
On the Generalization Power of Overfitted Two-Layer Neural Tangent Kernel Models

Peizhong Ju¹ Xiaojun Lin¹ Ness B. Shroff²

Abstract

In this paper, we study the generalization performance of min ℓ_2 -norm overfitting solutions for the neural tangent kernel (NTK) model of a two-layer neural network with ReLU activation that has no bias term. We show that, depending on the ground-truth function, the test error of overfitted NTK models exhibits characteristics that are different from the “double-descent” of other overparameterized linear models with simple Fourier or Gaussian features. Specifically, for a class of learnable functions, we provide a new upper bound of the generalization error that approaches a small limiting value, even when the number of neurons p approaches infinity. This limiting value further decreases with the number of training samples n . For functions outside of this class, we provide a lower bound on the generalization error that does not diminish to zero even when n and p are both large.

1. Introduction

Recently, there is significant interest in understanding why overparameterized deep neural networks (DNNs) can still generalize well (Zhang et al., 2017; Advani et al., 2020), which seems to defy the classical understanding of *bias-variance tradeoff* in statistical learning (Bishop, 2006; Hastie et al., 2009; Stein, 1956; James & Stein, 1992; LeCun et al., 1991; Tikhonov, 1943). Towards this direction, a recent line of study has focused on overparameterized linear models (Belkin et al., 2018b; 2019; Bartlett et al., 2020; Hastie et al., 2019; Muthukumar et al., 2019; Ju et al., 2020; Mei & Montanari, 2019). For linear models with simple features (e.g., Gaussian features and Fourier features) (Belkin et al., 2018b; 2019; Bartlett et al., 2020; Hastie et al., 2019;

Muthukumar et al., 2019; Ju et al., 2020), an interesting “double-descent” phenomenon has been observed. Thus, there is a region where the number of model parameters (or linear features) is larger than the number of samples (and thus overfitting occurs), but the generalization error actually decreases with the number of features. However, linear models with these simple features are still quite different from nonlinear neural networks. Thus, although such results provide some hint why overparameterization and overfitting may be harmless, it is still unclear whether similar conclusions apply to neural networks.

In this paper, we are interested in linear models based on the neural tangent kernel (NTK) (Jacot et al., 2018), which can be viewed as a useful intermediate step towards modeling nonlinear neural networks. Essentially, NTK can be seen as a linear approximation of neural networks when the weights of the neurons do not change much. Indeed, (Li & Liang, 2018; Du et al., 2018) have shown that, for a wide and fully-connected two-layer neural network, both the neuron weights and their activation patterns do not change much after gradient descent (GD) training with a sufficiently small step size. As a result, such a shallow and wide neural network is approximately linear in the weights when there are a sufficient number of neurons, which suggests the utility of the NTK model.

Despite its linearity, however, characterizing the double descent of such a NTK model remains elusive. The work in (Mei & Montanari, 2019) also studies the double-descent of a linear version of two-layer neural network. It uses the so-called “random-feature” model, where the bottom-layer weights are random and fixed, and only the top-layer weights are trained. (In comparison, the NTK model for such a two-layer neural network corresponds to training only the bottom-layer weights.) However, the setting there requires the number of neurons, the number of samples, and the data dimension to all grow proportionally to infinity. In contrast, we are interested in the setting where the number of samples is given, and the number of neurons is allowed to be much larger than the number of samples. As a consequence of the different setting, in (Mei & Montanari, 2019) eventually only *linear* ground-truth functions can be learned. (Similar settings are also studied in (d’Ascoli et al.,

¹School of Electrical and Computer Engineering, Purdue University, West Lafayette, Indiana, USA. ²Department of ECE and CSE, The Ohio State University, Columbus, Ohio, USA. Correspondence to: Xiaojun Lin <linx@purdue.edu>.

2020).) In contrast, we will show that far more complex functions can be learned in our setting. In a related work, (Ghorbani et al., 2019) shows that both the random-feature model and the NTK model can approximate highly *non-linear* ground-truth functions with a sufficient number of neurons. However, (Ghorbani et al., 2019) mainly studies the *expressiveness* of the models, and therefore does not explain why overfitting solutions can still generalize well. To the best of our knowledge, our work is the first to characterize the double-descent of overfitting solutions based on the NTK model.

Specifically, in this paper we study the generalization error of the min ℓ_2 -norm overfitting solution for a linear model based on the NTK of a two-layer neural network with ReLU activation that has no bias. Only the bottom-layer weights are trained. We are interested in min ℓ_2 -norm overfitting solutions because gradient descent (GD) can be shown to converge to such solutions while driving the training error to zero (Zhang et al., 2017) (see also Section 2). Given a class of ground truth functions (see details in Section 3), which we refer to as “learnable functions,” our main result (Theorem 1) provides an upper bound on the generalization error of the min ℓ_2 -norm overfitting solution for the two-layer NTK model with n samples and p neurons (for any finite p larger than a polynomial function of n). This upper bound confirms that the generalization error of the overfitting solution indeed exhibits descent in the overparameterized regime when p increases. Further, our upper bound can also account for the noise in the training samples.

Our results reveal several important insights. First, we find that the (double) descent of the overfitted two-layer NTK model is drastically different from that of linear models with simple Gaussian or Fourier features (Belkin et al., 2018b; 2019; Bartlett et al., 2020; Hastie et al., 2019; Muthukumar et al., 2019). Specifically, for linear models with simple features, when the number of features p increases, the generalization error will eventually grow again and approach the so-called “null risk” (Hastie et al., 2019), which is the error of a trivial model that predicts zero. In contrast, for the class of learnable functions described earlier, the generalization error of the overfitted NTK model will continue to descend as p grows to infinity, and will approach a limiting value that depends on the number of samples n . Further, when there is no noise, this limiting value will decrease to zero as the number of samples n increases. This difference is shown in Fig. 1(a). As p increases, the test mean-square-error (MSE) of min- ℓ_1 and min- ℓ_2 overfitting solutions for Fourier features (blue and red curves) eventually grow back to the null risk (the black dashed line), even though they exhibit a descent at smaller p . In contrast, the error of the overfitted NTK model continues to descend to a much lower level.

The second important insight is that the aforementioned behavior critically depends on the ground-truth function belonging to the class of “learnable functions.” Further, this class of learnable functions depend on the specific network architecture. For our NTK model (with RELU activation that has no bias), we precisely characterize this class of learnable functions. Specifically, for ground-truth functions that are outside the class of learnable functions, we show a lower bound on the generalization error that does not diminish to zero for any n and p (see Proposition 2 and Section 4). This difference is shown in Fig. 1(b), where we use an almost identical setting as Fig. 1(a), except a different ground-truth function. We can see in Fig. 1(b) that the test-error of the overfitted NTK model is always above the null risk and looks very different from that in Fig. 1(a). We note that whether certain functions are learnable or not critically depends on the specific structure of the NTK model, such as the choice of the activation unit. Recently, (Satpathi & Srikant, 2021) shows that all polynomials can be learned by 2-layer NTK model with ReLU activation that has a bias term, provided that the number of neurons p is sufficiently large. (See further discussions in Remark 2. However, (Satpathi & Srikant, 2021) does not characterize the descent of generalization errors as p increases.) This difference in the class of learnable functions between the two settings (ReLU with or without bias) also turns out to be consistent with the difference in the expressiveness of the neural networks. That is, shallow networks with biased-ReLU are known to be universal function approximators (Ji et al., 2019), while those without bias can only approximate the sum of linear functions and even functions (Ghorbani et al., 2019).

A closely related result to ours is the work in (Arora et al., 2019), which characterizes the generalization performance of wide two-layer neural networks whose bottom-layer weights are trained by gradient descent (GD) to overfit the training samples. In particular, our class of learnable functions almost coincides with that of (Arora et al., 2019). This is not surprising because, when the number of neurons is large, NTK becomes a close approximation of such two-layer neural networks. In that sense, the results in (Arora et al., 2019) are even more faithful in following the GD dynamics of the original two-layer network. However, the advantage of the NTK model is that it is easier to analyze. In particular, the results in this paper can quantify how the generalization error descends with p . In contrast, the results in (Arora et al., 2019) provide only a generalization bound that is independent of p (provided that p is sufficiently large), but do not quantify the descent behavior as p increases. Our numerical results in Fig. 1(a) suggest that, over a wide range of p , the descent behavior of the NTK model (the green curve) matches well with that of two-layer neural networks trained by gradient descent (the cyan curve). Thus, we believe that our results also provide guidance for the latter

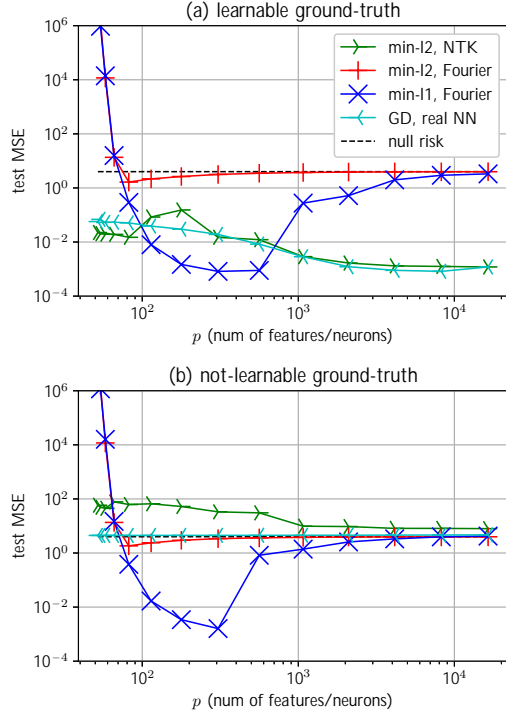


Figure 1. The test mean-square-error(MSE) vs. the number of features/neurons p for (a) learnable function and (b) not-learnable function when $n = 50$, $d = 2$, $k\epsilon k_2^2 = 0.01$. The corresponding ground-truth are (a) $f(\theta) = \sum_{k \in \{0,1,2,4\}g} (\sin(k\theta) + \cos(k\theta))$, and (b) $f(\theta) = \sum_{k \in \{3,5,7,9\}g} (\sin(k\theta) + \cos(k\theta))$. (Note that in 2-dimension every input \mathbf{x} on a unit circle can be represented by an angle $\theta \in [\pi, \pi]$. See the end of Section 4.) Every curve is the average of 9 random simulation runs. For GD on the real neural network (NN), we use the step size $1/\rho p$ and the number of training epochs is fixed at 2000.

model. The work in (Fiat et al., 2019) studies a different neural network architecture with gated ReLU, whose NTK model turns out to be the same as ours. However, similar to (Arora et al., 2019), the result in (Fiat et al., 2019) does not capture the speed of descent with respect to p either. Second, (Arora et al., 2019) only provides upper bounds on the generalization error. There is no corresponding lower bound to explain whether ground-truth functions outside a certain class are *not* learnable. Our result in Proposition 2 provides such a lower bound, and therefore more completely characterizes the class of learnable functions. (See further comparison in Remark 1 of Section 3 and Remark 3 of Section 5.) Another related work (Allen-Zhu et al., 2019) also characterizes the class of learnable functions for two-layer and three-layer networks. However, (Allen-Zhu et al., 2019) studies a training method that takes a new sample in every iteration, and thus does not overfit all training data. Finally, our paper studies generalization of NTK models for the regression setting, which is different from the classification

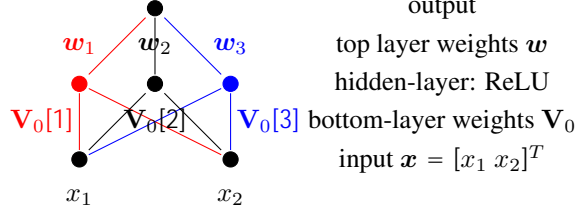


Figure 2. A two-layer neural network where $d = 2$, $p = 3$.

setting that assumes a separability condition, e.g., in (Ji & Telgarsky, 2019).

2. Problem Setup

We assume the following data model $y = f(\mathbf{x}) + \epsilon$, with the input $\mathbf{x} \in \mathbb{R}^d$, the output $y \in \mathbb{R}$, the noise $\epsilon \in \mathbb{R}$, and $f : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes the ground-truth function. Let (\mathbf{X}_i, y_i) , $i = 1, 2, \dots, n$ denote n training samples. We collect them as $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2 \ \dots \ \mathbf{X}_n] \in \mathbb{R}^{d \times n}$, $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_n]^T \in \mathbb{R}^n$, $\epsilon = [\epsilon_1 \ \epsilon_2 \ \dots \ \epsilon_n]^T \in \mathbb{R}^n$, and $\mathbf{F}(\mathbf{X}) = [f(\mathbf{X}_1) \ f(\mathbf{X}_2) \ \dots \ f(\mathbf{X}_n)]^T \in \mathbb{R}^n$. Then, the training samples can be written as $\mathbf{y} = \mathbf{F}(\mathbf{X}) + \epsilon$. After training (to be described below), we denote the trained model by the function \hat{f} . Then, for any new test data \mathbf{x} , we will calculate the test error by $|\hat{f}(\mathbf{x}) - f(\mathbf{x})|$, and the mean squared error (MSE) by $\mathbb{E}_{\mathbf{x}} [|\hat{f}(\mathbf{x}) - f(\mathbf{x})|^2]$.

For training, consider a fully-connected two-layer neural network with p neurons. Let $\mathbf{w}_j \in \mathbb{R}$ and $\mathbf{V}_0[j] \in \mathbb{R}^d$ denote the top-layer and bottom-layer weights, respectively, of the j -th neuron, $j = 1, 2, \dots, p$ (see Fig. 2). We collect them into $\mathbf{w} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_p]^T \in \mathbb{R}^p$, and $\mathbf{V}_0 = [\mathbf{V}_0[1]^T \ \mathbf{V}_0[2]^T \ \dots \ \mathbf{V}_0[p]^T]^T \in \mathbb{R}^{dp}$ (a column vector with dp elements). Note that with this notation, for any row or column vector \mathbf{v} with dp elements, $\mathbf{v}[j]$ denotes a (row/column) vector that consists of the $(jd + 1)$ -th to $(jd + d)$ -th elements of \mathbf{v} . We choose ReLU as the activation function for all neurons and there is no bias term in the ReLU activation function.

Now we are ready to introduce the NTK model (Jacot et al., 2018). We fix the top-layer weights \mathbf{w} , and let the initial bottom-layer weights \mathbf{V}_0 be randomly chosen. We then train only the bottom-layer weights. Let $\mathbf{V}_0 + \overline{\mathbf{V}}$ denote the bottom-layer weights after training. Thus, the change of the output after training is

$$\sum_{j=1}^n \mathbf{w}_j \mathbf{1}_{\hat{f}_{\mathbf{x}^T}(\mathbf{V}_0[j] + \overline{\mathbf{V}}[j]) > 0} (\mathbf{V}_0[j] + \overline{\mathbf{V}}[j])^T \mathbf{x} - \sum_{j=1}^n \mathbf{w}_j \mathbf{1}_{\hat{f}_{\mathbf{x}^T} \mathbf{V}_0[j] > 0} \mathbf{V}_0[j]^T \mathbf{x}.$$

In the NTK model, one assumes that $\overline{\mathbf{V}}$ is very small. As

a result, $\mathbf{1}_{f_{\mathbf{x}^T(\mathbf{V}_0[j] + \Delta \overline{\mathbf{V}}[j])} > 0g} = \mathbf{1}_{f_{\mathbf{x}^T \mathbf{V}_0[j]} > 0g}$ for most \mathbf{x} . Thus, the change of the output can be approximated by

$$\sum_{j=1}^n \mathbf{w}_j \mathbf{1}_{f_{\mathbf{x}^T \mathbf{V}_0[j]} > 0g} \overline{\mathbf{V}}[j]^T \mathbf{x} = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V},$$

where $\mathbf{V} \in \mathbb{R}^{dp}$ is given by $\mathbf{V}[j] := \mathbf{w}_j \overline{\mathbf{V}}[j]$, $j = 1, 2, \dots, p$, and $\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \in \mathbb{R}^{1 \times dp}$ is given by

$$\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}[j] := \mathbf{1}_{f_{\mathbf{x}^T \mathbf{V}_0[j]} > 0g} \mathbf{x}^T, \quad j = 1, 2, \dots, p. \quad (1)$$

In the NTK model, we assume that the output of the trained model is exactly given by Eq. (1), i.e.,

$$\hat{f}_{\Delta \mathbf{V}, \mathbf{V}_0}(\mathbf{x}) := \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V}. \quad (2)$$

In other words, the NTK model can be viewed as a linear approximation of the two-layer network when the change of the bottom-layer weights is small.

Define $\mathbf{H} \in \mathbb{R}^{n \times dp}$ such that its i -th row is $\mathbf{H}_i := \mathbf{h}_{\mathbf{V}_0, \mathbf{X}_i}$. Throughout the paper, we will focus on the following min- ℓ_2 -norm overfitting solution

$$\mathbf{V}^{\ell_2} := \arg \min_{\mathbf{V}} \|\mathbf{H}\mathbf{V} - \mathbf{y}\|_2, \quad \text{subject to } \mathbf{H}\mathbf{V} = \mathbf{y}.$$

Whenever \mathbf{V}^{ℓ_2} exists, it can be written in closed form as

$$\mathbf{V}^{\ell_2} = \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y}. \quad (3)$$

The reason that we are interested in \mathbf{V}^{ℓ_2} is that gradient descent (GD) or stochastic gradient descent (SGD) for the NTK model in Eq. (2) is known to converge to \mathbf{V}^{ℓ_2} (proven in Supplementary Material, Appendix B).

Using Eq. (2) and Eq. (3), the trained model is then

$$\hat{f}^{\ell_2}(\mathbf{x}) := \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V}^{\ell_2}. \quad (4)$$

In the rest of the paper, we will study the generalization error of Eq. (4).

We collect some assumptions. Define the unit sphere in \mathbb{R}^d as: $S^{d-1} := \{\mathbf{v} \in \mathbb{R}^d \mid \|\mathbf{v}\|_2 = 1\}$. Let $\mu(\cdot)$ denote the distribution of the input \mathbf{x} . Without loss of generality, we make the following assumptions: (i) the inputs \mathbf{x} are *i.i.d.* uniformly distributed in S^{d-1} , and the initial weights $\mathbf{V}_0[j]$'s are *i.i.d.* uniformly distributed in all directions in \mathbb{R}^d ; (ii) $p \gg n/d$ and $d \geq 2$; (iii) $\mathbf{X}_i \perp \mathbf{X}_j$ for any $i \neq j$, and $\mathbf{V}_0[k] \perp \mathbf{V}_0[l]$ for any $k \neq l$. We provide detailed justification of those assumptions in Supplementary Material, Appendix C.

3. Learnable Functions and Generalization Performance

We now show that the generalization performance of the overfitted NTK model in Eq. (4) crucially depends on the

ground-truth function $f(\cdot)$, where good generalization performance only occurs when the ground-truth function is ‘‘learnable.’’ Below, we first describe a candidate class of ground-truth functions, and explain why they may correspond to the class of ‘‘learnable functions.’’ Then, we will give an upper-bound on the generalization performance for this class of ground-truth functions. Finally, we will give a lower-bound on the generalization performance when the ground-truth functions are outside of this class.

We first define a set F^{ℓ_2} of ground-truth functions.

Definition 1. $F^{\ell_2} := \{f \stackrel{\text{a.e.}}{=} f_g \mid f_g(\mathbf{x}) = \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} g(\mathbf{z}) d\mu(\mathbf{z}), \text{ } kgk_1 < 1\}$.

Note that in Definition 1, $\stackrel{\text{a.e.}}{=}$ means two functions equals almost everywhere, and $kgk_1 := \int_{S^{d-1}} |g(\mathbf{z})| d\mu(\mathbf{z})$. The function $g(\mathbf{z})$ may be any finite-value function in $L^1(S^{d-1}; \mathbb{R})$. Further, we also allow $g(\mathbf{z})$ to contain (as components) Dirac δ -functions on S^{d-1} . Note that a δ -function $\delta_{\mathbf{z}_0}(\mathbf{z})$ has zero value for all $\mathbf{z} \in S^{d-1} \setminus \{\mathbf{z}_0\}$, but $\int_{S^{d-1}} \delta_{\mathbf{z}_0}(\mathbf{z}) d\mu(\mathbf{z}) = 1$. Thus, the function $g(\mathbf{z})$ may contain any sum of δ -functions and finite-value L^1 -functions.¹

To see why F^{ℓ_2} may correspond to the class of learnable functions, we can first examine what the learned function \hat{f}^{ℓ_2} in Eq. (4) should look like. Recall that $\mathbf{H}^T = [\mathbf{H}_1^T \dots \mathbf{H}_n^T]$. Thus, $\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T = \sum_{i=1}^n (\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}_i^T) \mathbf{e}_i^T$, where $\mathbf{e}_i \in \mathbb{R}^n$ denotes the i -th standard basis. Combining Eq. (3) and Eq. (4), we can see that the learned function in Eq. (4) is of the form

$$\begin{aligned} \hat{f}^{\ell_2}(\mathbf{x}) &= \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y} \\ &= \sum_{i=1}^n \left(\frac{1}{p} \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}_i^T \right) p \mathbf{e}_i^T (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y}. \end{aligned} \quad (5)$$

For all $\mathbf{x}, \mathbf{z} \in S^{d-1}$, define $C_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0} := \#\{j \in \{1, 2, \dots, p\} \mid \mathbf{z}^T \mathbf{V}_0[j] > 0, \mathbf{x}^T \mathbf{V}_0[j] > 0g\}$, and its cardinality is given by

$$|C_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0}| = \sum_{j=1}^p \mathbf{1}_{f_{\mathbf{z}^T \mathbf{V}_0[j]} > 0, \mathbf{x}^T \mathbf{V}_0[j] > 0g}. \quad (6)$$

Then, using Eq. (1), we can show $\frac{1}{p} \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}_i^T = \mathbf{x}^T \mathbf{X}_i \frac{j C_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p}$. It is not hard to show that

$$\frac{j C_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0}}{p} \approx \frac{\pi \arccos(\mathbf{x}^T \mathbf{z})}{2\pi}, \quad \text{as } p \rightarrow \infty. \quad (7)$$

¹Alternatively, we can also interpret $g(\mathbf{z})$ as a signed measure (Rao & Rao, 1983) on S^{d-1} . Then, δ -functions correspond to point masses, and the condition $kgk_1 < 1$ implies that the corresponding unsigned version of the measure on S^{d-1} is bounded.

where $\overset{P}{\rightarrow}$ denotes convergence in probability. (see Supplementary Material, Appendix D.5). Thus, if we let

$$g(\mathbf{z}) = \sum_{i=1}^n p e_i^T (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y} \delta_{\mathbf{x}_i}(\mathbf{z}), \quad (8)$$

then as $p \rightarrow \infty$, Eq. (5) should approach a function in F^{ℓ_2} . This explains why F^{ℓ_2} is a candidate class of ‘‘learnable functions.’’ However, note that the above discussion only addresses the *expressiveness* of the model. It is still unclear whether any function in F^{ℓ_2} can be learned with low generalization error. The following result provides the answer.

