqrt{(1-\eta)\lambda_0}}{d}, R^{'} = rac{CB^2\sqrt{L(0)}}{\sqrt{m}\lambda_0} \leq R.$$

From Lemma C.4, we can deduce that

$$m = \Omega\left(\frac{1}{1-\eta}\frac{d^6}{\lambda_0^3}\log^2\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right).$$

Proof of Corollary B.1. Similar as before, when $R' \leq R$ and $R'' \leq \frac{\sqrt{3\lambda_0}}{6}$, we have $\sigma_{min}(\boldsymbol{J}(t)) \geq \frac{\sqrt{3\lambda_0}}{3}$ and then $\lambda_{min}(\boldsymbol{H}(t)) \geq \frac{\lambda_0}{3}$ for $t = 0, \dots, k$.

$$= -\int_{0}^{1} \left\langle \frac{\partial \boldsymbol{u}(\boldsymbol{w}(s))}{\partial \boldsymbol{w}}, \eta \boldsymbol{J}(t)^{T} \boldsymbol{H}(t)^{-1} \boldsymbol{u}(\boldsymbol{w}(t)) \right\rangle ds$$

$$= -\int_{0}^{1} \left\langle \frac{\partial \boldsymbol{u}(\boldsymbol{w}(t))}{\partial \boldsymbol{w}}, \eta \boldsymbol{J}(t)^{T} \boldsymbol{H}(t)^{-1} \boldsymbol{u}(\boldsymbol{w}(t)) \right\rangle ds$$

$$+ \int_{0}^{1} \left\langle \frac{\partial \boldsymbol{u}(\boldsymbol{w}(t))}{\partial \boldsymbol{w}} - \frac{\partial \boldsymbol{u}(\boldsymbol{w}(s))}{\partial \boldsymbol{w}}, \eta \boldsymbol{J}(t)^{T} \boldsymbol{H}(t)^{-1} \boldsymbol{u}(\boldsymbol{w}(t)) \right\rangle ds$$

$$:= \boldsymbol{I}_{1}(t) + \boldsymbol{I}_{2}(t),$$
(76)

where the second equality is from the fundamental theorem of calculus and $\boldsymbol{w}(s) = s\boldsymbol{w}(t+1) + (1-s)\boldsymbol{w}(t) = s\boldsymbol{w}(t+1)$ $\boldsymbol{w}(t) - s\eta \boldsymbol{J}(t)^T \boldsymbol{H}(t)^{-1} \boldsymbol{u}(t).$

Note that
$$\frac{\partial \boldsymbol{u}(\boldsymbol{w}(t))}{\partial \boldsymbol{w}} = \boldsymbol{J}(t)$$
, thus $\boldsymbol{I}_1(t) = \eta \boldsymbol{u}(t)$. Plugging this into (76) yields that
 $\boldsymbol{u}(t+1) = (1-\eta)\boldsymbol{u}(t) + \boldsymbol{I}_2(t)$. (77)

Therefore, it remains only to bound $\|I_2(t)\|_2$.

