SHARP GENERALIZATION FOR NONPARAMETRIC RE GRESSION BY OVER-PARAMETERIZED NEURAL NET WORKS: A DISTRIBUTION-FREE ANALYSIS IN SPHER ICAL COVARIATE

Anonymous authors

008

009

010 011 012

013

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

038

039

040

041

042

043

Paper under double-blind review

ABSTRACT

Sharp generalization bound for neural networks trained by gradient descent (GD) is of central interest in statistical learning theory and deep learning. In this paper, we consider nonparametric regression by an over-parameterized two-layer NN trained by GD. We show that, if the neural network is trained by GD with early stopping, then the trained network renders a sharp rate of the nonparametric regression risk of $\mathcal{O}(\varepsilon_n^2)$, which is the same rate as that for the classical kernel regression trained by GD with early stopping, where ε_n is the critical population rate of the Neural Tangent Kernel (NTK) associated with the network and n is the size of the training data. It is remarked that our result does not require distributional assumptions on the covariate as long as the covariate lies on the unit sphere, in a strong contrast with many existing results which rely on specific distributions such as the spherical uniform data distribution or distributions satisfying certain restrictive conditions. As a special case of our general result, when the eigenvalues of the associated NTK decay at a rate of $\lambda_i \simeq j^{-\frac{a}{d-1}}$ for i > 1 which happens under certain distributional assumption such as the training features follow the spherical uniform distribution, we immediately obtain the minimax optimal rate of $\mathcal{O}(n^{-\frac{d}{2d-1}})$, which is the major results of several existing works in this direction. The neural network width in our general result is lower bounded by a function of only d and ε_n , and such width does not depend on the minimum eigenvalue of the empirical NTK matrix whose lower bound usually requires additional assumptions on the training data. Our results are built upon two significant technical results which are of independent interest. First, uniform convergence to the NTK is established during the training process by GD, so that we can have a nice decomposition of the neural network function at any step of the GD into a function in the Reproducing Kernel Hilbert Space associated with the NTK and an error function with a small L^{∞} -norm. Second, local Rademacher complexity is employed to tightly bound the Rademacher complexity of the function class comprising all the possible neural network functions obtained by GD. Our result formally fills the gap between training a classical kernel regression model and training an over-parameterized but finite-width neural network by GD for nonparametric regression without distributional assumptions about the spherical covariate.

044 045 046

047

1 INTRODUCTION

With the stunning success of deep learning in various areas of machine learning (LeCun et al., 2015), generalization analysis for neural networks is of central interest for statistical learning learning and deep learning. Considerable efforts have been made to analyze the optimization of deep neural networks showing that gradient descent (GD) and stochastic gradient descent (SGD) provably achieve vanishing training loss (Du et al., 2019b; Allen-Zhu et al., 2019b; Du et al., 2019a; Arora et al., 2019; Zou & Gu, 2019; Su & Yang, 2019). There are also extensive efforts devoted to generalization analysis of deep neural networks (DNNs) with algorithmic guarantees, that is, the generalization bounds

054 for neural networks trained by gradient descent or its variants. It has been shown that with sufficient 055 over-parameterization, that is, with enough number of neurons in hidden layers, the training dynam-056 ics of deep neural networks (DNNs) can be approximated by that of a kernel method with the kernel induced by the neural network architecture, termed the Neural Tangent Kernel (NTK), while other 058 studies such as (Yang & Hu, 2021) show that infinite-width neural networks can still learn features. The key idea of NTK based generalization analysis is that, for highly over-parameterized networks, the network weights almost remain around their random initialization. As a result, one can use the 060 first-order Taylor expansion around initialization to approximate the neural network functions and 061 analyze their generalization capability (Cao & Gu, 2019; Arora et al., 2019; Ghorbani et al., 2021). 062

- 063 Many existing works in generalization analysis of neural networks focus on clean data, but it is a 064 central problem in statistical learning that how neural networks can obtain sharp convergence rates for the risk of nonparametric regression where the observed data are corrupted by noise. Consider-065 able research has been conducted in this direction which shows that various types of DNNs achieve 066 optimal convergence rates for smooth (Yarotsky, 2017; Bauer & Kohler, 2019; Schmidt-Hieber, 067 2020; Jiao et al., 2023; Zhang & Wang, 2023) or non-smooth (Imaizumi & Fukumizu, 2019) tar-068 get functions for nonparametric regression. However, most of these works do not have algorithmic 069 guarantees, that is, the DNNs in these works are constructed specially to achieve optimal rates with no guarantees that an optimization algorithm, such as GD or its variants, can obtain such constructed 071 DNNs. To this end, efforts have been made in the literature to study the minimax optimal risk rates 072 for nonparametric regression with over-parameterized neural networks trained by GD with either 073 early stopping (Li et al., 2024) or ℓ^2 -regularization (Hu et al., 2021; Suh et al., 2022). However, 074 most existing works either require spherical uniform data distribution on the unit sphere (Hu et al., 075 2021; Suh et al., 2022) or certain restrictive conditions on the data distribution.
- 076 It remains an interesting and important question for the statistical learning and theoretical deep 077 learning literature that if an over-parameterized neural network trained by GD can achieve sharp 078 risk rates for nonparametric regression with milder assumptions or restrictions on the distribution 079 of the covariate, so that theoretical guarantees can be obtained for data in more practical scenar-080 ios. In this paper, we give a confirmative answer to this question. We present sharp risk rate for 081 nonparametric regression with an over-parameterized two-layer NN trained by GD with early stop-082 ping, which is distribution-free in spherical covariate. Throughout this paper, distribution-free in 083 spherical covariate means that there are no distributional assumptions about the covariate as long as the covariate lies on the unit sphere. Furthermore, our results give confirmative answers to certain 084 open questions or address particular concerns in the literature of training over-parameterized neural 085 networks by GD with early stopping for nonparametric regression with minimax optimal rates, such as the characterization of the stopping time in the early-stopping mechanism, the lower bound for 087 the network width, and the constant learning rate used in GD. Benefiting from our analysis which 088 is distribution-free in spherical covariate, our answers to these open questions or concerns do not require distributional assumptions about spherical covariate. Section 3 summarizes our main results with their significance and comparison to existing works.
- We organize this paper as follows. We first introduce the necessary notations in the remainder of this section. We then introduce in Section 2 the problem setup for nonparametric regression. Our main results are summarized in Section 3 and detailed in Section 5. The training algorithm for the over-parameterized two-layer neural network is introduced in Section 4. The roadmap of proofs is presented in Section 6.
- Notations. We use bold letters for matrices and vectors, and regular lower letter for scalars throughout this paper. The bold letter with a single superscript indicates the corresponding column of a 098 matrix, e.g., A_i is the *i*-th column of matrix A, and the bold letter with subscripts indicates the 099 corresponding element of a matrix or vector. We put an arrow on top of a letter with subscript if it 100 denotes a vector, e.g., $\vec{\mathbf{x}}_i$ denotes the *i*-th training feature. $\|\cdot\|_F$ and $\|\cdot\|_p$ denote the Frobenius norm 101 and the vector ℓ^p -norm or the matrix p-norm. [m:n] denotes all the natural numbers between m and 102 *n* inclusively, and [1: n] is also written as [n]. Var $[\cdot]$ denotes the variance of a random variable. I_n is 103 a $n \times n$ identity matrix. $\mathbb{I}_{\{E\}}$ is an indicator function which takes the value of 1 if event E happens, 104 or 0 otherwise. The complement of a set A is denoted by A^c , and |A| is the cardinality of the set 105 A. vec (·) denotes the vectorization of a matrix or a set of vectors, and tr (·) is the trace of a matrix. 106 We denote the unit sphere in *d*-dimensional Euclidean space by $\mathbb{S}^{d-1} \coloneqq \{\mathbf{x} : \mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|_2 = 1\}$. 107 Let $L^2(\mathbb{S}^{d-1},\mu)$ denote the space of square-integrable functions on \mathbb{S}^{d-1} with probability mea-

108 sure μ , and the inner product $\langle \cdot, \cdot \rangle_{\mu}$ and $\|\cdot\|_{\mu}^2$ are defined as $\langle f, g \rangle_{L^2} \coloneqq \int_{\mathbb{S}^{d-1}} f(x)g(x) d\mu(x)$ and 109 $\|f\|_{L^2}^2 \coloneqq \int_{\mathbb{S}^{d-1}} f^2(x) \mathrm{d}\mu(x) < \infty$. **B** $(\mathbf{x}; r)$ is the Euclidean closed ball centered at \mathbf{x} with radius 110 *r*. Given a function $g: \mathbb{S}^{d-1} \to \mathbb{R}$, its L^{∞} -norm is denoted by $||g||_{\infty} \coloneqq \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} |g(\mathbf{x})|$. L^{∞} is the function class whose elements almost surely have bounded L^{∞} -norm. $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and $||\cdot||_{\mathcal{H}}$ denote 111 112 the inner product and the norm in the Hilbert space \mathcal{H} . $a = \mathcal{O}(b)$ or $a \leq b$ indicates that there exists 113 a constant c > 0 such that $a \le cb$. \mathcal{O} indicates there are specific requirements in the constants of 114 the \mathcal{O} notation. a = o(b) and a = w(b) indicate that $\lim |a/b| = 0$ and $\lim |a/b| = \infty$, respectively. 115 $a \simeq b$ or $a = \Theta(b)$ denotes that there exists constants $c_1, c_2 > 0$ such that $c_1 b \le a \le c_2 b$. Through-116 out this paper we let the input space $\mathcal{X} = \mathbb{S}^{d-1}$, and Unif (\mathcal{X}) denotes the uniform distribution on 117 \mathcal{X} . The constants defined throughout this paper may change from line to line. For a Reproducing 118 Kernel Hilbert Space $\mathcal{H}, \mathcal{H}(\mu_0)$ denotes the ball centered at the origin with radius μ_0 in \mathcal{H} . We use 119 $\mathbb{E}_{P}[\cdot]$ to denote the expectation with respect to the distribution P.

2 PROBLEM SETUP

121 122 123

124

120

We introduce the problem setups for nonparametric regression in this section.

125 2.1 TWO-LAYER NEURAL NETWORK

We are given the training data $\{(\vec{\mathbf{x}}_i, y_i)\}_{i=1}^n$ where each data point is a tuple of feature vector $\vec{\mathbf{x}}_i \in \mathcal{X}$ and its response $y_i \in \mathbb{R}$. Throughout this paper we assume that no two training features 127 128 129 coincide, that is, $\overline{\mathbf{x}}_i \neq \overline{\mathbf{x}}_j$ for all $i, j \in [n]$ and $i \neq j$. We denote the training feature vectors by 130 $\mathbf{S} = \left\{ \vec{\mathbf{x}}_i \right\}_{i=1}^n$, and denote by P_n the empirical distribution over \mathbf{S} . All the responses are stacked 131 132 as a vector $\mathbf{y} = [y_1, \ldots, y_n]^\top \in \mathbb{R}^n$. The response y_i is given by $y_i = f^*(\mathbf{x}_i) + w_i$ for $i \in [n]$, where $\{w_i\}_{i=1}^n$ are i.i.d. sub-Gaussian random noise with mean 0 and variance proxy σ_0^2 , that is, $\mathbb{E} [\exp(\lambda w_i)] \le \exp(\lambda^2 \sigma_0^2/2)$ for any $\lambda \in \mathbb{R}$. f^* is the target function to be detailed later. We define 133 134 135 $\mathbf{y} \coloneqq [y_1, \dots, y_n], \mathbf{w} \coloneqq [w_1, \dots, w_n]^\top$, and use $f^*(\mathbf{S}) \coloneqq \left[f^*(\overrightarrow{\mathbf{x}}_1), \dots, f^*(\overrightarrow{\mathbf{x}}_n)\right]^\top$ to denote the 136 137 clean target labels. The feature vectors in S are drawn i.i.d. according to an underlying unknown 138 continuous data distribution P with μ being the probability measure for P. 139

We consider a two-layer NN (NN) in this paper whose mapping function is

141 142

143

144

150

151 152

154

where $\mathbf{x} \in \mathcal{X}$ is the input, $\sigma(\cdot) = \max\{\cdot, 0\}$ is the ReLU activation function, $\mathbf{W} = \left\{ \overrightarrow{\mathbf{w}}_r \right\}_{r=1}^m$

 $f(\mathbf{W}, \mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \sigma \left(\vec{\mathbf{w}}_r^{\top} \mathbf{x} \right),$

with $\vec{\mathbf{w}}_r \in \mathbb{R}^d$ for $r \in [m]$ denotes the weighting vectors in the first layer and m is the number of neurons. $\boldsymbol{a} = [a_1, \dots, a_m] \in \mathbb{R}^m$ denotes the weights of the second layer. Throughout this paper we also write \mathbf{W} as \mathbf{W}_S so as to indicate that the weighting vectors in \mathbf{W} are trained on the training features \mathbf{S} .

2.2 KERNEL AND KERNEL REGRESSION FOR NONPARAMETRIC REGRESSION

153 We define the kernel function

$$K(\mathbf{u}, \mathbf{v}) \coloneqq \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{2\pi} \left(\pi - \arccos \left\langle \mathbf{u}, \mathbf{v} \right\rangle \right), \quad \forall \mathbf{u}, \mathbf{v} \in \mathcal{X},$$
(2)

(1)

which is in fact the NTK associated with the two-layer NN (1), and K is a positive semi-definite (PSD) kernel. Let the gram matrix of K over the training data S be $\mathbf{K} \in \mathbb{R}^{n \times n}, \mathbf{K}_{ij} = K(\vec{\mathbf{x}}_i, \vec{\mathbf{x}}_j)$ for $i, j \in [n]$, and $\mathbf{K}_n := \mathbf{K}/n$ is the empirical NTK matrix. Let the eigendecomposition of \mathbf{K}_n be $\mathbf{K}_n = \mathbf{U}\Sigma\mathbf{U}^{\top}$ where U is a $n \times n$ orthogonal matrix, and Σ is a diagonal matrix with its diagonal elements $\{\widehat{\lambda}_i\}_{i=1}^n$ being eigenvalues of \mathbf{K}_n and sorted in a non-increasing order. It is proved in existing works, such as (Du et al., 2019b), that \mathbf{K}_n is non-singular, and it can be verified 162 that $\widehat{\lambda}_1 \in (0, 1/2)$. Let \mathcal{H}_K be the Reproducing Kernel Hilbert Space (RKHS) associated with K. Because K is continuous on the compact set $\mathcal{X} \times \mathcal{X}$, the integral operator $T_K \colon L^2(\mathcal{X}, \mu) \to L^2(\mathcal{X}, \mu), (T_K f)(\mathbf{x}) \coloneqq \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{x}') f(\mathbf{x}') d\mu(\mathbf{x}')$ is a positive, self-adjoint, and compact operator 163 164 on $L^2(\mathcal{X},\mu)$. By the spectral theorem, there is a countable orthonormal basis $\{e_j\}_{j>1} \subseteq L^2(\mathcal{X},\mu)$ 166 and $\{\lambda_j\}_{j\geq 1}$ with $\frac{1}{2} \geq \lambda_1 \geq \lambda_2 \geq \ldots > 0$ such that e_j is the eigenfunction of T_K with λ_j being the corresponding eigenvalue. That is, $T_K e_j = \lambda_j e_j, j \geq 1$. Let $\{\mu_\ell\}_{\ell\geq 1}$ be the distinct eigenvalues 167 168 associated with T_K , and let m_ℓ be the be the sum of multiplicity of the eigenvalue $\{\mu_{\ell'}\}_{\ell'=1}^{\ell}$. That is, 169 $m_{\ell'} - m_{\ell'-1}$ is the multiplicity of $\mu_{\ell'}$. It is well known that $\{v_j = \sqrt{\lambda_j} e_j\}_{j>1}$ is an orthonormal 170 basis of \mathcal{H}_K . For a positive constant μ_0 , we define $\mathcal{H}_K(\mu_0) \coloneqq \{f \in \mathcal{H}_K : \|f\|_{\mathcal{H}} \le \mu_0\}$ as the closed ball in \mathcal{H}_K centered at 0 with radius μ_0 . We note that $\mathcal{H}_K(\mu_0)$ is also specified by $\mathcal{H}_K(\mu_0) =$ 171 172 $\left\{ f \in L^2(\mathcal{X}, \mu) \colon f = \sum_{j=1}^{\infty} \beta_j e_j, \sum_{j=1}^{\infty} \beta_j^2 / \lambda_j \le \mu_0^2 \right\}.$ 173 174

The Task of Nonparametric Regression. With $f^* \in \mathcal{H}_K(\mu_0)$, the task of the analysis for nonparametric regression is to find an estimator \hat{f} from the training data $\left\{ \left(\vec{\mathbf{x}}_i, y_i \right) \right\}_{i=1}^n$ so that the risk $\mathbb{E}_P \left[\left(\hat{f} - f^* \right)^2 \right]$ can converge to 0 with a fast rate. In this work, we aim to establish a sharp rate of the risk where the over-parameterized neural network (1) trained by GD with early stopping serves as the estimator \hat{f} .

Sharp rate of the risk of nonparametric regression using classical kernel regression. The statistical learning literature has established rich results in the sharp convergence rates for the risk of 183 nonparametric kernel regression (Stone, 1985; Yang & Barron, 1999; Raskutti et al., 2014; Yuan & Zhou, 2016), with one representative result in (Raskutti et al., 2014) about kernel regression trained 185 by GD with early stopping. Let ε_n be the critical population rate of the PSD kernel K, which is 186 also referred to as the critical radius (Wainwright, 2019) of K. (Raskutti et al., 2014, Theorem 187 2) shows the following sharp bound for the nonparametric regression risk of a kernel regression 188 model trained by GD with early stopping when $f^* \in \mathcal{H}_K(\mu_0)$. That is, with probability at least 189 $1 - \Theta\left(\exp(-\Theta(n\varepsilon_n^2))\right),$ 190

- 191
- 192

 $\mathbb{E}_{P}\left[\left(f_{\widehat{T}} - f^{*}\right)^{2}\right] \lesssim \varepsilon_{n}^{2},\tag{3}$

where \hat{T} is the stopping time whose formal definition is deferred to Section 5.1, and $f_{\hat{T}}$ is the kernel regressor at the \hat{T} -th step of GD for the optimization problem of kernel regression. The risk bound (3) is rather sharp, since it is minimax optimal in several popular learning setups, such as the setup where the eigenvalues $\{\lambda_i\}_{i\geq 1}$ exhibit a certain polynomial decay. Such risk bound (3) also holds for a general PSD kernel rather than the NTK (2), and the risk bound (3) is also minimax optimal when the PSD kernel is low rank. It is also remarked that the risk bound (3) is distribution-free in the bounded covariate, that is, there are no distributional assumptions about the covariate when it is in a bounded input space. Interested readers are referred to (Raskutti et al., 2014) for more details.

The main result of this paper is that the over-parameterized two-layer NN (1) trained by GD with early stopping achieves the same order of risk rate as that in (3) with arbitrary continuous distribution of the spherical covariate, which are summarized in the next section.

204 205 206

3 SUMMARY OF MAIN RESULTS.

207 Our main results are summarized in this section.

208 First, Theorem 5.1 in Section 5.2 shows that the neural network (1) trained by GD with early stop-209 ping using Algorithm 1 enjoys a sharp rate of the nonparametric regression risk, $\mathcal{O}(\varepsilon_n^2)$, which is 210 the same as that for the classical kernel regression in (3). Such rate of nonparametric regression 211 risk in Theorem 5.1 is distribution-free in spherical covariate, and it immediately leads to minimax 212 optimal rates for certain special cases. For example, when the eigenvalues of the integral operator 213 associated with K has a particular polynomial eigenvalue decay rate (EDR), that is, $\lambda_i \simeq j^{-\frac{d}{d-1}}$ 214 for $j \ge 1$, then in this case $\varepsilon_n^2 \asymp n^{-\frac{d}{2d-1}}$ according to (Raskutti et al., 2014, Corollary 3), and 215 Theorem 5.1 renders the rate of the nonparametric regression risk of $\mathcal{O}(n^{-\frac{d}{2d-1}})$ which is minimax 216 Table 1: Comparison between our result and the existing works on the risk rates and assumptions 217 for nonparametric regression by training over-parameterized neural networks with algorithmic guarantees, and the listed results here are under a common and popular setup that $f^* \in \mathcal{H}_{\tilde{K}}$ and the responses $\{y_i\}_{i=1}^n$ are corrupted by i.i.d. Gaussian noise with zero mean and variance σ^2 . 218 219

220				
221	Existing Works and Our Result	Distributional Assumptions	Eigenvalue Decay Rate (EDR)	Rate of Nonparametric Regression Risk
	(Kuzborskij & Szepesvári, 2021, Theorem 2)	No	-	Not minimax optimal, $\sigma^2 + O(n^{\frac{-2}{2+d}})$
222	(Hu et al., 2021, Theorem 5.2), (Suh et al., 2022, Theorem 3.11)	P is Unif (\mathcal{X})	$\lambda_j \asymp j^{-\frac{d}{d-1}}$	minimax optimal, $O(n^{\frac{-d}{2d-1}})$
223	(Li et al., 2024, Proposition 13)	$\begin{array}{c} P \text{ satisfies} \\ \text{a restrictive condition:} \\ \text{the density } p(\mathbf{x}) \text{ for } \mathbf{x} \in \mathbb{R}^d \text{ satisfies} \\ p(x) \lesssim (1 + \ \mathbf{x}\ _2^2)^{-(d+2)/2}. \end{array}$	$\lambda_j \asymp j^{-rac{d}{d-1}}$	minimax optimal, $\mathcal{O}(n^{\frac{-d}{2d-1}})$
224				
225				
223	Our Result (Theorem 5.1)	No distributional assumption about P as long as $\mathcal{X} = \mathbb{S}^{d-1}$	No requirement for EDR	$\mathcal{O}(\varepsilon_n^2)$, which leads to the minimax
226				optimal rate $\mathcal{O}(n^{\frac{-a}{2d-1}})$ claimed in
227				(Figure 1 a., 2021 ; Sun et al., 2022) and (Li et al., 2024)
				as special cases.

228 229 230

231

232

233

> optimal for this special case (Stone, 1985; Yang & Barron, 1999; Yuan & Zhou, 2016). We refer to such EDR the polynomial EDR in the sequel. It is shown in (Bietti & Mairal, 2019; Bietti & Bach, 2021; Li et al., 2024) that the polynomial EDR holds for our NTK in (2) if $P = \text{Unif}(\mathcal{X})$, or P satisfies the distributional assumption for (Li et al., 2024, Proposition 13) in Table 1.

We remark that such a minimax optimal rate $\mathcal{O}(n^{-\frac{d}{2d-1}})$ is derived from Theorem 5.1 under the 235 special case of polynomial EDR, and this minimax optimal rate is also the major result of a series 236 of existing works in nonparametric regression by training over-parameterized neural networks (Hu 237 et al., 2021; Suh et al., 2022; Li et al., 2024) when the target function f^* belongs to $\mathcal{H}_{\tilde{K}}$, the RKHS 238 associated with the NTK \tilde{K} of the network in each particular existing work. We note that \tilde{K} is 239 the NTK of the network considered in a particular existing work which may not be the same as 240 our NTK in (2). We also note that one needs to set s = 1 in (Li et al., 2024, Proposition 13) 241 so that $f^* \in \mathcal{H}_{\tilde{K}}$, and in this case the risk rate for nonparametric regression in (Li et al., 2024, 242 Proposition 13) is $\mathcal{O}(n^{-\frac{d}{2d-1}})$. To the best of our knowledge, Theorem 5.1 presents the first sharp 243 risk rate for nonparametric regression which is distribution-free in spherical covariate, which is 244 closer to practical scenarios. In contrast, the minimax rates in (Hu et al., 2021; Suh et al., 2022) 245 require spherical uniform data distribution on \mathcal{X} . The recent work (Ko & Huo, 2024) also requires 246 certain distributional assumptions for the results about regression convergence rates which does not 247 have algorithmic guarantees. Although the minimax rate in another recent work (Li et al., 2024) 248 does not need the spherical uniform distribution, it still requires a restrictive condition on the data 249 distributions detailed in Table 1, and such condition is met by sub-Gaussian distributions. It is 250 under this condition that (Li et al., 2024) derives the polynomial EDR. Table 1 compares our work 251 to existing works for nonparametric regression with a common setup, that is, $f^* \in \mathcal{H}_{\tilde{K}}$ and the responses $\{y_i\}_{i=1}^n$ are corrupted by i.i.d. Gaussian noise. We further note that although the result in 252 (Kuzborskij & Szepesvári, 2021, Theorem 2) does not require distributional assumptions about the 253 covariate, its risk rate under this common setup is not minimax optimal due to the term σ^2 in the risk 254 bound. Furthermore, the other term $\mathcal{O}(n^{\frac{-2}{2+d}})$ in its risk bound suffers from the curse of dimension 255 with a slow rate to 0 for high-dimensional data. We also note that (Kuzborskij & Szepesvári, 2021, 256 Theorem 1) shows the minimax optimal rate of $\mathcal{O}(n^{-\frac{2}{2+d}})$, however, this rate is derived for the 257 noiseless case where the responses are not corrupted by noise. 258

259 Second, our results provide confirmative answers to several outstanding open questions or address 260 particular concerns in the existing literature about training over-parameterized neural networks for 261 nonparametric regression by GD with early stopping and sharp risk rates, which are detailed below. 262

Stopping time in the early-stopping mechanism. An open question raised in (Kuzborskij & 263 Szepesvári, 2021; Hu et al., 2021) is how to characterize the stopping time in the early-stopping 264 mechanism when training the over-parameterized network by GD. Let \hat{T} be the stopping time, (Li 265 et al., 2024, Proposition 13) shows that the stopping time should satisfy $\hat{T} \simeq n^{\frac{d}{2d-1}}$ under the distri-266 butional assumption in Table 1. In contrast, Theorem 5.1 provides a characterization of \widehat{T} showing 267 that $\widehat{T} \simeq \varepsilon_n^{-2}$, which is distribution-free in spherical covariate. Theorem 5.1 further suggests that 268 for each neural network function f_t obtained at the t-th step of GD with $t \approx \varepsilon_n^{-2}$, the sharp risk rate 269 of $\mathcal{O}(\varepsilon_n^2)$ is obtained.

270 Lower bound for the network width m. Our main result, Theorem 5.1, requires that the network 271 width m, which is the number of neurons in the first layer of the network, satisfies $m \gtrsim d^2/(\varepsilon_n^{16})$. 272 Such lower bound for m solely depends on d and ε_n . Under the polynomial EDR, Corollary 5.2, 273 which is a direct consequence of Theorem 5.1, shows that m should satisfy $m \gtrsim n^{\frac{16\alpha}{2\alpha+1}} d^2$ with 274 $\alpha = d/(2(d-1))$ (see (12)) so that GD with early stopping leads to the minimax rate of $\mathcal{O}(n^{-\frac{d}{2d-1}})$. 275 We remark that this is the first time that the lower bound for the network width m is specified only 276 in terms of n and d under the polynomial EDR with a minimax optimal risk rate for nonparamet-277 ric regression, which can be easily estimated from the training data. In contrast, under the same 278 polynomial EDR, all the existing works (Hu et al., 2021; Suh et al., 2022; Li et al., 2024) require 279 $m \gtrsim poly(n, 1/\lambda_n)$. The problem here is that one needs additional assumptions on the training 280 data (Bartlett et al., 2021; Nguyen et al., 2021) to find the lower bound for $\hat{\lambda}_n$, which is the minimal 281 eigenvalue of the empirical NTK matrix \mathbf{K}_n , to further estimate the lower bound for m using the 282 training data. 283

Corollary 5.2 also gives a competitive and smaller lower bound for the network width m than 284 some existing works which give explicit orders of the lower bound for m. For example, un-285 der the assumption of uniform spherical distribution, (Suh et al., 2022, Theorem 3.11) requires that $m/\log m \gtrsim L^{20}n^{24}$ where L is the number of layers of the DNN used in that work, and $m/\log m \gtrsim 2^{20}n^{24}$ even with L = 2 for the two-layer network (1) used in our work. Furthermore, 286 287 the proof of (Li et al., 2024, Proposition 13) suggests that $m \gtrsim n^{24} (\log m)^{12}$. Both lower bounds for m in (Suh et al., 2022, Theorem 3.11) and (Li et al., 2024, Proposition 13) are much larger than 288 289 our lower bound for m, $n^{\frac{16\alpha}{2\alpha+1}}d^2$, when $n \to \infty$ and d is fixed, which is the setup considered in 290 291 (Li et al., 2024). It is worthwhile to mention that (Suh et al., 2022; Li et al., 2024) use DNNs with 292 multiple layers for nonparametric regression. As shown in Table 1, through our careful analysis, a 293 shallow two-layer NN (1) exhibits the same minimax risk rate as its deeper counterpart under the same assumptions with much smaller network width. This observation further support the claim in (Bietti & Bach, 2021) that a shallow over-parameterized neural networks with ReLU activations 295 exhibit the same approximation properties as its deeper counterpart, in our nonparametric regression 296 setup. 297

298 **Training the network with learning rate** $\eta = \Theta(1)$. It is also worthwhile to mention that our 299 main result, Theorem 5.1, suggests that a constant learning rate $\eta = \Theta(1)$ can be used for GD when training the two-layer NN (1), which could lead to better empirical optimization performance 300 in practice. Some existing works in fact require an infinitesimal η . For example, (Li et al., 2024, 301 Proposition 13) is obtained by gradient flow where $\eta \to 0$ instead of the practical GD. Furthermore, 302 (Hu et al., 2021, Theorem 5.2) requires the learning rates for both the squared loss and the ℓ^2 -303 regularization term to have the order of $o(n^{-\frac{3d-1}{2d-1}}) \to 0$ as $n \to \infty$. We note that (Nitanda & 304 Suzuki, 2021) also employs constant learning rate in SGD to train neural networks. 305

More discussion about the literature. We herein provide more discussion about the results of this work and comparison to the existing relevant works with sharp rates for nonparametric regression. While this paper establishes sharp rate which is distribution-free in spherical covariate, such rate still depends on bounded input space ($\mathcal{X} = \mathbb{S}^{d-1}$) and the condition that the target function $f^* \in$ $\mathcal{H}_K(\mu_0)$. Some other existing works consider target function f^* not belonging to the RKHS ball centered at the origin with constant or low radius, such as (Haas et al., 2023; Bordelon et al., 2024). A more detailed discussion is deferred to Section B of the appendix.

