

A APPENDIX

A.1 PRELIMINARIES

A.1.1 NEURAL NETWORKS

Let us summarize all basic notations used in the NNs as follows:

1. Matrices are denoted by bold uppercase letters. For example, $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a real matrix of size $m \times n$ and \mathbf{A}^\top denotes the transpose of \mathbf{A} . $\|\mathbf{A}\|_F$ is the Frobenius norm of the matrix \mathbf{A} .

2. Vectors are denoted by bold lowercase letters. For example, $\mathbf{v} \in \mathbb{R}^n$ is a column vector of size n . Furthermore, denote $v(i)$ as the i -th elements of \mathbf{v} .

3. For a d -dimensional multi-index $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_d] \in \mathbb{N}^d$, we denote several related notations as follows:

$$\begin{aligned} (a) \quad |\boldsymbol{\alpha}| &= |\alpha_1| + |\alpha_2| + \dots + |\alpha_d|; \\ (b) \quad \mathbf{x}^\alpha &= x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}, \quad \mathbf{x} = [x_1, x_2, \dots, x_d]^\top; \end{aligned} \quad (19)$$

4. Assume $\mathbf{n} \in \mathbb{N}_+^n$, then $f(\mathbf{n}) = \mathcal{O}(g(\mathbf{n}))$ means that there exists positive C independent of \mathbf{n} , f, g such that $f(\mathbf{n}) \leq Cg(\mathbf{n})$ when all entries of \mathbf{n} go to $+\infty$.

5. Define $\sigma(x) = \max\{0, x\}$ and $\sigma_s(x) = (1-s)\text{Id}(x) + s\sigma(x)$ for $s > 0$. Two-layer NN structures are defined by:

$$\phi_{s_p}(\mathbf{x}; \boldsymbol{\theta}) := \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k \sigma_{s_p}(\boldsymbol{\omega}_k^\top \mathbf{x}). \quad (20)$$

6.

$$\mathcal{R}_{S, s_p}(\boldsymbol{\theta}) := \frac{1}{2n} \sum_{i=1}^n |f(\mathbf{x}_i) - \phi_{s_p}(\mathbf{x}_i; \boldsymbol{\theta})|^2, \quad (21)$$

it is assumed that the sequence $\{\mathbf{x}_i\}_{i=1}^n$ consists of independent and identically distributed (i.i.d.) random variables. These random variables are uniformly distributed within the hypercube $(0, 1)^d$, where d is the dimension of the input space.

A.1.2 RADEMACHER COMPLEXITY

In our further analysis, we will rely on the definition of Rademacher complexity and several lemmas related to it. Rademacher complexity is a fundamental concept in statistical learning theory and plays a crucial role in analyzing the performance of machine learning algorithms. It quantifies the complexity of a hypothesis class in terms of its ability to fit random noise in the data.

Definition 1 (Rademacher complexity [Anthony et al. \(1999\)](#)). Given a sample set $S = \{z_1, z_2, \dots, z_M\}$ on a domain \mathcal{Z} , and a class \mathcal{F} of real-valued functions defined on \mathcal{Z} , the empirical Rademacher complexity of \mathcal{F} in S is defined as

$$\text{Rad}_S(\mathcal{F}) := \frac{1}{M} \mathbf{E}_{\Sigma_M} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^M \tau_i f(z_i) \right],$$

where $\Sigma_M := \{\tau_1, \tau_2, \dots, \tau_M\}$ are independent random variables drawn from the Rademacher distribution, i.e., $\mathbf{P}(\tau_i = +1) = \mathbf{P}(\tau_i = -1) = \frac{1}{2}$ for $i = 1, 2, \dots, M$.

Lemma 3 (Rademacher complexity for linear predictors [Shalev-Shwartz & Ben-David \(2014\)](#)). Let $\Theta = \{\mathbf{w}_1, \dots, \mathbf{w}_m\} \in \mathbb{R}^d$. Let $\mathcal{G} = \{g(\mathbf{w}) = \mathbf{w}^\top \mathbf{x} : \|\mathbf{x}\|_1 \leq 1\}$ be the linear function class with parameter \mathbf{x} whose ℓ^1 norm is bounded by 1. Then

$$\text{Rad}_\Theta(\mathcal{G}) \leq \max_{1 \leq k \leq m} \|\mathbf{w}_k\|_\infty \sqrt{\frac{2 \log(2d)}{m}}.$$

Lemma 4 (Rademacher complexity and generalization gap [Shalev-Shwartz & Ben-David \(2014\)](#)). Suppose that f in \mathcal{F} are non-negative and uniformly bounded, i.e., for any $f \in \mathcal{F}$ and any $\mathbf{z} \in$

$\mathcal{Z}, 0 \leq f(\mathbf{z}) \leq B$. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the choice of n i.i.d. random samples $S = \{z_1, \dots, z_n\} \subset \mathcal{Z}$, we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbf{E}_{\mathbf{z}} f(\mathbf{z}) \right| \leq 2 \text{Rad}_S(\mathcal{F}) + 3B \sqrt{\frac{\log(4/\delta)}{2n}}.$$

A.2 PROOF OF THEOREM [1](#)

Before the proof, we need a lemma in the linear algebra.

Lemma 5. *Suppose \mathbf{A} and \mathbf{B} are strictly positive definite, we have that*

$$\lambda_{\min}(\mathbf{A} + \mathbf{B}) \geq \lambda_{\min}(\mathbf{A}) + \lambda_{\min}(\mathbf{B}). \quad (22)$$

Proof. Let λ_a be defined as $\lambda_{\min}(\mathbf{A})$ and λ_b as $\lambda_{\min}(\mathbf{B})$. Consequently, we can assert that $\mathbf{A} + \mathbf{B} - (\lambda_a + \lambda_b)\mathbf{I}$ possesses positive definiteness. If we designate λ as an eigenvalue of $\mathbf{A} + \mathbf{B}$, then it follows that $(\mathbf{A} + \mathbf{B})\mathbf{x}_* = \lambda\mathbf{x}_*$. This relationship can be expressed as:

$$(\mathbf{A} + \mathbf{B} - (\lambda_a + \lambda_b)\mathbf{I})\mathbf{x}_* = (\lambda - \lambda_a - \lambda_b)\mathbf{x}_*. \quad (23)$$

Consequently, we can deduce that $\lambda \geq \lambda_a + \lambda_b$, which further implies that $\lambda_{\min}(\mathbf{A} + \mathbf{B}) \geq \lambda_{\min}(\mathbf{A}) + \lambda_{\min}(\mathbf{B})$. \square

Proof of Theorem [1](#) For the case $s_p < 1$, let's start by considering the expression for the matrix $\mathbf{K}_p^{[\omega]}$ where

$$\mathbf{K}_p^{[\omega]} = (K_{ij,p}^{[\omega]})_{n \times n} = \left(\mathbf{E}_{(a,\omega)} a^2 \sigma'_{s_p}(\omega^\top \mathbf{x}_i) \sigma'_{s_p}(\omega^\top \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j \right)_{n \times n}. \quad (24)$$