$$\begin{aligned} \|\mathbf{I}_{2}(t)\|_{2} &= \left\| \int_{0}^{1} \left\langle \frac{\partial u(\boldsymbol{w}(t))}{\partial \boldsymbol{w}} - \frac{\partial u(\boldsymbol{w}(s))}{\partial \boldsymbol{w}}, \eta \mathbf{J}(t)^{T} \mathbf{H}(t)^{-1} u(\boldsymbol{w}(t)) \right\rangle ds \right\|_{2} \\ &\leq \int_{0}^{1} \|\mathbf{J}(\boldsymbol{w}(t)) - \mathbf{J}(\boldsymbol{w}(s))\|_{2} \|\eta \mathbf{J}(t)^{T} \mathbf{H}(t)^{-1} u(\boldsymbol{w}(t))\|_{2} ds \\ &\leq \eta \|\mathbf{J}(t)^{T} \mathbf{H}(t)^{-1}\|_{2} \|u(\boldsymbol{w}(t))\|_{2} \int_{0}^{1} \|\mathbf{J}(\boldsymbol{w}(t)) - \mathbf{J}(\boldsymbol{w}(s))\|_{2} ds \\ &\leq \eta \|\mathbf{J}(t)^{T} \mathbf{H}(t)^{-1}\|_{2} \|u(\boldsymbol{w}(t))\|_{2} \int_{0}^{1} \|\mathbf{J}(\boldsymbol{w}(t)) - \mathbf{J}(\boldsymbol{w}(s))\|_{2} ds \\ &\leq \eta \|\mathbf{J}(t)^{T} \mathbf{H}(t)^{-1}\|_{2} \|u(\boldsymbol{w}(t))\|_{2} \int_{0}^{1} \|\mathbf{J}(\boldsymbol{w}(t)) - \mathbf{J}(\boldsymbol{w}(s))\|_{2} ds \\ &\leq \eta \|\mathbf{J}(t)^{T} \mathbf{H}(t)^{-1}\|_{2} \|u(\boldsymbol{w}(t))\|_{2} \int_{0}^{1} \|\mathbf{J}(\boldsymbol{w}(t)) - \mathbf{J}(\boldsymbol{w}(s))\|_{2} ds \\ &\leq \frac{\eta \sqrt{L(t)}}{\sqrt{\lambda_{0}}} \int_{0}^{1} (\|\mathbf{J}(\boldsymbol{w}(t)) - \mathbf{J}(0)\|_{2} + \|\mathbf{J}(\boldsymbol{w}(s)) - \mathbf{J}(0)\|_{2}) ds \\ &\leq \frac{\eta \sqrt{L(t)}}{\sqrt{\lambda_{0}}} R'', \end{aligned}$$

where the last inequality follows from the fact that

$$\|\boldsymbol{w}_{r}(s) - \boldsymbol{w}_{r}(0)\|_{2} \le s \|\boldsymbol{w}_{r}(t+1) - \boldsymbol{w}_{r}(0)\|_{2} + (1-s)\|\boldsymbol{w}_{r}(t) - \boldsymbol{w}_{r}(0)\|_{2} \le R' \le R$$

and Lemma 4.5.

Plugging (78) into the recursion formula (77) yields that

$$\begin{aligned} \|\boldsymbol{u}(t+1)\|_{2}^{2} &= \|(1-\eta)\boldsymbol{u}(t) + \boldsymbol{I}_{2}(t)\|_{2}^{2} \\ &= (1-\eta)^{2}\|\boldsymbol{u}(t)\|_{2}^{2} + \|\boldsymbol{I}_{2}(t)\|_{2}^{2} + 2\langle(1-\eta)\boldsymbol{u}(t),\boldsymbol{I}_{2}(t)\rangle \\ &\leq (1-\eta)^{2}\|\boldsymbol{u}(t)\|_{2}^{2} + \|\boldsymbol{I}_{2}(t)\|_{2}^{2} + 2(1-\eta)\|\boldsymbol{u}(t)\|_{2}\|\boldsymbol{I}_{2}(t)\|_{2} \\ &\leq \left[(1-\eta)^{2} + \frac{C^{2}\eta^{2}(R'')^{2}}{\lambda_{0}} + 2(1-\eta)\frac{C\eta R''}{\sqrt{\lambda_{0}}}\right]\|\boldsymbol{u}(t)\|_{2}^{2}, \end{aligned}$$

where C is a universal constant.

1540 Then we can choose R'' such that

$$\|\boldsymbol{I}_{2}(t)\|_{2} \leq \frac{C\eta\sqrt{L(t)}R''}{\sqrt{\lambda_{0}}} \leq C_{1}\eta\sqrt{L(t)} = C_{1}\eta\sqrt{\boldsymbol{u}(t)},$$

1545 where C is a universal constant and C_1 is a constant to be determined.