313 314

315

316

4 TRAINING BY GRADIENT DESCENT AND PRECONDITIONED GRADIENT DESCENT

In the training process of our network (1), only W is optimized with a randomly initialized to ± 1 and then fixed. The following quadratic loss function is minimized during the training process:

321

$$L(\mathbf{W}) \coloneqq \frac{1}{2n} \sum_{i=1}^{n} \left(f(\mathbf{W}, \mathbf{\vec{x}}_{i}) - y_{i} \right)^{2}.$$
(4)

Algorithm	1 Training the	e Two-Layer NN
oy GD		

1: $\mathbf{W}(T) \leftarrow \text{Training-by-GD}(T, \mathbf{W}(0))$ 2: **input:** $T, \mathbf{W}(0)$ 3: **for** $t = 1, \dots, T$ **do** 4: Perform the *t*-th step of GD by (5) 5: **end for** 6: **return** $\mathbf{W}(T)$ $\begin{array}{ll} \textbf{324} \\ \textbf{325} \\ \textbf{326} \end{array} \quad \begin{array}{ll} \text{In the } (t+1) \text{-th step of GD with } t \geq 0, \text{ the weights of } \\ \textbf{the neural network, } \mathbf{W_S}, \text{ are updated by one-step of GD } \\ \textbf{through} \end{array}$

327 328

330

$$\operatorname{vec}\left(\mathbf{W}_{\mathbf{S}}(t+1)\right) - \operatorname{vec}\left(\mathbf{W}_{\mathbf{S}}(t)\right) = -\frac{\eta}{n} \mathbf{Z}_{\mathbf{S}}(t) (\widehat{\mathbf{y}}(t) - \mathbf{y}),$$
(5)

where $\mathbf{y}_i = y_i, \, \mathbf{\hat{y}}(t) \in \mathbb{R}^n$ with $[\mathbf{\hat{y}}(t)]_i = f(\mathbf{W}(t), \mathbf{x}_i)$. The notations with the subscripts **S** indicate the dependence on the training features **S**. We also denote $f(\mathbf{W}(t), \cdot)$ as $f_t(\cdot)$ as the neural network function with weighting vectors $\mathbf{W}(t)$ obtained after the *t*-th step of GD. We define $\mathbf{Z}_{\mathbf{S}}(t) \in \mathbb{R}^{md \times n}$ which is computed by

337

344 345 $(\mathbf{Z}_{\mathbf{S}}(t))_{[(r-1)d+1:rd]i} = \frac{1}{\sqrt{m}} \mathbb{I}_{\left\{\vec{\mathbf{w}}_{r}(t)^{\top}\vec{\mathbf{x}}_{i}\geq0\right\}} \vec{\mathbf{x}}_{i} a_{r}, i \in [n], r \in [m],$ (6) where $(\mathbf{Z}_{\mathbf{S}}(t))_{[(r-1)d+1:rd]i} \in \mathbb{R}^{d}$ is a vector with elements in the *i*-th column of $\mathbf{Z}_{\mathbf{S}}(t)$ with indices

where $(\mathbf{Z}_{\mathbf{S}}(t))_{[(r-1)d+1:rd]i} \in \mathbb{R}^{a}$ is a vector with elements in the *i*-th column of $\mathbf{Z}_{\mathbf{S}}(t)$ with indices in [(r-1)d+1:rd]. We employ the following particular symmetric random initialization so that $\hat{\mathbf{y}}(0) = \mathbf{0}$, which has been used in existing works such as (Chizat et al., 2019; Zhang et al., 2020). In our two-layer NN, *m* is even, $\{\vec{\mathbf{w}}_{2r'}(0)\}_{r'=1}^{m/2}$ and $\{a_{2r'}\}_{r'=1}^{m/2}$ are initialized randomly and independently according to

$$\vec{\mathbf{w}}_{2r'}(0) \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d), a_{2r'} \sim \operatorname{unif}\left(\{-1, 1\}\right), \quad \forall r' \in [m/2],$$
(7)

where $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes a Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, unif ($\{-1, 1\}$) 346 denotes a uniform distribution over $\{1, -1\}, 0 < \kappa \leq 1$ controls the magnitude of initialization. 347 We set $\vec{\mathbf{w}}_{2r'-1}(0) = \vec{\mathbf{w}}_{2r'}(0)$ and $a_{2r'-1} = -a_{2r}$ for all $r' \in [m/2]$. It then can be verified that 348 $\hat{\mathbf{y}}(0) = \mathbf{0}$, that is, the initial output of the two-layer network (1) is zero. Once randomly initialized, 349 a is fixed during the training. We use $\mathbf{W}(0)$ to denote the set of all the random weighting vectors 350 at initialization, that is, $\mathbf{W}(0) = \left\{ \vec{\mathbf{w}}_r(0) \right\}_{r=1}^m$. We run Algorithm 1 to train the two-layer NN by 351 352 GD, where T is the total number of steps for GD. Early stopping is enforced in Algorithm 1 through 353 a bounded T via $T \leq T$. 354

355 356

357

359 360

361 362

364

365 366 367

368

5 MAIN RESULTS

We present the definition of kernel complexity in this section, and then introduce the main results for nonparametric regression of this paper.

5.1 KERNEL COMPLEXITY

The local kernel complexity has been studied by (Bartlett et al., 2005; Koltchinskii, 2006; Mendelson, 2002). For the PSD kernel K, we define the empirical kernel complexity \hat{R}_K and the population kernel complexity R_K as

$$\widehat{R}_{K}(\varepsilon) \coloneqq \sqrt{\frac{1}{n} \sum_{i=1}^{n} \min\left\{\widehat{\lambda}_{i}, \varepsilon^{2}\right\}}, \quad R_{K}(\varepsilon) \coloneqq \sqrt{\frac{1}{n} \sum_{i=1}^{\infty} \min\left\{\lambda_{i}, \varepsilon^{2}\right\}}.$$
(8)

369 It can be verified that both $\sigma R_K(\varepsilon)$ and $\sigma \hat{R}_K(\varepsilon)$ are sub-root functions (Bartlett et al., 2005) in terms 370 of ε^2 . The formal definition of sub-root functions is deferred to Definition A.2 in the appendix. For 371 a given noise ratio σ , the critical empirical radius $\hat{\varepsilon}_n > 0$ is the smallest positive solution to the 372 inequality $\widehat{R}_K(\varepsilon) \leq \varepsilon^2/\sigma$, where $\widehat{\varepsilon}_n^2$ is the also the fixed point of $\sigma \widehat{R}_K(\varepsilon)$ as a function of ε^2 : 373 $\sigma \widehat{R}_K(\widehat{\varepsilon}_n) = \widehat{\varepsilon}_n^2$. Similarly, the critical population rate ε_n is defined to be the smallest positive 374 solution to the inequality $R_K(\varepsilon) \leq \varepsilon^2 / \sigma$, where ε_n^2 is the fixed point of $\sigma R_K(\varepsilon)$ as a function of ε^2 : 375 $\sigma R_K(\varepsilon_n) = \varepsilon_n^2$. In this paper we consider the case that $n\varepsilon_n^2 \to \infty$ as $n \to \infty$, which is also used 376 in standard analysis of nonparametric regression with minimax rates by kernel regression (Raskutti 377 et al., 2014).

Let $\eta_t \coloneqq \eta t$ for all $t \ge 0$, we then define the stopping time \widehat{T} as \widehat{T}

$$\widehat{T} := \min\left\{t \colon \widehat{R}_K(\sqrt{1/\eta_t}) > (\sigma\eta_t)^{-1}\right\} - 1.$$
(9)

The stopping time in fact limit the number of steps T in for Algorithm 1 as to be shown in Section 5.2, which in turn enforces the early stopping mechanism.

5.2 Results

380

382

383 384

386

387

388 389

390 391

392

393

394

395

397

Theorem 5.1. Let $c_T, c_t \in (0, 1]$ be arbitrary positive constants, and $c_T \hat{T} \leq T \leq \hat{T}$. Suppose $f^* \in \mathcal{H}_K(\mu_0)$, and *m* satisfies

$$m \gtrsim \frac{d^2}{\varepsilon_n^{16}},\tag{10}$$

and the neural network $f(\mathbf{W}(t), \cdot)$ is trained by GD using Algorithm 1 with the learning rate $\eta \in [1,2)$ and $T \leq \hat{T}$. Then for every $t \in [c_t T: T]$, with probability at least $1 - \exp(-\Theta(n)) - 7\exp(-\Theta(n\varepsilon_n^2)) - 2/n$ over the random noise w, the random training features S and the random initialization $\mathbf{W}(0)$, the stopping time satisfies $\hat{T} \approx \varepsilon_n^{-2}$, and $f(\mathbf{W}(t), \cdot) = f_t$ satisfies

$$\mathbb{E}_P\left[(f_t - f^*)^2\right] \lesssim \varepsilon_n^2. \tag{11}$$

Significance of Theorem 5.1 and comparison to existing works. To the best of our knowledge, 398 Theorem 5.1 is the first theoretical result which proves that over-parameterized neural network 399 trained by gradient descent with early stopping achieves sharp rate of $\mathcal{O}(\varepsilon_n^2)$, without distributional 400 assumption on the covariate as long as the input space \mathcal{X} is \mathbb{S}^{d-1} . To understand the sharpness of 401 the bound for the risk in (11), Corollary 5.2 shows that when the polynomial EDR holds, that is, 402 $\lambda_j \simeq j^{-\frac{d}{d-1}}$, then $\varepsilon_n^2 \simeq n^{-\frac{d}{2d-1}}$, and the rate of the risk is $\mathcal{O}(n^{-\frac{d}{2d-1}})$ which is minimax optimal under the polynomial EDR for $f^* \in \mathcal{H}_K(\mu_0)$ (Stone, 1985; Yang & Barron, 1999; Yuan & Zhou, 403 404 2016). (Bietti & Mairal, 2019; Bietti & Bach, 2021; Li et al., 2024) show that such polynomial EDR 405 holds for our NTK (2) if P is Unif (\mathcal{X}) , or P satisfies the distributional assumption for (Li et al., 406 2024, Proposition 13) in Table 1. The existing works (Hu et al., 2021; Suh et al., 2022; Li et al., 407 2024) prove the same minimax optimal rate for an over-parameterized neural network trained by GD 408 with either regularization or early stopping. However, it is remarked that such minimax optimal rates in these works are proved either for spherical uniform distribution on \mathcal{X} (Hu et al., 2021; Suh et al., 409 2022), or for distributions satisfying certain restrictive condition (Li et al., 2024). Table 1 compares 410 our result to existing works from the perspective of risk rates for nonparametric regression, required 411 distributional assumptions on the covariate, and the associated EDR. 412

413 We also emphasize that Theorem 5.1, for the first time, shows that the network width m required to 414 achieve the minimax rate can be quantized in terms of a well known quantity about the kernel K, 415 the critical population rate ε_n , and d in the manner of distribution-free in spherical covariate. More 416 discussions are referred to "Significance of Corollary 5.2".

Furthermore, Theorem 5.1, for the first time, gives an explicit characterization of the stopping time \hat{T} for training an over-parameterized neural network by GD with early stopping which is of the order $\hat{T} \simeq \varepsilon_n^{-2}$ and distribution-free in spherical covariate. This result suggests that t should be of the order $\Theta(\varepsilon_n^{-2})$ to ensure the sharp rate (11). Such result gives an order of the number of steps for GD when training the over-parameterized NN (1) so as to achieve the sharp risk bound $\mathcal{O}(\varepsilon_n^2)$. Under the polynomial EDR, the stopping time \hat{T} satisfies $\hat{T} \simeq n^{\frac{d}{2d-1}}$, which recovers the same result about the stopping time in (Li et al., 2024, Proposition 13).

425 When the polynomial EDR holds, we can apply Theorem 5.1 to obtain the following corollary.

426 Corollary 5.2 (Applying Theorem 5.1 to the special case of polynomial EDR). Suppose $\lambda_j \simeq j^{-2\alpha}$ 427 for $j \ge 1$ and $\alpha > 1/2$. Let $c_T, c_t \in (0, 1]$ be positive constants, and $c_T \hat{T} \le T \le \hat{T}$. Suppose m428 satisfies

429 430

$$m \gtrsim n^{\frac{16\alpha}{2\alpha+1}} d^2,\tag{12}$$

and the neural network $f(\mathbf{W}(t), \cdot)$ is trained by GD using Algorithm 1 with the learning rate $\eta \in [1, 2)$ and $T \leq \hat{T}$. Then for every $t \in [c_t T: T]$, with probability at least $1 - \exp(-\Theta(n)) - \Theta(n)$

 $\begin{array}{l} {}^{432}_{433} \\ {}^{433}_{434} \end{array} \qquad 7 \exp\left(-\Theta(n\varepsilon_n^2)\right) - 2/n \text{ over the random noise } \mathbf{w}, \text{ the random training features } \mathbf{S} \text{ and the random initialization } \mathbf{W}(0), \text{ the stopping time satisfies } \widehat{T} \asymp n^{\frac{d}{2d-1}}, \end{array}$

$$\mathbb{E}_P\left[(f_t - f^*)^2\right] \lesssim \left(\frac{1}{n}\right)^{\frac{2\alpha}{2\alpha+1}}.$$
(13)

Significance of Corollary 5.2. Corollary 5.2 shows that under the polynomial EDR, GD finds an over-parameterized neural network with minimax optimal rate of $\mathcal{O}(n^{-\frac{2\alpha}{2\alpha+1}}) = \mathcal{O}(n^{-\frac{d}{2d-1}})$, where $\alpha = d/(2(d-1))$, with a specific quantization of m in terms of only n and d in (12). In contrast, all the existing works (Hu et al., 2021; Suh et al., 2022; Li et al., 2024) require $m \gtrsim poly(n, 1/\hat{\lambda}_n)$, and additional assumptions on the training data (Bartlett et al., 2021; Nguyen et al., 2021) are required to bound $\hat{\lambda}_n$ from below so as to estimate the lower bound for m from the training data.

6 ROADMAP OF PROOFS

We present the roadmap of our theoretical results which lead to the main result, Theorem 5.1 in Section 5. We first present in the next subsection our results about the uniform convergence to the NTK (2) and more, which are crucial in the analysis of training dynamics by GD.

6.1 UNIFORM CONVERGENCE TO THE NTK AND MORE

We define functions

$$h(\mathbf{w}, \mathbf{x}, \mathbf{y}) \coloneqq \mathbf{x}^{\top} \mathbf{y} \mathbb{1}_{\{\mathbf{w}^{\top} \mathbf{x} \ge 0\}} \mathbb{1}_{\{\mathbf{w}^{\top} \mathbf{y} \ge 0\}}, \qquad \widehat{h}(\mathbf{W}, \mathbf{x}, \mathbf{y}) \coloneqq \frac{1}{m} \sum_{r=1}^{m} h(\vec{\mathbf{w}}_{r}, \mathbf{x}, \mathbf{y}), \qquad (14)$$

$$v_R(\mathbf{w}, \mathbf{x}) \coloneqq \mathbb{I}_{\{|\mathbf{w}^\top \mathbf{x}| \le R\}}, \qquad \qquad \widehat{v}_R(\mathbf{W}, \mathbf{x}) \coloneqq \frac{1}{m} \sum_{r=1}^m v_R(\vec{\mathbf{w}}_r, \mathbf{x}). \qquad (15)$$

Then we have the following theorem stating the uniform convergence of $\hat{h}(\mathbf{W}(0), \cdot, \cdot)$ to $K(\cdot, \cdot)$ and uniform convergence of $\hat{v}_R(\mathbf{W}(0), \mathbf{x})$ to $\frac{2R}{\sqrt{2\pi\kappa}}$ for a positive number $R \leq \eta T/\sqrt{m}$. While existing works such as (Li et al., 2024) also has uniform convergence results for over-parameterized neural network, our result does not depend on the Hölder continuity of the NTK.

Theorem 6.1. The following results hold with $m \ge \max\{d, n, 4\}$ and $m/\log m \ge d$.

(1) With probability at least
$$1 - 1/n$$
 over the random initialization $\mathbf{W}(0) = \left\{ \vec{\mathbf{w}}_r(0) \right\}_{r=1}^m$,

$$\sup_{\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{X}} \left| K(\mathbf{x}, \mathbf{y}) - \hat{h}(\mathbf{W}(0), \mathbf{x}, \mathbf{y}) \right| \le C_1(m, d, 1/n) \lesssim \sqrt{\frac{d \log m}{m}}.$$
(16)

(2) Suppose $m \ge (c_{\mathbf{u}}\eta T/R_0)^2$ for an arbitrary absolute positive constant $R_0 < \kappa$. Then with probability at least 1 - 1/n over the random initialization $\mathbf{W}(0) = \left\{ \vec{\mathbf{w}}_r(0) \right\}_{r=1}^m$,

$$\sup_{\mathbf{x}\in\mathcal{X}} |\widehat{v}_R(\mathbf{W}(0), \mathbf{x})| \le \frac{2R}{\sqrt{2\pi\kappa}} + C_2(m, R_0, 1/n) \lesssim \frac{\sqrt{d}}{m^{1/4}} + \frac{\eta T}{\sqrt{m}},\tag{17}$$

where $C_1(m, d, 1/n), C_2(m, R_0, 1/n)$ are two positive numbers depending on (m, d, n) and (m, R_0, n) , respectively, with their formal definitions deferred to (39) and (41) in Section C.2 of the appendix.

Proof. This theorem follows from Theorem C.1 and Theorem C.2 in Section C.2 of the appendix. 483 Note that $\hat{h}(\mathbf{W}, \mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{r=1}^{m} h(\vec{\mathbf{w}}_r, \mathbf{x}, \mathbf{y}) = \frac{1}{m/2} \sum_{r'=1}^{m/2} h(\vec{\mathbf{w}}_{2r}(0), \mathbf{x}, \mathbf{y})$, then part (1) directly fol-484 lows from Theorem C.1. Similarly, part (2) directly follows from Theorem C.2, and noting that $m \ge (c_{\mathbf{u}}\eta T/R_0)^2$ indicates $R \le R_0$. 486 Define

488

493

$$\mathcal{W}_0 \coloneqq \{ \mathbf{W}(0) \colon (\mathbf{16}), (\mathbf{17}) \text{ hold} \}$$

$$\tag{18}$$

be the set of all the good random initializations which satisfy (16) and (17) in Theorem 6.1. Theorem 6.1 shows that we have good random initialization with high probability, that is, $\Pr[\mathbf{W}(0) \in \mathcal{W}_0] \ge 1 - 2/n$. When $\mathbf{W}(0) \in \mathcal{W}_0$, the uniform convergence results, (16) and (17), hold with high probability, which is crucial for our main result in Theorem 5.1.

6.2 ROADMAP OF PROOFS

Because our main result, Theorem 5.1, is proved by Theorem C.10 and Theorem C.11 deferred to Section C.2, we illustrate in Figure 1, deferred to the appendix, the roadmap containing the intermediate theoretical results which lead to our main result, Theorem 5.1.

498 Summary of the technical approaches and novel results in the proofs. Theorem C.8 is the first 499 novel result in this work, showing that with high probability, the neural network function $f(\mathbf{W}(t), \cdot)$ 500 at step t of GD can be decomposed into two functions by $f(\mathbf{W}(t), \cdot) = f_t = h + e$, where $h \in \mathcal{H}_K$ is 501 a function in the RKHS associated with K with bounded \mathcal{H}_K -norm. The error function e has a small 502 L^{∞} -norm, that is, $\|e\|_{\infty} \leq w$ with w being a small number controlled by the network width m, that is, larger m leads to smaller w. Theorem C.10 is the second novel result, where we derive sharp and 504 novel bound for the nonparametric regression risk of the neural network function $f(\mathbf{W}(t), \cdot)$ in Theorem C.10, that is, $\mathbb{E}_P\left[(f_t - f^*)^2\right] - 2\mathbb{E}_{P_n}\left[(f_t - f^*)^2\right] \lesssim \varepsilon_n^2 + w$. To the best of our knowledge, Theorem C.10 is among the first in the literature to employ local Rademacher complexity so as to 505 506 obtain sharp rate for the risk of nonparametric regression which is distribution-free in spherical co-507 variate, and local Rademacher complexity is employed to tightly bound the Rademacher complexity 508 of the function class comprising all the possible neural network functions obtained by GD. 509

510

Novel proof strategy of this work. We remark that the proof strategy of our main result, Theo-511 rem 5.1, is significantly novel and different from the existing works in training over-parameterized 512 neural networks for nonparametric regression with minimax rates (Hu et al., 2021; Suh et al., 2022; 513 Li et al., 2024). In particular, the common proof strategy in these works uses the decomposition $f_t - f^* = (f_t - \hat{f}_t^{(\text{NTK})}) + (\hat{f}_t^{(\text{NTK})} - f^*)$ and then show that both $\left\| f_t - \hat{f}_t^{(\text{NTK})} \right\|_{L^2}$ and $\left\| \hat{f}_t^{(\text{NTK})} - f^* \right\|_{L^2}$ are bounded by certain minimax optimal rate, where $\hat{f}_t^{(\text{NTK})}$ is the kernel re-514 515 516 517 gressor obtained by either kernel ridge regression (Hu et al., 2021; Suh et al., 2022) or GD with 518 early stopping (Li et al., 2024). The remark after Theorem C.8 details a formulation of $\hat{f}_t^{(\text{NTK})}$. 519 $\left\| \hat{f}_t^{(\text{NTK})} - f^* \right\|_{L^2}$ is bounded by the minimax optimal rate under certain distributional assumptions 520 521 in the covariate, and this is one reason for the distributional assumptions about the covariate in existing works such as (Hu et al., 2021; Suh et al., 2022; Li et al., 2024). In a strong contrast, our analysis 522 does not rely on such decomposition of $f_t - f^*$. Instead of approximating f_t by $\hat{f}_t^{(\text{NTK})}$, we have a new decomposition of f_t by $f_t = h_t + e_t$ where f_t is approximated by h_t with e_t being the approxi-523 524 mation error. As suggested by the remark after Theorem C.8, we have $h_t = \hat{f}_t^{(\text{NTK})} + \hat{e}_2(\cdot, t)$ so that 525 $f_t = \hat{f}_t^{(\text{NTK})} + \hat{e}_2(\cdot, t) + e_t$. Our analysis only requires the network width m to be suitably large so that the \mathcal{H}_K -norm of $\hat{e}_2(\cdot, t)$ is bounded by a positive constant and $||e_t||_{\infty} \leq w$, while the common proof strategy in(Hu et al., 2021; Suh et al., 2022; Li et al., 2024) needs m to be sufficiently large so 526 527 528 that both $\|\widehat{e}_2(\cdot,t)\|_{\infty}$ and $\|e_t\|_{\infty}$ are bounded by an infinitesimal number (a minimax optimal rate 529 such as $\mathcal{O}(n^{-\frac{d}{2d-1}})$ and then $\left\|f_t - \hat{f}_t^{(\text{NTK})}\right\|_{L^2}$ is bounded by such minimax optimal rate. Detailed 530 531 in Section 3, such novel proof strategy leads to our sharp analysis, rendering a smaller lower bound 532 for m in our main result compared to some existing works. 533

534 7 CONCLUSION

In this paper, we show that an over-parameterized two-layer neural network trained by gradient descent (GD) with early stopping renders a sharp rate of the nonparametric regression risk with the order of $\Theta(\varepsilon_n^2)$ with ε_n being the critical population rate or the critical radius of the NTK, which is distribution-free in spherical covariate. We compare our results to the current state-of-the-art with a detailed roadmap of our technical approaches and results in our proofs.

540 REFERENCES

542

ized neural networks, going beyond two layers. In Hanna M. Wallach, Hugo Larochelle, Alina 543 Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), Advances in Neu-544 ral Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 6155–6166, 546 2019a. 547 548 Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. In International Conference on Machine Learning, volume 97 of Proceedings 549 of Machine Learning Research, pp. 242–252. PMLR, 2019b. 550 551 Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of op-552 timization and generalization for overparameterized two-layer neural networks. In International 553 Conference on Machine Learning, volume 97 of Proceedings of Machine Learning Research, pp. 554 322-332. PMLR, 2019. 555 Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. Ann. 556 Statist., 33(4):1497–1537, 08 2005. 558 Peter L. Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. 559 Acta Numerica, 30:87-201, 2021. doi: 10.1017/S0962492921000027. 560 Benedikt Bauer and Michael Kohler. On deep learning as a remedy for the curse of dimensionality 561 in nonparametric regression. Ann. Statist., 47(4):2261 - 2285, 2019. 562 563 Alberto Bietti and Francis R. Bach. Deep equals shallow for relu networks in kernel regimes. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 564 3-7, 2021. OpenReview.net, 2021. 565 566 Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. In Hanna M. 567 Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman 568 Garnett (eds.), Advances in Neural Information Processing Systems, pp. 12873–12884, 2019. 569 Blake Bordelon, Alexander B. Atanasov, and Cengiz Pehlevan. How feature learning can improve 570 neural scaling laws. CoRR, abs/2409.17858, 2024. doi: 10.48550/ARXIV.2409.17858. URL 571 https://doi.org/10.48550/arXiv.2409.17858. 572 573 Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and 574 deep neural networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence 575 d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), Advances in Neural Information Processing Systems, pp. 10835–10845, 2019. 576 577 Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. 578 Curran Associates Inc., Red Hook, NY, USA, 2019. 579 Simon S. Du, Jason D. Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds 580 global minima of deep neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), 581 International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning 582 Research, pp. 1675-1685. PMLR, 2019a. 583 584 Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes 585 over-parameterized neural networks. In International Conference on Learning Representations, 586 2019b. Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers 588 neural networks in high dimension. Ann. Statist., 49(2):1029 – 1054, 2021. 590 Moritz Haas, David Holzmüller, Ulrike von Luxburg, and Ingo Steinwart. Mind the spikes: Benign overfitting of kernels and neural networks in fixed dimension. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural In-592 formation Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023.

Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameter-

617

618

619

621

622

623

624

- Tianyang Hu, Wenjia Wang, Cong Lin, and Guang Cheng. Regularization matters: A nonparametric perspective on overparametrized neural network. In Arindam Banerjee and Kenji Fukumizu (eds.), *International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pp. 829–837. PMLR, 2021.
- Masaaki Imaizumi and Kenji Fukumizu. Deep neural networks learn non-smooth functions effectively. In Kamalika Chaudhuri and Masashi Sugiyama (eds.), *International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 869–878. PMLR, 2019.
- Yuling Jiao, Guohao Shen, Yuanyuan Lin, and Jian Huang. Deep nonparametric regression on
 approximate manifolds: Nonasymptotic error bounds with polynomial prefactors. *Ann. Statist.*,
 51(2):691 716, 2023.
- Hyunouk Ko and Xiaoming Huo. Universal consistency of wide and deep relu neural networks and minimax optimal convergence rates for kolmogorov-donoho optimal function classes. In *Fortyfirst International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27,* 2024. OpenReview.net, 2024.
- Vladimir Koltchinskii. Local rademacher complexities and oracle inequalities in risk minimization.
 Ann. Statist., 34(6):2593–2656, 12 2006.
- 613 Ilja Kuzborskij and Csaba Szepesvári. Nonparametric regression with shallow overparameterized
 614 neural networks trained by GD with early stopping. In Mikhail Belkin and Samory Kpotufe
 615 (eds.), *Conference on Learning Theory, COLT 2021, 15-19 August 2021, Boulder, Colorado,*616 USA, volume 134 of *Proceedings of Machine Learning Research*, pp. 2853–2890. PMLR, 2021.
 - B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302 1338, 2000.
- 44, 2015. Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521:436–444, 2015.
 - Michel. Ledoux. Probability in Banach Spaces [electronic resource] : Isoperimetry and Processes / by Michel Ledoux, Michel Talagrand. Classics in Mathematics. Springer Berlin Heidelberg, Berlin, Heidelberg, 1st ed. 1991. edition, 1991.
- Yicheng Li, Zixiong Yu, Guhan Chen, and Qian Lin. On the eigenvalue decay rates of a class of
 neural-network related kernel functions defined on general domains. *Journal of Machine Learning Research*, 25(82):1–47, 2024.
- Shahar Mendelson. Geometric parameters of kernel machines. In Jyrki Kivinen and Robert H.
 Sloan (eds.), *Conference on Computational Learning Theory*, volume 2375 of *Lecture Notes in Computer Science*, pp. 29–43. Springer, 2002.
- Quynh Nguyen, Marco Mondelli, and Guido F. Montúfar. Tight bounds on the smallest eigenvalue of the neural tangent kernel for deep relu networks. In Marina Meila and Tong Zhang (eds.), *International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8119–8129. PMLR, 2021.
- Atsushi Nitanda and Taiji Suzuki. Optimal rates for averaged stochastic gradient descent under neural tangent kernel regime. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- Garvesh Raskutti, Martin J. Wainwright, and Bin Yu. Early stopping and non-parametric regression: an optimal data-dependent stopping rule. *J. Mach. Learn. Res.*, 15(1):335–366, 2014.
- Lorenzo Rosasco, Mikhail Belkin, and Ernesto De Vito. On learning with integral operators. J.
 Mach. Learn. Res., 11:905–934, 2010.
- Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *Ann. Statist.*, 48(4):1875 1897, 2020.
- 647 Charles J. Stone. Additive Regression and Other Nonparametric Models. Ann. Statist., 13(2):689 705, 1985.

- Lili Su and Pengkun Yang. On learning over-parameterized neural networks: A functional approximation perspective. In *Advances in Neural Information Processing Systems*, pp. 2637–2646, 2019.
- Namjoon Suh, Hyunouk Ko, and Xiaoming Huo. A non-parametric regression viewpoint : Generalization of overparametrized deep RELU network under noisy observations. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. In Yonina C.
 Eldar and GittaEditors Kutyniok (eds.), *Compressed Sensing: Theory and Practice*, pp. 210–268.
 Cambridge University Press, 2012.
- Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.
- F. T. Wright. A Bound on Tail Probabilities for Quadratic Forms in Independent Random Variables
 Whose Distributions are not Necessarily Symmetric. *Ann. Probab.*, 1(6):1068 1070, 1973.
- Greg Yang and Edward J. Hu. Tensor programs IV: feature learning in infinite-width neural networks. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event,* volume 139 of *Proceedings of Machine Learning Research,* pp. 11727–11737. PMLR, 2021.
- Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *Ann. Statist.*, 27(5):1564 1599, 1999.
- Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94: 103–114, 2017.
- Ming Yuan and Ding-Xuan Zhou. Minimax optimal rates of estimation in high dimensional additive models. *Ann. Statist.*, 44(6):2564 2593, 2016.
- Kaiqi Zhang and Yu-Xiang Wang. Deep learning meets nonparametric regression: Are weight decayed dnns locally adaptive? In *International Conference on Learning Representations*. Open Review.net, 2023.
- Yaoyu Zhang, Zhi-Qin John Xu, Tao Luo, and Zheng Ma. A type of generalization error induced by initialization in deep neural networks. In Jianfeng Lu and Rachel A. Ward (eds.), *Proceedings of Mathematical and Scientific Machine Learning, MSML 2020, 20-24 July 2020, Virtual Conference / Princeton, NJ, USA*, volume 107 of *Proceedings of Machine Learning Research*, pp. 144–164. PMLR, 2020.
 - Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems*, pp. 2053–2062, 2019.
 - We present the basic mathematical results required in our proofs in Section A, then present proofs in the subsequent sections.
- 691 692 693

696

686

687

688

689 690

651

659

664

A MATHEMATICAL TOOLS

We introduce the basic definitions and mathematical results as the basic tools for the subsequent results in the next sections of this appendix.

697 Definition A.1. Let $\{\sigma_i\}_{i=1}^n$ be n i.i.d. random variables such that $\Pr[\sigma_i = 1] = \Pr[\sigma_i = -1] = \frac{1}{2}$. 698 The Rademacher complexity of a function class \mathcal{F} is defined as

700
701
$$\mathfrak{R}(\mathcal{F}) = \mathbb{E}_{\left\{\vec{\mathbf{x}}_{i}\right\}_{i=1}^{n}, \left\{\sigma_{i}\right\}_{i=1}^{n}} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} f(\vec{\mathbf{x}}_{i})\right].$$
(19)

The empirical Rademacher complexity is defined as

704

705 706

707

708 709 710

711 712

713

741

751 752 753

754 755

$$\widehat{\mathfrak{R}}(\mathcal{F}) = \mathbb{E}_{\{\sigma_i\}_{i=1}^n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\vec{\mathbf{x}}_i) \right],$$
(20)

For simplicity of notations, Rademacher complexity and empirical Rademacher complexity are also denoted by $\mathbb{E}\left[\sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}f(\vec{\mathbf{x}}_{i})\right]$ and $\mathbb{E}_{\sigma}\left[\sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}f(\vec{\mathbf{x}}_{i})\right]$ respectively.