For some $m \geq \left[1, \frac{\ln n}{\ln \frac{n}{2}}\right]$, define (recall that d is the dimension of \mathbf{x})

$$J_m(n, d) := 2^{1.5d+5.5} d^{0.5d} n^{(2+\frac{1}{m})(d-1)}. \quad (9)$$

Theorem 1. Assume a ground-truth function $f \stackrel{a.e.}{=} f_g \in F^{\ell_2}$ where $kgk_1 < 1^2$, $n \geq 2$, $m \geq \left[1, \frac{\ln n}{\ln \frac{n}{2}}\right]$, $d \leq n^4$, and $p \geq 6J_m(n, d) \ln \left(4n^{1+\frac{1}{m}}\right)$. Then, for any $q \geq \left[1, \frac{1}{1}\right]$ and for almost every $\mathbf{x} \in S^{d-1}$, we must have³

$$\begin{aligned} & \Pr_{\mathbf{v}_0, \mathbf{X}} \left\{ \underbrace{|\int f^{\ell_2}(\mathbf{x}) - f(\mathbf{x})|}_{\text{Term 1}} \leq n^{-\frac{1}{2}(1-\frac{1}{q})} \right. \\ & + \underbrace{\left(1 + \sqrt{J_m(n, d)n}\right) p^{-\frac{1}{2}(1-\frac{1}{q})}}_{\text{Term 2}} + \underbrace{\sqrt{J_m(n, d)n} k \epsilon k_2}_{\text{Term 3}}, \\ & \left. \text{for all } \epsilon \in \mathbb{R}^n \right\} \leq 2e^2 \left(\exp\left(\frac{p^{-\frac{1}{q}}}{8kgk_2^2}\right) \right) \\ & + \underbrace{\exp\left(\frac{p^{-\frac{1}{q}}}{8kgk_1^2}\right)}_{\text{Term 5}} + \underbrace{\exp\left(\frac{p^{-\frac{1}{q}}}{8nkgk_1^2}\right)}_{\text{Term 6}} + \underbrace{\frac{4}{p^{-\frac{1}{q}}}}_{\text{Term 7}}. \quad (10) \end{aligned}$$

To interpret Theorem 1, we can first focus on the noiseless case, where ϵ and Term 3 are zero. If we fix n and let $p \rightarrow \infty$, then Terms 2, 5, and 6 all approach zero. We can then conclude that, in the noiseless and heavily overparameterized setting ($p \rightarrow \infty$), the generalization error will converge to a small limiting value (Term 1) that depends only on n . Further, this limiting value (Term 1) will converge to zero (so do Terms 4 and 7) as $n \rightarrow \infty$, i.e.,

²The requirement of $kgk_1 < 1$ can be relaxed. We show in Supplementary Material, Appendix L that, even when g is a δ -function (so $kgk_1 = 1$), we can still have a similar result of Eq. (10) but Term 1 will have a slower speed of decay $O(n^{-\frac{1}{2(d-1)}(1-\frac{1}{q})})$ with respect to n instead of $O(n^{-\frac{1}{2}(1-\frac{1}{q})})$ shown in Eq. (10). Term 4 of Eq. (10) will also be different when g is a δ -function, but it still goes to zero when p and n are large.

³The notion \Pr_M in Eq. (10) emphasizes that randomness is in M .

when there are sufficiently many training samples. Finally, Theorem 1 holds even when there is noise.

The parameters of q and m can be tuned to make Eq. (10) sharper when n and p are large. For example, as we increase q , Term 1 will approach $n^{-0.5}$. Although a larger q makes Terms 4, 5, and 6 bigger, as long as n and p are sufficiently large, those terms will still be close to 0. Similarly, if we increase m , then $J_m(n, d)$ will approach the order of $n^{2(d-1)}$. As a result, Term 3 approaches the order of $n^{2d-0.5}$ times $k\epsilon k_2$ and the requirement $p \geq 6J_m(n, d) \ln \left(4n^{1+\frac{1}{m}}\right)$ approaches the order of $n^{2(d-1)} \ln n$.

Remark 1. We note that (Arora et al., 2019) shows that, for two-layer neural networks whose bottom-layer weights are trained by gradient descent, the generalization error for sufficiently large p has the following upper bound: for any $\zeta > 0$,

$$\begin{aligned} & \Pr \left\{ \mathbb{E}_{\mathbf{x}} |\int \hat{f}(\mathbf{x}) - f(\mathbf{x})| \leq \sqrt{\frac{2\mathbf{y}^T (\mathbf{H}^T)^{-1} \mathbf{y}}{n}} \right. \\ & \left. + O\left(\sqrt{\frac{\log \frac{n}{\zeta \min \text{eig}(\mathbf{H}^T)}}{n}}\right) \right\} \leq 1 - \zeta, \quad (11) \end{aligned}$$

where $\mathbf{H}^T = \lim_{p \rightarrow \infty} (\mathbf{H}\mathbf{H}^T/p) \in \mathbb{R}^{n \times n}$. For certain class of learnable functions (we will compare them with our F^{ℓ_2} in Section 4), the quantity $\mathbf{y}^T (\mathbf{H}^T)^{-1} \mathbf{y}$ is bounded. Thus, $\sqrt{\frac{2\mathbf{y}^T (\mathbf{H}^T)^{-1} \mathbf{y}}{n}}$ also decreases at the speed $1/\sqrt{n}$. The second $O(\cdot)$ -term in Eq. (11) contains the minimum eigenvalue of \mathbf{H}^T , which decreases with n . (Indeed, we show that this minimum eigenvalue is upper bounded by $O(n^{-\frac{1}{d-1}})$ in Supplementary Material, Appendix G.) Thus, Eq. (11) may decrease a little bit slower than $1/\sqrt{n}$, which is consistent with Term 1 in Eq. (10) (when q is large). Note that the term $2\mathbf{y}^T (\mathbf{H}^T)^{-1} \mathbf{y}$ in Eq. (11) captures how the complexity of the ground-truth function affects the generalization error. Similarly, the norm of $g(\cdot)$ also captures the impact⁴ of the complexity of the ground-truth function in Eq. (10). However, we caution that the GD solution in (Arora et al., 2019) is based on the original neural network, which is usually different from our $\min \ell_2$ -norm solution based on the NTK model (even though they are close for very large p). Thus, the two results may not be directly comparable.

Theorem 1 reveals several important insights on the generalization performance when the ground-truth function belongs to F^{ℓ_2} .

(i) Descent in the overparameterized region: When p increases, both sides of Eq. (10) decreases, suggesting that the test error of the overfitted NTK model decreases with

⁴Although Term 1 in Eq. (10) in its current form does not depend on $g(\cdot)$, it is possible to modify our proof so that the norm of $g(\cdot)$ also enters Term 1.

p . In Fig. 1(a), we choose a ground-truth function in F^{ℓ_2} (we will explain why this function is in F^{ℓ_2} later in Section 4). The test MSE of the aforementioned NTK model (green curve) confirms the overall trend⁵ of descent in the overparameterized region. We note that while (Arora et al., 2019) provides a generalization error upper-bound for large p (i.e., Eq. (11)), the upper bound there does not capture the dependency in p and thus does not predict this descent.

More importantly, we note a significant difference between the descent in Theorem 1 and that of min ℓ_2 -norm overfitting solutions for linear models with simple features (Belkin et al., 2018b; 2019; Bartlett et al., 2020; Hastie et al., 2019; Muthukumar et al., 2019; Liao et al., 2020; Jacot et al., 2020). For example, for linear models with Gaussian features, we can obtain (see, e.g., Theorem 2 of (Belkin et al., 2019)):

$$\text{MSE} = kf k_2^2 \left(1 - \frac{n}{p} \right) + \frac{\sigma^2 n}{p - n - 1}, \text{ for } p > n + 2 \quad (12)$$

where σ^2 denotes the variance of the noise. If we let $p \rightarrow \infty$ in Eq. (12), we can see that the MSE quickly approaches $kf k_2^2$, which is referred to as the “null risk” (Hastie et al., 2019), i.e., the MSE of a model that predicts zero. Note that the null-risk is at the level of the signal, and thus is quite large. In contrast, as $p \rightarrow \infty$, the test error of the NTK model converges to a value determined by n and ϵ (and is independent of the null risk). This difference is confirmed in Fig. 1(a), where the test MSE for the NTK model (green curve) is much lower than the null risk (the dashed line) when $p \rightarrow \infty$, while both the min ℓ_2 -norm (the red curve) and the min ℓ_1 -norm solutions (the blue curve) (Ju et al., 2020) with Fourier features rise to the null risk when $p \rightarrow \infty$. Finally, note that the descent in Theorem 1 requires p to increase much faster than n . Specifically, to keep Term 2 in Eq. (10) small, it suffices to let p increase a little bit faster than (n^{4d-1}) . This is again quite different from the descent shown in Eq. (12) and in other related work using Fourier and Gaussian features (Liao et al., 2020; Jacot et al., 2020), where p only needs to grow proportionally with n .

(ii) Speed of the descent: Since Theorem 1 holds for finite p , it also characterizes the speed of descent. In particular, Term 2 is proportional to $p^{-\frac{1}{2}}(1 - \frac{1}{q})$, which approaches $1/p^{\frac{1}{2}}$ when q is large. Again, such a speed of descent is not captured in (Arora et al., 2019). As we show in Fig. 1(a), the test error of the gradient descent solution under the original neural network (cyan curve) is usually quite close to that of

⁵This curve oscillates at the early stage when p is small. We suspect it is because, at small p , the convergence in Eq. (7) has not occurred yet, and thus the randomness in $\mathbf{V}_0[j]$ makes the simulation results more volatile.

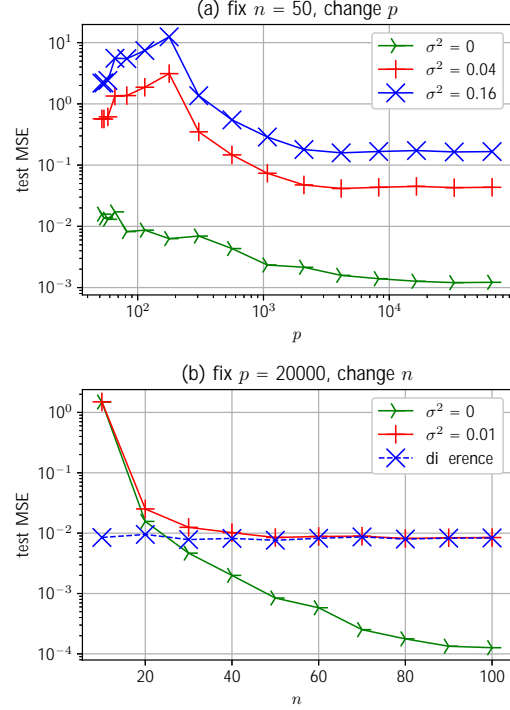


Figure 3. The test MSE of the overfitted NTK model for the same ground-truth function as Fig. 1(a). **(a)** We fix $n = 50$ and increase p for different noise level σ^2 . **(b)** We fix $p = 20000$ and increase n . All data points in this figure are the average of five random simulation runs.

the NTK model (green curve). Thus, our result provides useful guidance on how fast the generalization error descends with p for such neural networks.

(iii) The effect of noise: Term 3 in Eq. (10) characterizes the impact of the noise ϵ , which does not decrease or increase with p . Notice that this is again very different from Eq. (12), i.e., results of min ℓ_2 -norm overfitting solutions for simple features, where the noise term $\frac{\sigma^2 n}{p - n - 1} \rightarrow 0$ when $p \rightarrow \infty$. We use Fig. 3(a) to validate this insight. In Fig. 3(a), we fix $n = 50$ and plot curves of test MSE of NTK overfitting solution as p increases. We let the noise ϵ_i in the i -th training sample be *i.i.d.* Gaussian with zero mean and variance σ^2 . The green, red, and blue curves in Fig. 3(a) corresponds to the situation $\sigma^2 = 0$, $\sigma^2 = 0.04$, and $\sigma^2 = 0.16$, respectively. We can see that all three curves become flat when p is very large, and this phenomenon implies that the gap across different noise levels does not decrease when $p \rightarrow \infty$, which is in contrast to Eq. (12).

In Fig. 3(b), we instead fix $p = 20000$, and increase n . We plot the test MSE both for the noiseless setting (green curve) and for $\sigma^2 = 0.01$ (red curve). The difference between the two curves (dashed blue curve) then captures the impact of noise, which is related to Term 3 in Eq. (10). Somewhat

surprisingly, we find that the dashed blue curve is insensitive to n , which suggests that Term 3 in Eq. (10) may have room for improvement.

In summary, we have shown that any ground-truth function in F^{ℓ_2} leads to low generalization error for overfitted NTK models. It is then natural to ask what happens if the ground-truth function is not in F^{ℓ_2} . Let $\overline{F^{\ell_2}}$ denote the closure⁶ of F^{ℓ_2} , and $D(f, F^{\ell_2})$ denotes the L^2 -distance between f and F^{ℓ_2} (i.e., the infimum of the L^2 -distance from f to every function in F^{ℓ_2}).

Proposition 2. (i) For any given (\mathbf{X}, \mathbf{y}) , there exists a function $\hat{f}_p^{\ell_2} \in F^{\ell_2}$ such that, uniformly over all $\mathbf{x} \in S^{d-1}$, $\hat{f}_p^{\ell_2}(\mathbf{x}) \rightarrow \hat{f}^{\ell_2}(\mathbf{x})$ as $p \rightarrow \infty$. (ii) Consequently, if the ground-truth function $f \notin \overline{F^{\ell_2}}$ (or equivalently, $D(f, F^{\ell_2}) > 0$), then the MSE of $\hat{f}_p^{\ell_2}$ (with respect to the ground-truth function f) is at least $D(f, F^{\ell_2})$.

Intuitively, Proposition 2 (proven in Supplementary Material Appendix J) suggests that, if a ground-truth function is outside the closure of F^{ℓ_2} , then no matter how large n is, the test error of a NTK model with infinitely many neurons cannot be small (regardless whether or not the training samples contain noise). We validate this in Fig. 1(b), where a ground-truth function is chosen outside $\overline{F^{\ell_2}}$. The test MSE of NTK overfitting solutions (green curve) is above null risk (dashed black line) and thus is much higher compared with Fig. 1(a). We also plot the test MSE of the GD solution of the real neural network (cyan curve), which seems to show the same trend.

Comparing Theorem 1 and Proposition 2, we can clearly see that, all functions in F^{ℓ_2} are learnable by the overfitted NTK model, and all functions not in $\overline{F^{\ell_2}}$ are not.

4. What Exactly are the Functions in F^{ℓ_2} ?

Our expression for learnable functions in Definition 1 is still in an indirect form, i.e., through the unknown function $g(\cdot)$. In (Arora et al., 2019), the authors show that all functions of the form $(\mathbf{x}^T \mathbf{a})^l$, $l \in \{0, 1, 2, 4, 6, \dots\}$ are learnable by GD (assuming large p and small step size), for a similar 2-layer network with ReLU activation that has no bias. In the following, we will show that our learnable functions in Definition 1 also have a similar form. Further, we can show that any functions of the form $(\mathbf{x}^T \mathbf{a})^l$, $l \in \{3, 5, 7, \dots\}$ are not learnable. Our characterization uses an interesting connection to harmonics and filtering on S^{d-1} , which may be of independent interest.

Towards this end, we first note that the integral form in Def-

⁶We consider the normed space of all functions in $L^2(S^{d-1} \setminus \mathbb{R})$. Notice that although $g(\mathbf{z})$ in Definition 1 may not be in L^2 , f_g is always in L^2 . Specifically, $f_g(\mathbf{x})$ is bounded for every $\mathbf{x} \in S^{d-1}$ when $kgk_1 < 1$.

inition 1 can be viewed as a convolution on S^{d-1} (denoted by \sim). Specifically, for any $f_g \in F^{\ell_2}$, we can rewrite it as

$$f_g(\mathbf{x}) = g \sim h(\mathbf{x}) := \int_{\text{SO}(d)} g(\mathbf{S}\mathbf{e})h(\mathbf{S}^{-1}\mathbf{x})d\mathbf{S}, \quad (13)$$

$$h(\mathbf{x}) := \mathbf{x}^T \mathbf{e} \frac{\pi}{2\pi} \frac{\arccos(\mathbf{x}^T \mathbf{e})}{2\pi}, \quad (14)$$

where $\mathbf{e} := [0 \ 0 \ \dots \ 0 \ 1]^T \in \mathbb{R}^d$, and \mathbf{S} is a $d \times d$ orthogonal matrix that denotes a rotation in S^{d-1} , chosen from the set $\text{SO}(d)$ of all rotations. An important property of the convolution Eq. (13) is that it corresponds to multiplication in the frequency domain, similar to Fourier coefficients. To define such a transformation to the frequency domain, we use a set of hyper-spherical harmonics $\{Y_{\mathbf{K}}^l\}_{\mathbf{K}}^l$ (Vilenkin, 1968; Dokmanic & Petrinovic, 2009) when $d \geq 3$, which forms an orthonormal basis for functions on S^{d-1} . These harmonics are indexed by l and \mathbf{K} , where $\mathbf{K} = (k_1, k_2, \dots, k_{d-2})$ and $l = k_0 + k_1 + k_2 + \dots + k_{d-2} + 0$ (those k_i 's and l are all non-negative integers). Any function $f \in L^2(S^{d-1} \setminus \mathbb{R})$ (including even δ -functions (Li & Wong, 2013)) can be decomposed uniquely into these harmonics, i.e., $f(\mathbf{x}) = \sum_l \sum_{\mathbf{K}} c_f(l, \mathbf{K}) Y_{\mathbf{K}}^l(\mathbf{x})$, where $c_f(\cdot, \cdot)$ are projections of f onto the basis function. In Eq. (13), let $c_g(\cdot, \cdot)$ and $c_h(\cdot, \cdot)$ denote the coefficients corresponding to the decompositions of g and h , respectively. Then, we must have (Dokmanic & Petrinovic, 2009)

$$c_{f_g}(l, \mathbf{K}) = \gamma c_g(l, \mathbf{K})c_h(l, \mathbf{0}), \quad (15)$$

where γ is some normalization constant. Notice that in Eq. (15), the coefficient for h is $c_h(l, \mathbf{0})$ instead of $c_h(l, \mathbf{K})$, which is due to the intrinsic rotational symmetry of such convolution (Dokmanic & Petrinovic, 2009).

The above decomposition has an interesting ‘‘filtering’’ interpretation as follows. We can regard the function h as a ‘‘filter’’ or ‘‘channel,’’ while the function g as a transmitted ‘‘signal.’’ Then, the function f_g in Eq. (13) and Eq. (15) can be regarded as the received signal after g goes through the channel/filter h . Therefore, when coefficient $c_h(l, \mathbf{0})$ of h is non-zero, then the corresponding coefficient $c_{f_g}(l, \mathbf{K})$ for f_g can be any value (because we can arbitrarily choose g). In contrast, if a coefficient $c_h(l, \mathbf{0})$ of h is zero, then the corresponding coefficient $c_{f_g}(l, \mathbf{K})$ for f_g must also be zero for all \mathbf{K} .

Ideally, if h contains all ‘‘frequencies,’’ i.e., all coefficients $c_h(l, \mathbf{0})$ are non-zero, then f_g can also contain all ‘‘frequencies,’’ which means that F^{ℓ_2} can contain almost all functions. Unfortunately, this is not true for the function h given in Eq. (14). Specifically, using the harmonics defined in (Dokmanic & Petrinovic, 2009), the basis $Y_{\mathbf{0}}^l$ for $(l, \mathbf{0})$ turns out to have the form

$$Y_{\mathbf{0}}^l(\mathbf{x}) = \sum_{k=0}^{\lfloor \frac{l}{2} \rfloor} \binom{l}{k} a_{l,k} (\mathbf{x}^T \mathbf{e})^{l-2k}, \quad (16)$$

where $a_{l,k}$ are positive constants. Note that the expression Eq. (16) contains either only even powers of $\mathbf{x}^T \mathbf{e}$ (if l is even) or odd powers of $\mathbf{x}^T \mathbf{e}$ (if l is odd). Then, for the function h in Eq. (14), we have the following proposition (proven in Supplementary Material, Appendix K.4). We note that (Basri et al., 2019) has a similar harmonics analysis, where the expression of $c_h(l, \mathbf{0})$ is given. However, it is not obvious that the expression of $c_h(l, \mathbf{0})$ for all $l = 0, 1, 2, 4, 6, \dots$ given in (Basri et al., 2019) must be non-zero, which is made clear by Proposition 3 as follows.