1546 Thus, we can deduce that 1547

$$\begin{aligned} \|\boldsymbol{u}(t+1)\|_{2}^{2} &\leq \left[(1-\eta)^{2} + (C_{1}\eta)^{2} + 2(1-\eta)C_{1}\eta\right] \|\boldsymbol{u}(t)\|_{2}^{2} \\ &= \left[(1-\eta) + \eta(\eta C_{1}^{2} + 2(1-\eta)C_{1} + \eta - 1)\right] \|\boldsymbol{u}(t)\|_{2}^{2} \\ &\leq (1-\eta)\|\boldsymbol{u}(t)\|_{2}^{2}, \end{aligned}$$

where in the last inequality is due to that we can choose C_1 such that $\eta C_1^2 + 2(1-\eta)C_1 + \eta - 1 \le 0$.

Note that since $\eta \in (0, 1)$, the quadratic equation $\eta x^2 + 2(1 - \eta)x + \eta - 1 = 0$ has one negative root and one positive root, denoted as x_0 and x_1 respectively. Therefore, the condition $C_1 \le x_1$ is sufficient to satisfy the requirement. The explicit form of x_1 can be written as:

$$x_1 = \frac{2(\eta - 1) + \sqrt{4(1 - \eta)^2 - 4\eta(\eta - 1)}}{2\eta} = \frac{\sqrt{1 - \eta}}{1 + \sqrt{1 - \eta}} \ge \frac{\sqrt{1 - \eta}}{2}.$$

1561 Thus, $C_1 = \frac{\sqrt{1-\eta}}{2}$ is sufficient to satisfy that $\eta C_1^2 + 2(1-\eta)C_1 + \eta - 1 \le 0$.

1563 From this, we can deduce that

$$R^{''} \lesssim C_1 \sqrt{\lambda_0} \lesssim \sqrt{1-\eta} \sqrt{\lambda_0}$$

1566 Therefore, we can conclude that $\|\boldsymbol{u}(t)\|_2^2 \leq (1-\eta)^t \|\boldsymbol{u}(0)\|_2^2$ holds for $t = 0, \cdots, k$.

1	5	67
1	5	68

B.4. Proof of Corollary 4.9

Proof. In the proof of Theorem 4.7, we have proved that Condition 4 holds for all $t \in \mathbb{N}$. Thus, it is sufficient to prove that 1572 Condition 4 can lead to the conclusion in Corollary 4.9.

1573 Setting $\eta = 1$ in (77) yields that

$$\boldsymbol{u}(t+1) = \boldsymbol{I}_2(t).$$

1576 From (78), we have that

$$\|\boldsymbol{I}_{2}(t)\|_{2} \lesssim \frac{\sqrt{L(t)}}{\sqrt{\lambda_{0}}} \int_{0}^{1} \|\boldsymbol{J}(\boldsymbol{w}(t)) - \boldsymbol{J}(\boldsymbol{w}(s))\|_{2} ds.$$

$$\tag{79}$$

1579 Since $\boldsymbol{w}(s) = s\boldsymbol{w}(t+1) + (1-s)\boldsymbol{w}(t)$, then for any $r \in [m]$, we have $\|\boldsymbol{w}_r(s)\|_2 \le s\|\boldsymbol{w}_r(t+1)\|_2 + (1-s)\|\boldsymbol{w}_r(t)\|_2 \le B$. 1581 When $\sigma(\cdot)$ is smooth, we can deduce that for any $r \in [m]$,

$$\left\|\frac{\partial s_p(\boldsymbol{w}(s))}{\partial \boldsymbol{w}_r} - \frac{\partial s_p(\boldsymbol{w}(t))}{\partial \boldsymbol{w}_r}\right\|_2 \lesssim \frac{1}{\sqrt{n_1 m}} (B^2 + 1) \|\boldsymbol{w}_r(s) - \boldsymbol{w}_r(t)\|_2 \le \frac{1}{\sqrt{n_1 m}} (B^2 + 1) \|\boldsymbol{w}_r(t+1) - \boldsymbol{w}_r(t)\|_2$$