Given the derivative of the activation function:

$$\sigma'_{s_p}(x) = \begin{cases} 1, & x > 0 \\ (1 - s_p), & x < 0 \\ 0, & x = 0 \end{cases} \quad \sigma'_{s_{p+1}}(x) = \begin{cases} 1, & x > 0 \\ (1 - s_{p+1}), & x < 0 \\ 0, & x = 0 \end{cases}$$

we have

$$\sigma'_{s_{p+1}}(x) = \sigma'_{s_p}(x) + (s_p - s_{p+1})\sigma(-x) \quad (25)$$

$$\begin{aligned} & \mathbf{E}_{(a,\omega)} a^2 \sigma'_{s_{p+1}}(\omega^\top \mathbf{x}_i) \sigma'_{s_{p+1}}(\omega^\top \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j = \mathbf{E}_{(a,\omega)} a^2 \sigma'_{s_p}(\omega^\top \mathbf{x}_i) \sigma'_{s_p}(\omega^\top \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j \\ & - (s_p - s_{p+1}) \mathbf{E}_{(a,\omega)} a^2 \left[\sigma'(\omega^\top \cdot (-\mathbf{x}_i)) \sigma'_{s_p}(\omega^\top \mathbf{x}_j) (-\mathbf{x}_i) \cdot \mathbf{x}_j + \sigma'_{s_p}(\omega^\top \mathbf{x}_i) \sigma'(\omega^\top \cdot (-\mathbf{x}_j)) \mathbf{x}_i \cdot (-\mathbf{x}_j) \right] \\ & + (s_p - s_{p+1})^2 \mathbf{E}_{(a,\omega)} a^2 \sigma'(\omega^\top \cdot (-\mathbf{x}_i)) \sigma'(\omega^\top \cdot (-\mathbf{x}_j)) \mathbf{x}_i \cdot \mathbf{x}_j. \end{aligned} \quad (26)$$

Furthermore, since

$$\sigma'(x) = \frac{\sigma'_{s_p}(x) - s_p \sigma'_{s_p}(-x)}{1 - s_p^2}, \quad (27)$$

we have

$$\begin{aligned} & \sigma'(\omega^\top \cdot (-\mathbf{x}_i)) \sigma'_{s_p}(\omega^\top \mathbf{x}_j) (-\mathbf{x}_i) \cdot \mathbf{x}_j \\ & = \frac{1}{1 - s_p^2} \left[\sigma'_{s_p}(\omega^\top (-\mathbf{x}_i)) \sigma'_{s_p}(\omega^\top \mathbf{x}_j) (-\mathbf{x}_i) \cdot \mathbf{x}_j + s_p \sigma'_{s_p}(\omega^\top \mathbf{x}_i) \sigma'_{s_p}(\omega^\top \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j \right] \\ & \quad \sigma'_{s_p}(\omega^\top \mathbf{x}_i) \sigma'(\omega^\top (-\mathbf{x}_j)) \mathbf{x}_i \cdot (-\mathbf{x}_j) \\ & = \frac{1}{1 - s_p^2} \left[\sigma'_{s_p}(\omega^\top \mathbf{x}_i) \sigma'_{s_p}(\omega^\top (-\mathbf{x}_j)) (-\mathbf{x}_i) \cdot \mathbf{x}_j + s_p \sigma'_{s_p}(\omega^\top \mathbf{x}_i) \sigma'_{s_p}(\omega^\top \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j \right]. \end{aligned} \quad (28)$$

Therefore,

$$\mathbf{K}_{p+1}^{[\omega]} = \left(1 + \frac{2s_p(s_{p+1} - s_p)}{1 - s_p^2} \right) \mathbf{K}_p^{[\omega]} + \frac{s_{p+1} - s_p}{1 - s_p^2} (\mathbf{M}_p^{[\omega]} + \mathbf{H}_p^{[\omega]}) + (s_{p+1} - s_p)^2 \mathbf{T}_M^{[\omega]}. \quad (29)$$

When $s_p < 1$, with the initial condition $s_0 = 0$, we can establish the following inequalities based on Assumption [1](#) where $\mathbf{K}_0^{[\omega]}$, $\mathbf{M}_0^{[\omega]}$, $\mathbf{H}_0^{[\omega]}$ is strictly positive, and Lemma [5](#) holds:

$$\lambda_{\min}(\mathbf{K}_1^{[\omega]}) \geq 0. \quad (30)$$

The reason why $\mathbf{K}_0^{[\omega]}$, $\mathbf{M}_0^{[\omega]}$, $\mathbf{H}_0^{[\omega]}$ are positive definite matrices is indeed attributed to the fact that $\sigma'_0(x)$ is a constant function. Specifically, for $\mathbf{K}_0^{[\omega]}$, it can be represented as $(a(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n))^{\top} a(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$, which is inherently positive definite. Similar propositions can be derived for $\mathbf{M}_0^{[\omega]}$ and $\mathbf{H}_0^{[\omega]}$ based on the same principle.

Now, when $0 \leq s_p \leq s_{p+1}$ and $s_p < 1$ and Lemma [5](#) holds:

$$\lambda_{\min}(\mathbf{K}_{p+1}^{[\omega]}) \geq \lambda_{\min}(\mathbf{K}_p^{[\omega]}) \geq 0$$

due to Eqs. [\(29\)\(30\)](#).

For the case $s_p \geq 1$, we have that

$$\begin{aligned} & \mathbf{E}_{(a,\omega)} a^2 \sigma'_{s_{p+1}}(\omega^{\top} \mathbf{x}_i) \sigma'_{s_{p+1}}(\omega^{\top} \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j = \mathbf{E}_{(a,\omega)} a^2 \sigma'_{s_p}(\omega^{\top} \mathbf{x}_i) \sigma'_{s_p}(\omega^{\top} \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j \\ & - (s_p - s_{p+1}) \mathbf{E}_{(a,\omega)} a^2 \left[\sigma'(\omega^{\top} \cdot (-\mathbf{x}_i)) \sigma'_{s_p}(\omega^{\top} \mathbf{x}_j) (-\mathbf{x}_i) \cdot \mathbf{x}_j + \sigma'_{s_p}(\omega^{\top} \mathbf{x}_i) \sigma'(\omega^{\top} \cdot (-\mathbf{x}_j)) \mathbf{x}_i \cdot (-\mathbf{x}_j) \right] \\ & + (s_p - s_{p+1})^2 \mathbf{E}_{(a,\omega)} a^2 \sigma'(\omega^{\top} \cdot (-\mathbf{x}_i)) \sigma'(\omega^{\top} \cdot (-\mathbf{x}_j)) \mathbf{x}_i \cdot \mathbf{x}_j. \end{aligned} \quad (31)$$

Furthermore, since

$$\sigma'_{s_p}(x) = \sigma'(x) + (1 - s_p) \sigma'(-x), \quad (32)$$

we have

$$\begin{aligned} & \sigma'(\omega^{\top} \cdot (-\mathbf{x}_i)) \sigma'_{s_p}(\omega^{\top} \mathbf{x}_j) (-\mathbf{x}_i) \cdot \mathbf{x}_j \\ & = \sigma'(\omega^{\top} \cdot (-\mathbf{x}_i)) \sigma'(\omega^{\top} \mathbf{x}_j) (-\mathbf{x}_i) \cdot \mathbf{x}_j - (1 - s_p) \sigma'(\omega^{\top} \cdot (-\mathbf{x}_i)) \sigma'(\omega^{\top} (-\mathbf{x}_j)) (-\mathbf{x}_i) \cdot (-\mathbf{x}_j) \\ & \quad \sigma'_{s_p}(\omega^{\top} \mathbf{x}_i) \sigma'(\omega^{\top} (-\mathbf{x}_j)) \mathbf{x}_i \cdot (-\mathbf{x}_j) \\ & = \sigma'(\omega^{\top} \cdot \mathbf{x}_i) \sigma'(\omega^{\top} (-\mathbf{x}_j)) (-\mathbf{x}_i) \cdot \mathbf{x}_j - (1 - s_p) \sigma'(\omega^{\top} \cdot (-\mathbf{x}_i)) \sigma'(\omega^{\top} (-\mathbf{x}_j)) (-\mathbf{x}_i) \cdot (-\mathbf{x}_j). \end{aligned} \quad (33)$$