Proposition 3. $c_h(l, \mathbf{0})$ is zero for $l = 3, 5, 7, \dots$ and is non-zero for $l = 0, 1, 2, 4, 6, \dots$.

We are now ready to characterize what functions are in F^{ℓ_2} . By the form of Eq. (16), for any non-negative integer k , any even power $(\mathbf{x}^T \mathbf{e})^{2k}$ is a linear combination of $\frac{0}{0}, \frac{2}{0}, \dots, \frac{2k}{0}$, and any odd power $(\mathbf{x}^T \mathbf{e})^{2k+1}$ is a linear combination of $\frac{1}{0}, \frac{3}{0}, \dots, \frac{2k+1}{0}$. By Proposition 3, we thus conclude that any function $f_g(\mathbf{x}) = (\mathbf{x}^T \mathbf{e})^l$ where $l \geq 0, 1, 2, 4, 6, \dots$ can be written in the form of Eq. (15) in the frequency domain, and thus are in F^{ℓ_2} . In contrast, any function $f(\mathbf{x}) = (\mathbf{x}^T \mathbf{e})^l$ where $l \geq 3, 5, 7, \dots$ cannot be written in the form of Eq. (15), and are thus not in F^{ℓ_2} . Further, the ℓ_2 -norm of any latter function will also be equal to its distance to F^{ℓ_2} . Therefore, the generalization-error lower-bound in Proposition 2 will apply (with $D(f, F^{\ell_2}) = kf k_2$). Finally, by Eq. (13), F^{ℓ_2} is invariant under rotation and finite linear summation. Therefore, any finite sum of $(\mathbf{x}^T \mathbf{a})^l$, $l = 0, 1, 2, 4, 6, \dots$ must also belong to F^{ℓ_2} .

For the special case of $d = 2$, the input \mathbf{x} corresponds to an angle $\theta \in [\pi, \pi]$, and the above-mentioned harmonics become Fourier series $\sin(k\theta)$ and $\cos(k\theta)$, $k = 0, 1, \dots$. We can then get similar results that frequencies of $k \geq 0, 1, 2, 4, 6, \dots$ are learnable (while others are not), which explains the learnable and not-learnable functions in Fig. 1. Details can be found in Supplementary Material, Appendix K.5.

Remark 2. We caution that the above claim on non-learnable functions critically depends on the network architecture. That is, we assume throughout this paper that the ReLU activation has no bias. It is known from an expressiveness point of view that, using ReLU without bias, a shallow network can only approximate the sum of linear functions and even functions (Ghorbani et al., 2019). Thus, it is not surprising that other odd-power (but non-linear) polynomials cannot be learned. In contrast, by adding a bias, a shallow network using ReLU becomes a universal approximator (Ji et al., 2019). The recent work of (Satpathi & Srikant, 2021) shows that polynomials with all powers can be learned by the corresponding 2-layer NTK model. These results are consistent with ours because a ReLU activation function operating on $\mathbf{x} \in \mathbb{R}^{d-1}$ with a bias can be equivalently viewed as one

operating on a d -dimension input (with the last-dimension being fixed at $1/\sqrt{d}$) but with no bias. Even though only a subset of functions are learnable in the d -dimension space, when projected into a $(d-1)$ -dimension subspace, they may already span all functions. For example, one could write $(\mathbf{x}^T \mathbf{a})^3$ as a linear combination of $(\left[\frac{1}{\sqrt{d}} \right]^T \mathbf{b}_i)^{l_i}$, where $i \geq 1, 2, \dots, 5$, $[l_1, \dots, l_5] = [4, 4, 2, 1, 0]$, and $\mathbf{b}_i \in \mathbb{R}^d$ depends only on \mathbf{a} . (See Supplementary Material, Appendix K.6 for details.) It remains an interesting question whether similar difference arises for other network architectures (e.g., with more than 2 layers).

5. Proof Sketch of Theorem 1

In this section, we sketch the key steps to prove Theorem 1. Starting from Eq. (3), we have

$$\mathbf{V}^{\ell_2} = \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} (\mathbf{F}(\mathbf{X}) + \boldsymbol{\epsilon}). \quad (17)$$

For the learned model $\hat{f}^{\ell_2}(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}^T \mathbf{V}^{\ell_2}$ given in Eq. (4), the error for any test input \mathbf{x} is then

$$\begin{aligned} \hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) &= (\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}^T \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{F}(\mathbf{X}) - f(\mathbf{x})) \\ &\quad + \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}^T \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \boldsymbol{\epsilon}. \end{aligned} \quad (18)$$

In the classical ‘‘bias-variance’’ analysis with respect to MSE (Belkin et al., 2018a), the first term on the right-hand-side of Eq. (18) contributes to the bias and the second term contributes to the variance. We first quantify the second term (i.e., the variance) in the following proposition.

Proposition 4. For any $n \geq 2$, $m \geq \left\lceil 1, \frac{\ln n}{\ln \frac{n}{2}} \right\rceil$, $d \leq n^4$, if $p \geq 6J_m(n, d) \ln \left(4n^{1+\frac{1}{m}} \right)$, we must have $\Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ \left| \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}^T \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \boldsymbol{\epsilon} \right| \leq \sqrt{J_m(n, d) n \kappa \epsilon k_2}, \text{ for all } \boldsymbol{\epsilon} \in \mathbb{R}^n \right\} \geq 1 - \frac{2}{n}$.

The proof is in Supplementary Material Appendix F. Proposition 4 implies that, for fixed n and d , when $p \rightarrow \infty$, with high probability the variance will not exceed a certain factor of the noise $\kappa \epsilon k_2$. In other words, the variance will not go to infinity when $p \rightarrow \infty$. The main step in the proof is to lower bound $\min \text{eig}(\mathbf{H} \mathbf{H}^T) / p$, which is given by $1 / (J_m(n, d) n)$. Note that this is the main place where we used the assumption that \mathbf{x} is uniformly distributed. We expect that our main proof techniques can be generalized to other distributions (with a different expression of $J_m(n, d)$), which we leave for future work.

Remark 3. In the upper bound in (Arora et al., 2019) (i.e., Eq. (11)), any noise added to \mathbf{y} will at least contribute to the generalization upper bound Eq. (11) by a positive term $\boldsymbol{\epsilon}^T (\mathbf{H}^T)^{-1} \boldsymbol{\epsilon} / n$. Thus, their upper bound may also grow as $\min \text{eig}(\mathbf{H}^T)$ decreases. One of the contribution of Proposition 4 is to characterize this minimum eigenvalue.

We now bound the bias part. We first study the class of ground-truth functions that can be learned with fixed \mathbf{V}_0 . We refer to them as *pseudo ground-truth*, to differentiate them with the set F^{ℓ_2} of learnable functions for random \mathbf{V}_0 . They are defined with respect to the same $g(\cdot)$ function, so that we can later extend to the “real” ground-truth functions in F^{ℓ_2} when considering the randomness of \mathbf{V}_0 .

Definition 2. Given \mathbf{V}_0 , for any learnable ground-truth function $f_g \in F^{\ell_2}$ with the corresponding function $g(\cdot)$, define the corresponding *pseudo ground-truth* as

$$f_{\mathbf{V}_0}^g(\mathbf{x}) := \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{j C_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0}}{p} g(\mathbf{z}) d\mu(\mathbf{z}).$$

The reason that this class of functions may be the learnable functions for fixed \mathbf{V}_0 is similar to the discussions in Eq. (5) and Eq. (6). Indeed, using the same choice of $g(\mathbf{z})$ in Eq. (8), the learned function \hat{f}^{ℓ_2} in Eq. (5) at fixed \mathbf{V}_0 is always of the form in Definition 2.

The following proposition gives an upper bound of the generalization performance when the data model is based on the pseudo ground-truth and the NTK model uses exactly the same \mathbf{V}_0 .

Proposition 5. Assume fixed \mathbf{V}_0 (thus p and d are also fixed), there is no noise. If the ground-truth function is $f = f_{\mathbf{V}_0}^g$ in Definition 2 and $kgk_1 < 1$, then for any $\mathbf{x} \in S^{d-1}$ and $q \in [1, 1)$, we have $\Pr_{\mathbf{X}} \{ |j \hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x})| \leq n^{-\frac{1}{2}(1-\frac{1}{q})} \} \geq 1 - 2e^2 \exp\left(-\frac{q}{8kgk_1^2}\right)$.

The proof is in Supplementary Material, Appendix H. Note that both the threshold of the probability event and the upper bound coincide with Term 1 and Term 4, respectively, in Eq. (10). Here we sketch the proof of Proposition 5. Based on the definition of the pseudo ground-truth, we can rewrite $f_{\mathbf{V}_0}^g$ as $f_{\mathbf{V}_0}^g(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}^T \mathbf{V}$, where $\mathbf{V} \in \mathbb{R}^{dp}$ is given by, for all $j \in \{1, 2, \dots, pg\}$, $\mathbf{V}[j] = \int_{S^{d-1}} \mathbf{1}_{\mathbf{r}_{\mathbf{z}^T \mathbf{V}_0}[j] > 0} g(\mathbf{z}) \frac{g(\mathbf{z})}{p} d\mu(\mathbf{z})$. From Eq. (3) and Eq. (4), we can see that the learned model is $\hat{f}^{\ell_2}(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}^T \mathbf{P} \mathbf{V}$ where $\mathbf{P} := \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{H}$. Note that \mathbf{P} is an orthogonal projection to the row-space of \mathbf{H} . Further, it is easy to show that $k \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} k_2 \leq \frac{1}{\sqrt{p}}$. Thus, we have $|j \hat{f}^{\ell_2}(\mathbf{x}) - f_{\mathbf{V}_0}^g(\mathbf{x})| = |\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}^T (\mathbf{P} - \mathbf{I}) \mathbf{V}[j]| \leq \frac{1}{\sqrt{p}} k(\mathbf{P} - \mathbf{I}) \mathbf{V}[j] k_2$. The term $(\mathbf{P} - \mathbf{I}) \mathbf{V}$ can be interpreted as the distance from \mathbf{V} to the row-space of \mathbf{H} . Note that this distance is no greater than the distance between \mathbf{V} and any point in the row-space of \mathbf{H} . Thus, in order to get an upper bound on $k(\mathbf{P} - \mathbf{I}) \mathbf{V}[j] k_2$, we only need to find a vector $\mathbf{a} \in \mathbb{R}^n$ that makes $k \mathbf{V} - \mathbf{H}^T \mathbf{a} k_2$ as small as possible, especially when n is large. Our proof uses the vector \mathbf{a} such that its i -th element is $\mathbf{a}_i := \frac{g(\mathbf{X}_i)}{np}$. See Supplementary Material, Appendix H for the rest of the details.

The final step is to allow \mathbf{V}_0 to be random. Given any random \mathbf{V}_0 , any function $f_g \in F^{\ell_2}$ can be viewed as the summation of a pseudo ground-truth function (with the same $g(\cdot)$) and a difference term. This difference can be viewed as a special form of “noise”, and thus we can use Proposition 4 to quantify its impact. Further, the magnitude of this “noise” should decrease with p (because of Eq. (7)). Combining this argument with Proposition 5, we can then prove Theorem 1. See Supplementary Material, Appendix I for details.

6. Conclusions

In this paper, we studied the generalization performance of the min ℓ_2 -norm overfitting solution for a two-layer NTK model. We provide a precise characterization of the learnable ground-truth functions for such models, by providing a generalization upper bound for all functions in F^{ℓ_2} , and a generalization lower bound for all functions not in F^{ℓ_2} . We show that, while the test error of the overfitted NTK model also exhibits descent in the overparameterized regime, the descent behavior can be quite different from the double descent of linear models with simple features.

There are several interesting directions for future work. First, based on Fig. 3(b), our estimation of the effect of noise could be further improved. Second, it would be interesting to explore whether the methodology can be extended to NTK model for other neural networks, e.g., with different activation functions and with more than two layers.

Acknowledgements

This work is partially supported by an NSF sub-award via Duke University (IIS-1932630), by NSF grants CNS-1717493, CNS-1901057, and CNS-2007231, and by Office of Naval Research under Grant N00014-17-1-241. The authors would like to thank Professor R. Srikant at the University of Illinois at Urbana-Champaign and anonymous reviewers for their valuable comments and suggestions.

References

- Advani, M. S., Saxe, A. M., and Sompolinsky, H. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- Allen-Zhu, Z., Li, Y., and Liang, Y. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in neural information processing systems*, pp. 6158–6169, 2019.
- Arora, S., Du, S., Hu, W., Li, Z., and Wang, R. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pp. 322–332, 2019.

- Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 2020.
- Basri, R., Jacobs, D., Kasten, Y., and Kritchman, S. The convergence rate of neural networks for learned functions of different frequencies. *arXiv preprint arXiv:1906.00425*, 2019.
- Belkin, M., Hsu, D., Ma, S., and Mandal, S. Reconciling modern machine learning and the bias-variance trade-off. *arXiv preprint arXiv:1812.11118*, 2018a.
- Belkin, M., Ma, S., and Mandal, S. To understand deep learning we need to understand kernel learning. In *International Conference on Machine Learning*, pp. 541–549, 2018b.
- Belkin, M., Hsu, D., and Xu, J. Two models of double descent for weak features. *arXiv preprint arXiv:1903.07571*, 2019.
- Bishop, C. M. *Pattern recognition and machine learning*. Springer, 2006.
- Chaudhry, M. A., Qadir, A., Rafique, M., and Zubair, S. Extension of euler’s beta function. *Journal of computational and applied mathematics*, 78(1):19–32, 1997.
- Dokmanic, I. and Petrinovic, D. Convolution on the n -sphere with application to pdf modeling. *IEEE transactions on signal processing*, 58(3):1157–1170, 2009.
- Du, S. S., Zhai, X., Póczos, B., and Singh, A. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2018.
- Dutka, J. The incomplete beta function—a historical profile. *Archive for history of exact sciences*, pp. 11–29, 1981.
- d’Ascoli, S., Refinetti, M., Biroli, G., and Krzakala, F. Double trouble in double descent: Bias and variance (s) in the lazy regime. In *International Conference on Machine Learning*, pp. 2280–2290. PMLR, 2020.
- Fiat, J., Malach, E., and Shalev-Shwartz, S. Decoupling gating from linearity. *arXiv preprint arXiv:1906.05032*, 2019.
- Ghorbani, B., Mei, S., Misiakiewicz, T., and Montanari, A. Linearized two-layers neural networks in high dimension. *arXiv preprint arXiv:1904.12191*, 2019.
- Goemans, M. Chernoff bounds, and some applications. URL <http://math.mit.edu/~goemans/18310S15/chernoff-notes.pdf>, 2015.
- Hastie, T., Tibshirani, R., and Friedman, J. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- Hastie, T., Montanari, A., Rosset, S., and Tibshirani, R. J. Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
- Hayes, T. P. A large-deviation inequality for vector-valued martingales. *Combinatorics, Probability and Computing*, 2005.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- Jacot, A., Simsek, B., Spadaro, F., Hongler, C., and Gabriel, F. Implicit regularization of random feature models. In *International Conference on Machine Learning*, pp. 4631–4640. PMLR, 2020.
- James, W. and Stein, C. Estimation with quadratic loss. In *Breakthroughs in Statistics*, pp. 443–460. Springer, 1992.
- Ji, Z. and Telgarsky, M. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. *arXiv preprint arXiv:1909.12292*, 2019.
- Ji, Z., Telgarsky, M., and Xian, R. Neural tangent kernels, transportation mappings, and universal approximation. *arXiv preprint arXiv:1910.06956*, 2019.
- Ju, P., Lin, X., and Liu, J. Overfitting can be harmless for basis pursuit, but only to a degree. *Advances in Neural Information Processing Systems*, 33, 2020.
- LeCun, Y., Kanter, I., and Solla, S. A. Second order properties of error surfaces: Learning time and generalization. In *Advances in Neural Information Processing Systems*, pp. 918–924, 1991.
- Li, S. Concise formulas for the area and volume of a hyperspherical cap. *Asian Journal of Mathematics and Statistics*, 4(1):66–70, 2011.
- Li, Y. and Liang, Y. Learning overparameterized neural networks via stochastic gradient descent on structured data. *Advances in Neural Information Processing Systems*, 31: 8157–8166, 2018.
- Li, Y. and Wong, R. Integral and series representations of the dirac delta function. *arXiv preprint arXiv:1303.1943*, 2013.
- Liao, Z., Couillet, R., and Mahoney, M. W. A random matrix analysis of random fourier features: beyond the gaussian kernel, a precise phase transition, and the corresponding double descent. *arXiv preprint arXiv:2006.05013*, 2020.

- Mei, S. and Montanari, A. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- Muthukumar, V., Vodrahalli, K., and Sahai, A. Harmless interpolation of noisy data in regression. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pp. 2299–2303. IEEE, 2019.
- Rao, K. B. and Rao, M. B. *Theory of charges: a study of finitely additive measures*. Academic Press, 1983.
- Satpathi, S. and Srikant, R. The dynamics of gradient descent for overparametrized neural networks. In *3rd Annual Learning for Dynamics and Control Conference (LADC)*, 2021.
- Stein, C. Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. Technical report, Stanford University Stanford United States, 1956.
- Tikhonov, A. N. On the stability of inverse problems. In *Dokl. Akad. Nauk SSSR*, volume 39, pp. 195–198, 1943.
- Vilenkin, N. Y. Special functions and the theory of group representations. providence: American mathematical society. *sftp*, 1968.
- Wainwright, M. Uniform laws of large numbers, 2015. https://www.stat.berkeley.edu/~mhwain/stat210b/Chap4_Uniform_Feb4_2015.pdf, Accessed: Feb. 7, 2021.
- Wendel, J. G. A problem in geometric probability. *Math. Scand*, 11:109–111, 1962.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. In *5th International Conference on Learning Representations, ICLR 2017*, 2017.

A. Extra Notations

In addition to the notations that we have introduced in the main body of this paper, we need some extra notations that are used in the following appendices. The distribution of the initial weights $\mathbf{V}_0[j]$ is denoted by the probability density $\lambda(\cdot)$ on \mathbb{R}^d , and the directions of the initial weights (i.e., the normalized initial weights $\frac{\mathbf{V}_0[j]}{k\mathbf{V}_0[j]k_2}$) follows the probability density $\lambda(\cdot)$ on S^{d-1} . Let $\lambda_a(\cdot)$ be the Lebesgue measure on \mathbb{R}^a where the dimension a can be, e.g., $(d-1)$ and $(d-2)$.

Let $\text{Bin}(a, b)$ denote the binomial distribution, where a is the number of trials and b is the success probability. Let $I(\cdot, \cdot)$ denote the regularized incomplete beta function (Dutka, 1981). Let $B(\cdot, \cdot)$ denote the beta function (Chaudhry et al., 1997). Specifically,

$$B(x, y) := \int_0^1 t^{x-1}(1-t)^{y-1} dt, \quad (19)$$

$$I_x(a, b) := \frac{\int_0^x t^{a-1}(1-t)^{b-1} dt}{B(a, b)}. \quad (20)$$

Define a cap on a unit hyper-sphere S^{d-1} as the intersection of S^{d-1} with an open ball in \mathbb{R}^d centered at \mathbf{v} with radius r , i.e.,

$$B_v^r := \{\mathbf{v} \in S^{d-1} \mid \|\mathbf{v} - \mathbf{v}\|_2 < r\}. \quad (21)$$

Remark 4. For ease of exposition, we will sometimes neglect the subscript \mathbf{v} of B_v^r and use B^r instead, when the quantity that we are estimating only depends on r but not \mathbf{v} . For example, where we are interested in the area of B_v^r , it only depends on r but not \mathbf{v} . Thus, we write $\lambda_{d-1}(B^r)$ instead.

For any $\mathbf{x} \in \mathbb{R}^d$ such that $\mathbf{x}^T \mathbf{v} = 0$, define two halves of the cap B_v^r as

$$B_v^{r, \mathbf{x}, +} := \{\mathbf{v} \in B_v^r \mid \mathbf{x}^T \mathbf{v} > 0\}, \quad B_v^{r, \mathbf{x}, -} := \{\mathbf{v} \in B_v^r \mid \mathbf{x}^T \mathbf{v} < 0\}. \quad (22)$$

Define the set of directions of the initial weights $\mathbf{V}_0[j]$'s as

$$A_{\mathbf{V}_0} := \left\{ \frac{\mathbf{V}_0[j]}{k\mathbf{V}_0[j]k_2} \mid j \in \{1, 2, \dots, pg\} \right\}. \quad (23)$$

B. GD (gradient descent) Converges to Min ℓ_2 -Norm Solutions

We assume that the GD algorithm for minimizing the training MSE is given by

$$\mathbf{V}_{k+1}^{\text{GD}} = \mathbf{V}_k^{\text{GD}} - \gamma_k \sum_{i=1}^n (\mathbf{H}_i \mathbf{V}_k^{\text{GD}} - \mathbf{y}_i) \mathbf{H}_i^T, \quad (24)$$

where \mathbf{V}_k^{GD} denotes the solution in the k -th GD iteration ($\mathbf{V}_0^{\text{GD}} = \mathbf{0}$), and γ_k denotes the step size of the k -th iteration.