1586 and

$$\left\|\frac{\partial h_j(\boldsymbol{w}(s))}{\partial \boldsymbol{w}_r} - \frac{\partial h_j(\boldsymbol{w}(t))}{\partial \boldsymbol{w}_r}\right\|_2 \lesssim \frac{1}{\sqrt{n_1 m}} (B+1) \|\boldsymbol{w}_r(s) - \boldsymbol{w}_r(t)\|_2 \leq \frac{1}{\sqrt{n_1 m}} (B+1) \|\boldsymbol{w}_r(t+1) - \boldsymbol{w}_r(t)\|_2.$$

1591 From (75), we know that for any $r \in [m]$,

$$\|\boldsymbol{w}_r(t+1) - \boldsymbol{w}_r(t)\|_2 \lesssim \frac{B^2}{\sqrt{m\lambda_0}}\sqrt{L(t)}$$

Thus for any $s \in [0, 1]$, we have

1596
1597
$$\|J(w(s)) - J(w(t))\|_2^2$$

$$\leq \sum_{r=1}^{m} \left(\sum_{p=1}^{n_1} \left\| \frac{\partial s_p(\boldsymbol{w}(s))}{\partial \boldsymbol{w}_r} - \frac{\partial s_p(\boldsymbol{w}(t))}{\partial \boldsymbol{w}_r} \right\|_2^2 + \left\| \frac{\partial h_j(\boldsymbol{w}(s))}{\partial \boldsymbol{w}_r} - \frac{\partial h_j(\boldsymbol{w}(t))}{\partial \boldsymbol{w}_r} \right\|_2^2 \right)$$

 $\lesssim \frac{1}{m} \sum_{r=1}^{m} \left((B^4 + 1) \| \boldsymbol{w}_r(t+1) - \boldsymbol{w}_r(t) \|_2^2 + (B^2 + 1) \| \boldsymbol{w}_r(t+1) - \boldsymbol{w}_r(t) \|_2^2 \right)$ $\left(\frac{B^2}{\sqrt{m}\lambda_0}\sqrt{L(t)}\right)^2$.

Plugging this into (79), we have

$$\begin{split} \|\boldsymbol{I}_{2}(t)\|_{2} \lesssim \frac{\sqrt{L(t)}}{\sqrt{\lambda_{0}}} \int_{0}^{1} \|\boldsymbol{J}(\boldsymbol{w}(t)) - \boldsymbol{J}(\boldsymbol{w}(s))\|_{2} ds \\ \lesssim \frac{\sqrt{L(t)}}{\sqrt{\lambda_{0}}} \frac{B^{4}}{\sqrt{m\lambda_{0}}} \sqrt{L(t)} \end{split}$$

1614
1615
$$= \frac{B^4}{\sqrt{m\lambda_0^3}}L(t).$$

Combining with the fact $\boldsymbol{u}(t+1) = \boldsymbol{I}_2(t)$ yields that

$$\left\| \begin{pmatrix} \boldsymbol{s}(t+1) \\ \boldsymbol{h}(t+1) \end{pmatrix} \right\|_{2} \leq \frac{CB^{4}}{\sqrt{m\lambda_{0}^{3}}} \left\| \begin{pmatrix} \boldsymbol{s}(t) \\ \boldsymbol{h}(t) \end{pmatrix} \right\|_{2}^{2}$$

holds for $t \in \mathbb{N}$, where C is a universal constant.

In the proof above, we only require that $R^{'} \leq R$ and $R^{''} = CdR \leq \frac{\sqrt{3\lambda_0}}{6}$, leading to the requirement for m that

$$m = \Omega\left(\frac{d^6}{\lambda_0^3}\log^2\left(\frac{md}{\delta}\right)\log\left(\frac{n_1+n_2}{\delta}\right)\right).$$