Therefore,

$$\mathbf{K}_{p+1}^{[\omega]} = \mathbf{K}_p^{[\omega]} - (1 - s_p)(s_{p+1} - s_p)(\mathbf{M}_M^{[\omega]} + \mathbf{H}_M^{[\omega]}) + (s_{p+1} - s_p)(s_{p+1} - s_p + 2)\mathbf{T}_M^{[\omega]}.$$

When $0 \leq s_p \leq s_{p+1}$ and $s_p \geq 1$, we have that

$$\lambda_{\min}(\mathbf{K}_{p+1}^{[\omega]}) \geq \lambda_{\min}(\mathbf{K}_p^{[\omega]}) \geq 0$$

based on Assumption [1](#) as well as Lemma [5](#). Similar results can be derived for the Gram matrices with respect to the parameter a . \square

A.3 PROOFS IN t_1 ITERATION

Proof of Lemma [1](#) The proof can be found in [\(Luo et al., 2021, Lemma 9\)](#), for readable, we write the proof of this lemma here. Since $\mathbf{P}(|X| \leq B) \leq 2e^{-\frac{1}{2}B^2}$ if $X \sim \mathcal{N}(0, 1)$, we set $B = \sqrt{2 \log \frac{2m(d+1)}{\delta}}$ and obtain

$$\begin{aligned} \mathbf{P} \left(\max_{k \in [m]} \{ |a_k(0)|, \|\mathbf{w}_k(0)\|_{\infty} \} > B \right) &= \mathbf{P} \left(\max_{k \in [m], \alpha \in [d]} \{ |a_k(0)|, |(w_k(0))_{\alpha}| \} > B \right) \\ &= \mathbf{P} \left(\bigcup_{k=1}^m (|a_k(0)| > B) \cup \left(\bigcup_{\alpha=1}^d (|(w_k(0))_{\alpha}| > B) \right) \right) \\ &\leq \sum_{k=1}^m \mathbf{P}(|a_k(0)| > B) + \sum_{k=1}^m \sum_{\alpha=1}^d \mathbf{P}(|(w_k(0))_{\alpha}| > B) \\ &\leq 2me^{-\frac{1}{2}B^2} + 2mde^{-\frac{1}{2}B^2} \\ &= 2m(d+1)e^{-\frac{1}{2}B^2} \\ &= \delta. \end{aligned}$$

\square

Proof of Lemma 2 Let

$$\mathcal{G} := \{a\sigma_{s_1}(\boldsymbol{\omega}^\top \mathbf{x}), \mathbf{x} \in \Omega\} \quad (34)$$

and we have

$$|a(0)\sigma_{s_1}(\boldsymbol{\omega}^\top(0)\mathbf{x})| \leq 2d \log \frac{4m(d+1)}{\delta} =: B_1 \quad (35)$$

with probability at least $1 - \delta/2$ over the choice of $\boldsymbol{\theta}(0)$. Then we have

$$\begin{aligned} & \sup_{\mathbf{x} \in \Omega} \left| \frac{1}{m} \sum_{k=1}^m a_k(0)\sigma_{s_1}(\mathbf{w}_k(0) \cdot \mathbf{x}) \right| \\ &= \sup_{\mathbf{x} \in \Omega} \left| \frac{1}{m} \sum_{k=1}^m (a_k(0)\sigma_{s_1}(\mathbf{w}_k(0) \cdot \mathbf{x}) + B_1) - (\mathbf{E}_{(a,\mathbf{w})} a\sigma_{s_1}(\mathbf{w}^\top \mathbf{x}) + B_1) \right| \\ &\leq 2 \text{Rad}_{\boldsymbol{\theta}(0)}(\mathbf{G}) + 12d \left(\log \frac{4m(d+1)}{\delta} \right) \sqrt{\frac{2 \log(8/\delta)}{m}} \end{aligned}$$

with probability at least $1 - \delta$ over the choice of $\boldsymbol{\theta}(0)$. The Rademacher complexity can be estimated by

$$\begin{aligned} \text{Rad}_{\boldsymbol{\theta}(0)}(\mathbf{G}) &= \frac{1}{m} \mathbf{E}_\tau \left[\sup_{\mathbf{x} \in \Omega} \sum_{k=1}^m \tau_k a_k(0) \sigma(\mathbf{w}_k(0) \cdot \mathbf{x}) \right] \\ &\leq \frac{1}{m} \sqrt{2 \log \frac{4m(d+1)}{\delta}} \mathbf{E}_\tau \left[\sup_{\mathbf{x} \in \Omega} \sum_{k=1}^m \tau_k \mathbf{w}_k(0) \cdot \mathbf{x} \right] \\ &\leq \sqrt{2 \log \frac{4m(d+1)}{\delta}} \sqrt{2d \log \frac{4m(d+1)}{\delta}} \frac{\sqrt{d}}{\sqrt{m}} \\ &= \frac{2d \log \frac{4m(d+1)}{\delta}}{\sqrt{m}}, \end{aligned}$$

where the last inequality is a result of Lemma 1

Therefore, we have

$$\sup_{\mathbf{x} \in \Omega} |\phi_{s_1}(\mathbf{x}; \boldsymbol{\theta}(0))| \leq 2d \log \frac{4m(d+1)}{\delta} \left(2 + 6\sqrt{2 \log(8/\delta)} \right) \quad (36)$$

and

$$\mathcal{R}_{S, s_1}(\boldsymbol{\theta}(0)) \leq \frac{1}{2} \left[1 + 2d \log \frac{4m(d+1)}{\delta} \left(2 + 6\sqrt{2 \log(8/\delta)} \right) \right]^2. \quad (37)$$

□

Next we are going to proof Proposition 1 before that, we need the definition of sub-exponential random variables and sub-exponential Bernstein's inequality.

Definition 2 (Vershynin (2018)). *A random variable X is sub-exponential if and only if its sub-exponential norm is finite i.e.*

$$\|X\|_{\psi_1} := \inf\{s > 0 \mid \mathbf{E}_X[e^{|X|/s}] \leq 2.\} \quad (38)$$

Furthermore, the chi-square random variable X is a sub-exponential random variable and $C_{\psi, d} := \|X\|_{\psi_1}$.

Lemma 6. *Suppose that $\mathbf{w} \sim N(0, \mathbf{I}_d)$, $a \sim N(0, 1)$ and given $\mathbf{x}_i, \mathbf{x}_j \in \Omega$. Then we have*

(i) if $X := \sigma_{s_1}(\mathbf{w}^\top \mathbf{x}_i) \sigma_{s_1}(\mathbf{x} \cdot \mathbf{x}_j)$, then $\|X\|_{\psi_1} \leq dC_{\psi, d}$.

(ii) if $X := a^2 \sigma'_{s_1}(\mathbf{w}^\top \mathbf{x}_i) \sigma'_{s_1}(\mathbf{w}^\top \mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j$, then $\|X\|_{\psi_1} \leq dC_{\psi, d}$.

Proof. The proof is similar with (Luo et al., 2021 Lemma 14).