For data $\left\{\vec{\mathbf{x}}\right\}_{i=1}^{n}$ and a function class \mathcal{F} , we define the notation $R_n\mathcal{F}$ by $R_n\mathcal{F} := \sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i=1}^{n}\sigma_i f(\vec{\mathbf{x}}_i).$

Theorem A.1 ((Bartlett et al., 2005, Theorem 2.1)). Let \mathcal{X}, P be a probability space, $\left\{\vec{\mathbf{x}}_i\right\}_{i=1}^n$ be independent random variables distributed according to P. Let \mathcal{F} be a class of functions that map \mathcal{X} into [a, b]. Assume that there is some r > 0 such that for every $f \in \mathcal{F}, \operatorname{Var}\left[f(\vec{\mathbf{x}}_i)\right] \le r$. Then, for every x > 0, with probability at least $1 - e^{-x}$,

$$\sup_{f \in \mathcal{F}} \left(\mathbb{E}_P[f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_n}[f(\mathbf{x})] \right) \leq \inf_{\alpha > 0} \left(2(1+\alpha) \mathbb{E}_{\left\{ \overrightarrow{\mathbf{x}}_i \right\}_{i=1}^n, \left\{ \sigma_i \right\}_{i=1}^n} [R_n \mathcal{F}] + \sqrt{\frac{2rx}{n}} + (b-a) \left(\frac{1}{3} + \frac{1}{\alpha} \right) \frac{x}{n} \right),$$
(21)

and with probability at least $1 - 2e^{-x}$.

$$\sup_{f \in \mathcal{F}} \left(\mathbb{E}_P[f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim P_n}[f(\mathbf{x})] \right) \leq \inf_{\alpha \in (0,1)} \left(\frac{2(1+\alpha)}{1-\alpha} \mathbb{E}_{\{\sigma_i\}_{i=1}^n}[R_n \mathcal{F}] + \sqrt{\frac{2rx}{n}} + (b-a) \left(\frac{1}{3} + \frac{1}{\alpha} + \frac{1+\alpha}{2\alpha(1-\alpha)} \right) \frac{x}{n} \right).$$
(22)

733 P_n is the empirical distribution over $\left\{\vec{\mathbf{x}}_i\right\}_{i=1}^n$ with $\mathbb{E}_{\mathbf{x}\sim P_n}\left[f(\mathbf{x})\right] = \frac{1}{n}\sum_{i=1}^n f(\vec{\mathbf{x}}_i)$. Moreover, the 734 same results hold for $\sup_{f\in\mathcal{F}} \left(\mathbb{E}_{\mathbf{x}\sim P_n}[f(\mathbf{x})] - \mathbb{E}_P[f(\mathbf{x})]\right)$.

In addition, we have the contraction property for Rademacher complexity, which is due to Ledoux and Talagrand (Ledoux, 1991).
 The second second

Theorem A.2. Let ϕ be a contraction, that is, $|\phi(x) - \phi(y)| \le \mu |x - y|$ for $\mu > 0$. Then, for every function class \mathcal{F} ,

$$\mathbb{E}_{\{\sigma_i\}_{i=1}^n} \left[R_n \phi \circ \mathcal{F} \right] \le \mu \mathbb{E}_{\{\sigma_i\}_{i=1}^n} \left[R_n \mathcal{F} \right], \tag{23}$$

where $\phi \circ \mathcal{F}$ is the function class defined by $\phi \circ \mathcal{F} = \{\phi \circ f \colon f \in \mathcal{F}\}.$

743 744 745 *Definition* A.2 (Sub-root function,(Bartlett et al., 2005, Definition 3.1)). A function $\psi: [0, \infty) \rightarrow [0, \infty)$ is sub-root if it is nonnegative, nondecreasing and if $\frac{\psi(r)}{\sqrt{r}}$ is nonincreasing for r > 0.

Theorem A.3 ((Bartlett et al., 2005, Theorem 3.3)). Let \mathcal{F} be a class of functions with ranges in [*a*, *b*] and assume that there are some functional $T: \mathcal{F} \to \mathbb{R}+$ and some constant \overline{B} such that for every $f \in \mathcal{F}$, $\operatorname{Var}[f] \leq T(f) \leq \overline{BP}(f)$. Let ψ be a sub-root function and let r^* be the fixed point of ψ . Assume that ψ satisfies, for any $r \geq r^*$, $\psi(r) \geq \overline{BR}(\{f \in \mathcal{F}: T(f) \leq r\})$. Fix x > 0, then for any $K_0 > 1$, with probability at least $1 - e^{-x}$,

$$\forall f \in \mathcal{F}, \quad \mathbb{E}_{P}[f] \leq \frac{K_{0}}{K_{0}-1} \mathbb{E}_{P_{n}}[f] + \frac{704K_{0}}{\bar{B}}r^{*} + \frac{x\left(11(b-a) + 26\bar{B}K_{0}\right)}{n}$$

Also, with probability at least $1 - e^{-x}$,

$$\forall f \in \mathcal{F}, \quad \mathbb{E}_{P_n}\left[f\right] \leq \frac{K_0 + 1}{K_0} \mathbb{E}_P\left[f\right] + \frac{704K_0}{\bar{B}}r^* + \frac{x\left(11(b-a) + 26\bar{B}K_0\right)}{n}$$

Proposition A.4. Let \mathcal{F} be a class of functions with ranges in [0, b] for some positive constant b. Let ψ be a sub-root function such that for all $r \ge 0$, $\Re(\{f \in \mathcal{F} : \mathbb{E}_P [f(\mathbf{x})] \le r\}) \le \psi(r)$, and let r^* be the fixed point of ψ . Then for any $K_0 > 1$, with probability $1 - \exp(-x)$, every $f \in \mathcal{F}$ satisfies

$$\mathbb{E}_{P}\left[f\right] \le \frac{K_{0}}{K_{0}-1} \mathbb{E}_{P_{n}}\left[f\right] + \frac{704K_{0}}{b}r^{*} + \frac{x\left(11(b-a) + 26bK_{0}\right)}{n}.$$
(24)

B PROOFS FOR THEOREM 5.1 AND COROLLARY 5.2



Figure 1: Roadmap of major results leading to the main result, Theorem 5.1. The uniform convergence results in Theorem 6.1 are used in all the optimization results and Theorem C.8.

More discussion about the literature. We herein provide more discussion about the results of this work and comparison to the existing relevant works with sharp rates for nonparametric regression. While this paper establishes sharp rate which is distribution-free in spherical covariate, such rate still depends on bounded input space ($\mathcal{X} = \mathbb{S}^{d-1}$) and the condition that the target function $f^* \in$ $\mathcal{H}_K(\mu_0)$. Some other existing works consider target function f^* not belonging to the RKHS ball centered at the origin with constant or low radius, such as (Haas et al., 2023; Bordelon et al., 2024). We also note that in this work, only the first layer of an over-parameterized two-layer neural network is trained, while the weights of the second layer are randomly initialized and then fixed in the training process. In existing works such as (Hu et al., 2021; Suh et al., 2022; Allen-Zhu et al., 2019a), all the layers of a deep neural networks with more than two-layers are trained by GD or its variants. However, this work shows that only training the first layer still leads to sharp rate for nonparametric regression, which supports the claim in (Bietti & Bach, 2021) that a shallow over-parameterized neural networks with ReLU activations exhibit the same approximation properties as its deeper counterpart.

Proof of Theorem 5.1. We use Theorem C.10 and Theorem C.11 to prove this theorem.

First of all, it follows by Theorem C.11 that with probability at least $1 - \exp\left(-\Theta(n\hat{\varepsilon}_n^2)\right)$,

$$\mathbb{E}_{P_n}\left[(f_t - f^*)^2\right] \le \frac{3}{\eta t} \left(\frac{\mu_0^2}{2e} + 6\right)$$

Plugging such bound for $\mathbb{E}_{P_n} \left[(f_t - f^*)^2 \right]$ in (117) of Theorem C.10 leads to

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - \frac{6}{\eta t}\left(\frac{\mu_{0}^{2}}{2e} + 6\right) \le c_{0}'(\varepsilon_{n}^{2} + w).$$
(25)

Due to the definition of \widehat{T} and $\widehat{\varepsilon}_n^2$, we have

$$\hat{\varepsilon}_n^2 \le \frac{1}{\eta \hat{T}} \le \frac{2}{\eta (\hat{T}+1)} \le 2\hat{\varepsilon}_n^2.$$
(26)

⁸¹⁰ Lemma C.15 suggests that with probability at least $1 - 4\exp(-\Theta(n\varepsilon_n^2))$ over **S**, $\hat{\varepsilon}_n^2 \simeq \varepsilon_n^2$. Since $T \simeq \hat{T}$, for any $t \in [c_t T, T]$, we have

$$\frac{1}{\eta t} \asymp \frac{1}{\eta T} \asymp \frac{1}{\eta \widehat{T}} \asymp \widehat{\varepsilon}_n^2 \asymp \varepsilon_n^2.$$
(27)

We have $\Pr[\mathcal{W}_0] \ge 1 - 2/n$. Let $w = \varepsilon_n^2$, we now verify that $w \in (0, 1)$. Due to the definition of the fixed point, w > 0. Since $\sum_{i \ge 1} \lambda_i = \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{x}) d\mu(\mathbf{x}) = 1/2$, we have

$$0 < w = \frac{1}{n} \sum_{i \ge 1} \min\left\{\lambda_i, \varepsilon_n^2\right\} \le \frac{1}{n} \sum_{i \ge 1} \lambda_i \le \frac{1}{2n} < 1.$$

(11) then follows from (25) with $w = \varepsilon_n^2$, (27) and the union bound. The condition on m in (85) in Theorem C.10, together with $w = \varepsilon_n^2$ and (27) leads to the condition on m in (10). Furthermore, $\hat{T} \simeq \varepsilon_n^{-2}$ follows from (27) and $\eta = \Theta(1)$.

Proof of Corollary 5.2. We apply Theorem 5.1 to prove this corollary.

It is well known, such as (Raskutti et al., 2014, Corollary 3), that $\varepsilon_n^2 \approx n^{-\frac{2\alpha}{2\alpha+1}}$. It then can be verified by direct calculations that the condition on m, (10) in Theorem 5.1, is satisfied with the given condition (12). It then follows from (11) in Theorem 5.1 that $\mathbb{E}_P\left[(f_{\widehat{T}} - f^*)^2\right] \lesssim n^{-\frac{2\alpha}{2\alpha+1}}$.

C DETAILED PROOFS

Because Theorem 5.1 is proved by Theorem C.10 and Theorem C.11, in this section, we establish and prove all the theoretical results which lead to Theorem C.10 and Theorem C.11, along with the proof of Theorem C.10 and Theorem C.11.

C.1 BASIC DEFINITIONS

We introduce the following definitions for the proof of Theorem 5.2. We define

$$\mathbf{u}(t) \coloneqq \widehat{\mathbf{y}}(t) - \mathbf{y}.$$
(28)

Let $\tau \leq 1$ be a positive number, and $\varepsilon_0 \in (0, 1)$ is an arbitrary positive constant. For $t \geq 0$ and $T \geq 1$ we define the following quantities (or recall their definitions if defined before),

$$c_{\mathbf{u}} = \mu_0 / \min\left\{2, \sqrt{2e\eta}\right\} + \sigma + \tau + 1,$$

$$R = \frac{\eta c_{\mathbf{u}} T}{\sqrt{m}},$$
(29)

$$\mathcal{V}_t \coloneqq \left\{ \mathbf{v} \in \mathbb{R}^n \colon \mathbf{v} = -\left(\mathbf{I}_n - \eta \mathbf{K}_n\right)^t f^*(\mathbf{S}) \right\},\tag{30}$$

$$\mathcal{E}_{t,\tau} \coloneqq \left\{ \mathbf{e} \colon \mathbf{e} = \vec{\mathbf{e}}_1 + \vec{\mathbf{e}}_2 \in \mathbb{R}^n, \vec{\mathbf{e}}_1 = -\left(\mathbf{I}_n - \eta \mathbf{K}_n\right)^t \mathbf{w}, \left\|\vec{\mathbf{e}}_2\right\|_2 \le \sqrt{n\tau} \right\}.$$
 (31)

We define the set of neural network weights and the set of functions represented by the neural network during training as follows.

$$\mathcal{W}(\mathbf{S}, \mathbf{W}(0), T) \coloneqq \left\{ \mathbf{W} \colon \exists t \in [T] \text{ s.t. } \operatorname{vec}\left(\mathbf{W}\right) = \operatorname{vec}\left(\mathbf{W}(0)\right) - \sum_{t'=0}^{t-1} \frac{\eta}{n} \mathbf{Z}_{\mathbf{S}}(t') \mathbf{u}(t'), \right\}$$

$$\mathbf{u}(t') \in \mathbb{R}^{n}, \mathbf{u}(t') = \mathbf{v}(t') + \mathbf{e}(t'), \mathbf{v}(t') \in \mathcal{V}_{t'}, \mathbf{e}(t') \in \mathcal{E}_{t',\tau}, \text{ for all } t' \in [0, t-1] \right\}.$$
 (32)

867 $\mathcal{W}(\mathbf{S}, \mathbf{W}(0), T)$ is the set of weights of neural networks trained by GD on the training data **S** and random initialization $\mathbf{W}(0)$ with the preconditioner **M** generated by **Q** and the steps of GD no greater than *T*. The set of functions represented by the two-layer NN with weights in $\mathcal{W}(\mathbf{S}, \mathbf{W}(0), T)$ is then defined as

$$\mathcal{F}_{\mathrm{NN}}(\mathbf{S}, \mathbf{W}(0), T) \coloneqq \{ f_t = f(\mathbf{W}(t), \cdot) \colon \exists t \in [T], \mathbf{W}(t) \in \mathcal{W}(\mathbf{S}, \mathbf{W}(0), T) \} .$$
(33)

We define the function class $\mathcal{F}_{ext}(w,T)$ for any w > 0 as

$$\mathcal{F}_{\text{ext}}(w,T) \coloneqq \left\{ f \colon f = h + e, h \in \mathcal{H}_K(B_h), \|e\|_{\infty} \le w \right\},\tag{34}$$

where

$$B_h \coloneqq \mu_0 + 1 + \sqrt{2}. \tag{35}$$

C.2 THEOREM C.10, THEOREM C.11, AND THEIR PROOFS WITH RELATED THEORETICAL RESULTS

Theorem C.10 (repeat). Suppose $w \in (0, 1)$ and m satisfy

$$m\gtrsim \max\left\{\frac{(\eta T)^4\left(\sqrt{d}+1\right)^4}{w^4},(\eta T)^8d^2\right\},$$

and the neural network $f(\mathbf{W}(t), \cdot)$ is trained by GD in Algorithm 1 with the learning rate $\eta \in (0, 1/\hat{\lambda}_1)$ on random initialization $\mathbf{W}(0)$, and $T \leq \hat{T}$. Then for every $t \in [T]$, with probability at least $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\hat{\varepsilon}_n^2)) - \exp(-n\varepsilon_n^2) - 2/n$ over the random noise w, the random training features S and the random initialization $\mathbf{W}(0)$,

$$\mathbb{E}_{P}\left[(f_{t}-f^{*})^{2}\right]-2\mathbb{E}_{P_{n}}\left[(f_{t}-f^{*})^{2}\right]$$

$$\leq c_{0}\min_{0\leq Q\leq n}\left(\frac{B_{0}Q}{n}+w\left(\sqrt{\frac{Q}{n}}+1\right)+B_{h}\left(\frac{\sum\limits_{q=Q+1}^{\infty}\lambda_{q}}{n}\right)^{1/2}\right)^{2},$$

> Furthermore, with probability at least $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\hat{\varepsilon}_n^2)) - \exp(-n\varepsilon_n^2) - 2/n$ over the random noise w, the random training features S and the random initialization W(0),

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - 2\mathbb{E}_{P_{n}}\left[(f_{t} - f^{*})^{2}\right] \le c_{0}'(\varepsilon_{n}^{2} + w).$$
(36)

903 Here B_0, c_0, c'_0 are absolute positive constants depending on μ_0 , and c'_0 also depends on σ . **Theorem C.11 (repeat)**. Suppose the neural network trained after the *t*-th step of gradient descent, $f_t = f(\mathbf{W}(t), \cdot)$, satisfies $\mathbf{u}(t) = f_t(\mathbf{S}) - \mathbf{y} = \mathbf{v}(t) + \mathbf{e}(t)$ with $\mathbf{v}(t) \in \mathcal{V}_t$ and $\mathbf{e}(t) \in \mathcal{E}_{t,\tau}$ and $T \leq \widehat{T}$. If

$$\eta \in [1,2), \quad \tau \le \frac{1}{\eta T},$$

then for every $t \in [T]$, with probability at least $1 - \exp\left(-\Theta(n\hat{\varepsilon}_n^2)\right)$ over the random noise w, we have

$$\mathbb{E}_{P_n}\left[(f_t - f^*)^2\right] \le \frac{3}{\eta t} \left(\frac{\mu_0^2}{2e} + 3\right).$$

We have the following two theorems regarding the uniform convergence of $\hat{h}(\mathbf{W}(0), \cdot, \cdot)$ to $K(\cdot, \cdot)$ and the uniform convergence of $\hat{v}_R(\mathbf{W}(0), \cdot)$ to $\frac{2R}{\sqrt{2\pi\kappa}}$, which lay the foundation of the main results of this paper. The proofs are deferred to Section C.4. **Theorem C.1.** Let $\mathbf{W}(0) = \left\{ \vec{\mathbf{w}}_r(0) \right\}_{r=1}^m$, where each $\vec{\mathbf{w}}_r(0) \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)$ for $r \in [m]$. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over $\mathbf{W}(0)$,

$$\sup_{\mathbf{x}\in\mathcal{X},\mathbf{y}\in\mathcal{X}} \left| K(\mathbf{x},\mathbf{y}) - \hat{h}(\mathbf{W}(0),\mathbf{x},\mathbf{y}) \right| \le C_1(m,d,\delta),$$
(37)

where

$$C_1(m,d,\delta) \coloneqq \frac{1}{\sqrt{m}} \left(6(1+B\sqrt{d}) + \sqrt{2\log\frac{2(1+2m)^d}{\delta}} \right) + \frac{1}{m} \left(3 + \frac{7\log\frac{2(1+2m)^d}{\delta}}{3} \right), \quad (38)$$

and B is an absolute positive constant in Lemma C.19. In addition, when $m \ge \max\{d, n, 4\}$, $m/\log m \geq d$, and $\delta \asymp 1/n$,

$$C_1(m, d, \delta) \lesssim \sqrt{\frac{d\log m}{m}} + \frac{d\log m}{m} \lesssim \sqrt{\frac{d\log m}{m}}.$$
 (39)

Theorem C.2. Let $\mathbf{W}(0) = \left\{ \vec{\mathbf{w}}_r(0) \right\}_{r=1}^m$, where each $\vec{\mathbf{w}}_r(0) \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)$ for $r \in [m]$. Suppose $R \leq R_0$ for an arbitrary absolute positive constant $R_0 < \kappa$. B is an absolute positive constant in Lemma C.19. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over $\mathbf{W}(0)$,

> $\sup_{\mathbf{x}\in\mathcal{X}} \left| \widehat{v}_R(\mathbf{W}(0), \mathbf{x}) - \frac{2R}{\sqrt{2\pi\kappa}} \right| \le C_2(m, R_0, \delta),$ (40)

where

$$C_{2}(m, R_{0}, \delta) \coloneqq 3 \left((B\sqrt{d} + 1)\sqrt{m^{-\frac{1}{2}} + \frac{1}{m}} + \frac{\exp\left(-\frac{(\kappa^{2} - R_{0}^{2})^{2}}{4\kappa^{4}}m\right)}{\sqrt{m^{-\frac{1}{2}} + \frac{1}{m}}} \right) + \sqrt{\frac{2\log\frac{2(1+2m)^{d}}{\delta}}{m}} + \frac{7\log\frac{2(1+2m)^{d}}{\delta}}{3m}.$$
(41)

In addition, when $m \ge \max{\{d, n, 4\}}$, $m/\log m \ge d$, and $\delta \asymp 1/n$,

$$C_2(m,R_0,\delta) \lesssim \frac{\sqrt{d}}{m^{1/4}} + \sqrt{\frac{d\log m}{m}} + \frac{d\log m}{m} \lesssim \frac{\sqrt{d}}{m^{1/4}}$$

Lemma C.3. Suppose

$$m \gtrsim (\eta T)^4 (\sqrt{d} + \sqrt{\tau})^4 / \tau^4.$$
(42)

and the neural network $f(\mathbf{W}(t), \cdot)$ trained by gradient decent with the learning rate $\eta \in (0, 1/\hat{\lambda}_1)$ on random initialization $\mathbf{W}(0) \in \mathcal{W}_0$. Then with probability at least $1 - \exp\left(-\Theta(n)\right)$ over the random noise w, $\mathbf{W}(t) \in \mathcal{W}(\mathbf{S}, \mathbf{W}(0), T)$. Moreover, for all $t \in [0, T]$, $\mathbf{u}(t) = \mathbf{v}(t) + \mathbf{e}(t)$ where $\mathbf{u}(t) = \widehat{\mathbf{y}}(t) - \mathbf{y}, \mathbf{v}(t) \in \mathcal{V}_{K,t}, \mathbf{e}(t) \in \mathcal{E}_{K,t,\tau}, \text{ and } \|\mathbf{u}(t)\|_2 \leq c_{K,\mathbf{u}}\sqrt{n}.$

Proof. First, when $m \gtrsim (\eta T)^4 (\sqrt{d} + \sqrt{\tau})^4 / \tau^4$ with a proper constant, it can be verified that $\mathbf{E}_{m,\eta,\tau} \leq \tau \sqrt{n}/T$ where $\mathbf{E}_{m,\eta,\tau}$ is defined by (52) of Lemma C.5. Also, Theorem C.1 and Theo-rem C.2 hold when (42) holds. We then use mathematical induction to prove the lemma. We will first prove that $\mathbf{u}(t) = \mathbf{v}(t) + \mathbf{e}(t)$ where $\mathbf{v}(t) \in \mathcal{V}_t$, $\mathbf{e}(t) \in \mathcal{E}_{t,\tau}$, and $\|\mathbf{u}(t)\|_2 \leq c_{\mathbf{u}}\sqrt{n}$ for for all $t \in [0, T].$

966 for all
$$t \in [0, T]$$
, where $\sum_{t'=1}^{t} \cdot = 0$ for $t < 1$, and

When t = 0, we have

$$\mathbf{u}(0) = -\mathbf{y} = \mathbf{v}(0) + \mathbf{e}(0),\tag{43}$$

where $\mathbf{v}(0) \coloneqq -f^*(\mathbf{S}) = -(\mathbf{I} - \eta \mathbf{K}_n)^0 f^*(\mathbf{S}), \ \mathbf{e}(0) = -\mathbf{w} = \mathbf{e}_1(0) + \mathbf{e}_2(0) \ \text{with} \ \mathbf{e}_1(0) = -\mathbf{w} = \mathbf{e}_1(0) + \mathbf{e}_2(0) \ \mathbf{e}_1(0) = -\mathbf{e}_1(0) + \mathbf{e}_1(0) + \mathbf{e}_1(0) = -\mathbf{e}_1(0) + \mathbf{e}_1(0) + \mathbf{e}_1(0) = -\mathbf{e}_1(0) + \mathbf{e}_1(0) + \mathbf{e}_1(0) + \mathbf{e}_1(0) + \mathbf{e}_1(0) + \mathbf{e}_1(0) = -\mathbf{e}_1(0) + \mathbf{e}_1(0) + \mathbf{e}_1(0)$ $-(\mathbf{I} - \eta \mathbf{K}_n)^0 \mathbf{w}$ and $\mathbf{e}_2(0) = \mathbf{0}$. Therefore, $\mathbf{v}(0) \in \mathcal{V}_0$ and $\mathbf{e}(0) \in \mathcal{E}_{0,\tau}$. Also, it follows from the proof of Lemma C.4 that $\|\mathbf{u}(0)\|_2 \le c_{\mathbf{u}}$ with probability at least $1 - \exp(-\Theta(n))$ over the random noise w.

Suppose that for all $t_1 \in [0, t]$ with $t \in [0, T-1]$, $\mathbf{u}(t_1) = \mathbf{v}(t_1) + \mathbf{e}(t_1)$ where $\mathbf{v}(t_1) \in \mathcal{V}_{t_1}$, and $\mathbf{e}(t_1) = \overline{\mathbf{e}}_1(t_1) + \overline{\mathbf{e}}_2(t_1)$ with $\mathbf{v}(t_1) \in \mathcal{V}_{t_1}$ and $\mathbf{e}(t_1) \in \mathcal{E}_{t_1,\tau}$, and $\|\mathbf{u}(t_1)\|_2 \leq c_{\mathbf{u}}\sqrt{n}$ for all $t_1 \in [0, t]$

Then it follows from Lemma C.5 that the recursion $\mathbf{u}(t'+1) = (\mathbf{I} - \eta \mathbf{K}_n) \mathbf{u}(t') + \mathbf{E}(t'+1)$ holds for all $t' \in [0, t]$. As a result, we have

$$\mathbf{u}(t+1) = (\mathbf{I} - \eta \mathbf{K}_n) \mathbf{u}(t) + \mathbf{E}(t+1)$$

= $-(\mathbf{I} - \eta \mathbf{K}_n)^{t+1} f^*(\mathbf{S}) - (\mathbf{I} - \eta \mathbf{K}_n)^t \mathbf{w}$
+ $\sum_{t'=1}^t (\mathbf{I} - \eta \mathbf{K}_n)^{t-t'} \mathbf{E}(t')$
= $\mathbf{v}(t+1) + \mathbf{e}(t+1),$ (44)

where $\mathbf{v}(t+1)$ and $\mathbf{e}(t+1)$ are defined as

$$\mathbf{v}(t+1) \coloneqq -\left(\mathbf{I} - \eta \mathbf{K}_n\right)^{t+1} f^*(\mathbf{S}) \in \mathcal{V}_{t+1},\tag{45}$$

$$\mathbf{e}(t+1) \coloneqq \underbrace{-\left(\mathbf{I} - \eta \mathbf{K}_n\right)^{t+1} \mathbf{w}}_{\overrightarrow{\mathbf{e}}_1(t+1)} + \underbrace{\sum_{t'=1}^{t+1} \left(\mathbf{I} - \eta \mathbf{K}_n\right)^{t+1-t'} \mathbf{E}(t')}_{\overrightarrow{\mathbf{e}}_2(t+1)}.$$
(46)

We now prove the upper bound for $\vec{\mathbf{e}}_2(t+1)$. With $\eta \in (0, 1/\widehat{\lambda}_1)$, we have $\|\mathbf{I} - \eta \mathbf{K}_n\|_2 \in (0, 1)$. It follows that

$$\left\| \overrightarrow{\mathbf{e}}_{2}(t+1) \right\|_{2}$$

$$\leq \sum_{t'=1}^{t+1} \| \mathbf{I} - \eta \mathbf{K}_{n} \|_{2}^{t+1-t'} \| \mathbf{E}(t') \|_{2}$$

$$\leq \tau \sqrt{n}, \qquad (47)$$

where the last inequality follows from the fact that $\|\mathbf{E}(t)\|_2 \leq \mathbf{E}_{m,\eta,\tau} \leq \tau \sqrt{n}/T$ for all $t \in [T]$ and the induction hypothesis. It follows that $\mathbf{e}(t+1) \in \mathcal{E}_{t+1,\tau}$. Also, it follows from Lemma C.4 that

$$\begin{split} \|\mathbf{u}(t+1)\|_2 &\leq \|\mathbf{v}(t+1)\|_2 + \left\|\vec{\mathbf{e}}_1(t+1)\right\|_2 + \left\|\vec{\mathbf{e}}_2(t+1)\right\|_2 \\ &\leq \left(\frac{\mu_0}{\sqrt{2e\eta}} + \sigma + \tau + 1\right)\sqrt{n} = c_{\mathbf{u}}\sqrt{n}, \end{split}$$

This fact completes the induction step, which also completes the proof.

Lemma C.4. Let $t \in [T]$, $\mathbf{v} = -(\mathbf{I} - \eta \mathbf{K}_n)^t f^*(\mathbf{S})$, $\mathbf{e} = -(\mathbf{I} - \eta \mathbf{K}_n)^{t+1} \mathbf{w}$, and $\eta \in (0, 1/\widehat{\lambda}_1)$. Then with probability at least $1 - \exp(-\Theta(n))$ over the random noise \mathbf{w} ,

$$\|\mathbf{v}\|_{2} + \|\mathbf{e}\|_{2} \le \left(\frac{\mu_{0}}{\sqrt{2e\eta}} + \sigma + 1\right)\sqrt{n} \tag{48}$$

Proof. When $\mathbf{v} \in \mathcal{V}_t$ for $t \ge 1$, we have $\mathbf{v} = -(\mathbf{I} - \eta \mathbf{K}_n)^t f^*(\mathbf{S})$, and

1022
1023
1023
$$\|\mathbf{v}(t)\|_{2}^{2} = \sum_{i=1}^{n} \left(1 - \eta \widehat{\lambda}_{i}\right)^{2t} \left[\mathbf{U}^{\top} f^{*}(\mathbf{S})\right]_{i}^{2}$$

1024
1025
$$\stackrel{\text{(I)}}{=} \sum_{i=1}^{n} \frac{1}{2e\eta \widehat{\lambda}_{i} t} \left[\mathbf{U}^{\top} f^{*}(\mathbf{S}) \right]_{i}^{2}$$

$$\stackrel{\textcircled{0}}{\leq} \frac{n\mu_0^2}{2e\eta t}.$$
(49)

1029 Here ① follows Lemma C.13, ② follows by Lemma C.12.

Moreover, it follows from the concentration inequality about quadratic forms of sub-Gaussian ran dom variables in (Wright, 1973) that

$$\Pr\left[\left\|\mathbf{w}\right\|_{2}^{2} - \mathbb{E}\left[\left\|\mathbf{w}\right\|_{2}^{2}\right] > n\right] \le \exp\left(-\Theta(n)\right),\tag{50}$$

and
$$\mathbb{E}\left[\|\mathbf{w}\|_{2}\right] \leq \sqrt{\mathbb{E}\left[\|\mathbf{w}\|_{2}^{2}\right]} = \sqrt{n\sigma}$$
. Therefore, $\Pr\left[\|\mathbf{w}\|_{2} - \sqrt{n\sigma} > \sqrt{n}\right] \leq \exp\left(-\Theta(n)\right)$.