г	_	_	٦	
L				
L				
-				

C. Auxiliary Lemmas

Lemma C.1 (Theorem 3.1 in Kuchibhotla & Chakrabortty (2022)). If X_1, \dots, X_n are independent mean zero random variables with $||X_i||_{\psi_{\alpha}} < \infty$ for all $1 \le i \le n$ and some $\alpha > 0$, then for any vector $a = (a_1, \dots, a_n) \in \mathbb{R}^n$, the following holds true:

$$C(\alpha) := \max\{\sqrt{2}, 2^{1/\alpha}\} \begin{cases} \sqrt{8}(2\pi)^{1/4} e^{1/24} (e^{2/e}/\alpha)^{1/\alpha}, & \text{if } \alpha < 1, \\ 4e + 2(\log 2)^{1/\alpha}, & \text{if } \alpha \ge 1. \end{cases}$$

and for $\beta(\alpha) = \infty$ when $\alpha \leq 1$ and $\beta(\alpha) = \alpha/(\alpha - 1)$ when $\alpha > 1$,

$$L_n(\alpha) := \frac{4^{1/\alpha}}{\sqrt{2} \|b\|_2} \times \begin{cases} \|b\|_{\beta(\alpha)}, & \text{if } \alpha < 1, \\ 4e\|b\|_{\beta(\alpha)}/C(\alpha), & \text{if } \alpha \ge 1. \end{cases}$$

and $L_n^*(\alpha) = L_n(\alpha)C(\alpha) ||b||_2 / ||b||_{\beta(\alpha)}$.

In the following, we will provide some preliminary information about Orlicz norms. Let $f:[0,\infty) \to [0,\infty)$ be a non-decreasing function with f(0) = 0. The f-Orlicz norm of a real-valued random variable 1652 X is given by 1653 $||X||_f := \inf\{C > 0 : \mathbb{E}\left[f\left(\frac{|X|}{C}\right)\right] \le 1\}.$ 1654 1655 1656 If $||X||_{\psi_{\alpha}} < \infty$, we say that X is sub-Weibull of order $\alpha > 0$, where 1658 $\psi_{\alpha}(x) := e^{x^{\alpha}} - 1.$ 1659 1660 Note that when $\alpha \ge 1$, $\|\cdot\|_{\psi_{\alpha}}$ is a norm and when $0 < \alpha < 1$, $\|\cdot\|_{\psi_{\alpha}}$ is a quasi-norm. Moreover, since $(|a| + |b|)^{\alpha} \le 1$ 1661 $|a|^{\alpha}+|b|^{\alpha}$ holds for any $a,b\in\mathbb{R}$ and $0<\alpha<1,$ we can deduce that 1662 1663 $\mathbb{E}e^{\frac{|X+Y|^{\alpha}}{|C|^{\alpha}}} \leq \mathbb{E}e^{\frac{|X|^{\alpha}+|Y|^{\alpha}}{|C|^{\alpha}}} = \mathbb{E}e^{\frac{|X|^{\alpha}}{|C|^{\alpha}}}e^{\frac{|Y|^{\alpha}}{|C|^{\alpha}}} \leq \left(\mathbb{E}e^{\frac{2|X|^{\alpha}}{|C|^{\alpha}}}\right)^{1/2} \left(\mathbb{E}e^{\frac{2|Y|^{\alpha}}{|C|^{\alpha}}}\right)^{1/2}.