(i) $|X| \leq d\|\mathbf{w}\|_2^2 = dZ$ and

$$\begin{aligned} \|\mathbf{X}\|_{\psi_1} &= \inf \{s > 0 \mid \mathbf{E}_X \exp(|X|/s) \leq 2\} \\ &= \inf \{s > 0 \mid \mathbf{E}_{\mathbf{w}} \exp(|\sigma_{s_1}(\mathbf{w}^\top \mathbf{x}_i) \sigma_{s_1}(\mathbf{w}^\top \mathbf{x}_j)|/s) \leq 2\} \\ &\leq \inf \{s > 0 \mid \mathbf{E}_{\mathbf{w}} \exp(d\|\mathbf{w}\|_2^2/s) \leq 2\} \\ &= \inf \{s > 0 \mid \mathbf{E}_Z \exp(dZ/s) \leq 2\} \\ &= d \inf \{s > 0 \mid \mathbf{E}_Z \exp(|Z|/s) \leq 2\} \\ &= d \|\chi^2(d)\|_{\psi_1} \\ &\leq dC_{\psi,d} \end{aligned}$$

(ii) $|X| \leq d|a|^2 \leq dZ$ and $\|\mathbf{X}\|_{\psi_1} \leq dC_{\psi,d}$. \square

Theorem 3 (sub-exponential Bernstein's inequality (Vershynin (2018))). *Suppose that X_1, \dots, X_m are i.i.d. sub-exponential random variables with $\mathbf{E}X_1 = \mu$, then for any $s \geq 0$ we have*

$$\mathbf{P} \left(\left| \frac{1}{m} \sum_{k=1}^m X_k - \mu \right| \geq s \right) \leq 2 \exp \left(-C_0 m \min \left(\frac{s^2}{\|\mathbf{X}_1\|_{\psi_1}^2}, \frac{s}{\|\mathbf{X}_1\|_{\psi_1}} \right) \right),$$

where C_0 is an absolute constant.

Proof of Proposition 1. For any $\varepsilon > 0$, we define

$$\begin{aligned} \Omega_{ij,p}^{[a]} &:= \left\{ \boldsymbol{\theta}(0) \mid \left| G_{ij,p}^{[a]}(\boldsymbol{\theta}(0)) - K_{ij,p}^{[a]} \right| \leq \frac{\varepsilon}{n} \right\} \\ \Omega_{ij,p}^{[\omega]} &:= \left\{ \boldsymbol{\theta}(0) \mid \left| G_{ij,p}^{[\omega]}(\boldsymbol{\theta}(0)) - K_{ij,p}^{[\omega]} \right| \leq \frac{\varepsilon}{n} \right\}. \end{aligned} \quad (39)$$

Setting $\varepsilon \leq ndC_{\psi,d}$, by Theorem 3 and Lemma 6 we have

$$\begin{aligned} \mathbf{P}(\Omega_{ij,p}^{[a]}) &\geq 1 - 2 \exp \left(-\frac{mC_0\varepsilon^2}{n^2d^2C_{\psi,d}^2} \right), \\ \mathbf{P}(\Omega_{ij,p}^{[\omega]}) &\geq 1 - 2 \exp \left(-\frac{mC_0\varepsilon^2}{n^2d^2C_{\psi,d}^2} \right). \end{aligned} \quad (40)$$

Therefore, with probability at least $\left[1 - 2 \exp \left(-\frac{mC_0\varepsilon^2}{n^2d^2C_{\psi,d}^2} \right) \right]^{2n^2} \geq 1 - 4n^2 \exp \left(-\frac{mC_0\varepsilon^2}{n^2d^2C_{\psi,d}^2} \right)$ over the choice of $\boldsymbol{\theta}(0)$, we have

$$\begin{aligned} \left\| G_1^{[a]}(\boldsymbol{\theta}(0)) - K_1^{[a]} \right\|_F &\leq \varepsilon \\ \left\| G_1^{[p]}(\boldsymbol{\theta}(0)) - K_1^{[p]} \right\|_F &\leq \varepsilon. \end{aligned} \quad (41)$$

Hence by taking $\varepsilon = \frac{\lambda_1}{4}$ and $\delta = 4n^2 \exp \left(-\frac{mC_0\lambda_1^2}{16n^2d^2C_{\psi,d}^2} \right)$, where $\lambda_1 = \min\{\lambda_{a,1}, \lambda_{\omega,1}\}$

$$\begin{aligned} \lambda_{\min}(\mathbf{G}_1(\boldsymbol{\theta}(0))) &\geq \lambda_{\min}(\mathbf{G}_1^{[a]}(\boldsymbol{\theta}(0))) + \lambda_{\min}(\mathbf{G}_1^{[\omega]}(\boldsymbol{\theta}(0))) \\ &\geq \lambda_{a,1} + \lambda_{\omega,1} - \left\| G_1^{[a]}(\boldsymbol{\theta}(0)) - K_1^{[a]} \right\|_F - \left\| G_1^{[\omega]}(\boldsymbol{\theta}(0)) - K_1^{[\omega]} \right\|_F \\ &\geq \frac{3}{4}(\lambda_{a,1} + \lambda_{\omega,1}). \end{aligned} \quad (42)$$

\square

Proof of Proposition 2. Due to Proposition 1 and the definition of t_1^* , we have that for any $\delta \in (0, 1)$

$$\lambda_{\min}(\mathbf{G}_1(\boldsymbol{\theta})) \geq \frac{1}{2}(\lambda_{a,1} + \lambda_{\omega,1}) \quad (43)$$

with probability at least $1 - \delta$ over the choice of $\boldsymbol{\theta}(0)$.

As we know

$$G_{ij,1} = G_{ij,1}^{[a]} + G_{ij,1}^{[\omega]} = \sum_{k=1}^m \nabla_{a_k} \phi_{s_1}(\mathbf{x}_i; \boldsymbol{\theta}) \cdot \nabla_{a_k} \phi_{s_1}(\mathbf{x}_j; \boldsymbol{\theta}) + \frac{1}{m^2} \sum_{k=1}^m \nabla_{\omega_k} \phi_{s_1}(\mathbf{x}_i; \boldsymbol{\theta}) \cdot \nabla_{\omega_k} \phi_{s_1}(\mathbf{x}_j; \boldsymbol{\theta}) \quad (44)$$

and

$$\begin{cases} \frac{da_k(t)}{dt} = -\nabla_{a_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}) = -\frac{1}{n\sqrt{m}} \sum_{i=1}^n e_{i,1} \sigma_{s_p}(\mathbf{w}_k^\top \mathbf{x}_i) \\ \frac{d\omega_k(t)}{dt} = -\nabla_{\omega_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}) = -\frac{1}{n\sqrt{m}} \sum_{i=1}^n e_{i,1} a_i \sigma'_{s_p}(\mathbf{w}_k^\top \mathbf{x}_i) \mathbf{x}_i \end{cases}$$

where $e_{i,1} = |f(\mathbf{x}_i) - \phi_{s_1}(\mathbf{x}_i; \boldsymbol{\theta})|$.