Lemma 6. *If \mathbf{V}^{ℓ_2} exists and GD in Eq. (24) converges to zero-training loss (i.e., $\mathbf{H} \mathbf{V}_\gamma^{\text{GD}} = \mathbf{y}$), then $\mathbf{V}_\gamma^{\text{GD}} = \mathbf{V}^{\ell_2}$.*

Proof. Because $\mathbf{V}_0^{\text{GD}} = \mathbf{0}$ and Eq. (24), we know that \mathbf{V}_k^{GD} is in the row space of \mathbf{H} for any k . Thus, we can let $\mathbf{V}_\gamma^{\text{GD}} = \mathbf{H}^T \mathbf{a}$ where $\mathbf{a} \in \mathbb{R}^n$. When GD converges to zero training loss, we have $\mathbf{H} \mathbf{V}_\gamma^{\text{GD}} = \mathbf{y}$. Thus, we have $\mathbf{H} \mathbf{H}^T \mathbf{a} = \mathbf{y}$, which implies $\mathbf{a} = (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{y}$. Therefore, we must have $\mathbf{V}_\gamma^{\text{GD}} = \mathbf{H}^T \mathbf{a} = \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{y} = \mathbf{V}^{\ell_2}$. \square

C. Assumptions and Justifications

Because $\hat{f}_{\Delta \mathbf{V}, \mathbf{V}_0}(ax) = a \hat{f}_{\Delta \mathbf{V}, \mathbf{V}_0}(x)$ for any $a \in \mathbb{R}$, we can always do preprocessing to normalize the input \mathbf{x} . For simplicity, we focus on the simplest situation that the randomness for the inputs and the initial weights are uniform. Nonetheless, methods and results of this paper can be readily generalized to other continuous random variable distributions, which we leave for future work. We thus make the following Assumption 1.

Assumption 1. The input \mathbf{x} are uniformly distributed in S^{d-1} . The initial weights $\mathbf{V}_0[j]$'s are uniform in all directions. In other words, $\mu(\cdot)$ and $\lambda(\cdot)$ are both $\text{unif}(S^{d-1})$.

We study the overparameterized and overfitted setting, so in this paper we always assume $p \geq n/d$, i.e., the number of parameters pd is larger than or equal to the number of training samples n . The situation of $d = 1$ is relatively trivial, so we only consider the case $d \geq 2$. We then make Assumption 2.

Assumption 2. $p \geq n/d$ and $d \geq 2$.

If the input is a continuous random vector, then for any $i \neq j$, we have $\Pr f\mathbf{X}_i = \mathbf{X}_j g = 0$ and $\Pr f\mathbf{X}_i = -\mathbf{X}_j g = 0$ (because the probability that a continuous random variable equals to a given value is zero). Thus, $\Pr f\mathbf{X}_i \cdot \mathbf{X}_j g = 0$, and $\Pr f\mathbf{X}_i, \mathbf{X}_j g = 1$. Similarly, we can show that $\Pr f\mathbf{V}_0[k], \mathbf{V}_0[l] g = 1$. We thus make Assumption 3.

Assumption 3. $\mathbf{X}_i, \mathbf{X}_j$ for any $i \neq j$, and $\mathbf{V}_0[k], \mathbf{V}_0[l]$ for any $k \neq l$.

With these assumptions, the following lemma says that when p is large enough, with high probability \mathbf{H} has full row-rank (and thus \mathbf{V}^{ℓ_2} exists).

Lemma 7. $\lim_{p \rightarrow \infty} \Pr_{\mathbf{V}_0} \text{frank}(\mathbf{H}) = n \cdot \mathbf{X} g = 1$.

Proof. See Appendix E. □

D. Some Useful Supporting Results

Here we collect some useful lemmas that are needed for proofs in other appendices, many of which are estimations of certain quantities that we will use later.

D.1. Quantities related to the area of a cap on a hyper-sphere

The following lemma is introduced by (Li, 2011), which gives the area of a cap on a hyper-sphere with respect to the colatitude angle.

Lemma 8. Let $\phi \in [0, \frac{\pi}{2}]$ denote the colatitude angle of the smaller cap on S^{d-1} , then the area (in the measure of λ_{d-1}) of this hyper-spherical cap is

$$\frac{1}{2} \lambda_{d-1}(S^{d-1}) I_{\sin^2 \phi} \left(\frac{d-1}{2}, \frac{1}{2} \right).$$

The following lemma is another representation of the area of the cap with respect to the radius r (recall the definition of B^r in Eq. (21) and Remark 4).

Lemma 9. If $r \geq \frac{\sqrt{2}}{2}$, then we have

$$\lambda_{d-1}(B^r) = \frac{1}{2} \lambda_{d-1}(S^{d-1}) I_{r^2(1-\frac{r^2}{4})} \left(\frac{d-1}{2}, \frac{1}{2} \right).$$

Proof. Let ϕ denote the colatitude angle. By the law of cosines, we have

$$\cos \phi = 1 - \frac{r^2}{2}.$$

Thus, we have

$$\sin^2 \phi = 1 - \cos^2 \phi = 1 - \left(1 - \frac{r^2}{2} \right)^2 = r^2 \left(1 - \frac{r^2}{4} \right).$$

By Lemma 8, the result of this lemma thus follows. Notice that we require $r \geq \frac{\sqrt{2}}{2}$ to make sure that $\phi \in [0, \frac{\pi}{2}]$, which is required by Lemma 8. □

The area of a cap can be interpreted as the probability of the event that a uniformly-distributed random vector falls into that cap. We have the following lemma.

Lemma 10. *Suppose that a random vector $\mathbf{b} \in S^{d-1}$ follows uniform distribution in all directions. Given any $\mathbf{a} \in S^{d-1}$ and for any $c \in (0, 1)$, we have*

$$\Pr_{\mathbf{b}} \{j\mathbf{a}^T \mathbf{b}j > c\} = I_{1-c^2} \left(\frac{d-1}{2}, \frac{1}{2} \right).$$

Proof. Notice that $\{\mathbf{b} \mid \mathbf{a}^T \mathbf{b} > c\}$ is a hyper-spherical cap. Define its colatitude angle as ϕ . We have $\cos \phi = \mathbf{a}^T \mathbf{b} = c$. Thus, we have $\sin^2 \phi = 1 - c^2$. By Lemma 8, we then have

$$\lambda_{d-1}(\{\mathbf{b} \mid \mathbf{a}^T \mathbf{b} > c\}) = \frac{1}{2} \lambda_{d-1}(S^{d-1}) I_{1-c^2} \left(\frac{d-1}{2}, \frac{1}{2} \right).$$

Further, by symmetry, we have

$$\lambda_{d-1}(\{\mathbf{b} \mid j\mathbf{a}^T \mathbf{b}j > c\}) = 2\lambda_{d-1}(\{\mathbf{b} \mid \mathbf{a}^T \mathbf{b} > c\}) = \lambda_{d-1}(S^{d-1}) I_{1-c^2} \left(\frac{d-1}{2}, \frac{1}{2} \right).$$

Because \mathbf{b} follows uniform distribution in all directions, we have

$$\Pr_{\mathbf{b}} \{j\mathbf{a}^T \mathbf{b}j > c\} = \frac{\lambda_{d-1}(\{\mathbf{b} \mid j\mathbf{a}^T \mathbf{b}j > c\})}{\lambda_{d-1}(S^{d-1})} = I_{1-c^2} \left(\frac{d-1}{2}, \frac{1}{2} \right).$$

□

D.2. Estimation of certain norms

In this subsection, we will show $k\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}k_2 \leq \sqrt{p}$ in Lemma 11. We also upper bound the norm of the product of two matrices by the product of their norms in Lemma 12. At last, Lemma 13 states that if two vector differ a lot, then the sum of their norm cannot be too small.

Lemma 11. *$k\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}k_2 \leq \sqrt{p}$ for any $\mathbf{x} \in S^{d-1}$.*

Proof. This follows because the input \mathbf{x} is normalized. Specifically, by Eq. (1), we have

$$k\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}k_2 = \sqrt{\sum_{j=1}^p \|\mathbf{1}_{f_{\mathbf{x}^T \mathbf{V}_0[j]} > 0} \mathbf{x}^T\|_2^2} \leq \sqrt{p}. \quad (25)$$

□

Lemma 12. *If $\mathbf{C} = \mathbf{A}\mathbf{B}$, then $k\mathbf{C}k_2 \leq k\mathbf{A}k_2 k\mathbf{B}k_2$. Here \mathbf{A} , \mathbf{B} , and \mathbf{C} could be scalars, vectors, or matrices.*

Proof. This lemma directly follows the definition of matrix norm. □

Remark 5. Note that the (ℓ_2) matrix-norm (i.e., spectral norm) of a vector is exactly its ℓ_2 vector-norm (i.e., Euclidean norm)⁷. Therefore, when applying Lemma 12, we do not need to worry about whether \mathbf{A} , \mathbf{B} , and \mathbf{C} are matrices or vectors.

Lemma 13. *For any $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^d$, we have*

$$k\mathbf{v}_1k_2^2 + k\mathbf{v}_2k_2^2 \geq \frac{1}{2} k\mathbf{v}_1 - \mathbf{v}_2k_2^2.$$

⁷To see this, consider a (row or column) vector \mathbf{a} . The matrix norm of \mathbf{a} is

$$\max_{j|x|=1} k\mathbf{a}\mathbf{x}k_2 \text{ (when } \mathbf{a} \text{ is a column vector),}$$

$$\text{or } \max_{k\mathbf{x}k_2=1} k\mathbf{a}\mathbf{x}k_2 \text{ (when } \mathbf{a} \text{ is a row vector).}$$

In both cases, the value of the matrix-norm equals to $\sqrt{\sum a_i^2}$, which is exactly the ℓ_2 -norm (Euclidean norm) of \mathbf{a} .

Proof. It is easy to prove that $k \|\cdot\|_2^2$ is convex. Thus, we have

$$\begin{aligned} k\mathbf{v}_1 k_2^2 + k\mathbf{v}_2 k_2^2 &= k\mathbf{v}_1 k_2^2 + k \|\mathbf{v}_2\|_2^2 \\ &= 2 \left\| \frac{\mathbf{v}_1 + \mathbf{v}_2}{2} \right\|_2^2 \quad (\text{apply Jensen's inequality on the convex function } k \|\cdot\|_2^2) \\ &= \frac{1}{2} k\mathbf{v}_1 + \mathbf{v}_2 k_2^2. \end{aligned}$$

□

D.3. Estimates of certain tail probabilities

The following is the (restated) Corollary 5 of (Goemans, 2015).

Lemma 14. *If the random variable X follows $\text{Bino}(a, b)$, then for all $0 < \delta < 1$, we have*

$$\Pr\{X - ab > \delta ab\} < 2e^{-\delta^2 ab/3}.$$

The following lemma is the (restated) Theorem 1.8 of (Hayes, 2005).

Lemma 15 (Azuma–Hoeffding inequality for random vectors). *Let X_1, X_2, \dots, X_k be i.i.d. random vectors with zero mean (of the same dimension) in a real Euclidean space such that $\|X_i\|_2 \leq 1$ for all $i = 1, 2, \dots, k$. Then, for every $a > 0$,*

$$\Pr\left\{\left\|\sum_{i=1}^k X_i\right\|_2 \geq a\right\} < 2e^2 \exp\left(-\frac{a^2}{2k}\right).$$

In the following lemma, we use Azuma–Hoeffding inequality to upper bound the deviation of the empirical mean value of a bounded random vector from its expectation.

Lemma 16. *Let X_1, X_2, \dots, X_k be i.i.d. random vectors (of the same dimension) in a real Euclidean space such that $\|X_i\|_2 \leq U$ for all $i = 1, 2, \dots, k$. Then, for any $q \geq 1$,*

$$\Pr\left\{\left\|\frac{1}{k} \sum_{i=1}^k X_i - \mathbb{E} X_1\right\|_2 \geq k^{-\frac{1}{2q}}\right\} < 2e^2 \exp\left(-\frac{c_q}{8U^2}\right).$$

Proof. Because $\|X_i\|_2 \leq U$, we have $\|X_i - \mathbb{E} X_i\|_2 \leq U$. By triangle inequality, we have $\|X_i - \mathbb{E} X_i\|_2 \leq \|X_i\|_2 + \|\mathbb{E} X_i\|_2 \leq 2U$, i.e.,

$$\left\|\frac{X_i - \mathbb{E} X_i}{2U}\right\|_2 \leq 1. \quad (26)$$

We also have

$$\mathbb{E}\left[\frac{X_i - \mathbb{E} X_i}{2U}\right] = \frac{\mathbb{E} X_i - \mathbb{E} X_i}{2U} = \mathbf{0}. \quad (27)$$

We then have

$$\begin{aligned} &\Pr\left\{\left\|\frac{1}{k} \sum_{i=1}^k X_i - \mathbb{E} X_1\right\|_2 \geq k^{-\frac{1}{2q}}\right\} \\ &= \Pr\left\{\left\|\sum_{i=1}^k (X_i - \mathbb{E} X_i)\right\|_2 \geq k^{\frac{1}{2q} + \frac{1}{2}}\right\} \\ &= \Pr\left\{\left\|\sum_{i=1}^k \left(\frac{X_i - \mathbb{E} X_i}{2U}\right)\right\|_2 \geq \frac{k^{\frac{1}{2q} + \frac{1}{2}}}{2U}\right\} \\ &< 2e^2 \exp\left(-\frac{c_q}{8U^2}\right) \quad (\text{by Eqs. (26)(27) and letting } a = \frac{k^{\frac{1}{2q} + \frac{1}{2}}}{2U} \text{ in Lemma 15}). \end{aligned}$$

□

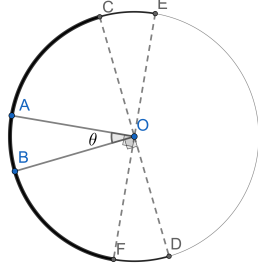


Figure 4. The arc \widehat{CBF} is $\frac{\pi - \theta}{2\pi}$ of the perimeter of the circle O .

D.4. Calculation of certain integrals

The following lemma calculates the ratio between the intersection area of two hyper-hemispheres and the area of the whole hyper-sphere.

Lemma 17.

$$\int_{S^{d-1}} \mathbf{1}_{\widehat{f}v \cdot z^T v > 0, \widehat{x}^T v > 0} d\lambda(v) = \frac{\pi \arccos(\widehat{x}^T z)}{2\pi}. \quad (28)$$

(Recall that $\lambda(\cdot)$ denotes the distribution of the normalized version of $\mathbf{V}_0[j]$ on S^{d-1} and is assumed to be uniform in all directions.)

Before we give the proof of Lemma 17, we give its geometric explanation.

Geometric explanation of Eq. (28): Indeed, since λ is uniform on S^{d-1} , the integral on the left-hand-side of Eq. (28) represents the probability that a random point falls into the intersection of two hyper-hemispheres that are represented by $\widehat{f}v \cdot z^T v > 0$ and $\widehat{x}^T v > 0$, respectively. We can calculate that probability by

$$\frac{\text{measure of a hyper-spherical lune with angle } \pi - \theta(z, \widehat{x})}{\text{measure of a unit hyper-sphere}} = \frac{\pi \arccos(\widehat{x}^T z)}{2\pi}, \quad (29)$$

where $\theta(\cdot, \cdot)$ denote the angle (in radians) between two vectors, which would lead to Eq. (28). To help readers understand Eq. (29), we give examples for 2D and 3D in Fig. 4 and Fig. 5, respectively. In the 2D case depicted in Fig. 4, \widehat{OA} denotes z , \widehat{OB} denotes \widehat{x} . Thus, the arc \widehat{EAF} denotes $\widehat{f}v \cdot z^T v > 0$, and the arc \widehat{CBD} denotes $\widehat{f}v \cdot \widehat{x}^T v > 0$. The intersection of \widehat{EAF} and \widehat{CBD} , i.e., the arc \widehat{CBF} , represents $\widehat{f}v \cdot z^T v > 0, \widehat{x}^T v > 0$. Notice that the angle of \widehat{CBF} equals $\pi - \theta$, where θ denotes the angle between z and \widehat{x} . Therefore, ratio of the length of \widehat{CBF} to the perimeter of the circle equals to $\frac{\angle \widehat{COF}}{2\pi} = \frac{\pi - \theta}{2\pi}$. Similarly, in the 3D case depicted in Fig. 5, the spherical lune \widehat{ICHF} denotes the intersection of the semi-sphere in the direction of \widehat{OA} and the semi-sphere in the direction of \widehat{OB} . We can see that the area of the spherical lune \widehat{ICHF} is still proportional to the angle $\angle \widehat{COF}$. Thus, we still have the result that the area of the spherical lune \widehat{ICHF} is $\frac{\pi - \theta}{2\pi}$ of the area of the whole sphere. The proof below, on the other hand, applies to arbitrary dimensions.

Proof. Due to symmetry, we know that the integral of Eq. (28) only depends on the angle between \widehat{x} and z . Thus, without loss of generality, we let

$$\widehat{x} = [\widehat{x}_1 \ \widehat{x}_2 \ \dots \ \widehat{x}_d] = [0 \ 0 \ \dots \ 0 \ 1 \ 0]^T, \quad z = [0 \ 0 \ \dots \ 0 \ \cos \theta \ \sin \theta]^T,$$

where

$$\theta = \arccos(\widehat{x}^T z) \in [0, \pi]. \quad (30)$$

Thus, for any $v = [v_1 \ v_2 \ \dots \ v_d]^T$ that makes $z^T v > 0$ and $\widehat{x}^T v > 0$, it only needs to satisfy

$$[\cos \theta \ \sin \theta] \begin{bmatrix} v_d \\ v_{d-1} \end{bmatrix} > 0, \quad [1 \ 0] \begin{bmatrix} v_d \\ v_{d-1} \end{bmatrix} > 0. \quad (31)$$

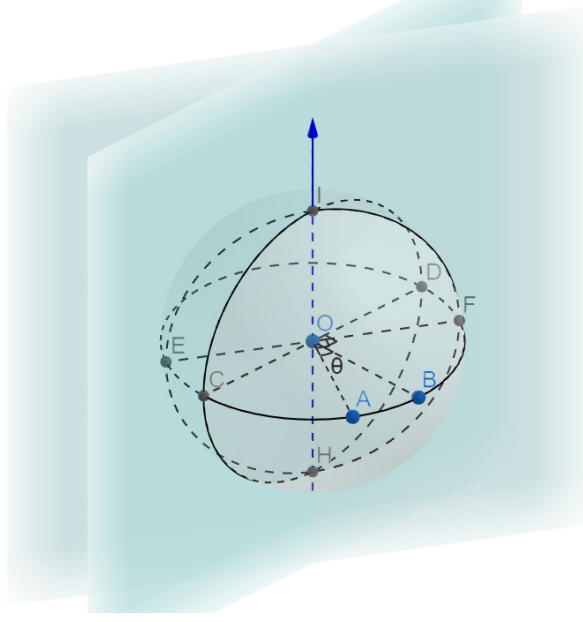


Figure 5. The area of the spherical lune ICHF is $\frac{\pi \cdot \theta}{2\pi}$ of the area of the whole sphere.

We compute the spherical coordinates $\varphi_{\mathbf{x}} = [\varphi_1^{\mathbf{x}} \ \varphi_2^{\mathbf{x}} \ \dots \ \varphi_{d-1}^{\mathbf{x}}]^T$ where $\varphi_1^{\mathbf{x}}, \dots, \varphi_{d-2}^{\mathbf{x}} \in [0, \pi]$ and $\varphi_{d-1}^{\mathbf{x}} \in [0, 2\pi)$ with the convention that

$$\begin{aligned} x_1 &= \cos(\varphi_1^{\mathbf{x}}), \\ x_2 &= \sin(\varphi_1^{\mathbf{x}}) \cos(\varphi_2^{\mathbf{x}}), \\ x_3 &= \sin(\varphi_1^{\mathbf{x}}) \sin(\varphi_2^{\mathbf{x}}) \cos(\varphi_3^{\mathbf{x}}), \\ &\vdots \\ x_{d-1} &= \sin(\varphi_1^{\mathbf{x}}) \sin(\varphi_2^{\mathbf{x}}) \dots \sin(\varphi_{d-2}^{\mathbf{x}}) \cos(\varphi_{d-1}^{\mathbf{x}}), \\ x_d &= \sin(\varphi_1^{\mathbf{x}}) \sin(\varphi_2^{\mathbf{x}}) \dots \sin(\varphi_{d-2}^{\mathbf{x}}) \sin(\varphi_{d-1}^{\mathbf{x}}). \end{aligned}$$

Thus, we have $\varphi_{\mathbf{x}} = [\pi/2 \ \pi/2 \ \dots \ \pi/2 \ 0]^T$. Similarly, the spherical coordinates for \mathbf{z} is $\varphi_{\mathbf{z}} = [\pi/2 \ \pi/2 \ \dots \ \pi/2 \ \theta]^T$. Let the spherical coordinates for \mathbf{v} be $\varphi_{\mathbf{v}} = [\varphi_1^{\mathbf{v}} \ \varphi_2^{\mathbf{v}} \ \dots \ \varphi_{d-1}^{\mathbf{v}}]^T$. Thus, Eq. (31) is equivalent to

$$\sin(\varphi_1^{\mathbf{v}}) \sin(\varphi_2^{\mathbf{v}}) \dots \sin(\varphi_{d-2}^{\mathbf{v}}) (\cos \theta \cos(\varphi_{d-1}^{\mathbf{v}}) + \sin \theta \sin(\varphi_{d-1}^{\mathbf{v}})) > 0, \quad (32)$$

$$\sin(\varphi_1^{\mathbf{v}}) \sin(\varphi_2^{\mathbf{v}}) \dots \sin(\varphi_{d-2}^{\mathbf{v}}) \cos(\varphi_{d-1}^{\mathbf{v}}) > 0. \quad (33)$$

Because $\varphi_1^{\mathbf{v}}, \dots, \varphi_{d-2}^{\mathbf{v}} \in [0, \pi]$ (by the convention of spherical coordinates), we have

$$\sin(\varphi_1^{\mathbf{v}}) \sin(\varphi_2^{\mathbf{v}}) \dots \sin(\varphi_{d-2}^{\mathbf{v}}) \geq 0.$$

Thus, for Eq. (32) and Eq. (33), we have

$$\cos(\theta - \varphi_{d-1}^{\mathbf{v}}) > 0, \quad \cos(\varphi_{d-1}^{\mathbf{v}}) > 0,$$

i.e., $\varphi_{d-1}^v \in (\pi/2, \pi/2) \setminus (\theta - \pi/2, \theta + \pi/2) \pmod{2\pi}$. We have

$$\begin{aligned}
 & \int_{S^{d-1}} \mathbf{1}_{fz^T v > 0, x^T v > 0} d\lambda(v) \\
 &= \frac{\int_{(\pi/2, \pi/2) \setminus (\theta - \pi/2, \theta + \pi/2)} \int_0^\pi \int_0^\pi \sin^{d-2}(\varphi_1) \sin^{d-3}(\varphi_2) \sin(\varphi_{d-2}) d\varphi_1 d\varphi_2 d\varphi_{d-1}}{\int_0^{2\pi} \int_0^\pi \int_0^\pi \sin^{d-2}(\varphi_1) \sin^{d-3}(\varphi_2) \sin(\varphi_{d-2}) d\varphi_1 d\varphi_2 d\varphi_{d-1}} \\
 &= \frac{\int_{(\pi/2, \pi/2) \setminus (\theta - \pi/2, \theta + \pi/2)} A d\varphi_{d-1}}{\int_0^{2\pi} A d\varphi_{d-1}} \\
 & \quad (\text{by defining } A := \int_0^\pi \int_0^\pi \sin^{d-2}(\varphi_1) \sin^{d-3}(\varphi_2) \sin(\varphi_{d-2}) d\varphi_1 d\varphi_2) \\
 &= \frac{\text{length of the interval } (\pi/2, \pi/2) \setminus (\theta - \pi/2, \theta + \pi/2)}{2\pi} \\
 &= \frac{\pi - \theta}{2\pi} \text{ (because } \theta \in [0, \pi] \text{ by Eq. (30))} \\
 &= \frac{\pi - \arccos(x^T z)}{2\pi} \text{ (by Eq. (30)).}
 \end{aligned}$$

The result of this lemma thus follows. □

D.5. Limits of $jC_{z,x}^{V_0}/p$ when $p \rightarrow \infty$

We introduce some notions given by (Wainwright, 2015).