As a result, we have

$$\|\mathbf{v}\|_{2} + \|\mathbf{e}\|_{2} \leq \sqrt{\frac{n\mu_{0}^{2}}{2e\eta}} + \|\mathbf{w}\|_{2} \leq \left(\frac{\mu_{0}}{\sqrt{2e\eta}} + \sigma + 1\right)\sqrt{n}.$$

1043 Lemma C.5. Let $0 < \eta < 1$, $0 \le t \le T - 1$ for $T \ge 1$, and suppose that $\|\widehat{\mathbf{y}}(t') - \mathbf{y}\|_2 \le c_{\mathbf{u}}\sqrt{n}$ 1044 holds for all $0 \le t' \le t$. Then

$$\widehat{\mathbf{y}}(t+1) - \mathbf{y} = (\mathbf{I} - \eta \mathbf{K}_n) \left(\widehat{\mathbf{y}}(t) - \mathbf{y}\right) + \mathbf{E}(t+1),$$
(51)

1046 where $\|\mathbf{E}(t+1)\|_2 \leq \mathbf{E}_{m,\eta,\tau}$, and $\mathbf{E}_{m,\eta,\tau}$ is defined by

$$\mathbf{E}_{m,\eta,\tau} \coloneqq \eta c_{\mathbf{u}} \sqrt{n} \left(3 \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_2(m/2, R_0, \delta) \right) + C_1(m/2, d, \delta) \right)$$
$$\lesssim \eta \sqrt{n} \left(\frac{\sqrt{d}}{m^{1/4}} + \frac{\eta T}{\sqrt{m}} \right). \tag{52}$$

Proof. Because $\|\widehat{\mathbf{y}}(t') - \mathbf{y}\|_2 \le \sqrt{n}c_{\mathbf{u}}$ holds for all $t' \in [0, t]$, by Lemma C.6, we have

$$\vec{\mathbf{w}}_r(t') - \vec{\mathbf{w}}_r(0) \Big\|_2 \le R, \quad \forall \, 0 \le t' \le t+1.$$
(53)

1057 Define two sets of indices 1058

$$E_{i,R} \coloneqq \left\{ r \in [m] \colon \left| \mathbf{w}_r(0)^\top \vec{\mathbf{x}}_i \right| > R \right\}, \quad \bar{E}_{i,R} \coloneqq [m] \setminus E_{i,R}.$$

We have

$$\begin{aligned} \widehat{\mathbf{y}}_{i}(t+1) - \widehat{\mathbf{y}}_{i}(t) &= \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_{r} \left(\sigma \left(\overrightarrow{\mathbf{w}}_{\mathbf{S},r}^{\top}(t+1) \overrightarrow{\mathbf{x}}_{i} \right) - \sigma \left(\overrightarrow{\mathbf{w}}_{\mathbf{S},r}^{\top}(t) \overrightarrow{\mathbf{x}}_{i} \right) \right) \\ &= \underbrace{\frac{1}{\sqrt{m}} \sum_{r \in E_{i,R}} a_{r} \left(\sigma \left(\overrightarrow{\mathbf{w}}_{\mathbf{S},r}^{\top}(t+1) \overrightarrow{\mathbf{x}}_{i} \right) - \sigma \left(\overrightarrow{\mathbf{w}}_{\mathbf{S},r}^{\top}(t) \overrightarrow{\mathbf{x}}_{i} \right) \right)}_{:=\mathbf{D}_{i}^{(1)}} \\ &+ \underbrace{\frac{1}{\sqrt{m}} \sum_{r \in \bar{E}_{i,R}} a_{r} \left(\sigma \left(\overrightarrow{\mathbf{w}}_{\mathbf{S},r}^{\top}(t+1) \overrightarrow{\mathbf{x}}_{i} \right) - \sigma \left(\overrightarrow{\mathbf{w}}_{\mathbf{S},r}^{\top}(t) \overrightarrow{\mathbf{x}}_{i} \right) \right)}_{:=\mathbf{E}_{i}^{(1)}} \\ &= \mathbf{D}_{i}^{(1)} + \mathbf{E}_{i}^{(1)}, \end{aligned}$$
(54)

and $\mathbf{D}^{(1)}, \mathbf{E}^{(1)} \in \mathbb{R}^n$ is a vector with their *i*-th element being $\mathbf{D}_i^{(1)}$ and $\mathbf{E}_i^{(1)}$ defined on the RHS of (54). Now we derive the upper bound for $\mathbf{E}_i^{(1)}$. For all $i \in [n]$ we have

1078
1079
$$\left| \mathbf{E}_{i}^{(1)} \right| = \left| \frac{1}{\sqrt{m}} \sum_{r \in \bar{E}_{i,R}} a_{r} \left(\sigma \left(\vec{\mathbf{w}}_{\mathbf{S},r}(t+1)^{\top} \vec{\mathbf{x}}_{i} \right) - \sigma \left(\vec{\mathbf{w}}_{\mathbf{S},r}(t)^{\top} \vec{\mathbf{x}}_{i} \right) \right) \right|$$

$$\begin{aligned} & \left| \begin{array}{l} 1080 \\ 1081 \\ 1081 \\ 1082 \\ 1082 \\ 1083 \\ 1084 \\ 1085 \\ 1086 \\ 1086 \\ 1086 \\ 1086 \\ 1086 \\ 1086 \\ 1088 \\ 1089 \\ 1089 \\ 1090 \\ 1090 \\ 1091 \\ 1092 \\ 1092 \\ 1093 \\ \end{aligned} \right| \left| \begin{array}{l} \frac{1}{\sqrt{m}} \sum_{r \in \bar{E}_{i,R}} \left\| \frac{\eta}{n} \left[\mathbf{Z}_{\mathbf{S}}(t) \right]_{[(r-1)d+1:rd]} \left(\hat{\mathbf{y}}(t) - \mathbf{y} \right) \right\|_{2} \\ \frac{\mathcal{Q}}{\sqrt{m}} \sum_{r \in \bar{E}_{i,R}} \frac{\eta}{\sqrt{m}} \\ \leq \eta c_{\mathbf{u}} \cdot \frac{\left| \bar{E}_{i,R} \right|}{m} . \end{aligned}$$

$$(55)$$

Here (1, 2) follow from (72) and (73) in the proof of Lemma C.6.

Let m be sufficiently large such that $R \leq R_0$ for the absolute positive constant $R_0 < \kappa$ specified in Theorem 6.1. Then it follows from Theorem C.2 that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over $\mathbf{W}(0)$,

$$\sup_{\mathbf{x}\in\mathcal{X}} \left| \widehat{v}_R(\mathbf{W}(0), \mathbf{x}) - \frac{2R}{\sqrt{2\pi\kappa}} \right| \le C_2(m/2, R_0, \delta),$$
(56)

where $\widehat{v}_R(\mathbf{W}(0), \mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \mathbb{I}_{\{\left| \vec{\mathbf{w}}_r(0)^\top \mathbf{x} \right| \le R\}}$, so that $\widehat{v}_R(\mathbf{W}(0), \vec{\mathbf{x}}_i) = \left| \overline{E}_{i,R} \right| / m$. It follows from (55), (56) and the induction hypothesis that

 $\left|\mathbf{E}_{i}^{(1)}\right| \leq \eta c_{\mathbf{u}} \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta)\right).$ (57)

It follows from (57) that $\left\|\mathbf{E}^{(1)}\right\|_2$ can be bounded by

$$\left\|\mathbf{E}^{(1)}\right\|_{2} \leq \eta c_{\mathbf{u}} \sqrt{n} \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta)\right).$$
(58)

 $\mathbf{D}_{i}^{(1)}$ on the RHS of (54) is expressed by

$$\begin{split} \mathbf{D}_{i}^{(1)} &= \frac{1}{\sqrt{m}} \sum_{r \in E_{i,R}} a_{r} \left(\sigma \left(\vec{\mathbf{w}}_{\mathbf{S},r}^{\top}(t+1) \vec{\mathbf{x}}_{i} \right) - \sigma \left(\vec{\mathbf{w}}_{\mathbf{S},r}^{\top}(t) \vec{\mathbf{x}}_{i} \right) \right) \\ &= \frac{1}{\sqrt{m}} \sum_{r \in E_{i,R}} a_{r} \mathbb{I}_{\left\{ \vec{\mathbf{w}}_{\mathbf{S},r}(t)^{\top} \vec{\mathbf{x}}_{i} \ge 0 \right\}} \left(\vec{\mathbf{w}}_{\mathbf{S},r}(t+1) - \vec{\mathbf{w}}_{\mathbf{S},r}(t) \right)^{\top} \vec{\mathbf{x}}_{i} \\ &= \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_{r} \mathbb{I}_{\left\{ \vec{\mathbf{w}}_{\mathbf{S},r}(t)^{\top} \vec{\mathbf{x}}_{i} \ge 0 \right\}} \left(-\frac{\eta}{n} \left[\mathbf{Z}_{\mathbf{S}}(t) \right]_{\left[(r-1)d:rd\right]} \left(\hat{\mathbf{y}}(t) - \mathbf{y} \right) \right)^{\top} \vec{\mathbf{x}}_{i} \end{split}$$

1124
1125
1126
1127
1127
1128

$$+\frac{1}{\sqrt{m}}\sum_{r\in\bar{E}_{i,R}}a_{r}\mathbb{I}_{\left\{\overrightarrow{\mathbf{w}}_{\mathbf{S},r}(t)^{\top}\overrightarrow{\mathbf{x}}_{i}\geq0\right\}}\left(\frac{\eta}{n}\left[\mathbf{Z}_{\mathbf{S}}(t)\right]_{\left[(r-1)d:rd\right]}\left(\widehat{\mathbf{y}}(t)-\mathbf{y}\right)\right)^{\top}\overrightarrow{\mathbf{x}}_{i}$$
1127

1127
1128
$$= -\frac{\eta}{n} \left[\mathbf{H}(t) \right]_i \left(\widehat{\mathbf{y}}(t) - \mathbf{y} \right)$$

1132
1133
$$\underbrace{\sqrt{m} \sum_{r \in \bar{E}_{i,R}} \{\mathbf{w}_{\mathbf{S},r}(t) \mid \mathbf{x}_i \ge 0\} \setminus n^{1-|\mathbf{S}| < r \leq 1}_{i=\mathbf{E}_i^{(2)}}}_{:=\mathbf{E}_i^{(2)}}$$

1136 where $\mathbf{H}(t) \in \mathbb{R}^{n \times n}$ is a matrix specified by

$$\mathbf{H}_{pq}(t) = \frac{\mathbf{\vec{x}}_{p}^{\top} \mathbf{\vec{x}}_{q}}{m} \sum_{r=1}^{m} \mathbb{I}_{\left\{\mathbf{\vec{w}}_{\mathbf{S},r}(t)^{\top} \mathbf{\vec{x}}_{p} \ge 0\right\}} \mathbb{I}_{\left\{\mathbf{\vec{w}}_{r}(0)^{\top} \mathbf{\vec{x}}_{q} \ge 0\right\}}, \quad \forall p \in [n], q \in [n].$$

1142 Let $\mathbf{D}^{(2)}, \mathbf{E}^{(2)} \in \mathbb{R}^n$ be a vector with their *i*-the element being $\mathbf{D}_i^{(2)}$ and $\mathbf{E}_i^{(2)}$ defined on the RHS 1143 of (59). $\mathbf{E}^{(2)}$ can be expressed by $\mathbf{E}^{(2)} = \frac{\eta}{n} \tilde{\mathbf{E}}^{(2)} (\hat{\mathbf{y}}(t) - \mathbf{y})$ with $\tilde{\mathbf{E}}^{(2)} \in \mathbb{R}^{n \times n}$ and

$$\tilde{\mathbf{E}}_{pq}^{(2)} = \frac{1}{m} \sum_{r \in \bar{E}_{i,R}} \mathrm{I\!I}_{\left\{\vec{\mathbf{w}}_{\mathbf{S},r}(t)^{\top} \vec{\mathbf{x}}_{p} \ge 0\right\}} \mathrm{I\!I}_{\left\{\vec{\mathbf{w}}_{r}(0)^{\top} \vec{\mathbf{q}}_{q} \ge 0\right\}} \vec{\mathbf{x}}_{q}^{\top} \vec{\mathbf{x}}_{p} \le \frac{1}{m} \sum_{r \in \bar{E}_{i,R}} 1 = \frac{\left|E_{i,R}\right|}{m}$$

for all $p \in [n], q \in [n]$. The spectral norm of $\tilde{\mathbf{E}}^{(2)}$ is bounded by

$$\left\|\tilde{\mathbf{E}}^{(2)}\right\|_{2} \leq \left\|\tilde{\mathbf{E}}^{(2)}\right\|_{\mathrm{F}} \leq n \frac{\left|\bar{E}_{i,R}\right|}{m} \stackrel{\textcircled{1}{=}}{\leq} n \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta)\right),\tag{60}$$

where ① follows from (56). Also, $\|\mathbf{H}(t)\|_2 \le \|\mathbf{H}(t)\|_F \le \sqrt{nN}$ for all $t \ge 0$. It follows from (60) that $\|\mathbf{E}^{(2)}\|_2$ can be bounded by

$$\left\| \mathbf{E}^{(2)} \right\|_{2} \leq \frac{\eta}{n} \left\| \tilde{\mathbf{E}}^{(2)} \right\|_{2} \left\| \mathbf{y}(t) - \mathbf{y} \right\|_{2}$$
$$\leq \eta c_{\mathbf{u}} \sqrt{n} \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta) \right).$$
(61)

1161 $\mathbf{D}_i^{(2)}$ on the RHS of (59) is expressed by

$$\mathbf{D}^{(2)} = -\frac{\eta}{n} \mathbf{H}(t) \left(\widehat{\mathbf{y}}(t) - \mathbf{y} \right)$$

1164
1165
1166

$$= \underbrace{-\frac{\eta}{n} \mathbf{K} \left(\widehat{\mathbf{y}}(t) - \mathbf{y} \right)}_{\mathbf{p}_{1}^{(3)}}$$

1167
1168
$$+ \frac{\eta}{n} \left(\mathbf{K} - \mathbf{H}(0) \right) \left(\widehat{\mathbf{y}}(t) - \mathbf{y} \right)$$

1169
$$:=\mathbf{E}^{(3)}$$

$$+\frac{\eta}{n} \left(\mathbf{H}(0) - \mathbf{H}(t) \right) \left(\widehat{\mathbf{y}}(t) - \mathbf{y} \right)$$

1172
1173

$$= \mathbf{D}^{(3)} + \mathbf{E}^{(4)}$$
.

1175 On the RHS of (62), $\mathbf{D}^{(3)}, \mathbf{E}^{(3)}, \mathbf{E}^{(4)} \in \mathbb{R}^n$ are vectors which are analyzed as follows. $\|\tilde{\mathbf{E}}^{(3)}\|_2$ is bounded by

$$\|\mathbf{K} - \mathbf{H}(0)\|_{2} \le \|\mathbf{K} - \mathbf{H}(0)\|_{F} \le nC_{1}(m/2, d, \delta),$$
(63)

(62)

where the last inequality holds with probability $1 - \delta$ over $\mathbf{W}(0)$ according to Theorem C.1.

1181 In order to bound $\mathbf{E}^{(4)}$, we first estimate the upper bound for $|\mathbf{H}_{ij}(t) - \mathbf{H}_{ij}(0)|$ for all $i, j \in [n]$. We note that

$$\mathbb{I}_{\left\{\mathbb{I}_{\left\{\overrightarrow{\mathbf{w}}_{\mathbf{S},r}(t)^{\top}\overrightarrow{\mathbf{x}}_{i}\right\}}\neq\mathbb{I}_{\left\{\mathbf{w}_{r}(0)^{\top}\overrightarrow{\mathbf{x}}_{i}\right\}}\right\}} \leq \mathbb{I}_{\left\{\left|\mathbf{w}_{r}(0)^{\top}\overrightarrow{\mathbf{x}}_{i}\right|\leq R\right\}} + \mathbb{I}_{\left\{\left\|\mathbf{w}_{\mathbf{S},r}(t)-\overrightarrow{\mathbf{w}}_{r}(0)\right\|_{2}>R\right\}}.$$
(64)

¹¹⁸⁶ It follows from (64) that

 $|\mathbf{H}_{ij}(t) - \mathbf{H}_{ij}(0)|$

$$\begin{aligned} & \underset{1180}{1180} \\ & \underset{1190}{1191} \end{aligned} = \left| \frac{\vec{\mathbf{x}}_i^\top \vec{\mathbf{x}}_j}{m} \sum_{r=1}^m \left(\mathrm{I}_{\left\{ \vec{\mathbf{w}}_{\mathbf{s},r}(t)^\top \vec{\mathbf{x}}_i \ge 0 \right\}} \mathrm{I}_{\left\{ \vec{\mathbf{w}}_r(0)^\top \vec{\mathbf{x}}_j \ge 0 \right\}} - \mathrm{I}_{\left\{ \mathbf{w}_r(0)^\top \vec{\mathbf{x}}_i \ge 0 \right\}} \mathrm{I}_{\left\{ \mathbf{w}_r(0)^\top \vec{\mathbf{x}}_j \ge 0 \right\}} \right) \right| \end{aligned}$$

$$\leq \frac{1}{m} \sum_{r=1}^{m} \mathbb{I}_{\left\{ \mathbf{I}_{\left\{ \overrightarrow{\mathbf{w}}_{\mathbf{S},r}(t)^{\top} \overrightarrow{\mathbf{x}}_{i} \geq 0 \right\}} \neq \mathbb{I}_{\left\{ \overrightarrow{\mathbf{w}}_{r}(0)^{\top} \overrightarrow{\mathbf{x}}_{i} \geq 0 \right\}} \right\}} \\
\leq \frac{1}{m} \sum_{r=1}^{m} \left(\mathbb{I}_{\left\{ \left| \overrightarrow{\mathbf{w}}_{r}(0)^{\top} \overrightarrow{\mathbf{x}}_{i} \right| \leq R \right\}} + \mathbb{I}_{\left\{ \left\| \mathbf{w}_{\mathbf{S},r}(t) - \overrightarrow{\mathbf{w}}_{r}(0) \right\|_{2} > R \right\}} \right) \\
\leq v_{R}(\mathbf{W}(0), \overrightarrow{\mathbf{x}}_{i}) \stackrel{\bigoplus}{\leq} \frac{2R}{\sqrt{2\pi}\kappa} + C_{2}(m/2, R_{0}, \delta),$$
(65)

where ① follows from (56).

It follows from (63) and (65) that $\|\mathbf{E}^{(3)}\|_2, \|\mathbf{E}^{(4)}\|_2$ are bounded by

$$\begin{aligned} \left\| \mathbf{E}^{(3)} \right\|_{2} &\leq \frac{\eta}{n} \| \mathbf{K} - \mathbf{H}(0) \|_{2} \| \widehat{\mathbf{y}}(t) - \mathbf{y} \|_{2} \\ &\leq \frac{\eta}{n} \cdot nC_{1}(m/2, d, \delta) \cdot \| \mathbf{y}(t) - \mathbf{y} \|_{2} \\ &\leq \eta c_{\mathbf{u}} \sqrt{n} C_{1}(m/2, d, \delta), \end{aligned}$$
(66)

$$\begin{aligned} \|\mathbf{E}^{(4)}\|_{2} &\leq \frac{\eta}{n} \|\mathbf{H}(0) - \mathbf{H}(t)\|_{2} \|\widehat{\mathbf{y}}(t) - \mathbf{y}\|_{2} \\ &\leq \frac{\eta}{n} \cdot n \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta)\right) \cdot \|\mathbf{y}(t) - \mathbf{y}\|_{2} \\ &\leq \eta c_{\mathbf{u}} \sqrt{n} \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta)\right). \end{aligned}$$

$$(67)$$

It follows from (59) and (62) that

$$\mathbf{D}_{i}^{(1)} = \mathbf{D}_{i}^{(3)} + \mathbf{E}_{i}^{(2)} + \mathbf{E}_{i}^{(3)} + \mathbf{E}_{i}^{(4)}.$$
(68)

It then follows from (54) that

$$\widehat{\mathbf{y}}_{i}(t+1) - \widehat{\mathbf{y}}_{i}(t) = \mathbf{D}_{i}^{(1)} + \mathbf{E}_{i}^{(1)}$$

$$= \mathbf{D}_{i}^{(3)} + \underbrace{\mathbf{E}_{i}^{(1)} + \mathbf{E}_{i}^{(2)} + \mathbf{E}_{i}^{(3)} + \mathbf{E}_{i}^{(4)}}_{:=\mathbf{E}_{i}}$$

$$= -\frac{\eta}{n} \mathbf{K} \left(\widehat{\mathbf{y}}(t) - \mathbf{y} \right) + \mathbf{E}_{i}, \tag{69}$$

where $\mathbf{E} \in \mathbb{R}^n$ with its *i*-th element being \mathbf{E}_i , and $\mathbf{E} = \mathbf{E}^{(1)} + \mathbf{E}^{(2)} + \mathbf{E}^{(3)} + \mathbf{E}^{(4)}$. It then follows from (58), (61), (66), and (67) that

$$\|\mathbf{E}\|_{2} \leq \eta c_{\mathbf{u}} \sqrt{n} \left(3 \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta) \right) + C_{1}(m/2, d, \delta) \right).$$
(70)

Finally, (69) can be rewritten as

$$\widehat{\mathbf{y}}(t+1) - \mathbf{y} = \left(\mathbf{I} - \frac{\eta}{n}\mathbf{K}\right)(\widehat{\mathbf{y}}(t) - \mathbf{y}) + \mathbf{E}(t+1),$$

which proves (51) with the upper bound for $||\mathbf{E}||_2$ in (70).

Lemma C.6. Suppose that $t \in [0, T-1]$ for $T \ge 1$, and $\|\widehat{\mathbf{y}}(t') - \mathbf{y}\|_2 \le \sqrt{n}c_{\mathbf{u}}$ holds for all $0 \leq t' \leq t$. Then $\|\bar{\mathbf{v}}\|$

$$\left\| \vec{\mathbf{w}}_{\mathbf{S},r}(t') - \vec{\mathbf{w}}_{r}(0) \right\|_{2} \le R, \quad \forall 0 \le t' \le t+1.$$
 (71)

Proof. Let $[\mathbf{Z}_{\mathbf{S}}(t)]_{[(r-1)d:rd]}$ denotes the submatrix of $\mathbf{Z}_{\mathbf{S}}(t)$ formed by the the rows of $\mathbf{Z}_{\mathbf{Q}}(t)$ with row indices in [(r-1)d: rd]. By the GD update rule we have for $t \in [0, T-1]$ that

$$\vec{\mathbf{w}}_{\mathbf{S},r}(t+1) - \vec{\mathbf{w}}_{\mathbf{S},r}(t) = -\frac{\eta}{n} \left[\mathbf{Z}_{\mathbf{S}}(t) \right]_{[(r-1)d:rd]} \left(\widehat{\mathbf{y}}(t) - \mathbf{y} \right),$$
(72)

We have $\left\| \left[\mathbf{Z}_{\mathbf{S}}(t) \right]_{[(r-1)d:rd]} \right\|_{2} \leq \sqrt{n/m}$. It then follows from (72) that

$$\left\| \overrightarrow{\mathbf{w}}_{\mathbf{S},r}(t+1) - \overrightarrow{\mathbf{w}}_{\mathbf{S},r}(t) \right\|_{2} \leq \frac{\eta}{n} \| \mathbf{Z}_{\mathbf{S}}(t) \|_{2} \| \widehat{\mathbf{y}}(t) - \mathbf{y} \|_{2} \leq \frac{\eta c_{\mathbf{u}}}{\sqrt{m}}.$$
(73)

Note that (71) trivially holds for t' = 0. For $t' \in [1, t + 1]$, it follows from (73) that

$$\begin{split} \left\| \vec{\mathbf{w}}_{\mathbf{S},r}(t') - \vec{\mathbf{w}}_{r}(0) \right\|_{2} &\leq \sum_{t''=0}^{t'-1} \left\| \vec{\mathbf{w}}_{\mathbf{S},r}(t''+1) - \vec{\mathbf{w}}_{\mathbf{S},r}(t'') \right\|_{2} \\ &\leq \frac{\eta}{\sqrt{m}} \sum_{t''=0}^{t'-1} c_{\mathbf{u}} \\ &\leq \frac{\eta c_{\mathbf{u}} T}{\sqrt{m}} = R, \end{split}$$
(74) which completes the proof.

which completes the proof.

Lemma C.7. Let $h(\cdot) = \sum_{t'=0}^{t-1} h(\cdot, t')$ for $t \in [T], T \leq \widehat{T}$ where $h(\cdot, t') = v(\cdot, t') + \widehat{e}(\cdot, t'),$ $v(\cdot, t') = \frac{\eta}{n} \sum_{i=1}^{n} K(\vec{\mathbf{x}}_j, \mathbf{x}) \mathbf{v}_j(t'),$ $\widehat{e}(\cdot,t') = \frac{\eta}{n} \sum_{i=1}^{n} K(\overrightarrow{\mathbf{x}}_{j},\mathbf{x}) \overrightarrow{\mathbf{e}}_{j}(t'),$

where $\mathbf{v}(t') \in \mathcal{V}_{t'}$, $\mathbf{e}(t') \in \mathcal{E}_{t',\tau}$ for all $0 \leq t' \leq t-1$. Suppose that $\tau \leq 1/(\eta T)$, then with probability at least $1 - \exp\left(-\Theta(n\widehat{\varepsilon}_n^2)\right)$ over the random noise w,

> $\|h\|_{\mathcal{H}_{K}} \leq B_{h} = \mu_{0} + 1 + \sqrt{2},$ (75)

and B_h is also defined in (35).

> *Proof.* We have $\mathbf{y} = f^*(\mathbf{S}) + \mathbf{w}, \mathbf{v}(t) = -(\mathbf{I} - \eta \mathbf{K}_n)^t f^*(\mathbf{S}), \mathbf{e}(t) = \overrightarrow{\mathbf{e}}_1(t) + \overrightarrow{\mathbf{e}}_2(t)$ with $\overrightarrow{\mathbf{e}}_1(t) = -(\mathbf{I} - \eta \mathbf{K}_n)^t f^*(\mathbf{S})$. $-(\mathbf{I} - \eta \mathbf{K}_n)^t \mathbf{w}, \| \vec{\mathbf{e}}_2(t) \|_2 \lesssim \sqrt{n\tau}.$ We define

$$\widehat{e}_1(\cdot,t) = \frac{\eta}{n} \sum_{j=1}^n K(\overrightarrow{\mathbf{x}}_j,\mathbf{x}) \left[\overrightarrow{\mathbf{e}}_1(t')\right]_j, \quad \widehat{e}_2(\cdot,t) = \frac{\eta}{n} \sum_{j=1}^n K(\overrightarrow{\mathbf{x}}_j,\mathbf{x}) \left[\overrightarrow{\mathbf{e}}_2(t')\right]_j,$$

Let Σ be the diagonal matrix containing eigenvalues of \mathbf{K}_n , we then have

$$\sum_{t'=0}^{t-1} v(\mathbf{x}, t') = \frac{\eta}{n} \sum_{j=1}^{n} \sum_{t'=0}^{t-1} \left[\left(\mathbf{I} - \eta \mathbf{K}_n \right)^{t'} f^*(\mathbf{S}) \right]_j K(\vec{\mathbf{x}}_j, \mathbf{x})$$
$$= \frac{\eta}{n} \sum_{j=1}^{n} \sum_{t'=0}^{t-1} \left[\mathbf{U} \left(\mathbf{I} - \eta \mathbf{\Sigma} \right)^{t'} \mathbf{U}^\top f^*(\mathbf{S}) \right]_j K(\vec{\mathbf{x}}_j, \mathbf{x}).$$
(76)

It follows from (76) that

1294
1295
$$\left\|\sum_{t'=0}^{t-1} v(\cdot, t')\right\|_{\mathcal{H}_K}^2$$

1296
1297
1298
$$= \frac{\eta^2}{n^2} f^*(\mathbf{S})^\top \mathbf{U} \sum_{t'=0}^{t-1} \left(\mathbf{I} - \eta \boldsymbol{\Sigma}\right)^{t'} \mathbf{U}^\top \mathbf{K} \mathbf{U} \sum_{t'=0}^{t-1} \left(\mathbf{I} - \eta \boldsymbol{\Sigma}\right)^{t'} \mathbf{U}^\top f^*(\mathbf{S})$$

 $\leq \frac{1}{n} \sum_{i=1}^{n} \frac{\left(1 - \left(1 - \eta \lambda_i\right)\right)}{\widehat{\lambda}_i} \left[\mathbf{U}^\top f^*(\mathbf{S})\right]_i^2$

1299
1300
1301
$$= \frac{1}{n} \left\| \eta \left(\mathbf{K}_{n} \right)^{1/2} \mathbf{U} \sum_{i=1}^{t-1} \left(\mathbf{I} - \eta \boldsymbol{\Sigma} \right)^{t'} \mathbf{U}^{\top} f^{*}(\mathbf{S}) \right\|^{2}$$

$$= \frac{1}{n} \left\| \eta \left(\mathbf{K}_n \right)^{1/2} \mathbf{U} \sum_{t'=0} \left(\mathbf{I} - \eta \mathbf{Z} \right)^{t'} \right\|_{t'=0}^{2}$$

where the last inequality follows from Lemma C.12.

 $\leq \mu_0^2$,

Similarly, we have

$$\left\|\sum_{t'=0}^{t-1}\widehat{e}_{1}(\cdot,t')\right\|_{\mathcal{H}_{K}}^{2} \leq \frac{1}{n}\sum_{i=1}^{n}\frac{\left(1-\left(1-\eta\widehat{\lambda}_{i}\right)^{t}\right)^{2}}{\widehat{\lambda}_{i}}\left[\mathbf{U}^{\top}\mathbf{w}\right]_{i}^{2}.$$
(78)

 $\overset{(1)}{\leq} \frac{\sigma^2}{n} \sum_{i=1}^{n} \min\left\{\frac{1}{\widehat{\lambda}_i}, \eta_t^2 \widehat{\lambda}_i\right\}$

 $\overset{\textcircled{0}}{\leq} \frac{\sigma^2 \eta_t}{n} \sum_{i=1}^n \min\left\{1, \eta_t \widehat{\lambda}_i\right\}$

 $= \frac{\sigma^2 \eta_t^2}{n} \sum_{i=1}^n \min\left\{\eta_t^{-1}, \widehat{\lambda}_i\right\}$

 $= \sigma^2 \eta_t^2 \widehat{R}_K^2(\sqrt{1/\eta_t}) < 1.$

 $\leq \frac{\sigma^2 \eta_t}{n} \sum_{i=1}^n \min\left\{\frac{1}{n_i \hat{\lambda}_i}, \eta_t \hat{\lambda}_i\right\}$

(77)

(80)

It then follows from the argument in the proof of (Raskutti et al., 2014, Lemma 9) that the RHS of (78) is bounded with high probability. We define a diagonal matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$ with \mathbf{R}_{ii} $(1 - (1 - \eta \hat{\lambda}_i)^t)^2 / \hat{\lambda}_i$ for $i \in [n]$. Then the RHS of (78) is $1/n \cdot \text{tr} (\mathbf{U} \mathbf{R} \mathbf{U}^\top \mathbf{w} \mathbf{w}^\top)$. It follows from (Wright, 1973) that

$$\Pr\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right) - \mathbb{E}\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right)\right] \ge u\right]$$
$$\le \exp\left(-c\min\left\{nu/\|\mathbf{R}\|_{2}, n^{2}u^{2}/\|\mathbf{R}\|_{F}^{2}\right\}\right)$$
(79)

 $\mathbb{E}\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right)\right] \leq \frac{\sigma^{2}}{n} \sum_{i=1}^{n} \frac{\left(1 - \left(1 - \eta\widehat{\lambda}_{i}\right)^{t}\right)^{2}}{\widehat{\lambda}_{i}}$

for all u > 0, and c is a positive constant. Recall that $\eta_t = \eta t$ for all $t \ge 0$, we have

Here (1) follows from the fact that $(1 - \eta \widehat{\lambda}_i)^t \ge \max \left\{ 0, 1 - t \eta \widehat{\lambda}_i \right\}$, and (2) follows from $\min \{a, b\} \leq \sqrt{ab} \text{ for any nonnegative numbers } a, b. \text{ Because } t \leq T \leq \widehat{T}, \text{ we have } \widehat{R}_K(\sqrt{1/\eta_t}) \leq T \leq \widehat{T}, \text{ for any nonnegative numbers } a, b \in \mathbb{R}$ $1/(\sigma \eta_t)$, so the last inequality holds.