Then finally we get that

$$\begin{aligned} \frac{d}{dt} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}(t)) &= \sum_{k=1}^m \left(\nabla_{a_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}) \frac{da_k(t)}{dt} + \nabla_{\omega_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}) \frac{d\omega_k(t)}{dt} \right) \\ &= -\sum_{k=1}^m (\nabla_{a_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}) \nabla_{a_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}) + \nabla_{\omega_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}) \nabla_{\omega_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta})) \\ &= -\frac{1}{n^2} \mathbf{e}_1^T \mathbf{G}_{ij,1}(\boldsymbol{\theta}(t)) \mathbf{e}_1 \\ &\leq -\frac{2}{n} \lambda_{\min}(\mathbf{G}_1(\boldsymbol{\theta})) \mathcal{R}_{S,s_1}(\boldsymbol{\theta}(t)) \\ &\leq -\frac{1}{n} (\lambda_{a,1} + \lambda_{\omega,1}) \mathcal{R}_{S,s_1}(\boldsymbol{\theta}(t)). \end{aligned} \quad (45)$$

Therefore,

$$\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(t)) \leq \mathcal{R}_{S,s_1}(\boldsymbol{\theta}(0)) \exp\left(-\frac{t}{n} (\lambda_{a,1} + \lambda_{\omega,1})\right). \quad (46)$$

□

A.4 PROOFS IN t_2 ITERATION

Proof of Theorem 3.3 For any $k \in [m]$, denote

$$\alpha(t) = \max_{k \in [m], s \in [0, t]} |a_k(s)|, \quad \omega(t) = \max_{k \in [m], s \in [0, t]} \|\mathbf{w}_k(s)\|_\infty$$

and we have

$$\left| \frac{da_k(t)}{dt} \right|^2 = |\nabla_{a_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta})|^2 = \left| \frac{1}{n\sqrt{m}} \sum_{i=1}^n e_{i,1} \sigma_{s_p}(\mathbf{w}_k^\top \mathbf{x}_i) \right|^2 \leq \frac{2d^2(\omega(t))^2 \mathcal{R}_{S,s_1}(\boldsymbol{\theta})}{m}. \quad (47)$$

Similarly, we have that

$$\left\| \frac{d\omega_k(t)}{dt} \right\|_\infty^2 \leq \frac{2d^2(\alpha(t))^2 \mathcal{R}_{S,s_1}(\boldsymbol{\theta})}{m}.$$

Due to the Proposition 2 we have

$$\begin{aligned} |a_k(t) - a_k(0)| &\leq \int_0^t |\nabla_{a_k} \mathcal{R}_{S,s_1}(\boldsymbol{\theta}(s))| ds \\ &\leq \frac{\sqrt{2d}}{\sqrt{m}} \int_0^t \omega(s) \sqrt{\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(s))} ds \\ &\leq \frac{\sqrt{2d}}{\sqrt{m}} \omega(t) \int_0^t \sqrt{\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(0))} \exp\left(-\frac{s}{2n} (\lambda_{a,1} + \lambda_{\omega,1})\right) ds \\ &\leq \frac{2\sqrt{2nd} \sqrt{\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(0))}}{\sqrt{m}(\lambda_{a,1} + \lambda_{\omega,1})} \omega(t), \end{aligned}$$

with probability at least $1 - \delta/2$ over the choice of $\boldsymbol{\theta}(0)$. Similarly, we have

$$\|\boldsymbol{\omega}_k(t) - \boldsymbol{\omega}_k(0)\|_\infty \leq \frac{2\sqrt{2}nd\sqrt{\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(0))}}{\sqrt{m}(\lambda_{a,1} + \lambda_{\boldsymbol{\omega},1})}\alpha(t). \quad (48)$$

Therefore, we have that

$$\begin{aligned} \alpha(t) &\leq \alpha(0) + \frac{1}{\sqrt{m}}\kappa\omega(t) \\ \omega(t) &\leq \omega(0) + \frac{1}{\sqrt{m}}\kappa\alpha(t) \end{aligned}$$

where $\kappa = \frac{2\sqrt{2}nd\sqrt{\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(0))}}{\lambda_{a,1} + \lambda_{\boldsymbol{\omega},1}}$. Therefore, when $m \geq \kappa^2$, we have

$$\max\{\alpha(t), \omega(t)\} \leq 2\alpha(0) + 2\omega(0).$$

Based on Lemma [I](#) with probability at least $1 - \delta/2$ over the choice of $\boldsymbol{\theta}(0)$ such that

$$\max_{k \in [m]} \{|a_k(0)|, \|\boldsymbol{\omega}_k(0)\|_\infty\} \leq \sqrt{2 \log \frac{4m(d+1)}{\delta}}. \quad (49)$$

Therefore, we have

$$\max_{k \in [m]} \{|a_k(t) - a_k(0)|, \|\boldsymbol{\omega}_k(t) - \boldsymbol{\omega}_k(0)\|_\infty\} \leq \frac{8\sqrt{2}nd\sqrt{\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(0))}}{\sqrt{m}(\lambda_{a,1} + \lambda_{\boldsymbol{\omega},1})} \sqrt{2 \log \frac{4m(d+1)}{\delta}} \quad (50)$$

with probability at least $1 - \delta$ over the choice of $\boldsymbol{\theta}(0)$. \square

Lemma 7. Suppose that $\boldsymbol{\omega} := \boldsymbol{\omega}(0) \sim N(0, \mathbf{I}_d)$, $a = a(0) \sim N(0, 1)$ and given $\mathbf{x}_i, \mathbf{x}_j \in \Omega$. If

$$m \geq \max \left\{ \frac{16n^2d^2C_{\psi,d}}{C_0\lambda^2} \log \frac{4n^2}{\delta}, \frac{8n^2d^2\mathcal{R}_{S,s_1}(\boldsymbol{\theta}(0))}{(\lambda_{a,1} + \lambda_{\boldsymbol{\omega},1})^2} \right\}$$

then with probability at least $1 - \delta$ over the choice of $\boldsymbol{\theta}(0)$, we have

(i) if $X := \sigma_{s_2}(\bar{\boldsymbol{\omega}}^\top(\boldsymbol{\omega})\mathbf{x}_i) \sigma_{s_2}(\bar{\boldsymbol{\omega}}^\top(\boldsymbol{\omega}) \cdot \mathbf{x}_j)$, then $\|X\|_{\psi_1} \leq 2dC_{\psi,d} + \frac{2d^2\psi(m)^2}{\log 2}$.

(ii) if $X := \bar{a}(a)^2 \sigma'_{s_2}(\bar{\boldsymbol{\omega}}^\top(\boldsymbol{\omega})\mathbf{x}_i) \sigma'_{s_2}(\bar{\boldsymbol{\omega}}^\top(\boldsymbol{\omega})\mathbf{x}_j) \mathbf{x}_i \cdot \mathbf{x}_j$, then $\|X\|_{\psi_1} \leq 2dC_{\psi,d} + \frac{2d^2\psi(m)^2}{\log 2}$.

Proof. (i)

$$|X| \leq d\|\bar{\boldsymbol{\omega}}(\boldsymbol{\omega})\|_2^2 \leq 2d\|\boldsymbol{\omega}\|_2^2 + 2d\|\bar{\boldsymbol{\omega}}(\boldsymbol{\omega}) - \boldsymbol{\omega}\|_2^2 \leq 2d|Z| + 2d^2\psi(m)^2$$

and

$$\begin{aligned} \|X\|_{\psi_1} &= \inf \{s > 0 \mid \mathbf{E}_X \exp(|X|/s) \leq 2\} \\ &= \inf \{s > 0 \mid \mathbf{E}_{\boldsymbol{\omega}} \exp(|\sigma_{s_2}(\bar{\boldsymbol{\omega}}^\top(\boldsymbol{\omega})\mathbf{x}_i) \sigma_{s_2}(\bar{\boldsymbol{\omega}}^\top(\boldsymbol{\omega}) \cdot \mathbf{x}_j)|/s) \leq 2\} \\ &\leq \inf \left\{ s > 0 \mid \mathbf{E}_{\boldsymbol{\omega}} \exp\left(\frac{2d|Z| + 2d^2\psi(m)^2}{s}\right) \leq 2 \right\} \\ &\leq \inf \{s > 0 \mid \mathbf{E}_Z \exp(2d|Z|/s) \leq 2\} + \inf \left\{ s > 0 \mid \mathbf{E}_{\boldsymbol{\omega}} \exp\left(\frac{2d^2\psi(m)^2}{s}\right) \leq 2 \right\} \\ &= 2d\|\chi^2(d)\|_{\psi_1} + \frac{2d^2\psi(m)^2}{\log 2} \\ &\leq 2dC_{\psi,d} + \frac{2d^2\psi(m)^2}{\log 2}. \end{aligned}$$