$ 1664 1665 1666 This implies that $||X + Y||_{\psi_{\alpha}} \le 2^{1/\alpha} \max\{||X||_{\psi_{\alpha}}, ||Y||_{\psi_{\alpha}}\} \le 2^{1/\alpha} (||X||_{\psi_{\alpha}} + ||Y||_{\psi_{\alpha}}).$ 1668 1669 Furthermore, for p, q > 0, we have $|||X|||_{\psi_p} = |||X|^{p/q}||_{\psi_q}^{q/p}$. And in the related proofs, we may frequently use the fact that 1670 for real-valued random variable $X \sim \mathcal{N}(0, 1)$, we have $\|X\|_{\psi_2} \leq \sqrt{6}$ and $\|X^2\|_{\psi_1} = \|X\|_{\psi_2}^2 \leq 6$. 1671 1672 **Lemma C.2.** If $||X||_{\psi_{\alpha}}, ||Y||_{\psi_{\beta}} < \infty$ with $\alpha, \beta > 0$, then we have $||XY||_{\psi_{\gamma}} \le ||X||_{\psi_{\alpha}} ||Y||_{\psi_{\beta}}$, where γ satisfies that 1673 1674 $\frac{1}{\gamma} = \frac{1}{\alpha} + \frac{1}{\beta}.$ 1675 1676 1677 *Proof.* Without loss of generality, we can assume that $||X||_{\psi_{\alpha}} = ||Y||_{\psi_{\beta}} = 1$. To prove this, let us use Young's inequality, 1678 which states that 1679 $xy \le \frac{x^p}{p} + \frac{y^q}{q}, for \ x, y \ge 0, p, q > 1.$ 1681 1682 Let $p = \alpha/\gamma, q = \beta/\gamma$, then 1683 $\mathbb{E}[\exp(|XY|^{\gamma})] \le \mathbb{E}\left[\exp\left(\frac{|X|^{\gamma p}}{n} + \frac{|Y|^{\gamma q}}{q}\right)\right]$ $= \mathbb{E}\left[\exp\left(\frac{|X|^{\alpha}}{p}\right)\exp\left(\frac{|Y|^{\beta}}{q}\right)\right]$ 1687 $\leq \mathbb{E}\left[\frac{\exp(|X|^{\alpha})}{p} + \frac{\exp(|Y|^{\beta})}{q}\right]$ 1690 $\leq \frac{2}{p} + \frac{2}{q}$ 1692 1693 = 2.where the first and second inequality follow from Young's inequality. From this, we have that $\|XY\|_{\psi_{\gamma}} \leq \|X\|_{\psi_{\alpha}} \|Y\|_{\psi_{\beta}}$. 1696 1698 **Lemma C.3** (Bernstein inequality, Theorem 3.1.7 in Giné & Nickl (2021)). Let X_i , $1 \le i \le n$ be independent centered 1699 random variables a.s. bounded by $c < \infty$ in absolute value. Set $\sigma^2 = 1/n \sum_{i=1}^n \mathbb{E}X_i^2$ and $S_n = 1/n \sum_{i=1}^n X_i$. Then, for 1700 all t > 0, u.