Moreover, we have the upper bounds for $\|\mathbf{R}\|_2$ and $\|\mathbf{R}\|_F$ as follows. First, we have

1348
1349
$$\|\mathbf{R}\|_{2} \leq \max_{i \in [n]} \frac{\left(1 - \left(1 - \eta \widehat{\lambda}_{i}\right)^{t}\right)^{2}}{\widehat{\lambda}_{i}}$$

1350
1351
$$\leq \min\left\{\frac{1}{\widehat{\lambda}_i}, \eta_t^2 \widehat{\lambda}_i\right\} \leq \eta_t.$$
(81)

We also have

$$\frac{1}{n} \|\mathbf{R}\|_{\mathrm{F}}^{2} = \frac{1}{n} \sum_{i=1}^{n} \frac{\left(1 - \left(1 - \eta \widehat{\lambda}_{i}\right)^{t}\right)^{4}}{(\widehat{\lambda}_{i})^{2}} \\
\leq \frac{\eta_{t}^{3}}{n} \sum_{i=1}^{n} \min\left\{\frac{1}{\eta_{t}^{3}\widehat{\lambda}_{i}^{2}}, \eta_{t}\widehat{\lambda}_{i}^{2}\right\} \\
\stackrel{(3)}{\leq} \frac{\eta_{t}^{3}}{n} \sum_{i=1}^{n} \min\left\{\widehat{\lambda}_{i}, \frac{1}{\eta_{t}}\right\} = \eta_{t}^{3}\widehat{R}_{K}^{2}(\sqrt{1/\eta_{t}}) \leq \frac{\eta_{t}}{\sigma^{2}},$$
(82)

where ③ follows from

$$\min\left\{\frac{1}{\eta_t^3\widehat{\lambda}_i^2}, \eta_t\widehat{\lambda}_i^2\right\} = \widehat{\lambda}_i \min\left\{\frac{1}{\eta_t^3\widehat{\lambda}_i^3}, \eta_t\widehat{\lambda}_i\right\} \le \widehat{\lambda}_i.$$

Combining (78)- (82) with u = 1 in (79), we have

1370
1371
$$\Pr\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right) - \mathbb{E}\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right)\right] \ge 1\right] \le \exp\left(-c\min\left\{n/\eta_t, n\sigma^2/\eta_t\right\}\right)$$
1372
$$\le \exp\left(-nc'/\eta_t\right) \le \exp\left(-c'n\widehat{\varepsilon}_n^2\right)$$

where $c' = c \min\{1, \sigma^2\}$, and the last inequality is due to the fact that $1/\eta_t \ge \hat{\varepsilon}_n^2$ since $t \le T \le \hat{T}$. It follows that with probability at least $1 - \exp\left(-\Theta(n\widehat{\varepsilon}_n^2)\right), \left\|\sum_{t'=0}^{t-1}\widehat{e}_1(\cdot, t')\right\|_{\mathcal{H}_K}^2 \leq 2.$

We now find the upper bound for $\left\|\sum_{t'=0}^{t-1} \widehat{e}_2(\cdot, t')\right\|_{\mathcal{H}_{K}}$. We have $\|\widehat{e}_{2}(\cdot,t')\|_{\mathcal{H}_{K}}^{2} \leq \frac{\eta^{2}}{n^{2}} \overrightarrow{\mathbf{e}}_{2}^{\top}(t') \mathbf{K} \overrightarrow{\mathbf{e}}_{2}(t')$ $\leq \eta^2 \widehat{\lambda}_1 \tau^2,$

so that

$$\left\|\sum_{t'=0}^{t-1} \widehat{e}_2(\cdot, t')\right\|_{\mathcal{H}_K} \leq \sum_{t'=0}^{t-1} \|\widehat{e}_2(\cdot, t')\|_{\mathcal{H}_K}$$
$$\leq T\eta \sqrt{\widehat{\lambda}_1} \tau \leq 1, \tag{83}$$

if $\tau \lesssim 1/(\eta T)$.

Finally, we have

$$\|h\|_{\mathcal{H}_{K}} \leq \left\|\sum_{t'=0}^{t-1} \widehat{v}(\cdot, t')\right\|_{\mathcal{H}_{K}} + \left\|\sum_{t'=0}^{t-1} \widehat{e}_{1}(\cdot, t')\right\|_{\mathcal{H}_{K}} + \left\|\sum_{t'=0}^{t-1} \widehat{e}_{2}(\cdot, t')\right\|_{\mathcal{H}_{K}}$$
$$\leq \mu_{0} + 1 + \sqrt{2} = B_{h}.$$

Theorem C.8. For every $t \in [T]$, let the neural network $f(\cdot) = f(\mathbf{W}(t), \cdot)$ be trained by gradient descent with the learning rate $\eta \in (0, 1/\hat{\lambda}_1)$ on the random initialization $\mathbf{W}(0) \in \mathcal{W}_0$ with $T \leq \hat{T}$. Then with probability at least $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\hat{\varepsilon}_n^2))$ over the random noise w, $f \in$ $\mathcal{F}_{NN}(\mathbf{S}, \mathbf{W}(0), T)$, and f can be decomposed by

f

$$= h + e \in \mathcal{F}_{\text{ext}}(w, T), \tag{84}$$

where $h \in \mathcal{H}_K(B_h)$ with B_h defined in (35), $e \in L^\infty$. When

$$m \gtrsim \max\left\{\frac{(\eta T)^4 \left(\sqrt{d}+1\right)^4}{w^4}, (\eta T)^8 d^2\right\},$$
(85)

1410 then

$$\|e\|_{\infty} \le w. \tag{86}$$

1413 In addition,

$$\|f\|_{\infty} \le \frac{B_h}{\sqrt{2}} + w. \tag{87}$$

Remark. We consider the kernel regression problem with the training loss $L(\alpha) = 1/2 \cdot \|\mathbf{K}_n \alpha - \mathbf{y}\|_2^2$. Letting $\beta = \mathbf{K}_n^{1/2} \alpha$ and then performing GD on β with this training loss and the learning rate η , it can be verified that the kernel regressor right after the *t*-th step of GD is

$$\widehat{f}_t^{(\text{NTK})} = \frac{\eta}{n} \sum_{t'=0}^{t-1} \sum_{i=1}^n K(\cdot, \mathbf{\vec{x}}_i) \boldsymbol{\alpha}_i^{(t')},$$
(88)

where $\alpha^{(t')} = (\mathbf{I}_n - \eta \mathbf{K}_n)^{t'} \mathbf{y}$. Following from the proof of Lemma C.6 and Theorem C.8, under the conditions of Theorem C.8 we have

$$h_t = \widehat{f}_t^{(\mathrm{NTK})} + \widehat{e}_2(\cdot, t)$$

where $\hat{e}_2(\cdot, t) = \frac{\eta}{n} \sum_{t'=0}^{t-1} \sum_{j=1}^n K(\cdot, \mathbf{x}_j) \left[\mathbf{e}_2(t') \right]_i$ and $\mathbf{e}_2(t')$ appears in the definition of $\mathcal{E}_{t,\tau}$ in (31). It is remarked that in our analysis, we approximate f_t by $h_t \in \mathcal{H}_K(B_h)$ with a small approximation error w, and we do not need to approximate f_t by the kernel regressor $\hat{f}_t^{(\text{NTK})}$ with a sufficiently small approximation error which is the common strategy used in existing works (Hu et al., 2021; Suh et al., 2022; Li et al., 2024). In fact, our analysis only requires m is suitably large so that the \mathcal{H}_K -norm of $\hat{e}_2(\cdot, t) = h_t - \hat{f}_t^{(\text{NTK})}$ is bounded by a positive constant rather than an infinitesimal number as $m \to \infty$, that is, $\|\hat{e}_2(\cdot, t)\|_{\mathcal{H}_K} \leq 1$, which is revealed by the proof of Lemma C.7.

 $\vec{\mathbf{w}}_r$ is expressed as

$$\vec{\mathbf{w}}_r = \vec{\mathbf{w}}_{\mathbf{S},r}(t) = \vec{\mathbf{w}}_r(0) - \sum_{t'=0}^{t-1} \frac{\eta}{n} \left[\mathbf{Z}_{\mathbf{S}}(t') \right]_{[(r-1)d:rd]} \mathbf{u}(t'), \tag{89}$$

1448 where the notation $\vec{\mathbf{w}}_{\mathbf{S},r}$ emphasizes that $\vec{\mathbf{w}}_r$ depends on the training data S.

We define the event

$$E_r(R) \coloneqq \left\{ \left| \vec{\mathbf{w}}_r(0)^\top \mathbf{x} \right| \le R \right\}, \quad r \in [m].$$

We now approximate $f(\mathbf{W}, \mathbf{x})$ by $g(\mathbf{x}) \coloneqq \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \mathbb{I}_{\left\{ \vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0 \right\}} \vec{\mathbf{w}}_r^\top \mathbf{x}$. We have

 $|f(\mathbf{W}, \mathbf{x}) - g(\mathbf{x})|$ 1456 $= \frac{1}{\sqrt{m}} \left| \sum_{r=1}^{m} a_r \sigma\left(\vec{\mathbf{w}}_r^\top \mathbf{x} \right) - \sum_{r=1}^{m} a_r \mathbb{I}_{\left\{ \vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0 \right\}} \vec{\mathbf{w}}_r^\top \mathbf{x} \right|$

$$\leq \frac{1}{\sqrt{m}} \sum_{r=1}^{m} \left| a_r \left(\mathbb{1}_{\{E_r(R)\}} + \mathbb{1}_{\{\bar{E}_r(R)\}} \right) \left(\sigma \left(\vec{\mathbf{w}}_r^\top \mathbf{x} \right) - \mathbb{1}_{\{\vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\}} \vec{\mathbf{w}}_r^\top \mathbf{x} \right) \right|$$

$$= 1 \qquad 1 \qquad m \qquad | \mathbf{x} \in [\mathbf{x}, \mathbf{x}]$$

1461
1462
$$= \frac{1}{\sqrt{m}} \sum_{r=1}^{m} \mathbb{1}_{\{E_r(R)\}} \left| \sigma \left(\overrightarrow{\mathbf{w}}_r^\top \mathbf{x} \right) - \mathbb{1}_{\{\overrightarrow{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\}} \overrightarrow{\mathbf{w}}_r^\top \mathbf{x} \right|$$
1463

$$= \frac{1}{\sqrt{m}} \sum_{r=1}^{m} \mathbb{I}_{\{E_r(R)\}} \left| \sigma \left(\overrightarrow{\mathbf{w}}_r^\top \mathbf{x} \right) - \sigma \left(\overrightarrow{\mathbf{w}}_r(0)^\top \mathbf{x} \right) - \mathbb{I}_{\{\overrightarrow{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\}} (\overrightarrow{\mathbf{w}}_r - \overrightarrow{\mathbf{w}}_r(0))^\top \mathbf{x} \right|$$

1466
1467
1468
$$\leq \frac{2R}{\sqrt{m}} \sum_{r=1}^{m} \mathbb{I}_{\{E_r(R)\}}.$$
(90)

1469 Plugging $R = \frac{\eta c_{\mathbf{u}} T}{\sqrt{m}}$ in (90), we have

$$|f(\mathbf{W}, \mathbf{x}) - g(\mathbf{x})| \leq \frac{2R}{\sqrt{m}} \sum_{r=1}^{m} \mathbb{1}_{\{E_r(R)\}}$$
$$= 2\eta c_{\mathbf{u}} T \cdot \frac{1}{m} \sum_{r=1}^{m} \mathbb{1}_{\{E_r(R)\}}$$
$$= 2\eta c_{\mathbf{u}} T \cdot \widehat{v}_R(\mathbf{W}(0), \mathbf{x})$$
$$\leq 2\eta c_{\mathbf{u}} T \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_2(m/2, R_0, \delta)\right).$$
(91)

Using (89), we can express $g(\mathbf{x})$ as

$$g(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \mathbf{I}_{\left\{\vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\right\}} \vec{\mathbf{w}}_r(0)^\top \mathbf{x}$$
$$- \sum_{t'=0}^{t-1} \frac{1}{\sqrt{m}} \sum_{r=1}^{m} \mathbf{I}_{\left\{\vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\right\}} \left(\frac{\eta}{n} \left[\mathbf{Z}_{\mathbf{S}}(t')\right]_{[(r-1)d:rd]} \mathbf{u}(t')\right)^\top \mathbf{x}$$
$$\underbrace{\bigoplus_{t'=0}^{t-1} \frac{\eta}{nm} \sum_{r=1}^{m} \mathbf{I}_{\left\{\vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\right\}} \sum_{j=1}^{n} \mathbf{I}_{\left\{\vec{\mathbf{w}}_r(t')^\top \vec{\mathbf{x}}_j \ge 0\right\}} \mathbf{u}_j(t') \vec{\mathbf{x}}_j^\top \mathbf{x}, \qquad (92)$$
$$= G_{t'}(\mathbf{x})}$$

where ① follows from the fact that $\frac{1}{\sqrt{m}} \sum_{r=1}^{m} a_r \mathbb{I}_{\{\vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\}} \vec{\mathbf{w}}_r(0)^\top \mathbf{x} = f(\mathbf{W}(0), \mathbf{x}) = 0$ due to the particular initialization of the two-layer NN. For each $G_{t'}$ in the RHS of (92), we have

$$G_{t'}(\mathbf{x}) \stackrel{\textcircled{0}}{=} \frac{\eta}{nm} \sum_{r=1}^{m} \mathrm{I}\!\!I_{\left\{\vec{\mathbf{w}}_{r}(0)^{\top}\mathbf{x}\geq0\right\}} \sum_{j=1}^{n} d_{t',r,j}\mathbf{u}_{j}(t')\vec{\mathbf{x}}_{j}^{\top}\mathbf{x}$$

$$+ \frac{\eta}{nm} \sum_{r=1}^{m} \mathrm{I}\!\!I_{\left\{\vec{\mathbf{w}}_{r}(0)^{\top}\mathbf{x}\geq0\right\}} \sum_{j=1}^{n} \mathrm{I}\!\!I_{\left\{\vec{\mathbf{w}}_{r}(0)^{\top}\vec{\mathbf{x}}\geq0\right\}} \mathbf{u}_{j}(t')\vec{\mathbf{x}}_{j}^{\top}\mathbf{x}$$

$$\stackrel{\textcircled{0}}{=} \frac{\eta}{n} \sum_{j=1}^{n} K(\mathbf{x},\vec{\mathbf{x}}_{j})\mathbf{u}_{j}(t') + \frac{\eta}{n} \sum_{j=1}^{n} q_{j}\mathbf{u}_{j}(t')$$

$$= E_{1}(\mathbf{x})$$

$$+ \underbrace{\frac{\eta}{nm} \sum_{r=1}^{m} \mathrm{I}\!\!I_{\left\{\vec{\mathbf{w}}_{r}(0)^{\top}\mathbf{x}\geq0\right\}} \sum_{j=1}^{n} d_{t',r,j}\mathbf{u}_{j}(t')\vec{\mathbf{x}}_{j}^{\top}\mathbf{x}}. \qquad (93)$$

$$= E_{2}(\mathbf{x})$$

1510 where 1511

$$d_{t',r,j} \coloneqq \mathrm{I}_{\left\{ \overrightarrow{\mathbf{w}}_r(t')^\top \overrightarrow{\mathbf{x}}_j \ge 0 \right\}} - \mathrm{I}_{\left\{ \overrightarrow{\mathbf{w}}_r(0)^\top \overrightarrow{\mathbf{x}}_j \ge 0 \right\}}$$

in 2, and and

$$q_j \coloneqq h(\mathbf{W}(0), \mathbf{x}_j, \mathbf{x}) - K(\mathbf{x}_j, \mathbf{x})$$

for all $j \in [n]$ in \mathfrak{G} .

We now analyze each term on the RHS of (93). Let $h(\cdot, t'): \mathcal{X} \to \mathbb{R}$ be defined by

$$h(\mathbf{x}, t') \coloneqq \frac{\eta}{n} \sum_{j=1}^{n} K(\mathbf{x}, \mathbf{\vec{x}}_j) \mathbf{u}_j(t'),$$

then $h(\cdot, t')$ is an element in the RKHS \mathcal{H}_K for each $t' \in [0, t-1]$. We further define

Ш

$$h(\cdot) \coloneqq \sum_{t'=0}^{t-1} h(\cdot, t'),$$
(94)

It follows from Theorem C.1 that, with probability $1 - \delta$ over $\mathbf{W}(0)$, $q_j \leq C_1(m/2, d, \delta)$ for all $j' \in [n]$ with $C_1(m/2, d, \delta)$ defined in (38). Moreover, $\|\mathbf{K}_{\mathbf{S},\mathbf{Q}}\|_2 \leq \sqrt{nN}, \mathbf{u}(t') \leq c_{\mathbf{u}}\sqrt{n}$ with high probability, so that we have

$$\|E_1\|_{\infty} = \left\| \frac{\eta}{n} \sum_{j=1}^n q_j \mathbf{u}_j(t') \right\|_{\infty} \le \frac{\eta}{n} \|\mathbf{u}(t')\|_2 \sqrt{n} C_1(m/2, d, \delta)$$
$$\le \eta c_{\mathbf{u}} C_1(m/2, d, \delta). \tag{95}$$

We now bound the last term on the RHS of (93). Define $\mathbf{X}' \in \mathbb{R}^{d \times n}$ with its *j*-column being $\mathbf{X}'_{j} = \frac{1}{m} \sum_{r=1}^{m} \mathbb{1}_{\left\{ \overrightarrow{\mathbf{w}}_{r}(0)^{\top} \mathbf{x} \ge 0 \right\}} d_{t',r,j} \overrightarrow{\mathbf{x}}_{j} \text{ for all } j \in [n], \text{ then } E_{2}(\mathbf{x}) = \frac{\eta}{n} \left(\mathbf{X}' \mathbf{u}(t') \right)^{\top} \mathbf{x}.$

We need to derive the upper bound for $\|\mathbf{X}'\|_2$. Because $\|\vec{\mathbf{w}}_r - \vec{\mathbf{w}}_r(0)\|_2 \leq R$, it follows that $\mathbb{I}_{\left\{\vec{\mathbf{w}}_{r}(t')^{\top}\vec{\mathbf{x}}_{j}\geq0\right\}} = \mathbb{I}_{\left\{\vec{\mathbf{w}}_{r}(0)^{\top}\vec{\mathbf{x}}_{j}\geq0\right\}} \text{ when } \left|\vec{\mathbf{w}}_{r}(0)^{\top}\mathbf{x}_{j'}'\right| > R \text{ for all } j'\in[n]. \text{ Therefore,}$

$$|d_{t',r,j'}| = \left| \mathbb{I}_{\left\{ \vec{\mathbf{w}}_r(t')^\top \vec{\mathbf{x}}_j \ge 0 \right\}} - \mathbb{I}_{\left\{ \vec{\mathbf{w}}_r(0)^\top \vec{\mathbf{x}}_j \ge 0 \right\}} \right| \le \mathbb{I}_{\left\{ \left| \vec{\mathbf{w}}_r(0)^\top \vec{\mathbf{x}}_j \right| \le R \right\}},$$

and it follows that

$$\frac{\left|\sum_{r=1}^{m} \mathrm{I\!I}_{\left\{\overrightarrow{\mathbf{w}}_{r}(0)^{\top}\overrightarrow{\mathbf{x}}_{i}\geq0\right\}}d_{t',r,j}\right|}{m} \leq \frac{\sum_{r=1}^{m} |d_{t',r,j}|}{m} \leq \frac{\sum_{r=1}^{m} \mathrm{I\!I}_{\left\{\left|\overrightarrow{\mathbf{w}}_{r}(0)^{\top}\overrightarrow{\mathbf{x}}_{j}\right|\leq R\right\}}}{m} = \widehat{v}_{R}(\mathbf{W}(0),\overrightarrow{\mathbf{x}}_{j})$$
$$\leq \frac{2R}{\sqrt{2\pi\kappa}} + C_{2}(m/2, R_{0}, \delta), \tag{96}$$

where \hat{v}_R is defined by (15), and the last inequality follows from Theorem C.2.

It follows from (96) that $\|\mathbf{X}'\|_2 \le \sqrt{n} \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_2(m/2, R_0, \delta)\right)$, and we have

$$\begin{aligned} \left\| E_2(\mathbf{x}) \right\|_{\infty} &\leq \frac{\eta}{n} \| \mathbf{X}' \|_2 \| \mathbf{u}(t') \|_2 \| \mathbf{x} \|_2 \\ &\leq \eta c_{\mathbf{u}} \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_2(m/2, R_0, \delta) \right). \end{aligned}$$
(97)

Combining (93), (95), and (97), for any $t' \in [0, t-1]$,

1563
$$\|G_{t'}(\mathbf{x}) - h(\mathbf{x}, t')\|_{\infty} \le \|E_1\|_{\infty} + \|E_2\|_{\infty}$$

1564
1565
$$\leq \eta c_{\mathbf{u}} \left(C_1(m/2, d, \delta) + \frac{2R}{\sqrt{2\pi\kappa}} + C_2(m/2, R_0, \delta) \right).$$
(98)

1566 Define $e(\cdot) = f(\mathbf{W}, \cdot) - h(\cdot)$, it then follows from (91), (92), and (98) that 1567 1568 1569 $\|e(\mathbf{x})\|_{\infty} \leq \|f(\mathbf{W}, \cdot) - g\|_{\infty} + \|g - h\|_{\infty}$ 1570 $\leq \|f(\mathbf{W}, \cdot) - g\|_{\infty} + \sum_{t=1}^{t-1} \|G_{t'} - h(\cdot, t')\|_{\Omega_{\varepsilon_0, \mathbf{Q}}}$ 1571 1572 1573 $\overset{\textcircled{0}}{\leq} 2\eta c_{\mathbf{u}} T \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_2(m/2, R_0, \delta) \right)$ 1574 1575 + $\eta c_{\mathbf{u}}T\left(C_1(m/2,d,\delta)+\frac{2R}{\sqrt{2\pi\kappa}}+C_2(m/2,R_0,\delta)\right)$ 1576 $\leq \eta c_{\mathbf{u}} T \left(C_1(m/2, d, \delta) + 3 \left(\frac{2R}{\sqrt{2\pi\kappa}} + C_2(m/2, R_0, \delta) \right) \right)$ 1579 1580 $\coloneqq \Delta_{m,n,N,c_x,\eta,\tau,\delta}.$ (99) 1581 We now give estimates for $\Delta_{m,n,N,c_x,\eta,\tau,\delta}$. Since $m \ge \max{\{d,n,4\}}$, we have $\sqrt{\frac{d \log m}{m}} \le \frac{\sqrt{d}}{m^{1/4}}$. As a result, 1585 $\Delta_{m,n,N,c_x,\eta,\tau,\delta} \lesssim \eta T \left(\frac{\sqrt{d}}{m^{1/4}} + \frac{\eta T}{\sqrt{m}} \right).$ 1587 1588 By direct calculations, for any w > 0, when 1589 1590 $m \gtrsim \frac{(\eta T)^4 \left(\sqrt{d}+1\right)^4}{4},$ 1591 1592 1593 we have $\Delta_{m,n,N,c_x,\eta,\tau,\delta} \leq w$. 1594 It follows from Lemma C.7 that with probability at least $1 - \exp\left(-\Theta(n\hat{\varepsilon}_n^2)\right)$ over the random noise 1595 w, 1596 1597 (100) $\|h\|_{\mathcal{H}_{K}} \leq B_{h},$ 1598 where B_h is defined in (35), and τ are required to satisfy 1599 1600 $\tau \lesssim 1/(\eta T).$ 1601 Lemma C.3 requires that $m \gtrsim (\eta T)^4 (\sqrt{d} + \sqrt{\tau})^4 / \tau^4$. As a result, we have 1602 1603 $m \gtrsim (\eta T)^8 d^2$. 1604 It also follows from the Cauchy-Schwarz inequality that $||h||_{\infty} \leq B_h/\sqrt{2}$. This together with (99) proves (87). 1607 1608 For B, w > 0, we define the function class 1609 1610 $\mathcal{F}(B, w) := \left\{ f \colon \exists h \in \mathcal{H}_K(B), \exists e \in L^\infty, \|e\|_\infty \le w \text{ s.t. } f = h + e \right\}.$ (101)1611 **Lemma C.9.** For every B, w > 0 every r > 0, 1612 1613 $\Re\left(\left\{f \in \mathcal{F}(B, w) \colon \mathbb{E}_P\left[f^2\right] \le r\right\}\right) \le \varphi_{B, w}(r),$ (102)1614 where 1615 1616 $\varphi_{B,w}(r) \coloneqq \min_{Q: Q \ge 0} \left((\sqrt{r} + w) \sqrt{\frac{Q}{n}} + B \left(\frac{\sum_{q=Q+1}^{\infty} \lambda_q}{n} \right)^{r} \right) + w.$ 1617 1618 (103)1619

Proof. We first decompose the Rademacher complexity of the function class $\{f \in \mathcal{F}(B, w) \colon \mathbb{E}_P [f^2] \leq r\}$ into two terms as follows:

$$\Re\left(\left\{f: f \in \mathcal{F}(B, w), \mathbb{E}_{P}\left[f^{2}\right] \leq r\right\}\right)$$

$$\leq \underbrace{\frac{1}{n}\mathbb{E}\left[\sup_{f \in \mathcal{F}(B, w): \mathbb{E}_{P}\left[f^{2}\right] \leq r} \sum_{i=1}^{n} \sigma_{i}h(\vec{\mathbf{x}}_{i})\right]}_{:=\mathcal{R}_{1}} + \underbrace{\frac{1}{n}\mathbb{E}\left[\sup_{f \in \mathcal{F}(B, w): \mathbb{E}_{P}\left[f^{2}\right] \leq r} \sum_{i=1}^{n} \sigma_{i}e(\vec{\mathbf{x}}_{i})\right]}_{:=\mathcal{R}_{2}}.$$
(104)

We now analyze the upper bounds for $\mathcal{R}_1, \mathcal{R}_2$ on the RHS of (104).

Derivation for the upper bound for \mathcal{R}_1 **.**

According to Definition 101 and Theorem C.8, for any $f \in \mathcal{F}(B,w)$, we have f = h + e with $h \in \mathcal{H}_K(B), e \in L^{\infty}, ||e||_{\infty} \le w.$

When $\mathbb{E}_P[f^2] \leq r$, it follows from the triangle inequality that $\|h\|_{L^2} \leq \|f\|_{L^2} + \|e\|_{L^2} \leq \sqrt{r} + C$ $w \coloneqq r_h$. We now consider $h \in \mathcal{H}_K(B)$ with $\|h\|_{L^2} \leq r_h$ in the remaining of this proof. We have

$$\sum_{i=1}^{n} \sigma_{i} f(\vec{\mathbf{x}}_{i}) = \sum_{i=1}^{n} \sigma_{i} \left(h(\vec{\mathbf{x}}_{i}) + e(\vec{\mathbf{x}}_{i}) \right)$$
$$= \left\langle h, \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}) \right\rangle_{\mathcal{H}_{K}} + \sum_{i=1}^{n} \sigma_{i} e(\vec{\mathbf{x}}_{i}).$$
(105)

Because $\{v_q\}_{q>1}$ is an orthonormal basis of \mathcal{H}_K , for any $0 \leq Q \leq n$, we further express the first term on the RHS of (105) as

$$\begin{aligned}
& \left\{ h, \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}) \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}) \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}) \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{i=1}^{Q} \sqrt{\lambda_{q}} \langle h, v_{q} \rangle_{\mathcal{H}_{K}} v_{q}, \sum_{q=1}^{Q} \frac{1}{\sqrt{\lambda_{q}}} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\rangle_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} v_{q} \right\}_{\mathcal{H}_{K}} \\
& \left\{ h, \sum_{q>Q} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{$$

Due to the fact that $h \in \mathcal{H}_K$, $h = \sum_{q=1}^{\infty} \beta_q^{(h)} v_q = \sum_{q=1}^{\infty} \sqrt{\lambda_q} \beta_q^{(h)} e_q$ with $v_q = \sqrt{\lambda_q} e_q$. Therefore,

$$\|h\|_{L^{2}}^{2} = \sum_{q=1}^{\infty} \lambda_{q} \beta_{q}^{(h)^{2}}, \text{ and}$$

$$\|h\|_{L^{2}}^{2} = \sum_{q=1}^{\infty} \lambda_{q} \beta_{q}^{(h)^{2}}, \text{ and}$$

$$\|\sum_{q=1}^{Q} \sqrt{\lambda_{q}} \langle h, v_{q} \rangle_{\mathcal{H}_{K}} v_{q} \|_{\mathcal{H}_{K}} = \left\|\sum_{q=1}^{Q} \sqrt{\lambda_{q}} \beta_{q}^{(h)} v_{q} \right\|_{\mathcal{H}_{K}}$$

$$= \sqrt{\sum_{q=1}^{Q} \lambda_{q} \beta_{q}^{(h)^{2}}} \leq \|h\|_{L^{2}} \leq r_{h}.$$
(107)
$$(107)$$

According to Mercer's Theorem, because the kernel K is continuous symmetric positive definite, it has the decomposition

$$K(\cdot, \vec{\mathbf{x}}_i) = \sum_{j=1}^{\infty} \lambda_j e_j(\cdot) e_j(\vec{\mathbf{x}}_i),$$

so that we have

$$\begin{cases} 1672\\ 1673 \end{cases} \left\langle \sum_{i=1}^{n} \sigma_{i} K(\cdot, \vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} = \left\langle \sum_{i=1}^{n} \sigma_{i} \sum_{j=1}^{\infty} \lambda_{j} e_{j} e_{j}(\vec{\mathbf{x}}_{i}), v_{q} \right\rangle_{\mathcal{H}_{K}} \end{cases}$$

$$= \left\langle \sum_{i=1}^{n} \sigma_{i} \sum_{j=1}^{\infty} \sqrt{\lambda_{j}} e_{j}(\vec{\mathbf{x}}_{i}) \cdot v_{j}, v_{q} \right\rangle_{\mathcal{H}_{K}}$$

$$= \sum_{i=1}^{n} \sigma_{i} \sqrt{\lambda_{q}} e_{q}(\vec{\mathbf{x}}_{i}).$$

$$(108)$$

Combining (106), (107), and (108), we have

 $\left\langle h, \sum_{i=1}^{n} \sigma_i K(\cdot, \vec{\mathbf{x}}_i) \right\rangle$

$$\leq \|h\|_{L^{2}} \left\| \sum_{q=1}^{Q} \sum_{i=1}^{n} \sigma_{i} e_{q}(\vec{\mathbf{x}}_{i}) v_{q} \right\|_{\mathcal{H}_{K}} + B \left\| \sum_{q=Q+1}^{\infty} \sum_{i=1}^{n} \sigma_{i} \sqrt{\lambda_{q}} e_{q}(\vec{\mathbf{x}}_{i}) v_{q} \right\|_{\mathcal{H}_{K}}$$

$$\leq r_{h} \sqrt{\sum_{i=1}^{Q} \left(\sum_{j=1}^{n} \sigma_{i} e_{q}(\vec{\mathbf{x}}_{i}) \right)^{2}} + B \sqrt{\sum_{i=1}^{\infty} \left(\sum_{j=1}^{n} \sigma_{i} \sqrt{\lambda_{q}} e_{q}(\vec{\mathbf{x}}_{i}) \right)^{2}}, \quad (109)$$

$$\leq r_h \sqrt{\sum_{q=1}^{Q} \left(\sum_{i=1}^{n} \sigma_i e_q(\vec{\mathbf{x}}_i)\right)^2} + B \sqrt{\sum_{q=Q+1}^{\infty} \left(\sum_{i=1}^{n} \sigma_i \sqrt{\lambda_q} e_q(\vec{\mathbf{x}}_i)\right)^2},\tag{109}$$

where ① is due to Cauchy-Schwarz inequality. Moreover, by Jensen's inequality we have

$$\mathbb{E}\left[\sqrt{\sum_{q=1}^{Q} \left(\sum_{i=1}^{n} \sigma_{i} e_{q}(\vec{\mathbf{x}}_{i})\right)^{2}}\right] \leq \sqrt{\mathbb{E}\left[\sum_{q=1}^{Q} \left(\sum_{i=1}^{n} \sigma_{i} e_{q}(\vec{\mathbf{x}}_{i})\right)^{2}\right]} \\ \leq \sqrt{\mathbb{E}\left[\sum_{q=1}^{Q} \sum_{i=1}^{n} e_{q}^{2}(\vec{\mathbf{x}}_{i})\right]} = \sqrt{nQ}.$$
 (110)

and similarly,

$$\mathbb{E}\left[\sqrt{\sum_{q=Q+1}^{\infty} \left(\sum_{i=1}^{n} \sigma_{i} \sqrt{\lambda_{q}} e_{q}(\vec{\mathbf{x}}_{i})\right)^{2}}\right] \leq \sqrt{\mathbb{E}\left[\sum_{q=Q+1}^{\infty} \lambda_{q} \sum_{i=1}^{n} e_{q}^{2}(\vec{\mathbf{x}}_{i})\right]} = \sqrt{n \sum_{q=Q+1}^{\infty} \lambda_{q}}.$$
 (111)

Since (109)-(111) hold for all $Q \ge 0$, it follows that

$$\mathbb{E}\left[\sup_{h\in\mathcal{H}_{K}(B),\|h\|_{L^{2}}\leq r_{h}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}h(\vec{\mathbf{x}}_{i})\right]\leq \min_{Q\colon Q\geq 0}\left(r_{h}\sqrt{nQ}+B\sqrt{n\sum_{q=Q+1}^{\infty}\lambda_{q}}\right).$$
 (112)

_

It follows from (104), (105), and (112) that

$$\begin{aligned}
 1719 \\
 1720 \\
 1721 \\
 1722 \\
 1723 \\
 1724 \\
 1725 \\
 1726 \\
 1727
 \end{aligned}

$$\mathbf{\mathcal{R}}_1 \leq \frac{1}{n} \mathbb{E} \left[\sup_{h \in \mathcal{H}_K(B), \|h\|_{L^2} \leq r_h} \sum_{i=1}^n \sigma_i h(\vec{\mathbf{x}}_i) \right] \\
 1729 \\
 1720 \\
 1721 \\
 1725 \\
 1726
 \end{aligned}$$

$$\mathbf{\mathcal{R}}_1 \leq \frac{1}{n} \mathbb{E} \left[\sup_{h \in \mathcal{H}_K(B), \|h\|_{L^2} \leq r_h} \sum_{i=1}^n \sigma_i h(\vec{\mathbf{x}}_i) \right] \\
 1725 \\
 1726
 \end{aligned}$$

$$(113)$$

$$\mathbf{\mathcal{R}}_1 \leq \frac{1}{n} \mathbb{E} \left[\sum_{h \in \mathcal{H}_K(B), \|h\|_{L^2} \leq r_h} \sum_{i=1}^n \sigma_i h(\vec{\mathbf{x}}_i) \right] \\
 1727
 \end{aligned}$$$$

Derivation for the upper bound for \mathcal{R}_2 **.**

1728 Because $\left|1/n\sum_{i=1}^{n}\sigma_{i}e(\vec{\mathbf{x}}_{i})\right| \leq w$ when $\|e\|_{\infty} \leq w$, we have 1730

1731

1732

1733 1734

$$\mathcal{R}_2 \leq \frac{1}{n} \mathbb{E} \left[\sup_{e \in L^{\infty} : \|e\|_{\infty} \leq w} \sum_{i=1}^n \sigma_i e(\vec{\mathbf{x}}_i) \right] \leq w.$$

 $\leq \min_{Q: Q \geq 0} \left(r_h \sqrt{\frac{Q}{n}} + B\left(\frac{\sum_{q=Q+1}^{\infty} \lambda_q}{n}\right)^{1/2} \right) + w.$

 $\Re\left(\left\{f\colon f\in\mathcal{F}(B,w),\mathbb{E}_P\left[f^2\right]\leq r\right\}\right)$

1735 It follows from (113) and (114) that

1737 1738

1736

1741 1742

Plugging r_h in the RHS of the above inequality completes the proof.