(ii) $|X| \leq d|a|^2 \leq 2d|Z| + 2d^2\psi(m)^2$ and $\|X\|_{\psi_1} \leq 2dC_{\psi,d} + \frac{2d^2\psi(m)^2}{\log 2}$. \square

To enhance simplicity and maintain consistent notation, we define:

$$C_{\psi,d,2} := 2C_{\psi,d} + \frac{2d\psi(m)^2}{\log 2}. \quad (51)$$

Proof of Proposition 4

$$\begin{aligned}\bar{k}_2^{[a]}(\mathbf{x}, \mathbf{x}') &:= \mathbf{E}_\omega \sigma_{s_2}(\bar{\omega}^\top(\omega)\mathbf{x}) \sigma_{s_2}(\bar{\omega}^\top(\omega)\mathbf{x}') \\ \bar{k}_2^{[\omega]}(\mathbf{x}, \mathbf{x}') &:= \mathbf{E}_{(a,\omega)} \bar{a}(a)^2 \sigma'_{s_2}(\bar{\omega}^\top(\omega)\mathbf{x}) \sigma'_{s_2}(\bar{\omega}^\top(\omega)\mathbf{x}') \mathbf{x} \cdot \mathbf{x}'.\end{aligned}\quad (52)$$

The Gram matrices, denoted as $\bar{\mathbf{K}}_2^{[a]}$ and $\bar{\mathbf{K}}_2^{[\omega]}$, corresponding to an infinite-width two-layer network with the activation function σ_{s_2} , can be expressed as follows:

$$\begin{aligned}\bar{K}_{ij,2}^{[a]} &= \bar{k}_2^{[a]}(\mathbf{x}_i, \mathbf{x}_j), \quad \bar{\mathbf{K}}_2^{[a]} = (\bar{K}_{ij,2}^{[a]})_{n \times n}, \\ \bar{K}_{ij,2}^{[\omega]} &= \bar{k}_2^{[\omega]}(\mathbf{x}_i, \mathbf{x}_j), \quad \bar{\mathbf{K}}_2^{[\omega]} = (\bar{K}_{ij,2}^{[\omega]})_{n \times n}.\end{aligned}\quad (53)$$

The proof can be divided into two main parts. The first part, seeks to establish that the difference between $\mathbf{K}_2^{[a]} + \mathbf{K}_2^{[\omega]}$ and $\bar{\mathbf{K}}_2^{[a]} + \bar{\mathbf{K}}_2^{[\omega]}$ is small. In this case, the proof draws upon Proposition 3, which underscores the potential for the error in $\|\boldsymbol{\theta}(0) - \boldsymbol{\theta}(t^*)\|_\infty$ to be highly negligible when m assumes a large value. The second part aims to demonstrate that the disparity between $\mathbf{G}(\boldsymbol{\theta}(t_1^*))$ and $\bar{\mathbf{K}}_2^{[a]} + \bar{\mathbf{K}}_2^{[\omega]}$ is minimal. This particular proof relies on the application of sub-exponential Bernstein's inequality as outlined in Vershynin (2018) (Theorem 3).

First of all, we prove that the difference between $\mathbf{K}_2^{[a]} + \mathbf{K}_2^{[\omega]}$ and $\bar{\mathbf{K}}_2^{[a]} + \bar{\mathbf{K}}_2^{[\omega]}$ is small. Due to

$$\begin{aligned}\left| \bar{k}_2^{[a]}(\mathbf{x}, \mathbf{x}') - k_2^{[a]}(\mathbf{x}, \mathbf{x}') \right| &\leq \mathbf{E}_\omega \left| \sigma_{s_2}(\bar{\omega}^\top(\omega)\mathbf{x}) \sigma_{s_2}(\bar{\omega}^\top(\omega)\mathbf{x}') - \sigma_{s_2}(\omega\mathbf{x}) \sigma_{s_2}(\omega\mathbf{x}') \right| \\ &\leq 2d \|\bar{\omega}^\top(\omega(0)) - \omega(0)\|_\infty \|\omega(0)\|_\infty \\ &\leq 2d\psi(m) \sqrt{2 \log \frac{4m(d+1)}{\delta}}\end{aligned}\quad (54)$$

with probability at least $1 - \delta$ over the choice of $\boldsymbol{\theta}(0)$. Therefore,

$$\|\mathbf{K}_2^{[a]} - \bar{\mathbf{K}}_2^{[a]}\|_F \leq 2n\psi(m) \sqrt{2 \log \frac{4m(d+1)}{\delta}}.\quad (55)$$

Similarly, we can obtain that

$$\|\mathbf{K}_2^{[\omega]} - \bar{\mathbf{K}}_2^{[\omega]}\|_F \leq 2n\psi(m) \sqrt{2 \log \frac{4m(d+1)}{\delta}}.\quad (56)$$

Set $\psi(m) \leq \frac{\min\{\lambda_{a,2}, \lambda_{\omega,2}\}}{16n \sqrt{2 \log \frac{4m(d+1)}{\delta}}}$, i.e.

$$m \geq n^4 \left(\frac{128\sqrt{2}d \sqrt{\mathcal{R}_{S, s_1}(\boldsymbol{\theta}(0))}}{(\lambda_{a,1} + \lambda_{\omega,1}) \min\{\lambda_{a,2}, \lambda_{\omega,2}\}} 2 \log \frac{4m(d+1)}{\delta} \right),$$

we have

$$\|\mathbf{K}_2^{[a]} - \bar{\mathbf{K}}_2^{[a]}\|_F, \|\mathbf{K}_2^{[\omega]} - \bar{\mathbf{K}}_2^{[\omega]}\|_F \leq \frac{1}{8} \min\{\lambda_{a,2}, \lambda_{\omega,2}\}.$$

Furthermore, by sub-exponential Bernstein's inequality as outlined in Vershynin (2018) (Theorem 3), for any $\varepsilon > 0$, we define

$$\begin{aligned}\Omega_{ij,2}^{[a]} &:= \left\{ \boldsymbol{\theta}(0) \mid \left| G_{ij,2}^{[a]}(\boldsymbol{\theta}(0)) - \bar{K}_{ij,2}^{[a]} \right| \leq \frac{\varepsilon}{n} \right\} \\ \Omega_{ij,2}^{[\omega]} &:= \left\{ \boldsymbol{\theta}(0) \mid \left| G_{ij,2}^{[\omega]}(\boldsymbol{\theta}(0)) - \bar{K}_{ij,2}^{[\omega]} \right| \leq \frac{\varepsilon}{n} \right\}.\end{aligned}\quad (57)$$

Setting $\varepsilon \leq ndC_{\psi,d,2}$, by Theorem 3 and Lemma 6, we have

$$\begin{aligned}\mathbf{P}(\Omega_{ij,2}^{[a]}) &\geq 1 - 2 \exp\left(-\frac{mC_0\varepsilon^2}{n^2d^2C_{\psi,d,2}}\right), \\ \mathbf{P}(\Omega_{ij,2}^{[\omega]}) &\geq 1 - 2 \exp\left(-\frac{mC_0\varepsilon^2}{n^2d^2C_{\psi,d,2}}\right).\end{aligned}\quad (58)$$