Glivenko-Cantelli class. Let F be a class of integrable real-valued functions with domain X , and let $X_1^k = fX_1, \dots, X_k \in \mathcal{G}$ be a collection of *i.i.d.* samples from some distribution P over X . Consider the random variable

$$kP_k - P_{k,\mathcal{F}} := \sup_{f \in \mathcal{F}} \left| \frac{1}{k} \sum_{i=1}^k f(X_i) - E[f] \right|,$$

which measures the maximum deviation (over the class F) between the sample average $\frac{1}{k} \sum_{i=1}^k f(X_i)$ and the population average $E[f] = E[f(X)]$. We say that F is a *Glivenko-Cantelli class* for P if $kP_k - P_{k,\mathcal{F}}$ converges to zero in probability as $k \rightarrow \infty$.

Polynomial discrimination. A class F of functions with domain X has polynomial discrimination of order $\nu \geq 1$ if for each positive integer k and collection $X_1^k = fX_1, \dots, X_k \in \mathcal{G}$ of k points in X , the set $F(X_1^k)$ has cardinality upper bounded by

$$\text{card}(F(X_1^k)) \leq (k+1)^\nu.$$

The following lemma is shown in Page 108 of (Wainwright, 2015).

Lemma 18. *Any bounded function class with polynomial discrimination is Glivenko-Cantelli.*

For our case, we care about the following value.

$$\left| \frac{jC_{z,x}^{V_0}}{p} - \frac{\pi - \arccos(x^T z)}{2\pi} \right| = \left| \frac{1}{p} \sum_{j=1}^p \mathbf{1}_{f_{x^T v_{0[j]} > 0, z^T v_{0[j]} > 0}} - E_{v \sim \lambda(\cdot)} [\mathbf{1}_{f_{x^T v > 0, z^T v > 0}}] \right| \text{ (by Lemma 17).}$$

In the language of Glivenko-Cantelli class, the function class F consists of functions $\mathbf{1}_{f_{x^T v > 0, z^T v > 0}}$ that map $v \in S^{d-1}$ to 0 or 1, where every $x \in S^{d-1}$ and $z \in S^{d-1}$ corresponds to a distinct function in F . According to Lemma 18, we need to calculate the order of the polynomial discrimination for this F . Towards this end, we need the following lemma, which can be derived from the quantity $Q_{n,N}$ in (Wendel, 1962) (which is the quantity $Q_{d,k}$ in the following lemma).

Lemma 19. Given $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k \in S^{d-1}$, the number of different values (i.e., the cardinality) of the set $\{(\mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_1} > 0g}, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_2} > 0g}, \dots, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_k} > 0g}) \mid \mathbf{x} \in S^{d-1}\}$ is at most $Q_{d,k}$, where

$$Q_{d,k} := \begin{cases} 2 \sum_{i=0}^{d-1} \binom{k-1}{i}, & \text{if } k > d, \\ 2^k, & \text{if } k \leq d. \end{cases}$$

Intuitively, Lemma 19 states the number of different regions that k hyper-planes through the origin (i.e., the kernel of the inner product with each \mathbf{v}_i) can cut S^{d-1} into, because all \mathbf{x} in one region corresponds to the same value of the tuple $(\mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_1} > 0g}, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_2} > 0g}, \dots, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_k} > 0g})$. For example, in the 2D case (i.e., $d = 2$), k diameters of a circle can at most cut the whole circle into $2k$ (which equals to $Q_{2,k}$) parts. Notice that if some \mathbf{v}_i 's are parallel (thus some diameters are overlapped), then the total number of different parts can only be smaller. That is why Lemma 19 states that the cardinality is ‘‘at most’’ $Q_{d,k}$.

The following lemma shows that the cardinality in Lemma 19 is polynomial in k .

Lemma 20. Recall the definition $Q_{d,k}$ in Lemma 19. For any integer $k \geq 1$ and $d \geq 2$, we must have $Q_{d,k} \leq (k+1)^{d+1}$.

Proof. When $k > d$, because $\binom{k-1}{i} \leq (k-1)^{d-1}$ when $i \leq d-1$, we have $Q_{d,k} = 2 \sum_{i=0}^{d-1} \binom{k-1}{i} \leq 2d(k+1)^{d-1} (k+1)^{d+1}$ (the last step uses $k \geq 1$ and $k > d$). When $k \leq d$, because $k \geq 1$, we have $Q_{d,k} = 2^k \leq (k+1)^k \leq (k+1)^d$. In summary, for any integer $k \geq 1$ and $d \geq 2$, the result $Q_{d,k} \leq (k+1)^{d+1}$ always holds. \square

We can now calculate the order of the polynomial discrimination for the function class F . Because

$$\begin{aligned} & \text{card} \left(\left\{ (\mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_1} > 0, \mathbf{z}^T \mathbf{v}_1 > 0g}, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_2} > 0, \mathbf{z}^T \mathbf{v}_2 > 0g}, \dots, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_k} > 0, \mathbf{z}^T \mathbf{v}_k > 0g}) \mid \mathbf{x} \in S^{d-1}, \mathbf{z} \in S^{d-1} \right\} \right) \\ & \text{card} \left(\left\{ (\mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_1} > 0g}, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_2} > 0g}, \dots, \mathbf{1}_{\mathbf{f}_{\mathbf{x}^T \mathbf{v}_k} > 0g}) \mid \mathbf{x} \in S^{d-1} \right\} \right) \\ & \text{card} \left(\left\{ (\mathbf{1}_{\mathbf{f}_{\mathbf{z}^T \mathbf{v}_1} > 0g}, \mathbf{1}_{\mathbf{f}_{\mathbf{z}^T \mathbf{v}_2} > 0g}, \dots, \mathbf{1}_{\mathbf{f}_{\mathbf{z}^T \mathbf{v}_k} > 0g}) \mid \mathbf{z} \in S^{d-1} \right\} \right), \end{aligned}$$

by Lemma 19 and Lemma 20, we know that

$$\text{card}(F(X_1^k)) \leq (Q_{d,k})^2 \leq (k+1)^{2(d+1)}.$$

(Here X_1^k means $f\mathbf{V}_0[1], \dots, \mathbf{V}_0[k]g$.)

Thus, F has polynomial discrimination with order at most $2(d+1)$. Notice that all functions in F is bounded because their outputs can only be 0 or 1. Therefore, by Lemma 18 (i.e., any bounded function class with polynomial discrimination is Glivenko-Cantelli), we know that F is Glivenko-Cantelli. In other words, we have shown the following lemma.

Lemma 21.

$$\sup_{\mathbf{x}, \mathbf{z} \in S^{d-1}} \left| \frac{jC_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0} j}{p} - \frac{\pi - \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} \right| \xrightarrow{p \rightarrow \infty} 0, \text{ as } p \rightarrow \infty. \quad (34)$$

E. Proof of Lemma 7 (H has full row-rank with high probability as $p \rightarrow \infty$)

In this section, we prove Lemma 7, i.e., the matrix \mathbf{H} has full row-rank with high probability when $p \rightarrow \infty$. We first introduce two useful lemmas as follows.

The following lemma states that, given \mathbf{X} (that satisfies Assumption 3) and $k \geq 1, 2, \dots, ng$, there always exists a vector $\mathbf{v} \in S^{d-1}$ that is only orthogonal to one training input \mathbf{X}_k but not orthogonal to other training inputs \mathbf{X}_i for all $i \neq k$. An intuitive explanation is that, because no training inputs are parallel (as stated in Assumption 3), the total set of vectors that are orthogonal to at least two training inputs is too small. That gives us many options to pick such a vector \mathbf{v} that is only orthogonal to one input but not others.

Lemma 22. For all $k \geq 1, 2, \dots, ng$ we have

$$T_k := \{ \mathbf{v} \in S^{d-1} \mid \mathbf{v}^T \mathbf{X}_k = 0, \mathbf{v}^T \mathbf{X}_i \neq 0, \text{ for all } i \geq 1, 2, \dots, ng, i \neq k \} \neq \emptyset.$$

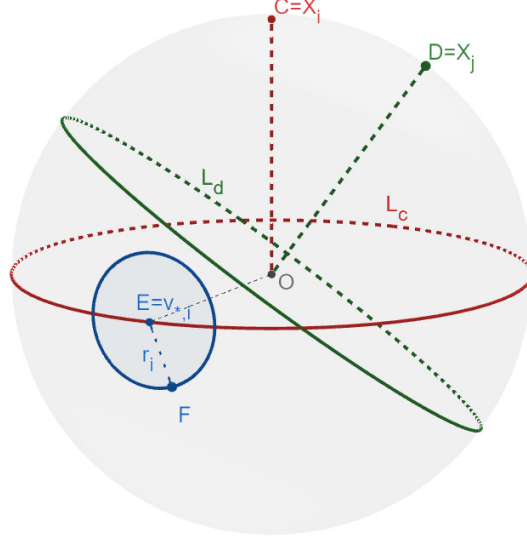


Figure 6. Geometric interpretation of $B_v^{r_i, \mathbf{X}_i}$ and $B_v^{r_i, \mathbf{X}_i}$ on a sphere (i.e., S^2).

Proof. We have

$$\begin{aligned} T_k &= S^{d-1} \setminus \ker(\mathbf{X}_k) \cap \left(\bigcup_{i \in \{1, 2, \dots, n\} \setminus \{k\}} \ker(\mathbf{X}_i) \right) \\ &= S^{d-1} \setminus \ker(\mathbf{X}_k) \cap \left(\bigcup_{i \in \{1, 2, \dots, n\} \setminus \{k\}} (S^{d-1} \setminus \ker(\mathbf{X}_k) \setminus \ker(\mathbf{X}_i)) \right). \end{aligned}$$

Because

$$\begin{aligned} \dim(S^{d-1} \setminus \ker(\mathbf{X}_k)) &= d-2, \\ \dim(S^{d-1} \setminus \ker(\mathbf{X}_k) \setminus \ker(\mathbf{X}_i)) &= d-3 \text{ for all } i \in \{1, 2, \dots, n\} \setminus \{k\} \text{ (because } \mathbf{X}_i \neq \mathbf{X}_k), \end{aligned} \quad (35)$$

we have

$$\begin{aligned} \lambda_{d-2}(S^{d-1} \setminus \ker(\mathbf{X}_k)) &= \lambda_{d-2}(S^{d-2}) > 0, \\ \lambda_{d-2}(S^{d-1} \setminus \ker(\mathbf{X}_k) \setminus \ker(\mathbf{X}_i)) &= 0 \text{ for all } i \in \{1, 2, \dots, n\} \setminus \{k\}. \end{aligned} \quad (36)$$

(When $d = 2$, the set in Eq. (35) is not defined. Nonetheless, Eq. (36) still holds when $d = 2$.) Thus, we have

$$\begin{aligned} \lambda_{d-2}(T_k) &= \lambda_{d-2}(S^{d-1} \setminus \ker(\mathbf{X}_k)) - \lambda_{d-2} \left(\bigcup_{i \in \{1, 2, \dots, n\} \setminus \{k\}} (S^{d-1} \setminus \ker(\mathbf{X}_k) \setminus \ker(\mathbf{X}_i)) \right) \\ &= \lambda_{d-2}(S^{d-1} \setminus \ker(\mathbf{X}_k)) - \sum_{i \in \{1, 2, \dots, n\} \setminus \{k\}} \lambda_{d-2}(S^{d-1} \setminus \ker(\mathbf{X}_k) \setminus \ker(\mathbf{X}_i)) \\ &= \lambda_{d-2}(S^{d-2}) \\ &> 0. \end{aligned}$$

Therefore, $T_k \neq \emptyset$. □

The following lemma plays an important role in answering whether \mathbf{H} has full row-rank. Further, it is also closely related to our estimation on the $\min \text{eig}(\mathbf{H}\mathbf{H}^T)$ later in Appendix F.

Lemma 23. Consider any $i \in \{1, 2, \dots, n\}$. For any $\mathbf{v}_i \in S^{d-1}$ satisfying $\mathbf{v}_i^T \mathbf{X}_i = 0$, we define

$$r_i := \min_{j \in \{1, 2, \dots, n\}, j \neq i} |\mathbf{v}_i^T \mathbf{X}_j|. \quad (37)$$

If there exist $k, l \in \{1, 2, \dots, p\}$ such that

$$\frac{\mathbf{V}_0[k]}{k\|\mathbf{V}_0[k]\|_{k_2}} \geq B_{\mathbf{v}_i, +}^{r_i, \mathbf{X}_i}, \quad \frac{\mathbf{V}_0[l]}{k\|\mathbf{V}_0[l]\|_{k_2}} \geq B_{\mathbf{v}_i, -}^{r_i, \mathbf{X}_i}, \quad (38)$$

then we must have

$$\mathbf{H}_j[k] = \mathbf{H}_j[l], \text{ for all } j \in \{1, 2, \dots, n\}, j \neq i, \quad (39)$$

$$\mathbf{H}_i[k] = \mathbf{X}_i^T, \quad (40)$$

$$\mathbf{H}_i[l] = \mathbf{0}. \quad (41)$$

(Notice that Eq. (38) implies $r_i > 0$.)

Remark 6. We first give an intuitive geometric interpretation of Lemma 23. In Fig. 6, the sphere centered at O denotes S^{d-1} , the vector OC denotes \mathbf{X}_i , the vector OD denotes one of other \mathbf{X}_j 's, the vector OE denotes \mathbf{v}_i , which is perpendicular to \mathbf{X}_i (i.e., $\mathbf{X}_i^T \mathbf{v}_i = 0$). The upper half of the cap E denotes $B_{\mathbf{v}_i, +}^{r_i, \mathbf{X}_i}$, the lower half of the cap E denotes $B_{\mathbf{v}_i, -}^{r_i, \mathbf{X}_i}$. The great circle L_c cuts the sphere into two semi-spheres. The semi-sphere in the direction of OC corresponds to all vectors \mathbf{v} on the sphere that have positive inner product with \mathbf{X}_i (i.e., $\mathbf{v}^T \mathbf{X}_i > 0$), and the semi-sphere in the opposite direction of OC corresponds to all vectors \mathbf{v} on the sphere that have negative inner product with \mathbf{X}_i (i.e., $\mathbf{v}^T \mathbf{X}_i < 0$). The great circle L_d is similar to the great circle L_c , but is perpendicular to the direction OD (i.e., \mathbf{X}_j). By choosing the radius of the cap E in Eq. (37), we can ensure that all great circles that are perpendicular to other \mathbf{X}_j 's do not pass the cap E. In other words, for the two semi-spheres cut by the great circle perpendicular to \mathbf{X}_j , $j \neq i$, the cap E must be contained in one of them. Therefore, vectors on the upper half of the cap E and the vectors on the lower half of the cap E must have the same sign when calculating the inner product with all \mathbf{X}_j 's, for all $j \neq i$.

Now, let us consider the meaning of Eq. (38) in this geometric setup depicted in Fig. 6. The expression $\frac{\mathbf{V}_0[k]}{k\|\mathbf{V}_0[k]\|_{k_2}} \geq B_{\mathbf{v}_i, +}^{r_i, \mathbf{X}_i}$ means that the direction of $\mathbf{V}_0[k]$ is in the upper half of the cap E. By the definition of $\mathbf{H}_i = \mathbf{h}_{\mathbf{V}_0, \mathbf{X}_i}$ in Eq. (1), we must then have $\mathbf{H}_i[k] = \mathbf{X}_i^T$. Similarly, the expression $\frac{\mathbf{V}_0[l]}{k\|\mathbf{V}_0[l]\|_{k_2}} \geq B_{\mathbf{v}_i, -}^{r_i, \mathbf{X}_i}$ means that the direction of $\mathbf{V}_0[l]$ is in the lower half of the cap E, and thus $\mathbf{H}_i[l] = \mathbf{0}$. Then, based on the discussions in the previous paragraph, we know that $\mathbf{V}_0[k]$ and $\mathbf{V}_0[l]$ has the same activation pattern under ReLU for all \mathbf{X}_j 's that $j \neq i$, which implies that $\mathbf{H}_j[k] = \mathbf{H}_j[l]$. These are precisely the conclusions in Eqs. (39)(40)(41).

Later in Appendix F, Lemma 23 plays an important role in estimating $\min_{\mathbf{a} \in S^{n-1}} k\mathbf{H}^T \mathbf{a} k_2^2$. To see this, let a_j denotes the j -th element of \mathbf{a} . By Eq. (39), we have $\sum_{j \in \{1, 2, \dots, n\}, j \neq i} ((\mathbf{H}^T \mathbf{a}_j)[k] - (\mathbf{H}^T \mathbf{a}_j)[l]) = \mathbf{0}$. By Eq. (40) and Eq. (41), we have $(\mathbf{H}^T \mathbf{a}_i)[k] - (\mathbf{H}^T \mathbf{a}_i)[l] = \mathbf{X}_i$. Combining them together, we have $(\mathbf{H}^T \mathbf{a})[k] - (\mathbf{H}^T \mathbf{a})[l] = a_i \mathbf{X}_i$. As long as a_i is not zero, then regardless values of other elements in \mathbf{a} , we always obtain that $\mathbf{H}^T \mathbf{a}$ is a non-zero vector. This implies $k\mathbf{H}^T \mathbf{a} k_2 > 0$, which will be useful for estimating $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$ in Appendix F.

Proof. By the definition of r_i , we have

$$|\mathbf{v}_i^T \mathbf{X}_j| \geq r_i > 0, \text{ for all } j \in \{1, 2, \dots, n\}, j \neq i. \quad (42)$$

For any $j \in \{1, 2, \dots, n\}, j \neq i$ and any $\mathbf{v} \in B_{\mathbf{v}_i, +}^{r_i, \mathbf{X}_i}$, since $k\|\mathbf{v} - \mathbf{v}_i\|_{k_2} < r_i$, we have

$$\begin{aligned} (\mathbf{v}^T \mathbf{X}_j)(\mathbf{v}_i^T \mathbf{X}_j) &= ((\mathbf{v} - \mathbf{v}_i)^T \mathbf{X}_j + \mathbf{v}_i^T \mathbf{X}_j)(\mathbf{v}_i^T \mathbf{X}_j) \\ &= (\mathbf{v}^T \mathbf{X}_j)^2 + (\mathbf{v}_i^T \mathbf{X}_j)((\mathbf{v} - \mathbf{v}_i)^T \mathbf{X}_j) \\ &= (\mathbf{v}^T \mathbf{X}_j)^2 - |\mathbf{v}_i^T \mathbf{X}_j| |(\mathbf{v} - \mathbf{v}_i)^T \mathbf{X}_j| \\ &= (\mathbf{v}^T \mathbf{X}_j)^2 - |\mathbf{v}_i^T \mathbf{X}_j| k\|\mathbf{v} - \mathbf{v}_i\|_{k_2} k\|\mathbf{X}_j\|_{k_2} \\ &> (\mathbf{v}^T \mathbf{X}_j)^2 - |\mathbf{v}_i^T \mathbf{X}_j| r_i \text{ (by Eq. (21))} \\ &= |\mathbf{v}_i^T \mathbf{X}_j| (|\mathbf{v}^T \mathbf{X}_j| - r_i) \\ &> 0 \text{ (by Eq. (42)).} \end{aligned}$$

Thus, for any $\mathbf{v}_1 \in B_{\mathbf{v},i}^{r_i}$, $\mathbf{v}_2 \in B_{\mathbf{v},i}^{r_i}$, $j \in \{1, 2\}$, $n, n \geq 1$, we have $(\mathbf{v}_1^T \mathbf{X}_j)(\mathbf{v}_2^T \mathbf{X}_j) > 0$ and $(\mathbf{v}_2^T \mathbf{X}_j)(\mathbf{v}_1^T \mathbf{X}_j) > 0$. It implies that

$$\text{sign}(\mathbf{v}_1^T \mathbf{X}_j) = \text{sign}(\mathbf{v}_2^T \mathbf{X}_j) = \text{sign}(\mathbf{v}_1^T \mathbf{X}_j). \quad (43)$$

By Eq. (38), we know that both $\mathbf{V}_0[k]$ and $\mathbf{V}_0[l]$ are in $B_{\mathbf{v},i}^{r_i}$. Applying Eq. (43), we have

$$\text{sign}(\mathbf{X}_j^T \mathbf{V}_0[k]) = \text{sign}(\mathbf{X}_j^T \mathbf{V}_0[l]), \text{ for all } j \in \{1, 2\}, n, n \geq 1.$$

Thus, by Eq. (1), we have

$$\mathbf{H}_j[k] = \mathbf{1}_{\{\mathbf{X}_j^T \mathbf{V}_0[k] > 0\}} \mathbf{X}_j^T = \mathbf{1}_{\{\mathbf{X}_j^T \mathbf{V}_0[l] > 0\}} \mathbf{X}_j^T = \mathbf{H}_j[l], \text{ for all } j \in \{1, 2\}, n, n \geq 1.$$

By Eq. (22), we have

$$\mathbf{X}_i^T \mathbf{V}_0[k] > 0, \mathbf{X}_i^T \mathbf{V}_0[l] < 0.$$

Thus, by Eq. (1), we have

$$\mathbf{H}_i[k] = \mathbf{1}_{\{\mathbf{X}_i^T \mathbf{V}_0[k] > 0\}} \mathbf{X}_i^T = \mathbf{X}_i^T, \quad \mathbf{H}_i[l] = \mathbf{1}_{\{\mathbf{X}_i^T \mathbf{V}_0[l] > 0\}} \mathbf{X}_i^T = \mathbf{0}.$$

□

Now, we are ready to prove Lemma 7.