(114)

1745 **Theorem C.10.** Suppose $w \in (0, 1)$ and m satisfy 1746

$$m \gtrsim \max\left\{\frac{(\eta T)^4 \left(\sqrt{d}+1\right)^4}{w^4}, (\eta T)^8 d^2\right\},$$
 (115)

and the neural network $f(\mathbf{W}(t), \cdot)$ is trained by GD in Algorithm 1 with the learning rate $\eta \in (0, 1/\hat{\lambda}_1)$ on random initialization $\mathbf{W}(0)$, and $T \leq \hat{T}$. Then for every $t \in [T]$, with probability at least $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\hat{\varepsilon}_n^2)) - \exp(-n\varepsilon_n^2) - 2/n$ over the random noise w, the random training features S and the random initialization $\mathbf{W}(0)$,

$$\mathbb{E}_P\left[(f_t - f^*)^2\right] - 2\mathbb{E}_{P_n}\left[(f_t - f^*)^2\right]$$

1756 1757 1758

1759 1760 1761

1765 1766

1768

$$\leq c_0 \min_{0 \leq Q \leq n} \left(\frac{B_0 Q}{n} + w \left(\sqrt{\frac{Q}{n}} + 1 \right) + B_h \left(\frac{\sum_{q=Q+1}^{\infty} \lambda_q}{n} \right)^{1/2} \right)^2, \tag{116}$$

Furthermore, with probability at least $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\hat{\varepsilon}_n^2)) - \exp(-n\varepsilon_n^2) - 2/n$ over the random noise w, the random training features **S** and the random initialization **W**(0),

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - 2\mathbb{E}_{P_{n}}\left[(f_{t} - f^{*})^{2}\right] \le c_{0}'(\varepsilon_{n}^{2} + w).$$
(117)

Here B_0, c_0, c'_0 are absolute positive constants depending on μ_0 , and c'_0 also depends on σ .

Proof. We first remark that the conditions on m, (115), is required by Lemma C.3 and Theorem C.8. It follows from Lemma C.3 and Theorem C.8 that for every $t \in [T]$, conditioned on an event Ω with probability at least $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\hat{\varepsilon}_n^2))$ over the random noise w, we have $\mathbf{W}(t) \in \mathcal{W}(\mathbf{S}, \mathbf{W}(0), T)$, and

1774

$$f(\mathbf{W}(t), \cdot) = f_t \in \mathcal{F}_{NN}(\mathbf{S}, \mathbf{W}(0), T)$$

1775 Moreover, conditioned on the event Ω , 1776

1777

1781

 $f_t \in \mathcal{F}_{\text{ext}}(\mathbf{Q}, w, T).$

We then derive the sharp upper bound for $\mathbb{E}_P\left[(f_t - f^*)^2\right]$ by applying Theorem A.3 to the function class

$$\mathcal{F} = \left\{ F = \left(f - f^* \right)^2 : f \in \mathcal{F}(B_h, w) \right\}.$$

 $\Re\left(\{F \in \mathcal{F} \colon T(F) \le r\}\right)$

Let $B_0 \coloneqq B_h/\sqrt{2} + 1 + \mu_0/\sqrt{2} \ge B_h/\sqrt{2} + w + \mu_0/\sqrt{2}$, then we have $||F||_{\infty} \le B_0^2$ with $F \in \mathcal{F}$, so that $\mathbb{E}_P [F^2] \le B_0^2 \mathbb{E}_P [F]$. Let $T(F) = B_0^2 \mathbb{E}_P [F]$ for $F \in \mathcal{F}$. Then $\operatorname{Var} [F] \le \mathbb{E}_P [F^2] \le T(F) = B_0^2 \mathbb{E}_P [F]$.

1786 We have

$$= \Re\left(\left\{(f-f^*)^2 \colon f \in \mathcal{F}(B_h, w), \mathbb{E}_P\left[(f-f^*)^2\right] \leq \frac{r}{B_0^2}\right\}\right)$$

$$\stackrel{\textcircled{0}}{\leq} 2B_0 \Re\left(\left\{f-f^* \colon f \in \mathcal{F}(B_h, w), \mathbb{E}_P\left[(f-f^*)^2\right] \leq \frac{r}{B_0^2}\right\}\right)$$

$$\stackrel{\textcircled{0}}{\leq} 4B_0 \Re\left(\left\{f \in \mathcal{F}(B_h, w) \colon \mathbb{E}_P\left[f^2\right] \leq \frac{r}{4B_0^2}\right\}\right),$$
(118)

where ① is due to the contraction property of Rademacher complexity in Theorem A.2. Since $f^* \in \mathcal{F}(B_h, w), f \in \mathcal{F}(B_h, w)$, we have $\frac{f-f^*}{2} \in \mathcal{F}(B_h, w)$ due to the fact that $\mathcal{F}(B_h, w)$ is symmetric and convex, and it follows that ② holds.

1799 It follows from (118) and Lemma C.9 that

$$B_0^2 \Re\left(\{F \in \mathcal{F} \colon T(F) \le r\}\right) \le 4B_0^3 \Re\left(\left\{f \colon f \in \mathcal{F}(B_h, w), \mathbb{E}_P\left[f^2\right] \le \frac{r}{4B_0^2}\right\}\right)$$
$$\le 4B_0^3 \varphi_{B_h, w}\left(\frac{r}{4B_0^2}\right) \coloneqq \psi(r).$$
(119)

 ψ defined as the RHS of (119) is a sub-root function since it is nonnegative, nondecreasing and $\frac{\psi(r)}{\sqrt{r}}$ is nonincreasing. Let r^* be the fixed point of ψ , and $0 \le r \le r^*$. It follows from (Bartlett et al., 2005, Lemma 3.2) that $0 \le r \le \psi(r) = 4B_0^3 \varphi\left(\frac{r}{4B_0^2}\right)$. Therefore, by the definition of φ in (103), for every $0 \le Q \le n$, we have

$$\frac{r}{4B_0^3} \le \left(\frac{\sqrt{r}}{2B_0} + w\right) \sqrt{\frac{Q}{n}} + B_h \left(\frac{\sum\limits_{q=Q+1}^\infty \lambda_q}{n}\right)^{1/2} + w.$$
(120)

1816 Solving the quadratic inequality (120) for r, we have

$$r \leq \frac{8B_0^4Q}{n} + 8B_0^3 \left(w\left(\sqrt{\frac{Q}{n}} + 1\right) + B_h\left(\frac{\sum\limits_{q=Q+1}^\infty \lambda_q}{n}\right)^{1/2} \right).$$
(121)

(121) holds for every $0 \le Q \le n$, so we have

$$r \le 8B_0^3 \min_{0 \le Q \le n} \left(\frac{B_0Q}{n} + w \left(\sqrt{\frac{Q}{n}} + 1 \right) + B_h \left(\frac{\sum_{q=Q+1}^\infty \lambda_q}{n} \right)^{1/2} \right).$$
(122)

1829 It then follows from (119) and Theorem A.3 that with probability at least $1 - \exp(-x)$ over the random training features S, 1831 (x-2) = 2 = 2

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - \frac{K_{0}}{K_{0} - 1} \mathbb{E}_{P_{n}}\left[(f_{t} - f^{*})^{2}\right] - \frac{x\left(11B_{0}^{2} + 26B_{0}^{2}K_{0}\right)}{n} \le \frac{704K_{0}}{B_{0}^{2}}r^{*}, \quad (123)$$

 or

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - 2\mathbb{E}_{P_{n}}\left[(f_{t} - f^{*})^{2}\right] \lesssim r^{*} + \frac{x}{n},$$
(124)

with $K_0 = 2$ in (123).

It follows from (122) and (124) that

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - 2\mathbb{E}_{P_{n}}\left[(f_{t} - f^{*})^{2}\right]$$

$$\lesssim \min_{0 \le Q \le n} \left(\frac{B_{0}Q}{n} + w\left(\sqrt{\frac{Q}{n}} + 1\right) + B_{h}\left(\frac{\sum_{q=Q+1}^{\infty}\lambda_{q}}{n}\right)^{1/2}\right) + \frac{x}{n}.$$

Let $x = n\varepsilon_n^2$ in the above inequality, and we note that the above argument requires Theorem C.8 which holds with probability at least $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\widehat{\varepsilon}_n^2))$ over the random noise w. Then (116) is proved combined with the facts that $\Pr[W_0] \ge 1 - 2/n$.

1849 We now prove (117). First, it follows from the definition of $\varphi_{B_h,w}$ in (103) that

$$\psi(r) = 4B_0^3 \varphi_{B_h,w} \left(\frac{r}{4B_0^2}\right)$$

$$= 4B_0^3 \min_{Q: Q \ge 0} \left(\left(\frac{\sqrt{r}}{2B_0} + w\right) \sqrt{\frac{Q}{n}} + B_h \left(\frac{\sum_{q=Q+1}^\infty \lambda_q}{n}\right)^{1/2} \right) + 4B_0^3 w$$

 $\leq 4B_0^3 B_h \min_{Q: Q \ge 0} \left(\sqrt{\frac{Qr}{n}} + \left(\frac{\sum_{q=Q+1}^\infty \lambda_q}{n} \right)^{1/2} \right) + 4B_0^3 w \left(\sqrt{\frac{Q}{n}} + 1 \right)$

$$\leq \frac{4\sqrt{2}B_0^2 B_h}{\sigma} \cdot \sigma R_K(\sqrt{r}) + 8B_0^3 w \coloneqq \psi_1(r),$$

where the last inequality follows from the Cauchy-Schwarz inequality. It can be verified that $\psi_1(r)$ is a sub-root function. Let the fixed point of $\psi_1(r)$ be r_1^* . Because the fixed point of $\sigma R_K(\sqrt{r})$ as a function of r is ε_n^2 , it follows from Lemma C.17 that

$$r_1^* \le \max\left\{\frac{32\sqrt{2}B_0^6 B_h^2}{\sigma^2}, 1\right\}\varepsilon_n^2 + 16B_0^3 w.$$
(125)

1/2

1871 It then follows from Theorem A.3 with $K_0 = 2$ that with probability at least $1 - \exp(-x)$,

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - 2\mathbb{E}_{P_{n}}\left[(f_{t} - f^{*})^{2}\right] \lesssim r_{1}^{*} + \frac{x}{n}.$$

1874 Letting $x = n\varepsilon_n^2$, then plugging the upper bound for r_1^* , (125), in the above inequality leads to

$$\mathbb{E}_{P}\left[(f_{t} - f^{*})^{2}\right] - 2\mathbb{E}_{P_{n}}\left[(f_{t} - f^{*})^{2}\right] \lesssim \varepsilon_{n}^{2} + 16B_{0}^{3}w.$$
(126)

Again, we note that the above argument requires Theorem C.8 which holds with probability at least 1878 $1 - \exp(-\Theta(n)) - \exp(-\Theta(n\hat{\varepsilon}_n^2))$ over the random noise w. Then (117) is proved with the fact 1879 that $\Pr[W_0] \ge 1 - 2/n$ and (126).

 Theorem C.11. Suppose the neural network trained after the *t*-th step of gradient descent, $f_t = f(\mathbf{W}(t), \cdot)$, satisfies $\mathbf{u}(t) = f_t(\mathbf{S}) - \mathbf{y} = \mathbf{v}(t) + \mathbf{e}(t)$ with $\mathbf{v}(t) \in \mathcal{V}_t$ and $\mathbf{e}(t) \in \mathcal{E}_{t,\tau}$ and $T \leq \widehat{T}$. If

$$\eta \in [1,2), \quad \tau \le \frac{1}{\eta T},\tag{127}$$

then for every $t \in [T]$, with probability at least $1 - \exp\left(-\Theta(n\hat{\varepsilon}_n^2)\right)$ over the random noise w, we have

$$\mathbb{E}_{P_n}\left[(f_t - f^*)^2\right] \le \frac{3}{\eta t} \left(\frac{\mu_0^2}{2e} + 3\right).$$
(128)

Proof. We have

$$f_t(\mathbf{S}) = f^*(\mathbf{S}) + \mathbf{w} + \mathbf{v}(t) + \mathbf{e}(t), \qquad (129)$$

where $\mathbf{v}(t) \in \mathcal{V}_t$, $\mathbf{e}(t) \in \mathcal{E}_{t,\tau}$, $\mathbf{e}(t) = \mathbf{e}_1(t) + \mathbf{e}_2(t)$ with $\mathbf{e}_1(t) = -(\mathbf{I}_n - \eta \mathbf{K}_n)^t$ w and $\|\mathbf{e}_2(t)\|_2 \leq \sqrt{n\tau}$. We have $\eta \lambda_1 \in (0,1)$ if $\eta \in [1,2)$. It follows from (129) that

$$\mathbb{E}_{P_n} \left[(f_t - f^*)^2 \right] = \frac{1}{n} \| f_t(\mathbf{S}) - f^*(\mathbf{S}) \|_2^2 = \frac{1}{n} \| \mathbf{v}(t) + \mathbf{w} + \mathbf{e}(t) \|_2^2$$

= $\frac{1}{n} \left\| - (\mathbf{I} - \eta \mathbf{K}_n)^t f^*(\mathbf{S}) + \left(\mathbf{I}_n - (\mathbf{I}_n - \eta \mathbf{K}_n)^t \right) \mathbf{w} + \vec{\mathbf{e}}_2(t) \right\|_2^2$
 $\stackrel{(1)}{\leq} \frac{3}{n} \sum_{i=1}^n \left(1 - \eta \widehat{\lambda}_i \right)^{2t} \left[\mathbf{U}^\top f^*(\mathbf{S}) \right]_i^2 + \frac{3}{n} \sum_{i=1}^n \left(1 - \left(1 - \eta \widehat{\lambda}_i \right)^t \right)^2 \left[\mathbf{U}^\top \mathbf{w} \right]_i^2 + \frac{3}{n} \left\| \vec{\mathbf{e}}_2(t) \right\|_2^2$

$$\overset{\textcircled{0}}{\leq} \frac{3\mu_{0}^{2}}{2e\eta t} + \frac{3}{n} \sum_{i=1}^{n} \left(1 - (1 - \eta\lambda_{i})^{t} \right)^{2} \left[\mathbf{U}^{\top} \mathbf{w} \right]_{i}^{2} + 3\tau^{2} \\
\leq \frac{3}{\eta t} \left(\frac{\mu_{0}^{2}}{2e} + \frac{1}{\eta} \right) + 3 \cdot \underbrace{\frac{1}{n} \sum_{i=1}^{n} \left(1 - (1 - \eta\lambda_{i})^{t} \right)^{2} \left[\mathbf{U}^{\top} \mathbf{w} \right]_{i}^{2}}_{:=E_{\varepsilon}} \\
\leq \frac{3}{\eta t} \left(\frac{\mu_{0}^{2}}{2e} + 2\widehat{\lambda}_{1} \right) + 3E_{\varepsilon} \leq \frac{3}{\eta t} \left(\frac{\mu_{0}^{2}}{2e} + 4 \right) + 3E_{\varepsilon}.$$
(130)

Here ① follows from the Cauchy-Schwarz inequality, ② follows from (49) in the proof of Lemma C.4. We then derive the upper bound for E_{ε} on the RHS of (130). We define the diagonal matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$ with $\mathbf{R}_{ii} = \left(1 - (1 - \eta \lambda_i)^t\right)^2$. Then we have

$$E_{\varepsilon} = 1/n \cdot \operatorname{tr} \left(\mathbf{U} \mathbf{R} \mathbf{U}^{\top} \mathbf{w} \mathbf{w}^{\top} \right)$$

1920 It follows from (Wright, 1973) that

$$\Pr\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right) - \mathbb{E}\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right)\right] \ge u\right]$$
$$\le \exp\left(-c\min\left\{nu/\|\mathbf{R}\|_{2}, n^{2}u^{2}/\|\mathbf{R}\|_{F}^{2}\right\}\right).$$
(131)

for all u > 0, and c is a positive constant. With $\eta_t = \eta t$ for all $t \ge 0$, we have

$$\mathbb{E}\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right)\right] \leq \frac{\sigma^{2}}{n} \sum_{i=1}^{n} \left(1 - \left(1 - \eta\widehat{\lambda}_{i}\right)^{t}\right)^{2}$$
$$\stackrel{\text{(I)}}{\leq} \frac{\sigma^{2}}{n} \sum_{i=1}^{n} \min\left\{1, \eta_{t}^{2}\widehat{\lambda}_{i}^{2}\right\}$$

 $\stackrel{\textcircled{0}}{\leq} \frac{\sigma^2 \eta_t}{n} \sum_{i=1}^n \min\left\{\frac{1}{\eta_t}, \widehat{\lambda}_i\right\} \\ = \sigma^2 \eta_t \widehat{R}_K^2(\sqrt{1/\eta_t}) \le \frac{1}{\eta_t}.$

 $\leq \frac{\sigma^2 \eta_t}{n} \sum_{i=1}^n \min\left\{\frac{1}{\eta_t}, \eta_t \widehat{\lambda}_i^2\right\}$

(132)

Here ① follows from the fact that $(1 - \eta \hat{\lambda}_i)^t \ge \max\left\{0, 1 - t\eta \hat{\lambda}_i\right\}$, and ② follows from $\min\left\{a, b\right\} \le \sqrt{ab}$ for any nonnegative numbers a, b. Because $t \le T \le \hat{T}$, we have $R_K(\sqrt{1/\eta_t}) \le 1/(\sigma\eta_t)$, so the last inequality holds.

Moreover, we have the upper bounds for $\|\mathbf{R}\|_2$ and $\|\mathbf{R}\|_F$ as follows. First, we have

1946
1947
1948
1949
1950

$$\|\mathbf{R}\|_{2} \leq \max_{i \in [n]} \left(1 - \left(1 - \eta \widehat{\lambda}_{i}\right)^{t}\right)^{2}$$

$$\leq \min\left\{1, \eta_{t}^{2} \widehat{\lambda}_{i}^{2}\right\} \leq 1.$$

We also have

$$\|\mathbf{R}\|_{\mathrm{F}}^{2} = \frac{1}{n} \sum_{i=1}^{n} \left(1 - \left(1 - \eta \widehat{\lambda}_{i} \right)^{t} \right)^{4}$$

$$\leq \frac{\eta_{t}}{n} \sum_{i=1}^{n} \min\left\{ \frac{1}{\eta_{t}}, \eta_{t}^{3} \widehat{\lambda}_{i}^{4} \right\}$$

$$\stackrel{(3)}{\leq} \frac{\eta_{t}}{n} \sum_{i=1}^{n} \min\left\{ \widehat{\lambda}_{i}, \frac{1}{\eta_{t}} \right\} = \eta_{t} \widehat{R}_{K}^{2}(\sqrt{1/\eta_{t}}) \leq \frac{1}{\sigma^{2} \eta_{t}}.$$
(134)

1961 If $1/\eta_t \leq \eta_t^3(\widehat{\lambda}_i)^4$, then $\min\left\{1/\eta_t, \eta_t^3(\widehat{\lambda}_i)^4\right\} = 1/\eta_t$. Otherwise, we have $\eta_t^4 \widehat{\lambda}_i^4 < 1$, so that 1962 $\eta_t \widehat{\lambda}_i < 1$ and it follows that $\min\left\{1/\eta_t, \eta_t^3(\widehat{\lambda}_i)^4\right\} \leq \eta_t^3 \widehat{\lambda}_i^4 \leq \widehat{\lambda}_i$. As a result, ③ holds.

¹⁹⁶⁴ Combining (131)- (134), we have

1966
$$\Pr\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right) - \mathbb{E}\left[1/n \cdot \operatorname{tr}\left(\mathbf{U}\mathbf{R}\mathbf{U}^{\top}\mathbf{w}\mathbf{w}^{\top}\right)\right] \ge u\right] \le \exp\left(-cn\min\left\{u, u^{2}\sigma^{2}\eta_{t}\right\}\right).$$
1967 Let $u = 1/n$ in the above inequality, we have

Let
$$u = 1/\eta_t$$
 in the above inequality, we have

 \overline{n}

$$\exp\left(-cn\min\left\{u, u^2\sigma^2\eta_t\right\}\right) = \exp\left(-c'n/\eta_t\right) \le \exp\left(-c'n\hat{\varepsilon}_n^2\right)$$

1971 where $c' = c \min\{1, \sigma^2\}$, and the last inequality is due to the fact that $1/\eta_t \ge \hat{\varepsilon}_n^2$ since $t \le T \le \hat{T}$. 1972 It follows that with probability at least $1 - \exp(-\Theta(n\hat{\varepsilon}_n^2))$,

$$E_{\varepsilon} \le u + \frac{1}{\eta_t} = \frac{2}{\eta_t}.$$
(135)

1976 It then follows from (130), (131)-(135) that

$$\mathbb{E}_{P_n}\left[(f_t - f^*)^2\right] \le \frac{3}{\eta t} \left(\frac{\mu_0^2}{2e} + 6\right)$$

holds with probability at least $1 - \exp\left(-c'n\widehat{\varepsilon}_n^2\right)$.

C.3 AUXILIARY RESULTS ABOUT REPRODUCING KERNEL HILBERT SPACES

Lemma C.12 (In the proof of (Raskutti et al., 2014, Lemma 8)). For any $f \in \mathcal{H}_K(\mu_0)$, we have

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\left[\mathbf{U}^{\top}f(\mathbf{S}')\right]_{i}^{2}}{\widehat{\lambda}_{i}} \le \mu_{0}^{2}.$$
(136)

Similarly, for $f \in \mathcal{H}_K(\mu_0)$, we have $\frac{1}{n} \sum_{i=1}^n \frac{\left[\mathbf{U}^\top f(\mathbf{S}')\right]_i^2}{\lambda_i} \le \mu_0^2$. Here, C.13. For any positive real number $a \in (0, 1)$ and natural number

Lemma C.13. For any positive real number $a \in (0, 1)$ and natural number t, we have

$$(1-a)^t \le e^{-ta} \le \frac{1}{eta}.$$
(137)

1997 Proof. The result follows from the facts that $\log(1-a) \leq a$ for $a \in (0,1)$ and $\sup_{u \in \mathbb{R}} ue^{-u} \leq 1/e$.

(133)

1998 Lemma C.14. ((Rosasco et al., 2010, Proposition 10)) With probability $1 - \delta$ over the training data S, for all $j \in [n]$,

$$\left|\lambda_j - \widehat{\lambda}_j\right| \le \sqrt{\frac{2\log\frac{2}{\delta}}{n}}.$$
(138)

Lemma C.15. With probability at least $1 - 2 \exp(-\Theta(n\varepsilon_n^2))$,

$$\varepsilon_n^2 \le c_1 \widehat{\varepsilon}_n^2. \tag{139}$$

Furthermore, with probability at least $1 - 2 \exp(-\Theta(n\varepsilon_n^2))$,

$$\widehat{\varepsilon}_n^2 \le c_1 \varepsilon_n^2. \tag{140}$$

Here c_1 is an absolute positive constant depending on σ .

Remark. Lemma C.15 shows that with probability at least $1 - 4 \exp(-\Theta(n\varepsilon_n^2))$, $\varepsilon_n^2 \simeq \hat{\varepsilon}_n^2$, which is also a fact used in kernel complexity or local Rademacher based analysis for kernel regression in the statistical learning literature. We herein provide a detailed proof to ensure the mathematical rigor of this paper.

Proof. Define function classes

$$\mathcal{F}_{t} := \left\{ f \in \mathcal{H}_{K} : \|f\|_{\mathcal{H}_{K}} \le 1, \|f\|_{L^{2}} \le t \right\}, \quad \widehat{\mathcal{F}}_{t} := \left\{ f \in \mathcal{H}_{K} : \|f\|_{\mathcal{H}_{K}} \le 1, \|f\|_{n} \le t \right\},$$

where $||f||_n^2 \coloneqq 1/n \cdot \sum_{i=1}^n f^2(\vec{\mathbf{x}}_i)$. Let $\mathcal{R}(t)$ be the Rademacher complexity of \mathcal{F}_t , that is,

$$\mathcal{R}(t) = \mathfrak{R}(\mathcal{F}_t) = \mathbb{E}_{\left\{\vec{\mathbf{x}}_i\right\}, \left\{\sigma_i\right\}} \left[\sup_{f \in \mathcal{F}_t} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\vec{\mathbf{x}}_i), \right]$$

and we will also write $\mathcal{R}(t) = \mathbb{E}\left[\sup_{f \in \mathcal{F}_t} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\vec{\mathbf{x}}_i), \right]$ for simplicity of notations. We let $\widehat{\mathcal{R}}(t)$ be the empirical Rademacher complexity of \mathcal{F}_t , that is,

$$\widehat{\mathcal{R}}(t) = \mathbb{E}_{\sigma} \left[\sup_{f \in \widehat{\mathcal{F}}_t} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\vec{\mathbf{x}}_i) \right]$$

By results of (Mendelson, 2002), there are universal constants c_{ℓ} and C_u with $0 < c_{\ell} < C_u$ such that when $t^2 \ge 1/n$, we have

$$c_{\ell}R_K(t) \le \mathcal{R}(t) \le C_uR_K(t), \quad c_{\ell}\widehat{R}_K(t) \le \widehat{\mathcal{R}}(t) \le C_u\widehat{R}_K(t).$$
 (141)

When $f \in \mathcal{F}_t$, $||f||_{\infty} \leq \tau_0 = \frac{1}{\sqrt{2}}$. It follows from Lemma C.16 that with probability at least $1 - \exp(-n\varepsilon_n^2)$,

$$\mathcal{F}_t \subseteq \left\{ f \in \mathcal{H}_K \colon \left\| f \right\|_{\mathcal{H}_K} \le 1, \left\| f \right\|_n \le \sqrt{c_2 t^2 + c_3 \varepsilon_n^2} \right\} \coloneqq \widehat{\mathcal{F}}_{\sqrt{c_2 t^2 + c_3 \varepsilon_n^2}}.$$
 (142)

Moreover, by the relation between Rademacher complexity and its empirical version in (Bartlett et al., 2005, Lemma A.4), for every x > 0, with probability at least $1 - \exp(-x)$,

$$\mathbb{E}\left[\sup_{f\in\widehat{\mathcal{F}}_{\sqrt{c_{2}t^{2}+c_{3}\varepsilon_{n}^{2}}}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}f(\vec{\mathbf{x}}_{i})\right] \leq 2\mathbb{E}_{\sigma}\left[\sup_{f\in\widehat{\mathcal{F}}_{\sqrt{c_{2}t^{2}+c_{3}\varepsilon_{n}^{2}}}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}f(\vec{\mathbf{x}}_{i})\right] + \frac{2\tau_{0}x}{n}.$$
 (143)

As a result,

2049
2050
2051
$$\mathcal{R}(t) \stackrel{\textcircled{1}}{\leq} \mathbb{E} \left[\sup_{f \in \widehat{\mathcal{F}}_{\sqrt{c_2 t^2 + c_3 \varepsilon_n^2}}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\vec{\mathbf{x}}_i) \right]$$

2052
2053
2054
2055
2056
2057

$$\overset{\textcircled{O}}{=} 2\mathbb{E}_{\sigma} \left[\sup_{f \in \widehat{\mathcal{F}}_{\sqrt{c_{2}t^{2} + c_{3}\varepsilon_{n}^{2}}}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} f(\vec{\mathbf{x}}_{i}) \right] + \frac{2\tau_{0}x}{n} = 2\widehat{\mathcal{R}}(\sqrt{c_{2}t^{2} + c_{3}\varepsilon_{n}^{2}}) + \frac{2\tau_{0}x}{n}.$$

Here ① follows from (142), and ② follows from (143). It follows from (141) and the above inequal-ity that

$$c_{\ell}/\sigma \cdot \sigma R_K(t) \le 2C_u/\sigma \cdot \sigma \widehat{R}_K(\sqrt{c_2 t^2 + c_3 \varepsilon_n^2}) + \frac{2\tau_0 x}{n}, \forall t^2 \ge 1/n.$$

Rewrite $R_K(t)$ as a function of $r = t^2$ as $R_K(t) = F_K(r)$. Similarly, $\widehat{R}_K(t) = \widehat{F}_K(r)$ with $r = t^2$. Then we have

$$\sigma F_K(r) \le \max\left\{2C_u/c_\ell, 1\right\} \cdot \sigma \widehat{F}_K(c_2r + c_3\varepsilon_n^2) + \frac{2\sigma\tau_0 x}{nc_\ell} \coloneqq G(r), \forall r \ge 1/n.$$
(144)

It can be verified that G(r) is a sub-root function, and let r_G^* be the fixed point of G. Let $x \ge 1$ $c_{\ell}/(2\sigma\tau_0)$, then $r_G^* \ge 1/n$. Moreover, $\sigma F_K(r)$ and $\sigma \widehat{F}_K(r)$ are sub-root functions, and they have fixed points ε_n^2 and $\widehat{\varepsilon}_n^2$, respectively. Set $r = r_G^* \ge 1/n$ in (144), we have

$$\sigma F_K(r_G^*) \le r_G^*,$$

and it follows from the above inequality and (Bartlett et al., 2005, Lemma 3.2) that $\varepsilon_n^2 \leq r_G^2$. Since $c_2 > 1$, it then follows from the properties about the fixed point of a sub-root function in Lemma C.17 that

$$\varepsilon_n^2 \le r_G^* \le \max\left\{2C_u/c_\ell, 1\right\}^2 \left(c_2\widehat{\varepsilon}_n^2 + \frac{2c_3\varepsilon_n^2}{c_2}\right) + \frac{4\sigma\tau_0 x}{nc_\ell}.$$

We can choose c_2 such that $c_2 > 2c_3 \max \{2C_u/c_\ell, 1\}^2$, then the above inequality indicates that

$$\varepsilon_n^2 \le c_{u,\ell} \widehat{\varepsilon}_n^2 + \frac{4\sigma \tau_0 x}{nc_\ell},$$

where $c_{u,\ell}$ is a constant depending on c_{ℓ} , C_u , c_2 , c_3 , and (139) is proved with $x = c' n \varepsilon_n^2$ where c' > 0 is a positive constant which is chosen such that $4c' \sigma \tau_0 / c_\ell < 1$.