$$P\left(S_n \ge \sqrt{\frac{2\sigma^2 t}{n} + \frac{ct}{3n}}\right) \le e^{-1}$$

1705	Lemma C.4. For $0 < \delta < 1$, with probability at least $1 - \delta$, we have that when $m \ge \log^2\left(\frac{n_1+n_2}{\delta}\right)$,
1706	$\parallel \langle \langle c \rangle \rangle \parallel^2 \langle c \rangle \langle c \rangle \rangle$
1707	$L(0) = \left\ \begin{pmatrix} \boldsymbol{s}(0) \\ \boldsymbol{l}(0) \end{pmatrix} \right\ _{l} = \mathcal{O} \left(d^2 \log \left(\frac{n_1 + n_2}{s} \right) \right).$
1708	$\ \left(\boldsymbol{n}(0) \right) \ _2 \left(\begin{array}{c} \delta \\ \delta \end{array} \right) \right)$
1710	
1710	<i>Proof.</i> Recall that for $p \in [n_1]$,
1711 1712	$1 \left[1 \sum_{m=1}^{m} \left(1 \sum_{m=1}^{m} \left($
1712	$s_p(0) = \frac{1}{\sqrt{m_1}} \left[\frac{1}{\sqrt{m_2}} \sum a_r \left(\sigma^{-}(\boldsymbol{w}_r(0)^T \boldsymbol{x}_p) w_{r0}(0) - \sigma^{-}(\boldsymbol{w}_r(0)^T \boldsymbol{x}_p) \ \boldsymbol{w}_{r1}(0) \ _2^2 \right) - f(\boldsymbol{x}_p) \right]$
1714	$\sqrt{n_1} \left[\sqrt{n_r} \right]_{r=1}$
1715	and for $j \in [n_2]$,
1716	$1 \begin{bmatrix} 1 & m \\ m \end{bmatrix} \begin{bmatrix} m \\ m \end{bmatrix} \begin{bmatrix} m \\ m \end{bmatrix}$
1717	$h_j(0) = rac{1}{\sqrt{n_2}} \left[rac{1}{\sqrt{m}} \sum a_r \sigma(oldsymbol{w}_r(0)^T oldsymbol{y}_j) - g(oldsymbol{y}_j) ight].$
1718	$\sqrt{r^2} \left[\sqrt{r^2} r^{-1} \right]$
1719	Then
1720	$I(0) = \sum_{n=1}^{n_1} \frac{1}{(1-n_1)^2} + \sum_{n=1}^{n_2} \frac{1}{(1-n_2)^2}$
1721	$L(0) = \sum_{j=1}^{\infty} \frac{1}{2} (s_p(0))^2 + \sum_{j=1}^{\infty} \frac{1}{2} (h_j(0))^2$
1722	p=1 $j=1$
1723	$= \frac{1}{2} \sum_{n=1}^{n_1} \left(\frac{1}{2} \sum_{n=1}^{m_1} \left(\frac{1}{2} \sum_{n=1}^{m_1} \left(\frac{1}{2} \sum_{n=1}^{m_1} \left(\frac{1}{2} \sum_{n=1}^{n_1} \left(\frac{1}{2} \sum_{n=1}^{$
1724	$ \leq \frac{1}{n_1} \sum_{r=1}^{n_1} \left(\frac{1}{\sqrt{m}} \sum_{r=1}^{n_1} a_r \left(\delta \left(w_r(0) \ x_p \right) w_{r0}(0) - \delta \left(w_r(0) \ x_p \right) \ w_{r1}(0) \ _2 \right) \right) + \frac{1}{n_1} \sum_{r=1}^{n_1} f \left(x_p \right) $
1726	p=1 ($r=1$) $p=1$
1720	$+\frac{1}{2}\sum_{n=1}^{n_2}\left(-\frac{1}{2}\sum_{n=1}^{m}a_n\sigma(m,(0)^Tm_n)\right)^2+\frac{1}{2}\sum_{n=1}^{n_2}a^2(m_n)$
1728	$+ \frac{1}{n_2} \sum_{i=1}^{r} \left(\frac{1}{\sqrt{m}} \sum_{r=1}^{r} a_r \sigma(\boldsymbol{w}_r(0) \mid \boldsymbol{y}_j) \right) + \frac{1}{n_2} \sum_{i=1}^{r} g(\boldsymbol{y}_j).$
1729	J-1 ($J-1$) $J-1$
1730	Note that $\left \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} \begin{pmatrix}$
1731	$\left a_r\left(\sigma\left(\boldsymbol{w}_r(0)^{T}\boldsymbol{x}_p\right)w_{r0}-\sigma\left(\boldsymbol{w}_r(0)^{T}\boldsymbol{x}_p\right)\ \boldsymbol{w}_{r1}(0)\ _2^2\right)\right \lesssim \ \boldsymbol{w}_r(0)\ _2^2 \boldsymbol{w}_r(0)^{T}\boldsymbol{x}_p $