Proof. We prove by contradiction. Suppose on the contrary that with some nonzero probability, the design matrix is not full row-rank as $p \neq 1$. Note that when the design matrix is not full row-rank, there exists a set of indices $I \subseteq \{1, \dots, n\}$ such that

$$\sum_{i \in I} b_i \mathbf{H}_i = \mathbf{0}, \quad b_i \neq 0 \text{ for all } i \in I. \quad (44)$$

The proof will be finished by two steps: 1) find an event \mathcal{J} that happens almost surely when $p \neq 1$; 2) prove this event \mathcal{J} contradicts Eq. (44).

Step 1:

Consider each $i \in \{1, 2, \dots, n\}$. By Lemma 22, we know that there exists a $\mathbf{v}_{i,+} \in S^{d-1}$ such that

$$\mathbf{v}_{i,+}^T \mathbf{X}_i = 0, \quad \mathbf{v}_{i,+}^T \mathbf{X}_j \neq 0, \text{ for all } j \in \{1, 2, \dots, n\}, n, n \geq 1. \quad (45)$$

Define

$$r_i = \min_{j \in \{1, 2, \dots, n\}, n, n \geq 1} |\mathbf{v}_{i,+}^T \mathbf{X}_j| > 0. \quad (46)$$

For all $i = 1, 2, \dots, n$, we define several events as follows.

$$\begin{aligned} \mathcal{J}_i &:= \left\{ A_{\mathbf{V}_0} \setminus B_{\mathbf{v}_{i,+}}^{r_i, \mathbf{X}_i} \neq \emptyset, A_{\mathbf{V}_0} \setminus B_{\mathbf{v}_{i,-}}^{r_i, \mathbf{X}_i} \neq \emptyset \right\}, \\ \mathcal{J}_{i,+} &:= \left\{ A_{\mathbf{V}_0} \setminus B_{\mathbf{v}_{i,+}}^{r_i, \mathbf{X}_i} \neq \emptyset \right\}, \\ \mathcal{J}_{i,-} &:= \left\{ A_{\mathbf{V}_0} \setminus B_{\mathbf{v}_{i,-}}^{r_i, \mathbf{X}_i} \neq \emptyset \right\}, \\ \mathcal{J} &:= \bigcap_{i=1}^n \mathcal{J}_i. \end{aligned}$$

(Recall the geometric interpretation in Remark 6. The events $\mathcal{J}_{i,+}$ and $\mathcal{J}_{i,-}$ mean that there exists $\mathbf{V}_0[j]/k\mathbf{V}_0[j]k_2$ in the upper half and the lower half of the cap E , respectively. The event $\mathcal{J}_i = \mathcal{J}_{i,+} \setminus \mathcal{J}_{i,-}$ means that there exist $\mathbf{V}_0[j]/k\mathbf{V}_0[j]k_2$

in both halves of the cap E. Finally, the event \mathcal{J} occurs when \mathcal{J}_i occurs for all i , although the vector $\mathbf{V}_0[j]/k\mathbf{V}_0[j]k_2$ that falls into the two halves may differ across i . As we will show later, whenever the event \mathcal{J} occurs, the matrix \mathbf{H} will have the full row-rank, which is why we are interested in the probability of the event \mathcal{J} .)

Those definitions implies that

$$\mathcal{J}_i^c = \mathcal{J}_{i,+}^c \cup \mathcal{J}_{i,-}^c \quad \text{for all } i = 1, 2, \dots, n, \quad (47)$$

$$\mathcal{J}^c = \bigcup_{i=1}^n \mathcal{J}_i^c. \quad (48)$$

Thus, we have

$$\begin{aligned} \Pr_{\mathbf{V}_0}[\mathcal{J}] &= 1 - \Pr_{\mathbf{V}_0}[\mathcal{J}^c] \\ &= 1 - \sum_{i=1}^n \Pr_{\mathbf{V}_0}[\mathcal{J}_i^c] \quad (\text{by Eq. (48) and the union bound}). \end{aligned} \quad (49)$$

For a fixed i , recall that by Eq. (46), we have $r_i > 0$. Because $B_{\mathbf{v},i,+}^{r_i, \mathbf{X}_i}$ and $B_{\mathbf{v},i,-}^{r_i, \mathbf{X}_i}$ are two halves of $B_{\mathbf{v},i}^{r_i, \mathbf{X}_i}$, we have

$$\lambda_{d-1}(B_{\mathbf{v},i,+}^{r_i, \mathbf{X}_i}) = \lambda_{d-1}(B_{\mathbf{v},i,-}^{r_i, \mathbf{X}_i}) = \frac{1}{2} \lambda_{d-1}(B_{\mathbf{v},i}^{r_i, \mathbf{X}_i}). \quad (50)$$

Therefore, we have

$$\begin{aligned} \Pr_{\mathbf{V}_0}[\mathcal{J}_i^c] &= \Pr_{\mathbf{V}_0}[\mathcal{J}_{i,+}^c] + \Pr_{\mathbf{V}_0}[\mathcal{J}_{i,-}^c] \quad (\text{by Eq. (47) and the union bound}) \\ &= \left(1 - \frac{\lambda_{d-1}(B_{\mathbf{v},i,+}^{r_i, \mathbf{X}_i})}{\lambda_{d-1}(S^{d-1})} \right)^p + \left(1 - \frac{\lambda_{d-1}(B_{\mathbf{v},i,-}^{r_i, \mathbf{X}_i})}{\lambda_{d-1}(S^{d-1})} \right)^p \\ &\quad (\text{all } \mathbf{V}_0[i]\text{'s are independent and Assumption 1}) \\ &= 2 \left(1 - \frac{\lambda_{d-1}(B_{\mathbf{v},i}^{r_i, \mathbf{X}_i})}{2\lambda_{d-1}(S^{d-1})} \right)^p \quad (\text{by Eq. (50)}). \end{aligned}$$

Notice that r_i is determined only by \mathbf{X} , and is independent of \mathbf{V}_0 and p . Therefore, we have

$$\lim_{p \uparrow \infty} \Pr_{\mathbf{V}_0}[\mathcal{J}_i^c] = 0. \quad (51)$$

Plugging Eq. (51) into Eq. (49), we have

$$\lim_{p \uparrow \infty} \Pr_{\mathbf{V}_0}[\mathcal{J}] = 1 \quad (\text{because } n \text{ is finite}).$$

Step 2:

To complete the proof, it remains to show that the event \mathcal{J} contradicts Eq. (44). Towards this end, we assume the event \mathcal{J} happens. By Eq. (44), we can pick one $i \geq l$. Further, by the definition of \mathcal{J} , there exists r_i such that $A_{\mathbf{V}_0} \setminus B_{\mathbf{v},i,+}^{r_i, \mathbf{X}_i} \neq \emptyset$ and $A_{\mathbf{V}_0} \setminus B_{\mathbf{v},i,-}^{r_i, \mathbf{X}_i} \neq \emptyset$. In other words, there must exist $k, l \geq l$, $k \neq l$, such that

$$\frac{\mathbf{V}_0[k]}{k\mathbf{V}_0[k]k_2} \geq B_{\mathbf{v},i,+}^{r_i, \mathbf{X}_i}, \quad \frac{\mathbf{V}_0[l]}{k\mathbf{V}_0[l]k_2} \geq B_{\mathbf{v},i,-}^{r_i, \mathbf{X}_i}.$$

By Lemma 23, we have

$$\mathbf{H}_j[k] = \mathbf{H}_j[l], \quad \text{for all } j \geq l, \quad n \geq n \text{ fixed}, \quad (52)$$

$$\mathbf{H}_i[k] = \mathbf{X}_i^T, \quad \mathbf{H}_i[l] = \mathbf{0}. \quad (53)$$

We now show that \mathbf{H} restricted to the columns corresponding to k and l cannot be linearly dependent. Specifically, we have

$$\begin{aligned}
 \sum_{j \geq l} b_j \mathbf{H}_j[k] &= b_i \mathbf{H}_i[k] + \sum_{j \geq l, j \neq i} b_j \mathbf{H}_j[k] \text{ (as we have picked } i \geq l \text{)} \\
 &= b_i \mathbf{H}_i[k] - b_j \mathbf{H}_i[l] + \sum_{j \geq l} b_j \mathbf{H}_j[l] \text{ (by Eq. (52))} \\
 &= b_i \mathbf{X}_i^T + \sum_{j \geq l} b_j \mathbf{H}_j[l] \text{ (by Eq. (53))} \\
 &\notin \sum_{j \geq l} b_j \mathbf{H}_j[l] \text{ (because } b_i \neq 0 \text{)}.
 \end{aligned}$$

This contradicts the assumption Eq. (44) that

$$\sum_{j \geq l} b_j \mathbf{H}_j[k] = \sum_{j \geq l} b_j \mathbf{H}_j[l] = \mathbf{0}.$$

The result thus follows. \square

F. Proof of Proposition 4 (the upper bound of the variance)

The following lemma shows the relationship between the variance term and $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$.

Lemma 24.

$$j\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon}_j \leq \frac{\rho_{\bar{p}} k \epsilon k_2}{\sqrt{\min \text{eig}(\mathbf{H}\mathbf{H}^T)}}.$$

Proof. We have

$$k\mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon} k_2 = \sqrt{(\mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon})^T \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon}} = \sqrt{\boldsymbol{\epsilon}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon}} \frac{k \epsilon k_2}{\sqrt{\min \text{eig}(\mathbf{H}\mathbf{H}^T)}}. \quad (54)$$

Thus, we have

$$\begin{aligned}
 &j\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon}_j \\
 &= k\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon} k_2 \text{ (}\ell_2\text{-norm of a number equals to its absolute value)} \\
 &k\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} k_2 \cdot k\mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon} k_2 \text{ (by Lemma 12)} \\
 &\frac{\rho_{\bar{p}} k \epsilon k_2}{\sqrt{\min \text{eig}(\mathbf{H}\mathbf{H}^T)}} \text{ (by Lemma 11 and Eq. (54)).}
 \end{aligned}$$

\square

The following lemma shows our estimation on $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$.

Lemma 25. For any $n \geq 2$, $m \geq \left[1, \frac{\ln n}{\ln \frac{n}{2}}\right]$, $d \leq n^4$, if $p \geq 6J_m(n, d) \ln\left(4n^{1+\frac{1}{m}}\right)$, we must have

$$\Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ \frac{\min \text{eig}(\mathbf{H}\mathbf{H}^T)}{p} \geq \frac{1}{J_m(n, d)n} \right\} \geq 1 - \frac{2}{n^{\frac{m}{2}}}.$$

Proposition 4 directly follows from Lemma 25 and Lemma 24.⁸

In rest of this section, we will show how to prove Lemma 25. The following lemma shows that, to estimate $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$, it is equivalent to estimate $\min_{\mathbf{a} \geq 2S^{n-1}} k\mathbf{H}^T \mathbf{a} k_2^2/p$.

⁸We can see that the key part during the proof of Proposition 4 is to estimate $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$. Lemma 25 shows a lower bound of $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$ which is almost (n^{1-2d}) when p is large. However, our estimation of this value may be loose. We will show a upper bound which is $O(n^{-\frac{1}{d-1}})$ (see Appendix G).

Lemma 26.

$$\min_{\mathbf{a} \in \mathbb{R}^{2S^n-1}} \text{eig}(\mathbf{H}\mathbf{H}^T) = \min_{\mathbf{a} \in \mathbb{R}^{2S^n-1}} \|\mathbf{H}^T \mathbf{a}\|_2^2.$$

Proof. Do the singular value decomposition (SVD) of \mathbf{H}^T as $\mathbf{H}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{W}^T$, where

$$\mathbf{\Sigma} \in \mathbb{R}^{(dp) \times n} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n).$$

By properties of singular values, we have

$$\min_{\mathbf{a} \in \mathbb{R}^{2S^n-1}} \|\mathbf{H}^T \mathbf{a}\|_2^2 = \min_{i \in \{1, 2, \dots, ng\}} \sigma_i^2.$$

We also have

$$\begin{aligned} \mathbf{H}\mathbf{H}^T &= \mathbf{W}\mathbf{\Sigma}^T \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T \\ &= \mathbf{W}\mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{W}^T \quad (\text{because } \mathbf{U}^T \mathbf{U} = \mathbf{I}) \\ &= \mathbf{W} \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2) \mathbf{W}^T. \end{aligned}$$

This equation is indeed the eigenvalue decomposition of $\mathbf{H}\mathbf{H}^T$, which implies that its eigenvalues are $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$. Thus, we have

$$\min_{\mathbf{a} \in \mathbb{R}^{2S^n-1}} \text{eig}(\mathbf{H}\mathbf{H}^T) = \min_{i \in \{1, 2, \dots, ng\}} \sigma_i^2 = \min_{\mathbf{a} \in \mathbb{R}^{2S^n-1}} \|\mathbf{H}^T \mathbf{a}\|_2^2.$$

□

Therefore, to finish the proof of Proposition 4, it only remains to estimate $\min_{\mathbf{a} \in \mathbb{R}^{2S^n-1}} \|\mathbf{H}^T \mathbf{a}\|_2^2$.

By Lemma 7 and its proof in Appendix E, we have already shown that $\mathbf{H}^T \mathbf{a}$ is not likely to be zero (i.e. $\min_{\mathbf{a} \in \mathbb{R}^{2S^n-1}} \|\mathbf{H}^T \mathbf{a}\|_2^2 > 0$) when $p \neq 1$. Here, we basically use the similar method as in Appendix E, but with more precise quantification.

Recall the definitions in Eqs. (21)(22)(23). For any $i \in \{1, 2, \dots, ng\}$, we choose one

$$\mathbf{v}_{\cdot, i} \in \mathbb{S}^{d-1} \text{ independently of } \mathbf{X}_j, j \neq i, \text{ such that } \mathbf{v}_{\cdot, i}^T \mathbf{X}_i = 0. \quad (55)$$

(Note that here, unlike in Eq. (45), we do not require $\mathbf{v}_{\cdot, i}^T \mathbf{X}_j \neq 0$ for all $j \neq i$. This is important as we would like \mathbf{X}_j to be independent of $\mathbf{v}_{\cdot, i}$ for all $j \neq i$.) Further, for any $0 < r_0 < 1$, we define

$$c_{r_0}^i := \min \left\{ j \in A_{\mathbf{V}_0} \setminus B_{\mathbf{v}_{\cdot, i}^{r_0}, \mathbf{X}_i}^+, j \in A_{\mathbf{V}_0} \setminus B_{\mathbf{v}_{\cdot, i}^{r_0}, \mathbf{X}_i}^- \right\}. \quad (56)$$

Then, we define

$$r_i := \min_{j \in \{1, 2, \dots, ng\}} |\mathbf{v}_{\cdot, i}^T \mathbf{X}_j|, \quad (57)$$

$$\hat{r} := \min_{i \in \{1, 2, \dots, ng\}} r_i. \quad (58)$$

(Note that here r_i or \hat{r} may be zero. Later we will show that they can be lower bounded with high probability.) Define

$$D_{\mathbf{X}} := \frac{\lambda_{d-1}(B^{\hat{r}})}{8n\lambda_{d-1}(S^{d-1})}. \quad (59)$$

Similar to Remark 6, these definitions have their geometric interpretation in Fig. 6. The value $c_{r_0}^i$ denotes the number of distinct pairs $\left(\frac{\mathbf{V}_0[k]}{\|\mathbf{V}_0[k]\|_2}, \frac{\mathbf{V}_0[l]}{\|\mathbf{V}_0[l]\|_2} \right)$ ⁹ such that $\frac{\mathbf{V}_0[k]}{\|\mathbf{V}_0[k]\|_2}$ is in the upper half of the cap E, and $\frac{\mathbf{V}_0[l]}{\|\mathbf{V}_0[l]\|_2}$ is in the lower half of the cap E. The quantities r_0, r_i , and \hat{r} can all be used as the radius of the cap E. The ratio $D_{\mathbf{X}}$ is proportional to the area of the cap E with radius \hat{r} (or equivalently, the probability that the normalized $\mathbf{V}_0[j]$ falls in the cap E).

The following lemma gives an estimation on $\|\mathbf{H}^T \mathbf{a}\|_2^2/p$ when \mathbf{X} is given. We put its proof in Appendix F.1.

⁹Here, ‘‘distinct’’ means that any normalized version of $\mathbf{V}_0[j]$ can appear at most in one pair.

Lemma 27. Given \mathbf{X} , we have

$$\Pr_{\mathbf{V}_0} \left\{ k\mathbf{H}^T \mathbf{a} k_2^2 \leq pD_{\mathbf{X}}, \text{ for all } \mathbf{a} \in S^{n-1} \right\} \geq 1 - 4ne^{-npD_{\mathbf{X}}/6}.$$

Notice that $D_{\mathbf{X}}$ only depends on \mathbf{X} and it may even be zero if $\hat{\rho}$ is zero. However, after we introduce the randomness of \mathbf{X} , we can show that $\hat{\rho}$ is lower bounded with high probability. We can then obtain the following lemma. We put its proof in Appendix F.2.

Define

$$C_d := \frac{2^{\rho} \bar{2}}{B\left(\frac{d-1}{2}, \frac{1}{2}\right)}, \quad (60)$$

$$D(n, d, \delta) := \frac{1}{16n} I_{\frac{\delta^2}{n^4 C_d^2}} \left(1 - \frac{\delta^2}{4n^4 C_d^2} \right) \left(\frac{d-1}{2}, \frac{1}{2} \right). \quad (61)$$

Lemma 28. For any $\delta \in \left(0, \frac{2}{\pi}\right]$, we have

$$\Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ k\mathbf{H}^T \mathbf{a} k_2^2 \leq pD(n, d, \delta), \text{ for all } \mathbf{a} \in S^{n-1} \right\} \geq 1 - 4ne^{-npD(n, d, \delta)/6} - \delta.$$

Notice that Lemma 28 is already very close to Lemma 25, and we put the final steps of the proof of Lemma 25 in Appendix F.3.

F.1. Proof of Lemma 27

Proof. Define events as follows.

$$\begin{aligned} \mathcal{J} &:= \left\{ k\mathbf{H}^T \mathbf{a} k_2^2 \leq pD_{\mathbf{X}}, \text{ for all } \mathbf{a} \in S^{n-1} \right\}, \\ \mathcal{J}_i &:= \left\{ \text{there exists } \mathbf{a} \in S^{n-1} \text{ that } i \in \arg \max_{j \in \{1, 2, \dots, n\}} |a_j|, \text{ and } k\mathbf{H}^T \mathbf{a} k_2^2 \leq pD_{\mathbf{X}} \right\}, \\ \mathcal{K}_i &:= \left\{ c_{r_i}^i \leq 2npD_{\mathbf{X}} \right\}, \text{ for } i = 1, 2, \dots, n. \end{aligned}$$

Those definitions directly imply that

$$\mathcal{J}^c = \bigcup_{i=1}^n \mathcal{J}_i. \quad (62)$$

Step 1: prove $\mathcal{J}_i \subseteq \mathcal{K}_i$

To show $\mathcal{J}_i \subseteq \mathcal{K}_i$, we only need to prove that \mathcal{J}_i implies \mathcal{K}_i . To that end, it suffices to show $k\mathbf{H}^T \mathbf{a} k_2^2 \geq \frac{c_{r_i}^i}{2n}$ for the vector \mathbf{a} defined in \mathcal{J}_i . Because $i \in \arg \max_{j=1}^n |a_j|$ and $k\mathbf{a} k_2 = 1$, we have

$$|a_i| \geq \frac{1}{n}. \quad (63)$$

By Eq. (56), we can construct $c_{r_i}^i$ pairs (k_j, l_j) for $j = 1, 2, \dots, c_{r_i}^i$ (all k_j 's are different and all l_j 's are different), such that

$$\frac{\mathbf{V}_0[k_j]}{k\mathbf{V}_0[k_j] k_2} \geq B_{\mathbf{v}, i, +}^{r_i, \mathbf{X}_i}, \quad \frac{\mathbf{V}_0[l_j]}{k\mathbf{V}_0[l_j] k_2} \geq B_{\mathbf{v}, i, -}^{r_i, \mathbf{X}_i}.$$

Thus, we have

$$\begin{aligned} (\mathbf{H}^T \mathbf{a})[k_j] - (\mathbf{H}^T \mathbf{a})[l_j] &= \sum_{k=1}^n a_k (\mathbf{H}_k[k_j] - \mathbf{H}_k[l_j]) \\ &= a_i (\mathbf{H}_i[k_j] - \mathbf{H}_i[l_j]) + \sum_{k \in \{1, 2, \dots, n\} \setminus \{i\}} a_k (\mathbf{H}_k[k_j] - \mathbf{H}_k[l_j]) \\ &= a_i \mathbf{X}_i \text{ (by Lemma 23)}. \end{aligned}$$

We then have

$$\begin{aligned} k(\mathbf{H}^T \mathbf{a})[k_j]k_2^2 + k(\mathbf{H}^T \mathbf{a})[l_j]k_2^2 & \frac{1}{2}ka_i\mathbf{X}_ik_2^2 \text{ (by Lemma 13)} \\ & \frac{1}{2n} \text{ (by Eq. (63)).} \end{aligned}$$

Further, we have

$$k\mathbf{H}^T \mathbf{a}k_2^2 = \sum_{j=1}^p k(\mathbf{H}^T \mathbf{a})[j]k_2^2 + \sum_{j=1}^{c_{r_i}^i} k(\mathbf{H}^T \mathbf{a})[k_j]k_2^2 + k(\mathbf{H}^T \mathbf{a})[l_j]k_2^2 = \frac{c_{r_i}^i}{2n}. \quad (64)$$

Clearly, if the event \mathcal{J}_i occurs, then $k\mathbf{H}\mathbf{a}k_2^2 \leq pD_{\mathbf{X}}$. Combining with Eq. (64), we then have $c_{r_i}^i \leq 2npD_{\mathbf{X}}$. In other words, the event \mathcal{K}_i must occur. Hence, we have shown that $\mathcal{J}_i \subseteq \mathcal{K}_i$.