Similarly, it follows from Lemma C.16 that with probability at least $1 - \exp(-n\varepsilon_n^2)$,

$$\widehat{\mathcal{F}}_t \subseteq \left\{ f \in \mathcal{H}_K \colon \|f\|_{\mathcal{H}_K} \le 1, \|f\|_{L^2} \le \sqrt{c_2 t^2 + c_3 \varepsilon_n^2} \right\} = \mathcal{F}_{\sqrt{c_2 t^2 + c_3 \varepsilon_n^2}}.$$
(145)

It follows from (Bartlett et al., 2005, Lemma A.4) again that for every x > 0, with probability at least $1 - \exp(-x)$,

$$\mathbb{E}_{\sigma}\left[\sup_{f\in\mathcal{F}_{\sqrt{c_{2}t^{2}+c_{3}\varepsilon_{n}^{2}}}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}f(\vec{\mathbf{x}}_{i})\right] \leq 2\mathbb{E}\left[\sup_{f\in\mathcal{F}_{\sqrt{c_{2}t^{2}+c_{3}\varepsilon_{n}^{2}}}}\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}f(\vec{\mathbf{x}}_{i})\right] + \frac{10\tau_{0}x}{12n}.$$
 (146)

As a result, we have

$$\widehat{\mathcal{R}}(t) \stackrel{\textcircled{1}}{\leq} \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}_{\sqrt{c_2 t^2 + c_3 \varepsilon_n^2}}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\vec{\mathbf{x}}_i) \right]$$

 $\overset{\textcircled{0}}{\leq} 2\mathbb{E} \left[\sup_{f \in \mathcal{F}_{\sqrt{c_2 t^2 + c_3 \varepsilon_n^2}}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(\vec{\mathbf{x}}_i) \right] + \frac{10\tau_0 x}{12n}$

$$= \int_{f \in \mathcal{F}_{\sqrt{c_2 t^2 + c_2 \varepsilon_2^2}}} n$$

2105
$$= 2\mathcal{R}(\sqrt{c_2t^2 + c_3\varepsilon_n^2}) + \frac{5\sqrt{2x}}{12n} \le 2C_u R_K(\sqrt{c_2t^2 + c_3\varepsilon_n^2}) + \frac{10\tau_0 x}{12n},$$

where ① follows from (145), and ② follows from (146). Using a similar argument for the proof of the first inequality in (139), we have

 $\widehat{\varepsilon}_n^2 \le r_G^* \le \max\left\{2C_u/c_\ell, 1\right\}^2 \left(c_2 + \frac{2c_3}{c_2}\right)\varepsilon_n^2 + \frac{10\tau_0 x}{12n},$

and the second inequality in (140) is approved with $x = \Theta(n\varepsilon_n^2)$.

Lemma C.16. Let K be a PSD kernel, then with probability at least $1 - \exp(-n\varepsilon_n^2)$,

$$|g||_{L^2}^2 \le c_2 ||g||_n^2 + c_3 \varepsilon_n^2, \quad \forall g \in \mathcal{H}_K(1).$$
(147)

Furthermore, with probability at least $1 - \exp(-n\varepsilon_n^2)$,

$$\left\|g\right\|_{n}^{2} \leq c_{2}\left\|g\right\|_{L^{2}}^{2} + c_{3}\varepsilon_{n}^{2}, \quad \forall g \in \mathcal{H}_{K}(1).$$

$$(148)$$

Here c_2, c_3 are positive constants with $c_2 > 1$.

2126 *Proof.* The results follow from Theorem A.1.

Lemma C.17. Suppose $\psi: [0, \infty) \to [0, \infty)$ is a sub-root function with the unique fixed point r^* .

Then the following properties hold.

2131 (1) Let $a \ge 0$, then $\psi(r) + a$ as a function of r is also a sub-root function with fixed point r_a^* , and 2132 $r^* \le r_a^* \le r^* + 2a$.

2133 2134 (2) Let $b \ge 1$, $c \ge 0$ then $\psi(br + c)$ as a function of r is also a sub-root function with fixed point r_b^* , 2135 and $r_b^* \le br^* + 2c/b$.

2136 (3) Let $b \ge 1$, then $\psi_b(r) = b\psi(r)$ is also a sub-root function with fixed point r_b^* , and $r_b^* \le b^2 r^*$. 2137 2138

2139 Proof. (1). Let $\psi_a(r) = \psi(r) + a$. It can be verified that $\psi_a(r)$ is a sub-root function because its **2140** nonnegative, nondecreasing and $\psi_a(r)/\sqrt{r}$ is nonincreasing. It follows from (Bartlett et al., 2005, **2141** Lemma 3.2) that ψ_a has unique fixed point denoted by r_a^* . Because $r^* = \psi(r^*) \le \psi(r^*) + a =$ **2142** $\psi_a(r^*)$, it follows from (Bartlett et al., 2005, Lemma 3.2) that $r^* \le r_a^*$. Furthermore, since

2143 2144

2145 2146

2109

2110 2111

2112 2113

2114 2115

2116 2117 2118

2121 2122

2124 2125

2127 2128

2130

$$\psi_a(r^* + 2a) = \psi(r^* + 2a) + a \le \psi(r^*)\sqrt{\frac{r^* + 2a}{r^*}} + a \le \sqrt{r^*(r^* + 2a)} + a \le r^* + 2a,$$

it follows from (Bartlett et al., 2005, Lemma 3.2) again that $r_a^* \le r^* + 2a$.

(2). Let $\psi_b(r) = \psi(br+c)$. It can be verified that $\psi_b(r)$ a sub-root function by checking the definition. Also, we have $\psi(b(br^*+2c/b)+c)/\sqrt{b(br^*+2c/b)+c} \le \psi(r^*)/\sqrt{r^*}$. It follows that

$$\psi_b\left(br^* + \frac{2c}{b}\right) = \psi\left(b\left(br^* + \frac{2c}{b}\right) + c\right) \le b\sqrt{\left(r^* + \frac{3c}{b^2}\right)r^*} \le b\left(r^* + \frac{3c}{2b^2}\right) \le br^* + \frac{2c}{b}.$$

2154 2155

2152 2153

Then it follows from (Bartlett et al., 2005, Lemma 3.2) that $r_b^* \leq br^* + 2c/b$.

(3). Let $\psi_b(r) = b\psi(r)$. It can be verified that $\psi_b(r)$ a sub-root function by checking the definition. Also, we have $\psi(b^2r^*)/\sqrt{b^2r^*} \le \psi(r^*)/\sqrt{r^*}$, so $\psi(b^2r^*) \le br^*$ and $\psi_b(b^2r^*) = b\psi(b^2r^*) \le b^2r^*$. Then it follows from (Bartlett et al., 2005, Lemma 3.2) that $r_b^* \le b^2r^*$.

2160 C.4 PROOFS OF THEOREM C.1 AND THEOREM C.2 2161

²¹⁶² We need the following definition of ε -net for the proof of Theorem C.1 and Theorem C.2.

2163 Definition C.1. (ε -net) Let (X, d) be a metric space and let $\varepsilon > 0$. A subset $N_{\varepsilon}(X, d)$ is called an ε -**2164** net of X if for every point $x \in X$, there exists some point $y \in N_{\varepsilon}(X, d)$ such that $d(x, y) \le \varepsilon$. The **2165** minimal cardinality of an ε -net of X, if finite, is denoted by $N(X, d, \varepsilon)$ and is called the covering **2166** number of X at scale ε .

Proof of Theorem C.1. First, we have $\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)} [h(\mathbf{w}, \mathbf{x}, \mathbf{y})] = K(\mathbf{x}, \mathbf{y})$. For any $\mathbf{x} \in \mathcal{X}$ and s > 0, define function class

2170 2171 2172

2167

$$\mathcal{H}_{\mathbf{x},s} \coloneqq \left\{ h(\cdot, \mathbf{x}', \mathbf{y}) \colon \mathbb{R}^d \to \mathbb{R} \colon \mathbf{x}' \in \mathbf{B} \left(\mathbf{x}; s \right) \cap \mathcal{X}, \mathbf{y} \in \mathcal{X} \right\}.$$
(149)

We first build an *s*-net for the unit sphere \mathcal{X} . By (Vershynin, 2012, Lemma 5.2), there exists an *s*-net $N_s(\mathcal{X}, \|\cdot\|_2)$ of \mathcal{X} such that $N(\mathcal{X}, \|\cdot\|_2, s) \le \left(1 + \frac{2}{s}\right)^d$.

2176 In the sequel, a function in the class $\mathcal{H}_{\mathbf{x},s}$ is also denoted as $h(\mathbf{w})$, omitting the presence of variables 2177 \mathbf{x}' and \mathbf{y} when no confusion arises. Let P_m be the empirical distribution over $\left\{ \vec{\mathbf{w}}_r(0) \right\}$ so that 2178 $\mathbb{E}_{\mathbf{w}\sim P_m}[h(\mathbf{w})] = \widehat{h}(\mathbf{W}(0), \mathbf{x}, \mathbf{y}).$ Given $\mathbf{x} \in N(\mathcal{X}, s)$, we aim to estimate the upper bound for 2179 the supremum of empirical process $\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)} [h(\mathbf{w})] - \mathbb{E}_{\mathbf{w} \sim P_m} [h(\mathbf{w})]$ when function h ranges 2180 over the function class $\mathcal{H}_{\mathbf{x},s}$. To this end, we apply Theorem A.1 to the function class $\mathcal{H}_{\mathbf{x},s}$ with 2181 $\mathbf{W}(0) = \left\{ \overrightarrow{\mathbf{w}}_{r}(0) \right\}_{r=1}^{m}$. It can be verified that $h \in [0,1]$ for any $h \in \mathcal{H}_{\mathbf{x},s}$. It follows that we can 2182 2183 set a = 0, b = 1 in Theorem A.1. With probability at least $1 - 2e^{-x}$ over the random initialization 2184 W(0),2185

2186

$$\begin{aligned}
& \sum_{h \in \mathcal{H}_{\mathbf{x},s}} \sup_{h \in \mathcal{H}_{\mathbf{x},s}} \left| \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0},\kappa^{2}\mathbf{I}_{d})} \left[h(\mathbf{w}) \right] - \mathbb{E}_{\mathbf{w} \sim P_{m}} \left[h(\mathbf{w}) \right] \right| \\
& \sum_{h \in \mathcal{H}_{\mathbf{x},s}} \sup_{h \in \mathcal{H}_{\mathbf{x},s}} \left| \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0},\kappa^{2}\mathbf{I}_{d})} \left[h(\mathbf{w}) \right] - \mathbb{E}_{\mathbf{w} \sim P_{m}} \left[h(\mathbf{w}) \right] \right| \\
& \sum_{\alpha \in (0,1)} \left(2(1+\alpha) \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{h \in \mathcal{H}_{\mathbf{x},s}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} h(\vec{\mathbf{w}}_{r}(0)) \right] + \sqrt{\frac{2rx}{m}} + (b-a) \left(\frac{1}{3} + \frac{1}{\alpha} \right) \frac{x}{m} \right) \\
& (150)
\end{aligned}$$

2193 where $\{\sigma\}_{r=1}^{m}$ are i.i.d. Rademacher random variables taking values of ± 1 with equal probability. 2194 It can be verified that $\operatorname{Var}[h] \leq \mathbb{E}_{\mathbf{w}} \left[h(\mathbf{w}, \mathbf{x}', \mathbf{y})^2\right] \leq 1$. Setting $\alpha = \frac{1}{2}$ in (150), it follows that with 2196 probability at least $1 - \delta$,

2197 2198

> 2199 2200 2201

$$\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s)\cap\mathcal{X},\mathbf{y}\in\mathcal{X}} \left| K(\mathbf{x}',\mathbf{y}) - \widehat{h}(\mathbf{W}(0),\mathbf{x}',\mathbf{y}) \right| \le 3\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) + \sqrt{\frac{2\log\frac{2}{\delta}}{m}} + \frac{7\log\frac{2}{\delta}}{3m}.$$
 (151)

Here $\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) = \mathbb{E}_{\mathbf{W}(0),\{\sigma_r\}_{r=1}^m} \left[\sup_{h \in \mathcal{H}_{\mathbf{x},s}} \frac{1}{m} \sum_{r=1}^m \sigma_r h(\vec{\mathbf{w}}_r(0)) \right]$ is the Rademacher complexity of the function class $\mathcal{H}_{\mathbf{x},s}$. By Lemma C.18, $\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) \leq \frac{1}{\sqrt{m}} + B\sqrt{ds}(s+1) + \sqrt{s} + s$. Plugging such upper bound for $\mathcal{R}(\mathcal{H}_{\mathbf{x},s})$ in (151), we have

$$\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s)\cap\mathcal{X},\mathbf{y}\in\mathcal{X}} \left| K(\mathbf{x}',\mathbf{y}) - \hat{h}(\mathbf{W}(0),\mathbf{x}',\mathbf{y}) \right| \\ \leq 3\left(\frac{1}{\sqrt{m}} + B\sqrt{ds}(s+1) + \sqrt{s} + s\right) + \sqrt{\frac{2\log\frac{2}{\delta}}{m}} + \frac{7\log\frac{2}{\delta}}{3m}.$$
(152)

Setting $s = \frac{1}{m}$, we have

$$\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s)\cap\mathcal{X},\mathbf{y}\in\mathcal{X}} \left| K(\mathbf{x}',\mathbf{y}) - \hat{h}(\mathbf{W}(0),\mathbf{x}',\mathbf{y}) \right| \\
\leq 3\left(\frac{1}{\sqrt{m}} + \frac{B\sqrt{d}\left(1+\frac{1}{m}\right)}{\sqrt{m}} + \frac{1}{\sqrt{m}} + \frac{1}{m} \right) + \sqrt{\frac{2\log\frac{2}{\delta}}{m}} + \frac{7\log\frac{2}{\delta}}{3m} \\
\leq \frac{1}{\sqrt{m}} \left(6(1+B\sqrt{d}) + \sqrt{2\log\frac{2}{\delta}} \right) + \frac{1}{m} \left(3 + \frac{7\log\frac{2}{\delta}}{3}\right).$$
(153)

By union bound, with probability at least $1 - (1 + 2m)^d \delta$ over $\mathbf{W}(0)$, (153) holds for arbitrary $\mathbf{x} \in N(\mathcal{X}, s)$. In this case, for any $\mathbf{x}' \in \mathcal{X}, \mathbf{y} \in \mathcal{X}$, there exists $\mathbf{x} \in N_s(\mathcal{X}, \|\cdot\|_2)$ such that $\|\mathbf{x}' - \mathbf{x}\|_2 \leq s$, so that $\mathbf{x}' \in \mathbf{B}(\mathbf{x}; s) \cap \mathcal{X}$, and (153) holds. Changing the notation \mathbf{x}' to \mathbf{x} , the conclusion is proved.

Lemma C.18. Let $\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) \coloneqq \mathbb{E}_{\mathbf{W}(0),\{\sigma_r\}_{r=1}^m} \left[\sup_{h \in \mathcal{H}_{\mathbf{x},s}} \frac{1}{m} \sum_{r=1}^m \sigma_r h(\vec{\mathbf{w}}_r(0)) \right]$ be the Rademacher complexity of the function class $\mathcal{H}_{\mathbf{x},s}$, B is a positive constant. Then

$$\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) \le \frac{1}{\sqrt{m}} + B\sqrt{ds}(s+1) + \sqrt{s} + s.$$
(154)

Proof. We have

$$\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) = \mathbb{E}_{\mathbf{W}(0),\{\sigma_r\}_{r=1}^m} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s), \mathbf{y} \in \mathcal{X}} \frac{1}{m} \sum_{r=1}^m \sigma_r h(\vec{\mathbf{w}}_r(0), \mathbf{x}', \mathbf{y}) \right]$$

$$\leq \mathcal{R}_1 + \mathcal{R}_2, \tag{155}$$

where

$$\mathcal{R}_{1} = \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s), \mathbf{y} \in \mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right],$$
$$\mathcal{R}_{2} = \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s), \mathbf{y} \in \mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \left(h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}', \mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right) \right].$$
(156)

Here (155) follows from the subadditivity of superemum and the fact that $\sum_{r=1}^{m} \sigma_r h(\vec{\mathbf{w}}_r(0), \mathbf{x}', \mathbf{y}) =$ $\sum_{r=1}^{m} \sigma_r h(\vec{\mathbf{w}}_r(0), \mathbf{x}, \mathbf{y}) + \sum_{r=1}^{m} \sigma_r \Big(h(\vec{\mathbf{w}}_r(0), \mathbf{x}', \mathbf{y}) - h(\vec{\mathbf{w}}_r(0), \mathbf{x}, \mathbf{y}) \Big).$

Now we bound \mathcal{R}_1 and \mathcal{R}_2 separately. For \mathcal{R}_1 , we have

$$\mathcal{R}_{1} = \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s), \mathbf{y} \in \mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right]$$

$$\stackrel{\textcircled{1}}{=} \mathbb{E}_{\mathbf{W}(0)} \left[\mathbb{E}_{\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{y} \in \mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \mathbf{x}^{\top} \mathbf{y} \mathbb{I}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{x} \geq 0\}} \mathbb{I}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{y} \geq 0\}} \right] \right]$$

$$= \mathbb{E}_{\mathbf{W}(0)} \left[\mathbb{E}_{\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{y} \in \mathcal{X}} \frac{1}{m} \mathbf{y}^{\top} \mathbb{I}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{y} \geq 0\}} \left(\sum_{r=1}^{m} \sigma_{r} \mathbf{x} \mathbb{I}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{x} \geq 0\}} \right) \right] \right]$$

$$\stackrel{\textcircled{2}}{\leq} \mathbb{E}_{\mathbf{W}(0)} \left[\mathbb{E}_{\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{y} \in \mathcal{X}} \frac{1}{m} \|\mathbf{y}\|_{2} \left| \mathbb{I}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{y} \geq 0\}} \right| \left\| \sum_{r=1}^{m} \sigma_{r} \mathbf{x} \mathbb{I}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{x} \geq 0\}} \right\|_{2} \right] \right]$$

$$\stackrel{\texttt{(3)}}{\leq} \mathbb{E}_{\mathbf{W}(0)} \left[\mathbb{E}_{\{\sigma_r\}_{r=1}^m} \left[\frac{1}{m} \left\| \sum_{r=1}^m \sigma_r \mathbf{x} \mathbf{I}_{\{\vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\}} \right\|_2 \right] \right] \\ \stackrel{\texttt{(4)}}{=} \mathbb{E}_{\mathbf{W}(0)} \left[\mathbb{E}_{\{\sigma_r\}_{r=1}^m} \left[\frac{1}{m} \sqrt{\left(\sum_{r=1}^m \sigma_r \mathbf{I}_{\{\vec{\mathbf{w}}_r(0)^\top \mathbf{x} \ge 0\}} \right)^2} \right] \right]$$

Г

$$\begin{split} & \stackrel{\scriptstyle{\Theta}}{=} \mathbb{E}_{\mathbf{W}(0)} \left[\mathbb{E}_{\{\sigma_{r}\}_{r=1}^{m}} \left[\frac{1}{m} \sqrt{\left(\sum_{r=1}^{m} \sigma_{r} \mathbb{1}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{x} \ge 0\}} \right)} \right] \right] \\ & \stackrel{\scriptstyle{\Theta}}{\leq} \mathbb{E}_{\mathbf{W}(0)} \left[\frac{1}{m} \sqrt{\mathbb{E}_{\{\sigma_{r}\}_{r=1}^{m}}} \left[\left(\sum_{r\in[m],r'\in[m]}^{m} \sigma_{r} \sigma_{r'} \mathbb{1}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{x} \ge 0\}} \mathbb{1}_{\{\vec{\mathbf{w}}_{r'}(0)^{\top} \mathbf{x} \ge 0\}} \right) \right] \right] \\ & = \mathbb{E}_{\mathbf{W}(0)} \left[\frac{1}{m} \sqrt{\mathbb{E}_{\{\sigma_{r}\}_{r=1}^{m}}} \left[\left(\sum_{r\in[m],r'\in[m]}^{m} \sigma_{r} \sigma_{r'} \mathbb{1}_{\{\vec{\mathbf{w}}_{r}(0)^{\top} \mathbf{x} \ge 0\}} \mathbb{1}_{\{\vec{\mathbf{w}}_{r'}(0)^{\top} \mathbf{x} \ge 0\}} \right) \right] \right] \\ & \stackrel{\scriptstyle{\Theta}}{\leq} \mathbb{E}_{\mathbf{W}(0)} \left[\frac{1}{m} \cdot \sqrt{m} \right] = \frac{1}{\sqrt{m}}. \end{split}$$

$$\tag{157}$$

In (157), (1) is due to the fact that the operand of the supremum operator does not depend on \mathbf{x}' and the Fubini Theorem. 2 follows from the Cauchy-Schwarz inequality. 3 is due to the fact that $\|\mathbf{y}\|_2 = 1, \mathbb{I}_{\left\{\vec{\mathbf{w}}_r(0)^\top \mathbf{y} \ge 0\right\}} \in \{0, 1\}.$ (4) follows from $\|\mathbf{x}\|_2 = 1$, and (5) is due to the Jensen's inequality. (6) follows from the property of Rademacher variable, that is, $\mathbb{E}_{\{\sigma_r\}_{r=1}^m}[\sigma_r\sigma_{r'}] = \mathbb{I}_{\{r=r'\}}$, and the fact that $\mathbb{I}_{\left\{ \vec{\mathbf{w}}_{r}(0)^{\top}\mathbf{x}\geq0\right\} }\mathbb{I}_{\left\{ \vec{\mathbf{w}}_{r'}(0)^{\top}\mathbf{x}\geq0\right\} }\in\{0,1\}.$

For \mathcal{R}_2 , we first define

$$Q\coloneqq \frac{1}{m}\sum_{r=1}^m \mathrm{I\!I}_{\left\{\mathbf{1}_{\left\{\mathbf{x}'^\top \overrightarrow{\mathbf{w}}_r(0)\geq 0\right\}}\neq \mathrm{I\!I}_{\left\{\mathbf{x}^\top \overrightarrow{\mathbf{w}}_r(0)\geq 0\right\}}\right\}},$$

which is the average number of weights in $\mathbf{W}(0)$ whose inner products with x and x' have different signs. Our observation is that, if $|\mathbf{x}^{\top} \vec{\mathbf{w}}_r(0)| > s \|\vec{\mathbf{w}}_r(0)\|_2$, then $\mathbf{x}^{\top} \vec{\mathbf{w}}_r(0)$ has the same sign as $\mathbf{x}^{\prime \top} \overrightarrow{\mathbf{w}}_{r}(0)$. To see this, by the Cauchy-Schwarz inequality,

$$\left\|\mathbf{x}^{\prime\top}\vec{\mathbf{w}}_{r}(0) - \mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right\| \leq \left\|\mathbf{x}^{\prime} - \mathbf{x}\right\|_{2} \left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} \leq s \left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2},$$
(158)

then we have $\mathbf{x}^{\top} \vec{\mathbf{w}}_r(0) > s \left\| \vec{\mathbf{w}}_r(0) \right\|_2 \Rightarrow \mathbf{x}'^{\top} \vec{\mathbf{w}}_r(0) \geq \mathbf{x}^{\top} \vec{\mathbf{w}}_r(0) - s \left\| \vec{\mathbf{w}}_r(0) \right\|_2 > 0$, and $\mathbf{x}^{\top} \vec{\mathbf{w}}_r(0) < -s \left\| \vec{\mathbf{w}}_r(0) \right\|_2 \Rightarrow \mathbf{x}'^{\top} \vec{\mathbf{w}}_r(0) \leq \mathbf{x}^{\top} \vec{\mathbf{w}}_r(0) + s \left\| \vec{\mathbf{w}}_r(0) \right\|_2 < 0.$

As a result, $Q \leq \frac{1}{m} \sum_{r=1}^{m} \mathbb{I}_{\left\{ \left| \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \right| \leq s \left\| \vec{\mathbf{w}}_{r}(0) \right\|_{2} \right\}}$, and it follows that

$$\mathbb{E}_{\mathbf{W}(0)}\left[Q\right] \leq \mathbb{E}_{\mathbf{W}(0)}\left[\frac{1}{m}\sum_{r=1}^{m}\mathbb{I}_{\left\{\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|\leq s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}\right\}}\right] = \Pr\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right\|\leq s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}\right]$$
$$= \Pr\left[\frac{\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right\|}{\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}}\leq s\right], \quad (159)$$

where the last equality holds because each $\vec{\mathbf{w}}_r(0), r \in [m]$, follows a continuous Gaussian distribu-tion. By Lemma C.19, $\Pr\left[\frac{\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|}{\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}} \le s\right] \le B\sqrt{ds}$ for an absolute positive constant *B*. According to this inequality and (159), it follows that

$$\mathbb{E}_{\mathbf{W}(0)}\left[Q\right] \le B\sqrt{ds}.\tag{160}$$

By Markov's inequality, we have

$$\Pr\left[Q \ge \sqrt{s}\right] \le B\sqrt{ds},\tag{161}$$

where the probability is with respect to the probability measure space of $\mathbf{W}(0)$. Let A be the event that $Q \ge \sqrt{s}$. We denote by Ω_s the subset of the probability measure space of $\mathbf{W}(0)$ such that A happens, then $\Pr[\Omega_s] \le B\sqrt{ds}$. Now we aim to bound \mathcal{R}_2 by estimating its bound on Ω_s and its complement. First, we have

$$\mathcal{R}_{2} = \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s), \mathbf{y} \in \mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \left(h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}', \mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right) \right]$$
$$= \mathbb{E}_{\mathbf{W}(0) \in \Omega_{s},\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s), \mathbf{y} \in \mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \left(h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}', \mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right) \right]$$

$$+\underbrace{\mathbb{E}_{\mathbf{W}(0)\notin\Omega_{s},\{\sigma_{r}\}_{r=1}^{m}}\left[\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s),\mathbf{y}\in\mathcal{X}}\frac{1}{m}\sum_{r=1}^{m}\sigma_{r}\left(h(\mathbf{w}_{r}(0),\mathbf{x}',\mathbf{y})-h(\mathbf{w}_{r}(0),\mathbf{x},\mathbf{y})\right)\right]}_{\mathcal{R}_{22}},$$
(162)

where we used the convention that $\mathbb{E}_{\mathbf{W}(0)\in A}[\cdot] = \mathbb{E}_{\mathbf{W}(0)}\left[\mathbb{1}_{\{\mathbf{W}(0)\in A\}}\times\cdot\right]$. Now we estimate the upper bound for \mathcal{R}_{21} and \mathcal{R}_{22} separately. Let $I = \left\{r \in [m]: \mathbb{1}_{\{\mathbf{x}'^{\top}\vec{\mathbf{w}}_{r}(0)\geq 0\}} \neq \mathbb{1}_{\{\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\geq 0\}}\right\}$. When $\mathbf{W}(0) \notin \Omega_{s}$, we have $Q < \sqrt{s}$. In this case, it follows that $|I| \leq m\sqrt{s}$. Moreover, when $r \in I$, either $\mathbb{1}_{\{\mathbf{x}'^{\top}\vec{\mathbf{w}}_{r}(0)\geq 0\}} = 0$ or $\mathbb{1}_{\{\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\geq 0\}} = 0$. As a result,

$$\begin{aligned} \left| h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}', \mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right| \\ &= \left| \mathbf{x}'^{\top} \mathbf{y} \mathbf{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \mathbf{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} - \mathbf{x}^{\top} \mathbf{y} \mathbf{I}_{\left\{ \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \mathbf{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right| \\ &\leq \max \left\{ \mathbf{x}'^{\top} \mathbf{y} \mathbf{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \mathbf{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}}, \mathbf{x}^{\top} \mathbf{y} \mathbf{I}_{\left\{ \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \mathbf{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right\} \\ &\leq \max \left\{ \mathbf{x}'^{\top} \mathbf{y}, \mathbf{x}^{\top} \mathbf{y} \right\} \le 1. \end{aligned}$$
(163)

When $r \in [m] \setminus I$, we have

$$\begin{aligned} \left| h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}', \mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right| \\ &= \left| \mathbf{x}'^{\top} \mathbf{y} \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \mathbb{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} - \mathbf{x}^{\top} \mathbf{y} \mathbb{I}_{\left\{ \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \mathbb{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right| \\ &= \left| \left(\mathbf{x}' \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} - \mathbf{x} \mathbb{I}_{\left\{ \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right)^{\top} \mathbf{y} \mathbb{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right| \\ &\stackrel{(1)}{\leq} \left\| \mathbf{x}' \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} - \mathbf{x} \mathbb{I}_{\left\{ \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right\|_{2} \| \mathbf{y} \|_{2} \left| \mathbb{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right| \\ &\stackrel{(2)}{\leq} \left\| \mathbf{x}' \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} - \mathbf{x} \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} + \mathbf{x} \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} - \mathbf{x} \mathbb{I}_{\left\{ \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \ge 0 \right\}} \right\|_{2} \end{aligned}$$

$$\begin{array}{l} 2376\\ 2377\\ 2378\\ 2379 \end{array} \leq \|\mathbf{x}' - \mathbf{x}\|_2 \left| \mathbb{I}_{\left\{ \mathbf{x}'^\top \vec{\mathbf{w}}_r(0) \ge 0 \right\}} \right| + \|\mathbf{x}\|_2 \left| \mathbb{I}_{\left\{ \mathbf{x}'^\top \vec{\mathbf{w}}_r(0) \ge 0 \right\}} - \mathbb{I}_{\left\{ \mathbf{x}^\top \vec{\mathbf{w}}_r(0) \ge 0 \right\}} \right|$$

$$\begin{array}{l} (164) \\ (164) \end{array}$$

where ① follows from the Cauchy-Schwarz inequality, ② is due to the fact that $\left| \mathbb{I}_{\left\{ \mathbf{y}^{\top} \vec{\mathbf{w}}_{r}(0) \geq 0 \right\}} \right| \in \left\{ 0, 1 \right\}$ and $\left\| \mathbf{y} \right\|_{2} = 1$. ③ follows from $\mathbf{x}' \in \mathbf{B}(\mathbf{x}; s)$, $\left| \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \geq 0 \right\}} \right| \in \{0, 1\}$, and $\left\| \mathbf{y} \right\|_{2} = \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \geq 0 \right\}} = \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \geq 0 \right\}}$ because $r \notin I$.