1732	and $ a_r \sigma(\boldsymbol{w}_r(0)^T \boldsymbol{y}_i) \leq \boldsymbol{w}_r(0) _2^2 \boldsymbol{w}_r(0)^T \boldsymbol{y}_i .$
1733	$\begin{array}{cccccccccccccccccccccccccccccccccccc$
1734	Since $\ \ \boldsymbol{w}_r(0)\ _2^2\ _{\psi_1} = \mathcal{O}(d)$ and $\ \boldsymbol{w}_r(0)^T \boldsymbol{y}_j\ _{\psi_2}, \ \boldsymbol{w}_r(0)^T \boldsymbol{x}_p\ _{\psi_2} = \mathcal{O}(1)$, from Lemma C.2, we have that
1726	$\ \ \boldsymbol{w}_{-}(0)\ _{2}^{2}\ \boldsymbol{w}_{-}(0)^{T}\boldsymbol{x}_{-}\ \ _{t} = \mathcal{O}(d)\ \ \boldsymbol{w}_{-}(0)^{T}\boldsymbol{u}_{t}\ \ _{t} = \mathcal{O}(d)$
1730	$\ \ \omega_{1}(0)\ _{2}\ \omega_{1}(0)-\omega_{p}\ \ \psi_{\frac{2}{3}}-\varepsilon(\omega),\ \omega_{1}(0)-y_{\frac{2}{3}}\ \ \psi_{\frac{2}{3}}-\varepsilon(\omega).$
1738	Applying Lemma C 1 with that for final $u \in [u_1]$ and $i \in [u_2]$ with match bility at least 1 $0 e^{-t}$
1739	Applying Lemma C.1 yields that for fixed $p \in [n_1]$ and $j \in [n_2]$ with probability at least $1 - 2e^{-s}$,
1740	$\left 1 \sum_{m=1}^{m} \left(\frac{1}{2} \left(-\frac{1}{2} \right)^{T} \right) - \frac{1}{2} \left(-\frac{1}{2} \right)^{T} \right = \left \frac{1}{2} \left(-\frac{1}{2} \right)^{T} \right = $
1741	$\left\ \frac{1}{\sqrt{m}} \sum a_r \left(\sigma \left(\boldsymbol{w}_r(0)^T \boldsymbol{x}_p \right) w_{r0}(0) - \sigma \left(\boldsymbol{w}_r(0)^T \boldsymbol{x}_p \right) \ \boldsymbol{w}_{r1}(0) \ _2^2 \right) \right\ \lesssim d\sqrt{t} + \frac{1}{\sqrt{m}} t^{\frac{1}{2}}$
1742	$ \mathbf{v} ^{r=1}$
1743	and with probability at least $1 - 2e^{-t}$,
1744	
1745	$\left \frac{1}{\sqrt{r}} \sum a_r \sigma(\boldsymbol{w}_r(0)^T \boldsymbol{y}_j) \right \lesssim d\sqrt{t} + \frac{a}{\sqrt{r}} t^{\frac{3}{2}}.$
1740	$\left \sqrt{m}\sum_{r=1}^{m}\right $
1747	
1749	Then taking a union bound for all $p \in [n_1]$ and $j \in [n_2]$ with $2(n_1 + n_2)e^{-t} = \delta$ yields that
1750	$(d_{a})^2$
1751	$L(0) \lesssim \left(d\sqrt{t} + \frac{a}{\sqrt{m}} t^{\frac{3}{2}} \right)$
1752	$\sqrt{\frac{\sqrt{11}}{12.3}}$
1753	$\lesssim d^2t + rac{d^2t^2}{d^2t}$
1754	m ((m + m) 1 (m + m))
1755	$=d^2\left(\log\left(\frac{n_1+n_2}{s}\right)+\frac{1}{m}\log^3\left(\frac{n_1+n_2}{s}\right)\right)$
1756	$\langle \langle o \rangle m = \langle o \rangle \rangle$
1/5/	$\lesssim d^2 \log \left(\frac{n_1 + n_2}{2} \right)$
1750	\sim \circ \langle δ $/$
1137	

1760	since $m > \log^2(n_1+n_2)$
1761	since $m \ge \log \left(\frac{-\delta}{\delta}\right)$.
1761	
1762	
1703	
1764	
1765	
1766	
1767	
1768	
1769	
1770	
1771	
1772	
1773	
1774	
1775	
1776	
1777	
1778	
1779	
1780	
1781	
1782	
1783	
1784	
1785	
1786	
1787	
1788	
1789	
1790	
1791	
1792	
1793	
1794	
1795	
1796	
1797	
1798	
1799	
1800	
1801	
1802	
1803	
1804	
1805	
1806	
1807	
1808	
1809	
1810	
1811	
1812	
1813	
1814	