Therefore, with probability at least $\left[1 - 2 \exp\left(-\frac{mC_0\varepsilon^2}{n^2 d^2 C_{\psi,d,2}^2}\right)\right]^{2n^2} \geq 1 - 4n^2 \exp\left(-\frac{mC_0\varepsilon^2}{n^2 d^2 C_{\psi,d,2}^2}\right)$ over the choice of $\boldsymbol{\theta}(0)$, we have

$$\begin{aligned} \left\|G_2^{[a]}(\boldsymbol{\theta}(0)) - \bar{K}_2^{[a]}\right\|_F &\leq \varepsilon \\ \left\|G_2^{[p]}(\boldsymbol{\theta}(0)) - \bar{K}_2^{[p]}\right\|_F &\leq \varepsilon. \end{aligned} \quad (59)$$

Hence by taking $\varepsilon = \frac{1}{8} \min\{\lambda_{a,2}, \lambda_{\omega,2}\}$ and $\delta = 4n^2 \exp\left(-\frac{mC_0\lambda_1^2}{16n^2 d^2 C_{\psi,d,2}^2}\right)$, we obtain that

$$\begin{aligned} \lambda_{\min}(\mathbf{G}_2(\boldsymbol{\theta}(t_1^*))) &\geq \lambda_{\min}\left(\mathbf{G}_2^{[a]}(\boldsymbol{\theta}(t_1^*))\right) + \lambda_{\min}\left(\mathbf{G}_2^{[\omega]}(\boldsymbol{\theta}(t_1^*))\right) \\ &\geq \lambda_{a,1} + \lambda_{\omega,1} - \left\|\mathbf{G}_2^{[a]}(\boldsymbol{\theta}(t_1^*)) - \bar{\mathbf{K}}_2^{[a]}\right\|_F - \left\|\mathbf{G}_2^{[\omega]}(\boldsymbol{\theta}(t_1^*)) - \bar{\mathbf{K}}_2^{[\omega]}\right\|_F \\ &\quad - \left\|\bar{\mathbf{K}}_2^{[a]} - \bar{\mathbf{K}}_2^{[a]}\right\|_F - \left\|\bar{\mathbf{K}}_2^{[\omega]} - \bar{\mathbf{K}}_2^{[\omega]}\right\|_F \\ &\geq \frac{3}{4}(\lambda_{a,2} + \lambda_{\omega,2}). \end{aligned} \quad (60)$$

□

Proof of Proposition 5 Due to Proposition 4 and the definition of t_2^* , we have that for any $\delta \in (0, 1)$

$$\lambda_{\min}(\mathbf{G}_2(\boldsymbol{\theta}(t))) \geq \frac{1}{2}(\lambda_{a,1} + \lambda_{\omega,1}) \quad (61)$$

for any $t \in [t_1^*, t_2^*]$ with probability at least $1 - \delta$ over the choice of $\boldsymbol{\theta}(0)$.

As we know

$$G_{ij,2} = G_{ij,2}^{[a]} + G_{ij,2}^{[\omega]} = \sum_{k=1}^m \nabla_{a_k} \phi_{s_2}(\mathbf{x}_i; \boldsymbol{\theta}) \cdot \nabla_{a_k} \phi_{s_2}(\mathbf{x}_j; \boldsymbol{\theta}) + \frac{1}{m^2} \sum_{k=1}^m \nabla_{\omega_k} \phi_{s_2}(\mathbf{x}_i; \boldsymbol{\theta}) \cdot \nabla_{\omega_k} \phi_{s_2}(\mathbf{x}_j; \boldsymbol{\theta}) \quad (62)$$

and

$$\begin{cases} \frac{da_k(t)}{dt} = -\nabla_{a_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}) = -\frac{1}{n\sqrt{m}} \sum_{i=1}^n e_{i,2} \sigma_{s_p}(\mathbf{w}_k^\top \mathbf{x}_i) \\ \frac{d\omega_k(t)}{dt} = -\nabla_{\omega_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}) = -\frac{1}{n\sqrt{m}} \sum_{i=1}^n e_{i,2} a_i \sigma'_{s_p}(\mathbf{w}_k^\top \mathbf{x}_i) \mathbf{x}_i \end{cases}$$

where $e_{i,2} = |f(\mathbf{x}_i) - \phi_{s_2}(\mathbf{x}_i; \boldsymbol{\theta})|$.

Then finally we get that

$$\begin{aligned} \frac{d}{dt} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}(t)) &= \sum_{k=1}^m \left(\nabla_{a_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}) \frac{da_k(t)}{dt} + \nabla_{\omega_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}) \frac{d\omega_k(t)}{dt} \right) \\ &= - \sum_{k=1}^m (\nabla_{a_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}) \nabla_{a_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}) + \nabla_{\omega_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta}) \nabla_{\omega_k} \mathcal{R}_{S,s_2}(\boldsymbol{\theta})) \\ &= -\frac{1}{n^2} \mathbf{e}_2^T \mathbf{G}_{ij,2}(\boldsymbol{\theta}(t)) \mathbf{e}_2 \\ &\leq -\frac{2}{n} \lambda_{\min}(\mathbf{G}_2(\boldsymbol{\theta})) \mathcal{R}_{S,s_2}(\boldsymbol{\theta}(t)) \\ &\leq -\frac{1}{n} (\lambda_{a,2} + \lambda_{\omega,2}) \mathcal{R}_{S,s_2}(\boldsymbol{\theta}(t)). \end{aligned} \quad (63)$$

Therefore,

$$\mathcal{R}_{S,s_2}(\boldsymbol{\theta}(t)) \leq \mathcal{R}_{S,s_2}(\boldsymbol{\theta}(t_1^*)) \exp\left(-\frac{t-t_1^*}{n} (\lambda_{a,2} + \lambda_{\omega,2})\right). \quad (64)$$

□

A.5 EXPERIMENTAL DETAILS FOR THE HRTA

A.5.1 FUNCTION APPROXIMATION USING SUPERVISED LEARNING

Example 1 (Approximating $\sin(2\pi x)$). *In the first example, our goal is to approximate the function $\sin(2\pi x)$ within the interval $[0, 1]$ using two-layer neural networks (NNs) and the HRTA. We will provide a detailed explanation of the training process for the case of $s = 0.5$, which corresponds to the homotopy training case. The training process is divided into two steps:*

1. *In the first step, we employ the following approximation function:*

$$\phi_{\frac{1}{2}}(x; \boldsymbol{\theta}) := \frac{1}{\sqrt{1000}} \sum_{k=1}^{1000} a_k \sigma_{\frac{1}{2}}(\omega_k x) \quad (65)$$

to approximate the function $\sin(2\pi x)$. Here, $\sigma_{\frac{1}{2}}(x) = \frac{1}{2} \text{Id}(x) + \frac{1}{2} \sigma(x)$, and the initial values of the parameters are drawn from a normal distribution $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. We select random sample points (or grid points) $\{x_i\}_{i=1}^{100}$, which are uniformly distributed in the interval $[0, 1]$. The loss function in this step is defined as

$$\mathcal{R}_{S, \frac{1}{2}}(\boldsymbol{\theta}) := \frac{1}{200} \sum_{i=1}^{100} |f(x_i) - \phi_{\frac{1}{2}}(x_i; \boldsymbol{\theta})|^2. \quad (66)$$

Therefore, we employ the Adam optimizer to train this model over 3000 steps to complete the first step of the process.