By (163) and (164), we have

$$\frac{1}{m}\sum_{r=1}^{m}\sigma_{r}\left(h(\vec{\mathbf{w}}_{r}(0),\mathbf{x}',\mathbf{y})-h(\vec{\mathbf{w}}_{r}(0),\mathbf{x},\mathbf{y})\right) \\
=\frac{1}{m}\sum_{r\in I}\sigma_{r}\left(h(\vec{\mathbf{w}}_{r}(0),\mathbf{x}',\mathbf{y})-h(\vec{\mathbf{w}}_{r}(0),\mathbf{x},\mathbf{y})\right) + \frac{1}{m}\sum_{r\in[m]\setminus I}\sigma_{r}\left(h(\vec{\mathbf{w}}_{r}(0),\mathbf{x}',\mathbf{y})-h(\vec{\mathbf{w}}_{r}(0),\mathbf{x},\mathbf{y})\right) \\
\leq\frac{1}{m}\sum_{r\in I}\left|h(\vec{\mathbf{w}}_{r}(0),\mathbf{x}',\mathbf{y})-h(\vec{\mathbf{w}}_{r}(0),\mathbf{x},\mathbf{y})\right| + \frac{1}{m}\sum_{r\in[m]\setminus I}\left|h(\vec{\mathbf{w}}_{r}(0),\mathbf{x}',\mathbf{y})-h(\vec{\mathbf{w}}_{r}(0),\mathbf{x},\mathbf{y})\right| \\
\stackrel{(0)}{=}\frac{m\sqrt{s}}{m} + \frac{m-m\sqrt{s}}{m}s \leq \sqrt{s} + s,$$
(165)

where ① uses the bounds in (163) and (164).

Using (165), we now estimate the upper bound for \mathcal{R}_{22} by

$$\mathcal{R}_{22} = \mathbb{E}_{\mathbf{W}(0)\notin\Omega_s,\{\sigma_r\}_{r=1}^m} \left[\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s),\mathbf{y}\in\mathcal{X}} \frac{1}{m} \sum_{r=1}^m \sigma_r \left(h(\vec{\mathbf{w}}_r(0),\mathbf{x}',\mathbf{y}) - h(\vec{\mathbf{w}}_r(0),\mathbf{x},\mathbf{y}) \right) \right]$$

$$\leq \mathbb{E}_{\mathbf{W}(0)\notin\Omega_s,\{\sigma_r\}_{r=1}^m} \left[\sqrt{s} + s \right] \leq \sqrt{s} + s.$$
(166)

When $\mathbf{W}(0) \in \Omega_s$, by the second last inequality of (164), we have

$$\left| h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}', \mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0), \mathbf{x}, \mathbf{y}) \right|$$

$$\leq \|\mathbf{x}' - \mathbf{x}\|_{2} \left| \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \geq 0 \right\}} \right| + \|\mathbf{x}\|_{2} \left| \mathbb{I}_{\left\{ \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \geq 0 \right\}} - \mathbb{I}_{\left\{ \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \geq 0 \right\}} \right| \leq s + 1.$$
(167)

According to (167), for \mathcal{R}_{21} , we have

$$\mathcal{R}_{21} = \mathbb{E}_{\mathbf{W}(0)\in\Omega_{s},\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s),\mathbf{y}\in\mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \left(h(\vec{\mathbf{w}}_{r}(0),\mathbf{x}',\mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0),\mathbf{x},\mathbf{y}) \right) \right]$$

$$\leq \mathbb{E}_{\mathbf{W}(0)\in\Omega_{s},\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s),\mathbf{y}\in\mathcal{X}} \frac{1}{m} \sum_{r=1}^{m} \left| \sigma_{r} \left(h(\vec{\mathbf{w}}_{r}(0),\mathbf{x}',\mathbf{y}) - h(\vec{\mathbf{w}}_{r}(0),\mathbf{x},\mathbf{y}) \right) \right| \right]$$

$$\stackrel{(1)}{\leq} \mathbb{E}_{\mathbf{W}(0)\in\Omega_{s},\{\sigma_{r}\}_{r=1}^{m}} \left[s+1 \right] = (s+1) \Pr\left[\Omega_{s}\right] \leq B\sqrt{ds}(s+1)$$
(168)

2425 Combining (162), (166), and (168), we have the upper bound for \mathcal{R}_2 as

$$\mathcal{R}_2 = \mathcal{R}_{21} + \mathcal{R}_{22} \le B\sqrt{ds}(s+1) + \sqrt{s} + s.$$
(169)

Plugging (157) and (169) in (155), we have

2440 2441

2446 2447

2448

2449 2450

2456 2457

2458

 $\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) \le \mathcal{R}_1 + \mathcal{R}_2 \le \frac{1}{\sqrt{m}} + B\sqrt{ds}(s+1) + \sqrt{s} + s.$ (170)

Lemma C.19. Let $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)$ with $\kappa > 0$. Then for any $\varepsilon \in (0, 1)$ and $\mathbf{x} \in \mathcal{X}$, **Pr** $\left[\frac{|\mathbf{x}^\top \mathbf{w}|}{\|\mathbf{w}\|_2} \le \varepsilon\right] \le B\sqrt{d\varepsilon}$ where *B* is an absolute positive constant.

Remark. In fact, B can be set to $2(2\pi)^{-1/2}$ when $d \to \infty$.

2442 *Proof.* Let $z = \frac{\mathbf{x}^{\top} \mathbf{w}}{\|\mathbf{w}\|_2}$. It can be verified that $z^2 \sim z_1$ where z_1 is a random variable following 2443 the Beta distribution $\text{Beta}(\frac{1}{2}, \frac{d-1}{2})$. Therefore, the distribution of z has the following continuous 2444 probability density function p_z with respect to the Lebesgue measure,

$$p_z(x) = (1 - x^2)^{\frac{d-3}{2}} \mathbb{I}_{\{|x| \le 1\}} / B',$$
(171)

where $B' = \int_{-1}^{1} (1-x^2)^{\frac{d-3}{2}} dx$ is the normalization factor. It can be verified by standard calculation that $1/B' \leq \frac{B\sqrt{d}}{2}$ for an absolute positive constant B.

Because $1 - x^2 \le 1$ over $x \in [-1, 1]$, we have $B' \le 1$. In addition,

$$\Pr\left[\frac{|\mathbf{x}^{\top}\mathbf{w}|}{\|\mathbf{w}\|_{2}} \le \varepsilon\right] = \Pr\left[-\varepsilon \le z \le \varepsilon\right] = \frac{1}{B'} \int_{-\varepsilon}^{\varepsilon} (1-x^{2})^{\frac{d-3}{2}} \mathrm{d}x \le B\sqrt{d}\varepsilon, \quad (172)$$

where the last inequality is due to the fact that $1 - x^2 \le 1$ for $x \in [-\varepsilon, \varepsilon]$ with $\varepsilon \in (0, 1)$.

Proof of Theorem C.2. We follow the same proof strategy as that for Theorem C.1.

First, we have $\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)} [v_R(\mathbf{w}, \mathbf{x})] = \Pr[|\mathbf{w}^\top \mathbf{x}| \le R]$. For any $\mathbf{x} \in \mathcal{X}$ and s > 0, define function class

2461 2462 2463

2464

$$\mathcal{V}_{\mathbf{x},s} \coloneqq \left\{ v_R(\cdot, \mathbf{x}') \colon \mathbb{R}^d \to \mathbb{R} \colon \mathbf{x}' \in \mathbf{B}\left(\mathbf{x}; s\right) \cap \in \mathcal{X} \right\}.$$
(173)

2465 We first build an *s*-net for the unit sphere \mathcal{X} . By (Vershynin, 2012, Lemma 5.2), there exists an *s*-net 2466 $N_s(\mathcal{X}, \|\cdot\|_2)$ of \mathcal{X} such that $N(\mathcal{X}, \|\cdot\|_2, s) \le (1 + \frac{2}{s})^d$.

In the sequel, a function in the class $\mathcal{V}_{\mathbf{x}}$ is also denoted as $v_R(\mathbf{w})$, omitting the presence of \mathbf{x} when no confusion arises. Let P_m be the empirical distribution over $\left\{ \vec{\mathbf{w}}_r(0) \right\}$ and $\mathbb{E}_{\mathbf{w}\sim P_m} \left[v_R(\mathbf{w}) \right] =$ $\hat{v}_R(\mathbf{W}(0), \mathbf{x})$.

Given $\mathbf{x} \in N_s(\mathcal{X}, \|\cdot\|_2)$, we aim to estimate the upper bound for the supremum of empirical process $\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)} [v_R(\mathbf{w})] - \mathbb{E}_{\mathbf{w} \sim P_m} [v_R(\mathbf{w})]$ when function v_R ranges over the function class $\mathcal{V}_{\mathbf{x},s}$. To this end, we apply Theorem A.1 to the function class $\mathcal{V}_{\mathbf{x},s}$ with $\mathbf{W}(0) = \left\{ \vec{\mathbf{w}}_r(0) \right\}_{r=1}^m$. It can be verified that $v_R \in [0, 1]$ for any $v_R \in \mathcal{V}_{\mathbf{x},s}$. It follows that we can set a = 0, b = 1 in Theorem A.1. Setting $\alpha = \frac{1}{2}$ in Theorem A.1, then with probability at least $1 - 2e^{-x}$ over the random initialization $\mathbf{W}(0)$,

2478 2479

$$\sup_{v_{R}\in\mathcal{V}_{\mathbf{x},s}} \left| \mathbb{E}_{\mathbf{w}\sim\mathcal{N}(\mathbf{0},\kappa^{2}\mathbf{I}_{d})} \left[v_{R}(\mathbf{w}) \right] - \mathbb{E}_{\mathbf{w}\sim P_{m}} \left[v_{R}(\mathbf{w}) \right] \right| \\
\leq \inf_{\alpha\in(0,1)} \left(3\mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{v_{R}\in\mathcal{V}_{\mathbf{x},s}} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} v_{R}(\vec{\mathbf{w}}_{r}(0)) \right] + \sqrt{\frac{2rx}{m}} + \frac{7(b-a)x}{3m} \right), \quad (174)$$

where $\{\sigma\}_{r=1}^{m}$ are i.i.d. Rademacher random variables taking values of ± 1 with equal probability. It can be verified that $\operatorname{Var}[v_R] \leq \mathbb{E}_{\mathbf{w}} \left[v_R(\mathbf{w}, \mathbf{x})^2 \right] \leq 1$, so r can be set to 1. It follows that with probability at least $1 - \delta$,

$$\sup_{\mathbf{x}'\in\mathbf{B}(\mathbf{x};s)\cap\mathcal{X}} \left| \widehat{v}_R(\mathbf{W}(0),\mathbf{x}') - \Pr\left[\left| \mathbf{w}^\top \mathbf{x}' \right| \le R \right] \right| \le 3\mathcal{R}(\mathcal{H}_{\mathbf{x},s}) + \sqrt{\frac{2\log\frac{2}{\delta}}{m}} + \frac{7\log\frac{2}{\delta}}{3m}.$$
 (175)

Here $\mathcal{R}(\mathcal{V}_{\mathbf{x},s}) = \mathbb{E}_{\mathbf{W}(0),\{\sigma_r\}_{r=1}^m} \left[\sup_{v_R \in \mathcal{V}_{\mathbf{x},s}} \frac{1}{m} \sum_{r=1}^m \sigma_r v_R(\vec{\mathbf{w}}_r(0)) \right]$ is the Rademacher complexity of the function class $\mathcal{V}_{\mathbf{x},s}$. By Lemma C.21, $\mathcal{R}(\mathcal{V}_{\mathbf{x},s}) \leq (B\sqrt{d}+1)\sqrt{m^{-\frac{1}{2}}+s} + \frac{\exp\left(-\frac{(\kappa^2 - R_0^2)^2}{4\kappa^4}m\right)}{\sqrt{m^{-\frac{1}{2}}+s}}$.

Plugging such upper bound for $\mathcal{R}(\mathcal{V}_{\mathbf{x},s})$ in (175), we have

 $\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s) \cap \mathcal{X}} \left| \widehat{v}_R(\mathbf{W}(0), \mathbf{x}') - \Pr\left[\left| \mathbf{w}^\top \mathbf{x}' \right| \le R \right] \right|$

$$\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s) \cap \mathcal{X}} \left| \widehat{v}_{R}(\mathbf{W}(0), \mathbf{x}') - \Pr\left[\left| \mathbf{w}^{\top} \mathbf{x}' \right| \le R \right] \right| \\ \le 3 \left((B\sqrt{d}+1)\sqrt{m^{-\frac{1}{2}} + s} + \frac{\exp\left(-\frac{(\kappa^{2} - R_{0}^{2})^{2}}{4\kappa^{4}}m \right)}{\sqrt{m^{-\frac{1}{2}} + s}} \right) + \sqrt{\frac{2\log\frac{2}{\delta}}{m}} + \frac{7\log\frac{2}{\delta}}{3m}.$$
(176)

Setting $s = \frac{1}{m}$, we have

2506 2507 2508

2499

2512 2513

2514

2515

2516 By union bound, with probability at least $1 - (1 + 2m)^d \delta$ over $\mathbf{W}(0)$, (177) holds for arbitrary 2517 $\mathbf{x} \in N(\mathcal{X}, s)$. In this case, for any $\mathbf{x}' \in \mathcal{X}$, there exists $\mathbf{x} \in N(\mathcal{X}, s)$ such that $\|\mathbf{x}' - \mathbf{x}\|_2 \leq s$, so that $\mathbf{x}' \in \mathbf{B}(\mathbf{x}; s) \cap \mathcal{X}$, and (177) holds.

 $\leq 3\left((B\sqrt{d}+1)\sqrt{m^{-\frac{1}{2}}+\frac{1}{m}}+\frac{\exp\left(-\frac{(\kappa^{2}-R_{0}^{2})^{2}}{4\kappa^{4}}m\right)}{\sqrt{m^{-\frac{1}{2}}+\frac{1}{m}}}\right)+\sqrt{\frac{2\log\frac{2}{\delta}}{m}}+\frac{7\log\frac{2}{\delta}}{3m}$

Note that $\Pr\left[\left|\mathbf{w}^{\top}\mathbf{x}'\right| \le R\right] \le \frac{2R}{\sqrt{2\pi\kappa}}$ for any $\mathbf{x}' \in \mathcal{X}$, changing the notation \mathbf{x}' to \mathbf{x} completes the proof.

2522 2523

(177)

Lemma C.20. Let $\mathbf{w} \in \mathbb{R}^d$ be a Gaussian random vector distribute according to $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}_d)$. Then $\Pr[\|\mathbf{w}\|_2 \ge R'] \ge 1 - \exp\left(-\left(\frac{\sqrt{m}}{2} - \frac{{R'}^2}{2\sqrt{m}\kappa^2}\right)^2\right)$ for any R' > 0. Then $\Pr[\|\mathbf{w}\|_2 \ge R'] \ge 1 - \exp\left(-\left(\frac{\sqrt{m}}{2} - \frac{{R'}^2}{2\sqrt{m}\kappa^2}\right)^2\right)$ for any R' > 0.

2528 Proof. Let $X = \frac{\|\mathbf{w}\|_2^2}{\kappa^2}$, then X follows the chi-square distribution with m degrees of freedom, that 2529 is, $X \sim \chi^2(m)$. By (Laurent & Massart, 2000, Lemma 1), we have the following concentration 2530 inequalities for any x > 0, 2531 $\mathbb{P}\left[[X = \chi^2 - \chi$

$$\Pr\left[X - m \ge 2\sqrt{mx} + 2x\right] \le \exp(-x), \Pr\left[m - X \ge 2\sqrt{mx}\right] \le \exp(-x).$$
(178)

2532 2533

Setting $x = \left(\frac{\sqrt{m}}{2} - \frac{R'^2}{2\sqrt{m}\kappa^2}\right)^2$ in the second inequality in (178), we have $m - 2\sqrt{mx} = \frac{R'^2}{\kappa^2}$ and 2535

2536
2537
$$\Pr\left[X \ge \frac{{R'}^2}{\kappa^2}\right] \ge 1 - \exp(-x).$$
 (179)

It follows from (179) that

$$\Pr[\|\mathbf{w}\|_{2} \ge R'] \ge 1 - \exp(-x) = 1 - \exp\left(-\left(\frac{\sqrt{m}}{2} - \frac{{R'}^{2}}{2\sqrt{m}\kappa^{2}}\right)^{2}\right), \quad (180)$$

which completes the proof.

Lemma C.21. Suppose $R \leq R_0$ for an absolute positive constant $R_0 < \kappa$. Let $\mathcal{R}(\mathcal{V}_{\mathbf{x},s}) \coloneqq$ $\mathbb{E}_{\mathbf{W}(0),\{\sigma_r\}_{r=1}^m} \left[\sup_{v_R \in \mathcal{V}_{\mathbf{x},s}} \frac{1}{m} \sum_{r=1}^m \sigma_r v_R(\vec{\mathbf{w}}_r(0)) \right]$ be the Rademacher complexity of the function class $\mathcal{V}_{\mathbf{x},s}$. Then

 $\mathcal{R}(\mathcal{V}_{\mathbf{x},s}) \le (B\sqrt{d}+1)\sqrt{m^{-\frac{1}{2}}+s} + \frac{\exp\left(-\frac{(\kappa^2 - R_0^2)^2}{4\kappa^4}m\right)}{\sqrt{m^{-\frac{1}{2}}+s}},$ (181)

where B is a positive constant.

 \mathcal{R}

Proof. We have

$$(\mathcal{V}_{\mathbf{x},s}) = \mathbb{E}_{\mathbf{W}(0),\{\sigma_r\}_{r=1}^m} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^m \sigma_r v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') \right]$$

$$\leq \mathcal{R}_1 + \mathcal{R}_2, \tag{182}$$

where

$$\mathcal{R}_{1} = \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} v_{R}(\vec{\mathbf{w}}_{r}(0), \mathbf{x}) \right],$$
$$\mathcal{R}_{2} = \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \left(v_{R}(\vec{\mathbf{w}}_{r}(0), \mathbf{x}') - v_{R}(\vec{\mathbf{w}}_{r}(0), \mathbf{x}) \right) \right].$$
(183)

Here (182) follows from the subadditivity of superemum and the fact that $\sum_{r=1}^{m} \sigma_r v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') =$ $\sum_{r=1}^{m} \sigma_r v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}) + \sum_{r=1}^{m} \sigma_r \Big(v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') - v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}) \Big).$

Now we bound \mathcal{R}_1 and \mathcal{R}_2 separately. For \mathcal{R}_1 , we have

$$\mathcal{R}_{1} = \mathbb{E}_{\mathbf{W}(0), \{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} v_{R}(\mathbf{w}_{r}(0), \mathbf{x}) \right] = 0.$$
(184)

For \mathcal{R}_2 , we first define

$$Q = \frac{1}{m} \sum_{r=1}^{m} \mathbb{I}_{\left\{ \mathbb{I}_{\left\{ \left| \mathbf{x}'^{\top} \overrightarrow{\mathbf{w}}_{r}(0) \right| \leq R \right\}} \neq \mathbb{I}_{\left\{ \left| \mathbf{x}^{\top} \overrightarrow{\mathbf{w}}_{r}(0) \right| \leq R \right\}} \right\}},$$

which is the number of weights in W(0) whose inner products with x and x' have different signs. Note that if $\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right\| - R$ > $s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}$, then $\mathbb{I}_{\left\{\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|\leq R\right\}} = \mathbb{I}_{\left\{\left|\mathbf{x}^{\prime\top}\vec{\mathbf{w}}_{r}(0)\right|\leq R\right\}}$. To see this, by the Cauchy-Schwarz inequality,

 $\begin{aligned} \left|\mathbf{x}'^{\top}\vec{\mathbf{w}}_{r}(0)-\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right| \leq \|\mathbf{x}'-\mathbf{x}\|_{2}\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} \leq s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}, \quad (185) \end{aligned}$ then we have $\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right| - R > s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} \Rightarrow \left|\mathbf{x}'^{\top}\vec{\mathbf{w}}_{r}(0)\right| - R \geq \left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right| - s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} - R > 0, \\ \text{and } \left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right| - R < -s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} \Rightarrow \left|\mathbf{x}'^{\top}\vec{\mathbf{w}}_{r}(0)\right| - R \leq \left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right| + s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} - R < 0. \end{aligned}$ As a result, $Q \leq \frac{1}{m}\sum_{r=1}^{m} \mathbbmmatrix_{1}^{m} [\left|\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right| - R\right| \leq s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} \right|^{2}. \end{aligned}$ For any fixed $r \in [m]$, by Lemma C.20, Pr $\left[\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2} \geq R'\right] \geq 1 - \exp\left(-\left(\frac{\sqrt{m}}{2} - \frac{R'^{2}}{2\sqrt{m\kappa^{2}}}\right)^{2}\right)$ holds for any $R' \geq 0.$ Set $R' = \sqrt{m}R_{0}$ for the constant $R_{0} < \kappa.$ Because $R \leq R_{0}$, it follows that $\Pr\left[\frac{R}{\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}} \leq m^{-\frac{1}{2}}\right] \geq \Pr\left[\frac{R_{0}}{\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}} \leq m^{-\frac{1}{2}}\right] \geq 1 - \exp\left(-\frac{(\kappa^{2} - R_{0}^{2})^{2}}{4\kappa^{4}}m\right).$ (186)
Due to the fact that $\mathbbmmatrix_{1}^{T}\left|\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right| \leq m^{-\frac{1}{2}}\right] \leq \mathbbmmatrix_{r}^{T}\left|\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right| \leq m^{-\frac{1}{2}}\left[\mathbbmmatrix_{r}^{T}\left|\mathbf{w}_{r}(0)\right|_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{r}^{T}\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{r}^{T}\right|_{2}^{T}\left[\mathbbmmatrix_{r}^{T}\left(\mathbf{x}^{\top}\mathbf{w}_{r}(0)\right]_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right|_{2}^{T}\left[\mathbbmmatrix_{r}^{T}\vec{\mathbf{w}}_{r}(0)\right]_{2}^{T}\right] + \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\mathbf{w}_{r}(0)\right|_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right] + \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right\|_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right] \\ \leq \mathbbmmatrix_{r}^{T}\left[\left\|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|_{2}^{T}\right]$

$$\leq \Pr\left[\frac{1}{\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}} \leq m^{-\frac{1}{2}} + s\right] + \exp\left(-\frac{(\kappa^{2} - R_{0}^{2})^{2}}{4\kappa^{4}}m\right) \\
\stackrel{(3)}{\leq} B\sqrt{d}(m^{-\frac{1}{2}} + s) + \exp\left(-\frac{(\kappa^{2} - R_{0}^{2})^{2}}{4\kappa^{4}}m\right),$$
(187)

where we used the convention that $\mathbb{E}_{\vec{\mathbf{w}}_{r}(0)\in A}[\cdot] = \mathbb{E}_{\vec{\mathbf{w}}_{r}(0)}\left[\mathbb{I}_{\{A\}}\times\cdot\right]$ in ① with A being an event. ② is due to (186). By Lemma C.19, $\Pr\left[\frac{\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|}{\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}}\leq m^{-\frac{1}{2}}+s\right]\leq B\sqrt{d}(m^{-\frac{1}{2}}+s)$ for an absolute constant B, so ③ holds.

According to (187), we have

$$\mathbb{E}_{\mathbf{W}(0)}\left[Q\right] \leq \mathbb{E}_{\mathbf{W}(0)}\left[\frac{1}{m}\sum_{r=1}^{m}\mathbb{I}_{\left\{\left|\left|\mathbf{x}^{\top}\vec{\mathbf{w}}_{r}(0)\right|-R\right|\leq s\left\|\vec{\mathbf{w}}_{r}(0)\right\|_{2}\right\}}\right]$$
$$\leq B\sqrt{d}(m^{-\frac{1}{2}}+s) + \exp\left(-\frac{(\kappa^{2}-R_{0}^{2})^{2}}{4\kappa^{4}}m\right).$$
(188)

2642 Define $s' := m^{-\frac{1}{2}} + s$. By Markov's inequality, we have

$$\Pr\left[Q \ge \sqrt{s'}\right] \le B\sqrt{ds'} + \frac{\exp\left(-\frac{(\kappa^2 - R_0^2)^2}{4\kappa^4}m\right)}{\sqrt{s'}},\tag{189}$$

where the probability is with respect to the probability measure space of $\mathbf{W}(0)$. Now we aim to bound \mathcal{R}_2 by estimating its bound on Ω_s and its complement. First, we have

$$\mathcal{R}_{2} = \mathbb{E}_{\mathbf{W}(0),\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \left(v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}') - v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}) \right) \right]$$

$$= \mathbb{E}_{\mathbf{W}(0): Q \geq \sqrt{s'},\{\sigma_{r}\}_{r=1}^{m}} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^{m} \sigma_{r} \left(v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}') - v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}) \right) \right]$$

$$\mathcal{R}_{21}$$

$$\mathcal{R}_{22}$$

$$\mathcal{R}_{22}$$

$$\mathcal{R}_{22}$$

$$\mathcal{R}_{22}$$

$$\mathcal{R}_{21}$$

$$\mathcal{R}_{22}$$

Now we estimate the upper bound for \mathcal{R}_{22} and \mathcal{R}_{21} separately. Let

$$I = \left\{ r \in [m] \colon \mathbb{I}_{\left\{ \left| \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \right| \leq R \right\}} \neq \mathbb{I}_{\left\{ \left| \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \right| \leq R \right\}} \right\}.$$

2666 When $Q < \sqrt{s'}$, $|I| \leq m\sqrt{s'}$. Moreover, when $r \in I$, either $\mathbb{I}_{\{|\mathbf{x}'^{\top} \vec{\mathbf{w}}_r(0)| \leq R\}} = 0$ or 2667 $\mathbb{I}_{\{|\mathbf{x}^{\top} \vec{\mathbf{w}}_r(0)| \leq R\}} = 0$. As a result,

$$\begin{aligned} \left| v_{R}(\vec{\mathbf{w}}_{r}(0), \mathbf{x}') - v_{R}(\vec{\mathbf{w}}_{r}(0), \mathbf{x}) \right| &= \left| \mathbb{I}_{\left\{ \left| \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \right| \leq R \right\}} - \mathbb{I}_{\left\{ \left| \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \right| \leq R \right\}} \right| \\ &\leq \max \left\{ \mathbb{I}_{\left\{ \left| \mathbf{x}'^{\top} \vec{\mathbf{w}}_{r}(0) \right| \leq R \right\}}, \mathbb{I}_{\left\{ \left| \mathbf{x}^{\top} \vec{\mathbf{w}}_{r}(0) \right| \leq R \right\}} \right\} \\ &\leq 1. \end{aligned}$$

$$(191)$$

When $r \in [m] \setminus I$, we have

$$\left| v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') - v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}) \right| = \left| \mathbb{I}_{\left\{ \left| \mathbf{x}'^\top \vec{\mathbf{w}}_r(0) \right| \le R \right\}} - \mathbb{I}_{\left\{ \left| \mathbf{x}^\top \vec{\mathbf{w}}_r(0) \right| \le R \right\}} \right| = 0.$$
(192)

By (191) and (192), we have

2683
2684
2685
$$\frac{1}{m}\sum_{r=1}^{m}\sigma_{r}\left(v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}')-v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x})\right)$$
2686
$$\frac{1}{m}\sum_{r\in I}\sigma_{r}\left(v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}')-v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x})\right)+\frac{1}{m}\sum_{r\in[m]\setminus I}\sigma_{r}\left(v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}')-v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x})\right)$$
2690
$$\leq \frac{1}{m}\sum_{r\in I}\left|v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}')-v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x})\right|+\frac{1}{m}\sum_{r\in[m]\setminus I}\left|v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}')-v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x})\right|$$
2692
$$(193)$$
2695
$$= 1-\frac{2}{m}\sum_{r\in I}\left|v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}')-v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x})\right|+\frac{1}{m}\sum_{r\in[m]\setminus I}\left|v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x}')-v_{R}(\vec{\mathbf{w}}_{r}(0),\mathbf{x})\right|$$
2691
$$(193)$$

where ① uses the bounds in (191) and (192).

2697 Using (193), we now estimate the upper bound for \mathcal{R}_{22} by

$$\mathcal{R}_{22} = \mathbb{E}_{\mathbf{W}(0): Q < \sqrt{s'}, \{\sigma_r\}_{r=1}^m} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^m \sigma_r \left(v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') - v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}) \right) \right]$$

$$\leq \mathbb{E}_{\mathbf{W}(0): Q < \sqrt{s'}, \{\sigma_r\}_{r=1}^m} \left[\sqrt{s'} \right] = \sqrt{s'}.$$
(194)

When $Q \ge \sqrt{s'}$, by (191), we still have $\left| v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') - v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}) \right| \le 1$. For \mathcal{R}_{21} , we have

$$\mathcal{R}_{21} = \mathbb{E}_{\mathbf{W}(0): Q \ge \sqrt{s'}, \{\sigma_r\}_{r=1}^m} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^m \sigma_r \left(v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') - v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}) \right) \right]$$

$$\leq \mathbb{E}_{\mathbf{W}(0): Q \ge \sqrt{s'}, \{\sigma_r\}_{r=1}^m} \left[\sup_{\mathbf{x}' \in \mathbf{B}(\mathbf{x};s)} \frac{1}{m} \sum_{r=1}^m \left| \sigma_r \left(v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}') - v_R(\vec{\mathbf{w}}_r(0), \mathbf{x}) \right) \right| \right]$$

$$\leq \mathbb{E}_{\mathbf{W}(0): Q \ge \sqrt{s'}, \{\sigma_r\}_{r=1}^m} \left[1 \right] = \Pr\left[Q \ge \sqrt{s'} \right] \le B\sqrt{ds'} + \frac{\exp\left(-\frac{(\kappa^2 - R_0^2)^2}{4\kappa^4} m \right)}{\sqrt{s'}}, \quad (195)$$

where the last inequality is due to (189). Combining (190), (194), and (195), we have the upper bound for \mathcal{R}_2 as

$$\mathcal{R}_{2} = \mathcal{R}_{21} + \mathcal{R}_{22} \le (B\sqrt{d} + 1)\sqrt{s'} + \frac{\exp\left(-\frac{(\kappa^{2} - R_{0}^{2})^{2}}{4\kappa^{4}}m\right)}{\sqrt{s'}}.$$
(196)

2722 Plugging (184) and (196) in (182), we have

$$\mathcal{R}(\mathcal{V}_{\mathbf{x},s}) \le \mathcal{R}_1 + \mathcal{R}_2 \le (B\sqrt{d} + 1)\sqrt{s'} + \frac{\exp\left(-\frac{(\kappa^2 - R_0^2)^2}{4\kappa^4}m\right)}{\sqrt{s'}},\tag{197}$$

which completes the proof.



Figure

Figure 2: Illustration of the test loss by GD

2754 D SIMULATION STUDY

We present simulation results for GD in this section. We randomly sample *n* points $\left\{ \vec{\mathbf{x}}_i \right\}_{i=1}^n$ as a i.i.d. sample of random variables distributed uniformly on the unit sphere in \mathbb{R}^{50} . n ranges within [100, 1000] with a step size of 100. We set the target function to $f^*(\mathbf{x}) = \mathbf{s}^\top \mathbf{x}$ where $\mathbf{s} \sim \text{Unif}(\mathcal{X})$ is randomly sampled. We also uniformly and independenly sample 1000 points on the unit sphere in \mathbb{R}^{50} as the test data. We train the two-layer NN (1) using either GD by Algorithm 1 or GD by Algoirthm 1 with $m \asymp n^2$ on a NVIDIA A100 GPU card with a learning rate $\eta = 0.1$, and report the test loss in Figure 2. It can be observed that early-stopping is always helpful in training neural networks with better generalization, as the test loss initially decreases and then increases with over-training. Figure 2 illustrates the test loss with respect to the steps (or epochs) of GD for n = 100, 500, 1000. For each n in [100, 1000] with a step size of 100, we find the step of GD where minimum test loss is achieved, denoted by \hat{t}_n which is the empirical early stopping time. We note that the theoretically predicted early stopping time is $\hat{\varepsilon}_n = n^{-d/(2d-1)}$, and we compute the ratio of early stopping time for each n by $\hat{t}_n/\hat{\varepsilon}_n$. Such ratios for different values of n are illustrated in the bottom right figure of Figure 2. It is observed that the ratio of early stopping time is roughly stable and distributed between [8, 10], suggesting that predicted early stopping time is empirically proportional to the empirical early stopping time.