Step 2: estimate the probability of \mathcal{K}_i

For all $j \in \{1, 2, \dots, p\}$, because $\mathbf{V}_0[j]$ is uniformly distributed in all directions, for any fixed $0 < r_0 < 1$, we have

$$\Pr_{\mathbf{V}_0} \left\{ \frac{\mathbf{V}_0[j]}{k\mathbf{V}_0[j]k_2} \geq B_{\mathbf{v}, i, +}^{r_0, \mathbf{X}_i} \mid i \right\} = \frac{\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{2\lambda_{d-1}(S^{d-1})}.$$

Thus, $jA_{\mathbf{V}_0} \setminus B_{\mathbf{v}, i, +}^{r_0, \mathbf{X}_i}$ follows the distribution $\text{Bino} \left(p, \frac{\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{2\lambda_{d-1}(S^{d-1})} \right)$ given i and \mathbf{X} . By Lemma 14 (with $\delta = \frac{1}{2}$), we have

$$\Pr_{\mathbf{V}_0} \left\{ jA_{\mathbf{V}_0} \setminus B_{\mathbf{v}, i, +}^{r_0, \mathbf{X}_i} < \frac{p\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{4\lambda_{d-1}(S^{d-1})} \mid i \right\} \leq 2 \exp \left(- \frac{p\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{48\lambda_{d-1}(S^{d-1})} \right). \quad (65)$$

Similarly, we have

$$\Pr_{\mathbf{V}_0} \left\{ jA_{\mathbf{V}_0} \setminus B_{\mathbf{v}, i, -}^{r_0, \mathbf{X}_i} < \frac{p\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{4\lambda_{d-1}(S^{d-1})} \mid i \right\} \leq 2 \exp \left(- \frac{p\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{48\lambda_{d-1}(S^{d-1})} \right). \quad (66)$$

By plugging Eq. (65) and Eq. (66) into Eq. (56) and applying the union bound, we have

$$\Pr_{\mathbf{V}_0} \left\{ c_{r_0}^i < \frac{p\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{4\lambda_{d-1}(S^{d-1})} \mid i \right\} \leq 4 \exp \left(- \frac{p\lambda_{d-1}(B_{\mathbf{v}}^{r_0})}{48\lambda_{d-1}(S^{d-1})} \right).$$

By letting $r_0 = \hat{r}$ and by Eq. (59), we thus have

$$\Pr_{\mathbf{V}_0} \left\{ c_{r_i}^i \leq 2npD_{\mathbf{X}} \mid i \right\} \leq 4 \exp \left(- \frac{1}{6} npD_{\mathbf{X}} \right),$$

i.e.,

$$\Pr_{\mathbf{V}_0}[\mathcal{K}_i] \geq 4 \exp \left(- \frac{1}{6} npD_{\mathbf{X}} \right), \text{ for all } i = 1, 2, \dots, n. \quad (67)$$

Step3: estimate the probability of \mathcal{J}

We have

$$\begin{aligned} \Pr_{\mathbf{V}_0}[\mathcal{J}^c] & \sum_{i=1}^n \Pr_{\mathbf{V}_0}[\mathcal{J}_i^c] \text{ (by Eq. (62) and the union bound)} \\ & \sum_{i=1}^n \Pr_{\mathbf{V}_0}[\mathcal{K}_i^c] \text{ (by } \mathcal{J}_i \subseteq \mathcal{K}_i \text{ proven in Step 1)} \\ & 4n \exp \left(- \frac{1}{6} npD_{\mathbf{X}} \right) \text{ (by Eq. (67)).} \end{aligned}$$

Thus, we have

$$\Pr_{\mathbf{V}_0}[\mathcal{J}] = 1 - \Pr_{\mathbf{V}_0}[\mathcal{J}^c] \geq 1 - 4n \exp \left(- \frac{1}{6} npD_{\mathbf{X}} \right).$$

The result of this lemma thus follows. \square

F.2. Proof of Lemma 28

Based on Lemma 27, it remains to estimate $\hat{\rho}$, which will then allow us to bound $D_{\mathbf{X}}$. Towards this end, we need a few lemmas to estimate $B\left(\frac{d-1}{2}, \frac{1}{2}\right)$ and $I_x\left(\frac{d-1}{2}, \frac{1}{2}\right)$.

Lemma 29. For any $x \geq \mathbb{R}$, we must have $x + 1 \leq e^x$.

Proof. Consider a function $g(x) = e^x - x - 1$. It remains to show that $g(x) \geq 0$ for all x . We have $g'(x) = e^x - 1$. In other words, $g'(x) \leq 0$ when $x \leq 0$, and $g'(x) \geq 0$ when $x \geq 0$. Thus, $g(x)$ is monotone decreasing when $x \leq 0$, and is monotone increasing when $x \geq 0$. Hence, we know that $g(x)$ achieves its minimum value at $x = 0$, i.e., $g(x) \geq g(0) = 0$ for any x . The conclusion of this lemma thus follows. \square

Lemma 30. For any $d \geq 5$, we have

$$\left(1 - \frac{1}{d-3}\right)^{d-3} \geq \frac{1}{e^2}.$$

Proof. By letting $x = \frac{1}{d-4}$ in Lemma 29, we have

$$\frac{d-3}{d-4} = \frac{1}{d-4} + 1 \leq \exp\left(\frac{1}{d-4}\right),$$

i.e.,

$$\frac{d-4}{d-3} \geq \exp\left(-\frac{1}{d-4}\right). \quad (68)$$

Thus, we have

$$\begin{aligned} \left(1 - \frac{1}{d-3}\right)^{d-3} &= \left(\frac{d-4}{d-3}\right)^{d-3} \\ &\geq \exp\left(-\frac{d-3}{d-4}\right) \\ &= \exp\left(-1 - \frac{1}{d-4}\right) \\ &\geq \exp(-2) \text{ (because } \exp(\cdot) \text{ is monotone increasing and } d \geq 5). \end{aligned}$$

\square

Lemma 31. For any $d \geq 5$, we must have

$$\frac{2}{e} \sqrt{\frac{1}{d-3}} \geq \frac{1}{d}$$

Proof. Because $1 - \frac{4}{e^2} \approx 0.46 > 0.6$, we have

$$\begin{aligned} &\frac{3}{5} > 1 - \frac{4}{e^2} \\ \Rightarrow &\frac{3}{d} > 1 - \frac{4}{e^2} \text{ (because } d \geq 5) \\ \Rightarrow &1 - \frac{3}{d} < \frac{4}{e^2} \\ \Rightarrow &\frac{d-3}{d} < \frac{4}{e^2} \\ \Rightarrow &\frac{4}{e^2} > \frac{d-3}{d} > 1 \\ \Rightarrow &\frac{2}{e} \sqrt{\frac{1}{d-3}} \geq \frac{1}{d}. \end{aligned}$$

\square

Lemma 32.

$$B\left(\frac{d-1}{2}, \frac{1}{2}\right) \geq \left[\frac{1}{d}, \pi\right].$$

Further, if $d \geq 5$, we have

$$B\left(\frac{d-1}{2}, \frac{1}{2}\right) \geq \left[\frac{1}{d}, \frac{4}{d-3}\right].$$

Proof. When $d = 2$, we have $B\left(\frac{d-1}{2}, \frac{1}{2}\right) = \pi$. When $d = 3$, we have $B\left(\frac{d-1}{2}, \frac{1}{2}\right) = 2$. When $d = 4$, we have $B\left(\frac{d-1}{2}, \frac{1}{2}\right) \approx 1.57$. It is easy to verify that the statement of the lemma holds for $d = 2, 3$, and 4. It remains to validate the case of $d \geq 5$. We first prove the lower bound. For any $m \in (0, 1)$, we have

$$\begin{aligned} B\left(\frac{d-1}{2}, \frac{1}{2}\right) &= \int_0^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \\ &\geq \int_m^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \quad (\text{because } t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} \geq 0) \\ &\geq m^{\frac{d-3}{2}} \int_m^1 (1-t)^{\frac{1}{2}} dt \\ &\quad (\text{because } t^{\frac{d-3}{2}} \text{ is monotone increasing with respect to } t \text{ when } d \geq 5) \\ &= m^{\frac{d-3}{2}} \left(2 \sqrt{1-t} \Big|_m^1 \right) \\ &= m^{\frac{d-3}{2}} \frac{2}{\sqrt{1-m}}. \end{aligned}$$

By letting $m = 1 - \frac{1}{d-3}$, we thus have

$$B\left(\frac{d-1}{2}, \frac{1}{2}\right) \geq \left(1 - \frac{1}{d-3}\right)^{\frac{d-3}{2}} 2\sqrt{\frac{1}{d-3}}.$$

Then, applying Lemma 30, we have

$$B\left(\frac{d-1}{2}, \frac{1}{2}\right) \geq \frac{2}{e} \sqrt{\frac{1}{d-3}}.$$

Thus, by Lemma 31, we have

$$B\left(\frac{d-1}{2}, \frac{1}{2}\right) \geq \frac{1}{d}.$$

Now we prove the upper bound. For any $m \in (0, 1)$, we have

$$\begin{aligned} B\left(\frac{d-1}{2}, \frac{1}{2}\right) &= \int_0^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \\ &= \int_0^m t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt + \int_m^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \\ &\leq \int_0^m t^{\frac{d-3}{2}} (1-m)^{\frac{1}{2}} dt + \int_m^1 (1-t)^{\frac{1}{2}} dt \\ &= \frac{2}{d-1} m^{\frac{d-1}{2}} (1-m)^{\frac{1}{2}} + 2 \sqrt{1-m} \\ &\leq \frac{2}{d-1} (1-m)^{\frac{1}{2}} + 2 \sqrt{1-m} \quad (\text{because } m < 1 \text{ and } d \geq 5). \end{aligned}$$

By letting $m = 1 - \frac{1}{d-3}$, we thus have

$$B\left(\frac{d-1}{2}, \frac{1}{2}\right) \frac{2^{\frac{d-1}{3}}}{d-1} + \frac{2}{d-3} \frac{4}{d-3}.$$

Notice that $\frac{4}{d-3} = 2^{\frac{2}{d-3}} < \pi$. The result of this lemma thus follows. □

Lemma 33. Recall C_d is defined in Eq. (60). If $d \geq n^4$ and $\delta \geq 1$, then

$$\left(1 - \frac{\delta^2}{4n^4 C_d^2}\right)^{\frac{d-1}{2}} \geq \frac{1}{2}.$$

Proof. We have

$$\begin{aligned} \left(1 - \frac{\delta^2}{4n^4 C_d^2}\right)^{\frac{d-1}{2}} &= \left(1 - \frac{\delta^2}{4n^4 C_d^2}\right)^{d-1} \\ &\geq 1 - \frac{(d-1)\delta^2}{4n^4 C_d^2} \quad (\text{by Bernoulli's inequality } (1+x)^a \geq 1+ax) \\ &= 1 - \frac{(d-1) \left(B\left(\frac{d-1}{2}, \frac{1}{2}\right)\right)^2}{4n^4 \cdot 8} \quad (\text{by } \delta \geq 1 \text{ and Eq. (60)}) \\ &\geq 1 - \frac{(d-1)\pi^2}{32n^4} \quad (\text{by Lemma 32}) \\ &\geq 1 - \frac{d}{n^4} \cdot \frac{\pi^2}{32} \\ &\geq \frac{1}{2} \quad (\text{because } n^4 \geq d \text{ and } \pi \leq 4). \end{aligned}$$

□

Lemma 34. For any $\delta \geq \left(0, \frac{2}{\pi}\right]$, we must have $\frac{\delta}{n^2 C_d} \geq \frac{1}{2}$.

Proof. Because Eq. (60), $\delta \geq \frac{2}{\pi}$, and $n \geq 1$, this lemma directly follows by Lemma 32. □

Lemma 35. For any $x \geq [0, 1]$, we must have

$$I_x\left(\frac{d-1}{2}, \frac{1}{2}\right) \geq \frac{C_d}{2(d-1)} x^{\frac{d-1}{2}},$$

and

$$\lim_{x \downarrow 0} \frac{I_x\left(\frac{d-1}{2}, \frac{1}{2}\right)}{x^{\frac{d-1}{2}}} = \frac{C_d}{2(d-1)}.$$

Proof. we have

$$\begin{aligned}
 I_x \left(\frac{d-1}{2}, \frac{1}{2} \right) &= \frac{\int_0^x t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt}{B \left(\frac{d-1}{2}, \frac{1}{2} \right)} \\
 &= \frac{C_d}{2 \rho^{\frac{d-1}{2}}} \int_0^x t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \text{ (by Eq. (60))} \\
 &\geq \left[\frac{C_d}{2 \rho^{\frac{d-1}{2}}} \int_0^x t^{\frac{d-3}{2}} dt, \frac{C_d}{2 \rho^{\frac{d-1}{2}} \rho^{\frac{1}{2}} x} \int_0^x t^{\frac{d-3}{2}} dt \right] \\
 &\quad \text{(because } (1-t)^{1/2} \geq \left[1, \rho^{\frac{1}{2}} \frac{1}{x} \right]) \\
 &\geq \left[\rho^{\frac{d-1}{2}} \frac{C_d}{2(d-1)} x^{\frac{d-1}{2}}, \rho^{\frac{d-1}{2}} \frac{C_d}{2(d-1)} \rho^{\frac{1}{2}} \frac{1}{x} x^{\frac{d-1}{2}} \right].
 \end{aligned}$$

Thus, we have

$$\frac{I_x \left(\frac{d-1}{2}, \frac{1}{2} \right)}{x^{\frac{d-1}{2}}} \geq \left[\rho^{\frac{d-1}{2}} \frac{C_d}{2(d-1)}, \rho^{\frac{d-1}{2}} \frac{C_d}{2(d-1)} \rho^{\frac{1}{2}} \frac{1}{x} \right],$$

which implies

$$\lim_{x \downarrow 0} \frac{I_x \left(\frac{d-1}{2}, \frac{1}{2} \right)}{x^{\frac{d-1}{2}}} = \rho^{\frac{d-1}{2}} \frac{C_d}{2(d-1)}.$$

□

Lemma 36. For any $x \in [\frac{1}{2}, 1)$ and for any $d \in \mathbb{N}, 3 \leq d$, we have

$$I_x \left(\frac{d-1}{2}, \frac{1}{2} \right) \geq \frac{2\sqrt{2(1-x)}}{B \left(\frac{d-1}{2}, \frac{1}{2} \right)}.$$

We also have

$$\lim_{(1-x) \downarrow 0^+} \frac{1}{\rho^{\frac{d-1}{2}}} \frac{I_x \left(\frac{d-1}{2}, \frac{1}{2} \right)}{1-x} = \frac{2}{B \left(\frac{d-1}{2}, \frac{1}{2} \right)}.$$

Proof. By the definition of regularized incomplete beta function in Eq. (20), we have

$$I_x \left(\frac{d-1}{2}, \frac{1}{2} \right) = \frac{\int_0^x t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt}{B \left(\frac{d-1}{2}, \frac{1}{2} \right)} = 1 - \frac{\int_x^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt}{B \left(\frac{d-1}{2}, \frac{1}{2} \right)}.$$

Thus, it remains to show that

$$\int_x^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \leq 2\sqrt{2(1-x)}, \text{ and} \tag{69}$$

$$\lim_{(1-x) \downarrow 0^+} \frac{\int_x^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt}{\rho^{\frac{d-1}{2}} \frac{1}{x}} = 2. \tag{70}$$

First, we prove Eq. (69). Case 1: $d = 2$. We have

$$\begin{aligned}
 & \int_x^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \\
 &= \int_x^1 t^{\frac{1}{2}} (1-t)^{\frac{1}{2}} dt \\
 & \stackrel{1}{=} \int_x^1 (1-t)^{\frac{1}{2}} dt \text{ (because } t^{\frac{1}{2}} \text{ is monotone decreasing in } [x, 1]) \\
 &= 2\sqrt{\frac{1-x}{x}} \\
 &= 2\sqrt{2(1-x)} \text{ (because } x = \frac{1}{2}).
 \end{aligned}$$

Case 2: $d \geq 3$. Then $t^{\frac{d-3}{2}}$ is monotone increasing in $[x, 1]$. Thus, we have

$$\int_x^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt \geq \int_x^1 (1-t)^{\frac{1}{2}} dt = 2\sqrt{1-x} - 2\sqrt{2(1-x)}.$$

To conclude, for all $d \geq 2, 3$, Eq. (69) holds.

Second, we prove Eq. (70). We have

$$\begin{aligned}
 & \frac{\int_x^1 t^{\frac{d-3}{2}} (1-t)^{\frac{1}{2}} dt}{\sqrt{\frac{1-x}{x}}} \geq 2 \left[\frac{\min\{1, x^{\frac{d-3}{2}} g\} \int_x^1 (1-t)^{\frac{1}{2}} dt}{\sqrt{\frac{1-x}{x}}}, \frac{\max\{1, x^{\frac{d-3}{2}} g\} \int_x^1 (1-t)^{\frac{1}{2}} dt}{\sqrt{\frac{1-x}{x}}} \right] \\
 &= \left[2 \min\{1, x^{\frac{d-3}{2}} g\}, 2 \max\{1, x^{\frac{d-3}{2}} g\} \right].
 \end{aligned}$$

Since $\lim_{x \rightarrow 1} x^{\frac{d-3}{2}} = 1$, Eq. (70) thus follows. \square

Now we are ready to prove Lemma 28.