2. *In the second step, we employ the following approximation function:*

$$\phi(x; \boldsymbol{\theta}) := \frac{1}{\sqrt{1000}} \sum_{k=1}^{1000} a_k \sigma(\omega_k x) \quad (67)$$

to approximate the function $\sin(2\pi x)$. Here the initial values of the parameters are the results in the first step. The loss function in this step is defined as

$$\mathcal{R}_S(\boldsymbol{\theta}) := \frac{1}{200} \sum_{i=1}^{100} |f(x_i) - \phi(x_i; \boldsymbol{\theta})|^2. \quad (68)$$

Therefore, we employ the Adam optimizer to train this model over 13000 steps to complete the second step of the process and finish the training.

For the purpose of comparison, we employ a traditional method with the following approximation function:

$$\phi(x; \boldsymbol{\theta}) := \frac{1}{\sqrt{1000}} \sum_{k=1}^{1000} a_k \sigma(\omega_k x) \quad (69)$$

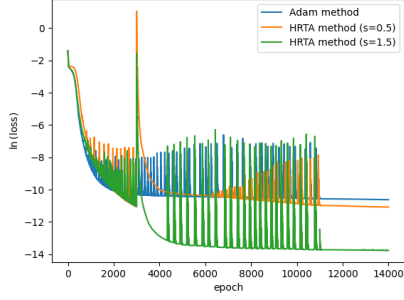
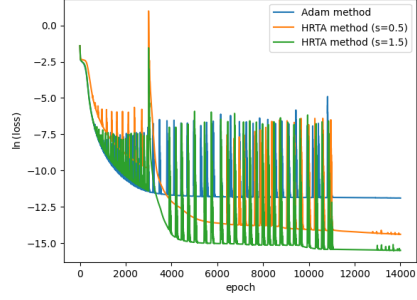
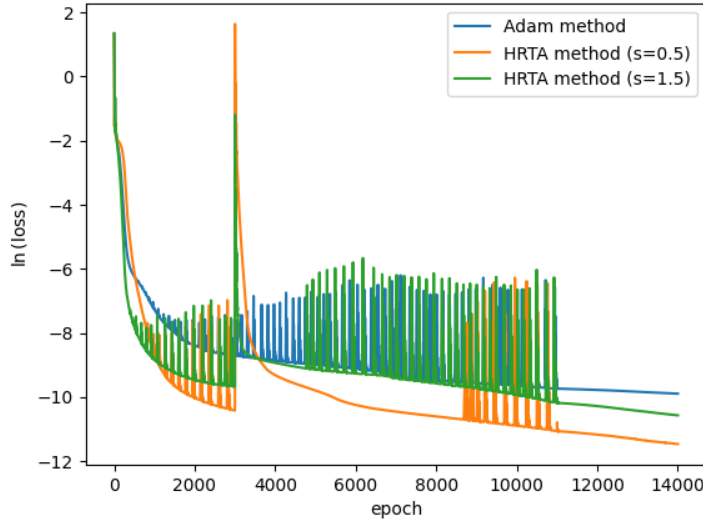
to approximate the function $\sin(2\pi x)$. Here, the initial values of the parameters are sampled from a normal distribution $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. We select the same random sample points (or grid points) $x_{i=1}^{100}$ as used in the HRTA. The loss function in this step is defined as

$$\mathcal{R}_S(\boldsymbol{\theta}) := \frac{1}{200} \sum_{i=1}^{100} |f(x_i) - \phi(x_i; \boldsymbol{\theta})|^2. \quad (70)$$

Therefore, we employ the Adam optimizer to train this model over 16000 steps to complete the training.

In addition, we conducted experiments with neural networks that were not highly overparameterized, containing only 200 and 400 nodes. The results are illustrated in the following figures:

Example 2 (Approximating $\sin(2\pi(x_1 + x_2 + x_3))$). *The training methods in Example 1 and this current scenario share the same structure. The only difference is that in this case, all instances of ω and x used in Example 1 have been extended to three dimensions. In Figure 3 we demonstrate that HRTA is effective in a highly overparameterized scenario, comprising 125 sample points with 1000 nodes. Additionally, we illustrate that HRTA remains effective in a scenario with less overparameterization, involving 400 nodes and 400 sample points. The results are presented below Figure 8*

Figure 6: Approximation for $\sin(2\pi x)$ by NNs with 200 nodesFigure 7: Approximation for $\sin(2\pi x)$ by NNs with 400 nodesFigure 8: Approximation for $\sin(2\pi(x_1 + x_2 + x_3))$ with less overparameterization

A.5.2 SOLVING PARTIAL DIFFERENTIAL EQUATIONS BY DEEP RITZ METHOD [YU & E \(2018\)](#)

Example 3. In this example, we aim to solve the Poisson equation given by:

$$\begin{cases} -\Delta u(x_1, x_2) = \pi^2 [\cos(\pi x_1) + \cos(\pi x_2)] & \text{in } \Omega, \\ \frac{\partial u}{\partial \nu} = 0 & \text{on } \partial\Omega, \end{cases}$$

by homotopy relaxation training methods, where Ω is a domain within the interval $[0, 1]^2$. The exact solution to this equation is denoted as $u^*(x_1, x_2) = \cos(\pi x_1) + \cos(\pi x_2)$.

1. In the first step, we employ the following approximation function:

$$\bar{\phi}(\mathbf{x}; \boldsymbol{\theta}) := \frac{1}{\sqrt{1000}} \sum_{k=1}^{1000} a_k \bar{\sigma}(\boldsymbol{\omega}_k \mathbf{x}) \quad (71)$$

to solve Poisson equations. Here, $\bar{\sigma}(x) = \frac{1}{2} \text{ReLU}^2(x)$, and the initial values of the parameters are drawn from a normal distribution $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. We select random sample points (or grid points) $\{x_i\}_{i=1}^{400}$, which are uniformly distributed in the interval $[0, 1]^2$. As per [\(Lu et al., 2021\)](#) Proposition 1), the loss function in the Deep Ritz method for solving this Poisson equation is indeed given by:

$$\mathcal{R}_{S, \frac{1}{2}}(\boldsymbol{\theta}) := \frac{1}{800} \sum_{i=1}^{400} [|u^*(\mathbf{x}_i) - \bar{\phi}(\mathbf{x}_i; \boldsymbol{\theta})|^2 + |\nabla u^*(\mathbf{x}_i) - \nabla \bar{\phi}(\mathbf{x}_i; \boldsymbol{\theta})|^2]. \quad (72)$$

This loss function captures the discrepancy between the exact solution $u^*(\mathbf{x}_i)$ and the network's output $\bar{\phi}(\mathbf{x}_i; \boldsymbol{\theta})$, as well as the gradient of the exact solution and the gradient of the network's output, for each sampled point \mathbf{x}_i . Therefore, we employ the Adam optimizer to train this model over 16000 steps to complete the step.

2. In the second step, we employ the following approximation function:

$$\bar{\phi}_{\frac{3}{2}}(x; \boldsymbol{\theta}) := \frac{1}{\sqrt{1000}} \sum_{k=1}^{1000} a_k \bar{\sigma}(\omega_k x) \quad (73)$$

to solve Poisson equations. Here the initial values of the parameters are the results in the first step. The loss function in this step is defined as

$$\mathcal{R}_S(\boldsymbol{\theta}) := \frac{1}{800} \sum_{i=1}^{400} \left[|u^*(\mathbf{x}_i) - \bar{\phi}_{\frac{3}{2}}(\mathbf{x}_i; \boldsymbol{\theta})|^2 + |\nabla u^*(\mathbf{x}_i) - \nabla \bar{\phi}_{\frac{3}{2}}(\mathbf{x}_i; \boldsymbol{\theta})|^2 \right]. \quad (74)$$

Therefore, we employ the Adam optimizer to train this model over 13000 steps to complete the step.