Recall \mathbf{v}_i defined in Eq. (55). For any $b \geq \left(0, \frac{1}{2}\right]$, we have, for \mathbf{x} independent of \mathbf{v}_i and with distribution μ ,

$$\begin{aligned}
 \Pr_{\mathbf{x}, \mu} \left\{ \mathbf{v}_i^T \mathbf{x} \geq b \right\} &= I_{1-b^2} \left(\frac{d-1}{2}, \frac{1}{2} \right) \text{ (because Lemma 10)} \\
 &\geq 1 - \frac{2\sqrt{2(1-(1-b^2))}}{B\left(\frac{d-1}{2}, \frac{1}{2}\right)} \text{ (by Lemma 36)} \\
 &= 1 - C_d b \text{ (by the definition of } C_d \text{ in Eq. (60)).}
 \end{aligned} \tag{71}$$

Since each of the $\mathbf{X}_j, j \neq i$, is independent of \mathbf{v}_i , we have

$$\begin{aligned}
 & \Pr_{\mathbf{X}} \left\{ \min_{j \in \{1, 2, \dots, n\} \setminus \{i\}} \mathbf{v}_i^T \mathbf{X}_j \geq b \right\} \\
 &= \left(\Pr_{\mathbf{x}, \mu} \left\{ \mathbf{v}_i^T \mathbf{x} \geq b \right\} \right)^{n-1} \text{ (because each } \mathbf{X}_j, j \neq i, \text{ is i.i.d. and independent of } \mathbf{v}_i) \\
 &= (1 - C_d b)^{n-1} \text{ (by Eq. (71))} \\
 &\geq 1 - (n-1)C_d b \text{ (by Bernoulli's inequality)} \\
 &\geq 1 - nC_d b.
 \end{aligned}$$

Or, equivalently,

$$\Pr_{\mathbf{X}} \left\{ \min_{i \in \{1, 2, \dots, n\} \setminus \{j\}} \mathbf{v}_i^T \mathbf{X}_i < b \right\} \leq nC_d b. \tag{72}$$

Recall the definition of r_i and \hat{r} in Eqs. (57)(58). Thus, we then have

$$\begin{aligned}
 & \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ \hat{r} < \frac{\delta}{n^2 C_d} \right\} \\
 & n \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ r_i < \frac{\delta}{n^2 C_d} \right\} \text{ (by Eq. (58) and the union bound)} \\
 & = n \Pr_{\mathbf{X}} \left\{ r_i < \frac{\delta}{n^2 C_d} \right\} \text{ (because } r \text{ is independent of } \mathbf{V}_0) \\
 & = n \Pr_{\mathbf{X}} \left\{ \min_{j \in \{1, 2, \dots, n\}} |\mathbf{v}_{\cdot, i}^T \mathbf{X}_j| < \frac{\delta}{n^2 C_d} \right\} \text{ (by Eq. (57))} \\
 & \leq n \Pr_{\mathbf{X}} \left\{ \frac{\delta}{n^2 C_d} < \frac{\delta}{n^2 C_d} \right\} \text{ (by letting } b = \frac{\delta}{n^2 C_d} \text{ in Eq. (72) and } b \leq \frac{1}{2} \text{ because of Lemma 34)} \\
 & = \delta.
 \end{aligned} \tag{73}$$

By Lemma 9 and Eq. (61), we have

$$\lambda_{d-1}(B_{n^2 C_d}^{\frac{\delta}{n^2}}) = \frac{1}{2} \lambda_{d-1}(S^{d-1}) I_{\frac{\delta^2}{n^4 C_d^2}}(1 - \frac{\delta^2}{4n^4 C_d^2}) \left(\frac{d-1}{2}, \frac{1}{2} \right) = 8n \lambda_{d-1}(S^{d-1}) D(n, d, \delta).$$

Thus, we have

$$D_{\mathbf{X}} = D(n, d, \delta), \text{ when } \hat{r} \leq \frac{\delta}{n^2 C_d}. \tag{74}$$

By Eq. (59) and Eq. (74), we have

$$D_{\mathbf{X}} = D(n, d, \delta), \text{ when } \hat{r} \leq \frac{\delta}{n^2 C_d}.$$

Notice that \hat{r} only depends on \mathbf{X} and is independent of \mathbf{V}_0 . By Lemma 27, for any \mathbf{X} that makes $\hat{r} \leq \frac{\delta}{n^2 C_d}$, we must have

$$\Pr_{\mathbf{V}_0} \left\{ k \mathbf{H}^T \mathbf{a} k_2^2 \leq p D(n, d, \delta), \text{ for all } \mathbf{a} \geq S^{n-1} \right\} \geq 1 - 4n e^{-np D(n, d, \delta)/6}.$$

In other words,

$$\Pr_{\mathbf{V}_0} \left\{ k \mathbf{H}^T \mathbf{a} k_2^2 \leq p D(n, d, \delta), \text{ for all } \mathbf{a} \geq S^{n-1} \mid \text{any given } \mathbf{X} \text{ such that } \hat{r} \leq \frac{\delta}{n^2 C_d} \right\} \geq 1 - 4n e^{-np D(n, d, \delta)/6}.$$

We thus have

$$\Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ k \mathbf{H}^T \mathbf{a} k_2^2 \leq p D(n, d, \delta), \text{ for all } \mathbf{a} \geq S^{n-1} \mid \hat{r} \leq \frac{\delta}{n^2 C_d} \right\} \geq 1 - 4n e^{-np D(n, d, \delta)/6}. \tag{75}$$

Thus, we have

$$\begin{aligned}
 & \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ k \mathbf{H}^T \mathbf{a} k_2^2 \leq p D(n, d, \delta), \text{ for all } \mathbf{a} \geq S^{n-1} \right\} \\
 & \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ \hat{r} \leq \frac{\delta}{n^2 C_d}, \text{ and } k \mathbf{H}^T \mathbf{a} k_2^2 \leq p D(n, d, \delta), \text{ for all } \mathbf{a} \geq S^{n-1} \right\} \\
 & = \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ k \mathbf{H}^T \mathbf{a} k_2^2 \leq p D(n, d, \delta), \text{ for all } \mathbf{a} \geq S^{n-1} \mid \hat{r} \leq \frac{\delta}{n^2 C_d} \right\} \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ \hat{r} \leq \frac{\delta}{n^2 C_d} \right\} \\
 & \geq (1 - 4n e^{-np D(n, d, \delta)/6}) (1 - \delta) \text{ (by Eq. (73) and Eq. (75))} \\
 & \geq 1 - 4n e^{-np D(n, d, \delta)/6} - \delta.
 \end{aligned}$$

The result of this lemma thus follows.

E.3. Proof of Lemma 25

Based on Lemma 28, it only remains to estimate $D(n, d, \delta)$. We start with a lemma.

Lemma 37. *If $\delta \geq 1$ and $d \leq n^4$, we must have*

$$D(n, d, \delta) \leq 2^{1.5d} 5.5^d n^{-2d+1} \delta^{d-1}. \quad (76)$$

For any given $\delta \geq 0$ and d , we must have

$$\lim_{n \rightarrow \infty} \frac{D(n, d, \delta)}{n^{2d-1}} = 2^{1.5d} 1.5 \left(B \left(\frac{d-1}{2}, \frac{1}{2} \right) \right)^{d-2} \frac{1}{d-1} \delta^{d-1}.$$

Proof. We start from

$$\begin{aligned} \frac{1}{(d-1)C_d^{d-2}} &= \frac{\left(B \left(\frac{d-1}{2}, \frac{1}{2} \right) \right)^{d-2}}{(d-1) \left(\frac{\sqrt{2}}{2} \right)^{d-2}} \quad (\text{by Eq. (60)}) \\ &= \frac{1}{(d-1) d^{\frac{d-1}{2}} \left(\frac{\sqrt{2}}{2} \right)^{d-2}} \quad (\text{by Lemma 32}) \\ &= \frac{1}{d^{\frac{d}{2}} \left(\frac{\sqrt{2}}{2} \right)^d} \\ &= (8d)^{-\frac{d}{2}}. \end{aligned} \quad (77)$$

Thus, we have

$$\begin{aligned} D(n, d, \delta) &= \frac{1}{16n} \frac{C_d}{2(d-1)} \left(\frac{\delta^2}{n^4 C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}} \quad (\text{by Eq. (61) and Lemma 35}) \\ &= \frac{1}{16} \frac{1}{\sqrt{2}} \frac{1}{(d-1)C_d^{d-2}} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right)^{\frac{d-1}{2}} \frac{\delta^{d-1}}{n^{2d-1}} \\ &= \frac{1}{32} \frac{1}{\sqrt{2}} (8d)^{-\frac{d}{2}} \frac{\delta^{d-1}}{n^{2d-1}} \quad (\text{by Lemma 33 and Eq. (77)}) \\ &= 2^{1.5d} 5.5^d n^{-2d+1} \delta^{d-1}. \end{aligned}$$

For any given d and $\delta \geq 0$, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{D(n, d, \delta)}{n^{2d-1}} &= \lim_{n \rightarrow \infty} \frac{1}{16n^{2d-2}} I_{\frac{\delta^2}{n^4 C_d^2}} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \left(\frac{d-1}{2}, \frac{1}{2} \right) \quad (\text{by Eq. (61)}) \\ &= \lim_{n \rightarrow \infty} \frac{\left(\frac{\delta^2}{n^4 C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}}{16n^{2d-2}} \frac{I_{\frac{\delta^2}{n^4 C_d^2}} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \left(\frac{d-1}{2}, \frac{1}{2} \right)}{\left(\frac{\delta^2}{n^4 C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}} \\ &= \frac{1}{16} \lim_{n \rightarrow \infty} \frac{\left(\frac{\delta^2}{C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}}{\left(\frac{\delta^2}{n^4 C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}} \frac{I_{\frac{\delta^2}{n^4 C_d^2}} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \left(\frac{d-1}{2}, \frac{1}{2} \right)}{\left(\frac{\delta^2}{n^4 C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}} \\ &= \frac{1}{16} \lim_{n \rightarrow \infty} \frac{\left(\frac{\delta^2}{C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}}{\left(\frac{\delta^2}{n^4 C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}} \lim_{n \rightarrow \infty} \frac{I_{\frac{\delta^2}{n^4 C_d^2}} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \left(\frac{d-1}{2}, \frac{1}{2} \right)}{\left(\frac{\delta^2}{n^4 C_d^2} \left(1 + \frac{\delta^2}{4n^4 C_d^2} \right) \right)^{\frac{d-1}{2}}} \\ &= \frac{1}{16} \frac{\delta^{d-1}}{C_d^{d-1}} \frac{C_d}{2(d-1)} \quad (\text{by Lemma 35}) \\ &= 2^{1.5d} 1.5 \left(B \left(\frac{d-1}{2}, \frac{1}{2} \right) \right)^{d-2} \frac{1}{d-1} \delta^{d-1} \quad (\text{by Eq. (60)}). \end{aligned}$$

□

Now we are ready to finish our proof of Lemma 25.

We have

$$\begin{aligned} D(n, d, \delta) \Big|_{\delta = \frac{1}{m} \frac{1}{n}} &= \frac{1}{2^{1.5d+5.5} d^{0.5d} n^{2d-1} n^{\frac{d-1}{m}}} \text{ (by Eq. (76))} \\ &= \frac{1}{2^{1.5d+5.5} d^{0.5d} n^{(2+\frac{1}{m})(d-1)}} \\ &= \frac{1}{J_m(n, d)n} \text{ (by Eq. (9)).} \end{aligned}$$

Thus, when $p \geq 6J_m(n, d) \ln \left(4n^{1+\frac{1}{m}} \right)$, we have

$$1 - 4ne^{-npD(n, d, \delta)/6} \Big|_{\delta = \frac{1}{m} \frac{1}{n}} \geq 1 - \frac{2}{\pi n}.$$

Then, we have

$$m \geq \left[1, \frac{\ln n}{\ln \frac{\pi}{2}} \right] \Rightarrow \left(\frac{\pi}{2} \right)^m \geq n \Rightarrow n^{\frac{1}{m}} \leq \frac{\pi}{2} \Rightarrow \frac{1}{m} \leq \frac{2}{\pi} \Rightarrow \delta \leq \frac{2}{\pi}.$$

By Lemma 26 and Lemma 28, the conclusion of Lemma 25 thus follows.

G. Upper bound of $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$

By Lemma 26, to get an upper bound of $\min \text{eig}(\mathbf{H}\mathbf{H}^T)/p$, it is equivalent to get an upper bound of $\min_{\mathbf{a} \in S^{n-1}} k\mathbf{H}^T \mathbf{a} k_2^2/p$. To that end, we only need to construct a vector \mathbf{a} and calculate the value of $k\mathbf{H}^T \mathbf{a} k_2^2/p$, which automatically becomes an upper bound $\min_{\mathbf{a} \in S^{n-1}} k\mathbf{H}^T \mathbf{a} k_2^2/p$.

The following lemma shows that, for given \mathbf{X} , if two input training data \mathbf{X}_i and \mathbf{X}_k are close to each other, then $\min_{\mathbf{a} \in S^{n-1}} k\mathbf{H}^T \mathbf{a} k_2^2/p$ is unlikely to be large.

Lemma 38. *If there exist \mathbf{X}_i and \mathbf{X}_k such that $i \neq k$ and $\theta := \arccos(\mathbf{X}_i^T \mathbf{X}_k)$, then*

$$\Pr \left\{ \min_{\mathbf{a} \in S^{n-1}} k\mathbf{H}^T \mathbf{a} k_2^2 \leq \frac{3p\theta^2}{8} + \frac{3p\theta}{4\pi} \right\} \leq 2 \exp \left(-\frac{p}{24} \right) + 2 \exp \left(-\frac{p\theta}{12} \right).$$

Intuitively, Lemma 38 is true because, when \mathbf{X}_i and \mathbf{X}_k are similar, \mathbf{H}_i and \mathbf{H}_k (the i -th and k -th row of \mathbf{H} , respectively) will also likely be similar, i.e., $k\mathbf{H}_i - \mathbf{H}_k k_2$ is not likely to be large. Thus, we can construct \mathbf{a} such that $\mathbf{H}^T \mathbf{a}$ is proportional to $\mathbf{H}_i - \mathbf{H}_k$, which will lead to the result of Lemma 38. We put the proof of Lemma 38 in Appendix G.1.

The next step is to estimate such difference between \mathbf{X}_i and \mathbf{X}_k (or equivalently, the angle θ between them). We have the following lemma.

Lemma 39. *When $n \geq \pi(d-1)$, there must exist two different \mathbf{X}_i 's such that the angle between them is at most*

$$\theta = \pi \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right)^{\frac{1}{d-1}} n^{-\frac{1}{d-1}}.$$

Lemma 39 is intuitive because S^{d-1} has limited area. When there are many \mathbf{X}_i 's on S^{d-1} , there must exist at least two \mathbf{X}_i 's that are relatively close. We put the proof of Lemma 39 in Appendix G.2.

Finally, we have the following lemma.

Lemma 40. When $n \gg \pi(d-1)$, we have

$$\begin{aligned} \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \frac{\min \text{eig}(\mathbf{H}\mathbf{H}^T)}{p} \geq \frac{3\pi^2}{8} \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right)^{\frac{2}{d-1}} n^{-\frac{2}{d-1}} \right. \\ \left. + \frac{3}{4} \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right)^{\frac{1}{d-1}} n^{-\frac{1}{d-1}} \right\} \\ \geq 1 - 2 \exp\left(-\frac{p}{24}\right) - 2 \exp\left(-\frac{p}{12} \pi \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right)^{\frac{1}{d-1}} n^{-\frac{1}{d-1}}\right). \end{aligned}$$

Proof. This lemma directly follows by combining Lemma 26, Lemma 38, and Lemma 39. \square

By Lemma 40, we can conclude that when p is much larger than $n^{\frac{1}{d-1}}$, $\frac{\min \text{eig}(\mathbf{H}\mathbf{H}^T)}{p} = O(n^{-\frac{1}{d-1}})$ with high probability.

G.1. Proof of Lemma 38

We first prove a useful lemma.

Lemma 41. For any $\varphi \in [0, 2\pi]$, we must have $\sin \varphi \leq \varphi$. For any $\varphi \in [0, \pi/2]$, we must have $\varphi \leq \frac{\pi}{2} \sin \varphi$.

Proof. To prove the first part of the lemma, note that

$$\frac{d(\varphi - \sin \varphi)}{d\varphi} = 1 - \cos \varphi \geq 0.$$

Thus, the function $(\varphi - \sin \varphi)$ is monotone increasing with respect to $\varphi \in [0, 2\pi]$. Thus, we have

$$\min_{\varphi \in [0, 2\pi]} (\varphi - \sin \varphi) = (\varphi - \sin \varphi)|_{\varphi=0} = 0.$$

In other words, we have $\sin \varphi \leq \varphi$ for any $\varphi \in [0, 2\pi]$.

To prove the second part of the lemma, note that when $\varphi \in [0, \pi/2]$, we have

$$\frac{d^2(\varphi - \frac{\pi}{2} \sin \varphi)}{d\varphi^2} = \frac{\pi}{2} \sin \varphi \geq 0.$$

Thus, the function $\varphi - \frac{\pi}{2} \sin \varphi$ is convex with respect to $\varphi \in [0, \pi/2]$. Because the maximum of a convex function must be attained at the endpoint of the domain interval, we have

$$\max_{\varphi \in [0, \pi/2]} (\varphi - \frac{\pi}{2} \sin \varphi) = \max_{\varphi \in \{0, \pi/2\}} (\varphi - \frac{\pi}{2} \sin \varphi) = 0.$$

Thus, we have $\varphi \leq \frac{\pi}{2} \sin \varphi$ for any $\varphi \in [0, \pi/2]$. \square

Now we are ready to prove Lemma 38.

Proof. Through the proof, we fix \mathbf{X}_i and \mathbf{X}_k , and only consider the randomness of \mathbf{V}_0 . Because θ is the angle between \mathbf{X}_i and \mathbf{X}_k and because of Assumption 1, we have

$$\begin{aligned} \|\mathbf{X}_i - \mathbf{X}_k\|_2 &= 2 \sin \frac{\theta}{2} \\ &\geq 2 \frac{\theta}{2} \quad (\text{by Lemma 41}) \\ &= \theta. \end{aligned} \tag{78}$$

Let $\mathbf{a} = \frac{1}{\sqrt{2}}(\mathbf{e}_i \quad \mathbf{e}_k)$, where \mathbf{e}_q denotes the q -th standard basis vector, $q = 1, 2, \dots, n$. Then, we have

$$\begin{aligned}
 k\mathbf{H}^T \mathbf{a} k_2^2 &= \frac{1}{2} k\mathbf{H}_i^T \quad \mathbf{H}_k^T k_2^2 \\
 &= \frac{1}{2} \sum_{j=1}^p \left\| \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0g} \mathbf{X}_i \quad \mathbf{1}_{f(\mathbf{X}_k^T \mathbf{V}_0[j]) > 0g} \mathbf{X}_k \right\|_2^2 \quad (\text{by Eq. (1)}) \\
 &= \frac{1}{2} \sum_{j=1}^p \left(\mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0, \mathbf{X}_k^T \mathbf{V}_0[j] > 0g} k\mathbf{X}_i \quad \mathbf{X}_k k_2^2 + \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j])(\mathbf{X}_k^T \mathbf{V}_0[j]) < 0g} \right) \quad (\text{by } k\mathbf{X}_i k_2^2 = k\mathbf{X}_k k_2^2 = 1) \\
 &= \frac{1}{2} \sum_{j=1}^p \left(\mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0, \mathbf{X}_k^T \mathbf{V}_0[j] > 0g} \theta^2 + \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j])(\mathbf{X}_k^T \mathbf{V}_0[j]) < 0g} \right) \quad (\text{by Eq. (78)}) \\
 &= \frac{\theta^2}{2} \sum_{j=1}^p \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0g} + \frac{1}{2} \sum_{j=1}^p \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j])(\mathbf{X}_k^T \mathbf{V}_0[j]) < 0g}. \tag{79}
 \end{aligned}$$

Since \mathbf{X}_i is fixed and the direction of $\mathbf{V}_0[j]$ is uniformly distributed, we have $\Pr_{\mathbf{V}_0} f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0g = \frac{1}{2}$ and

$$\begin{aligned}
 \Pr_{\mathbf{V}_0} f(\mathbf{X}_i^T \mathbf{V}_0[j])(\mathbf{X}_k^T \mathbf{V}_0[j]) < 0g &= 2 \Pr_{\mathbf{V}_0} f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0, \mathbf{X}_k^T \mathbf{V}_0[j] < 0g \\
 &= 2 \Pr_{\mathbf{V}_0} f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0, \mathbf{X}_k^T \mathbf{V}_0[j] > 0g \\
 &= 2 \int_{S^{d-1}} \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{v}) > 0, \mathbf{X}_k^T \mathbf{v} > 0g} d\lambda(\mathbf{v}) \\
 &= 2 \frac{\pi}{2\pi} \left(\frac{\pi}{2} - \theta \right) \quad (\text{by Lemma 17}) \\
 &= \frac{\theta}{\pi}.
 \end{aligned}$$

Thus, based on the randomness of \mathbf{V}_0 , when \mathbf{X} are given, we have

$$\begin{aligned}
 \sum_{j=1}^p \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0g} &\sim \text{Bino} \left(p, \frac{1}{2} \right), \\
 \sum_{j=1}^p \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j])(\mathbf{X}_k^T \mathbf{V}_0[j]) < 0g} &\sim \text{Bino} \left(p, \frac{\theta}{\pi} \right).
 \end{aligned}$$

By letting $\delta = \frac{1}{2}$, $a = p$, $b = \frac{1}{2}$ in Lemma 14, we then have

$$\Pr_{\mathbf{V}_0} \left\{ \sum_{j=1}^p \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j]) > 0g} \leq \frac{3p}{4} \right\} \leq 2 \exp \left(-\frac{p}{24} \right), \tag{80}$$

$$\Pr_{\mathbf{V}_0} \left\{ \sum_{j=1}^p \mathbf{1}_{f(\mathbf{X}_i^T \mathbf{V}_0[j])(\mathbf{X}_k^T \mathbf{V}_0[j]) < 0g} \geq \frac{3p\theta}{2\pi} \right\} \leq 2 \exp \left(-\frac{p\theta}{12\pi} \right). \tag{81}$$

Thus, we have

$$\begin{aligned}
 & \Pr_{\mathbf{V}_0} \left\{ k\mathbf{H}^T \mathbf{a} k_2^2 \leq \frac{3p\theta^2}{8} + \frac{3p\theta}{4\pi} \right\} \\
 & \Pr_{\mathbf{V}_0} \left\{ \frac{\theta^2}{2} \sum_{j=1}^p \mathbf{1}_{\tilde{r}(\mathbf{x}_i^T \mathbf{v}_{0[j]}) > 0g} + \frac{1}{2} \sum_{j=1}^p \mathbf{1}_{\tilde{r}(\mathbf{x}_i^T \mathbf{v}_{0[j]}) (\mathbf{x}_k^T \mathbf{v}_{0[j]}) < 0g} \leq \frac{3p\theta^2}{8} + \frac{3p\theta}{4\pi} \right\} \\
 & \text{(by Eq. (79))} \\
 & \Pr_{\mathbf{V}_0} \left\{ \left\{ \sum_{j=1}^p \mathbf{1}_{\tilde{r}(\mathbf{x}_i^T \mathbf{v}_{0[j]}) > 0g} > \frac{3p}{4} \right\} \wedge \left\{ \sum_{j=1}^p \mathbf{1}_{\tilde{r}(\mathbf{x}_i^T \mathbf{v}_{0[j]}) (\mathbf{x}_k^T \mathbf{v}_{0[j]}) < 0g} \leq \frac{3p\theta}{2\pi} \right\} \right\} \\
 & \Pr_{\mathbf{V}_0} \left\{ \sum_{j=1}^p \mathbf{1}_{\tilde{r}(\mathbf{x}_i^T \mathbf{v}_{0[j]}) > 0g} > \frac{3p}{4} \right\} + \Pr_{\mathbf{V}_0} \left\{ \sum_{j=1}^p \mathbf{1}_{\tilde{r}(\mathbf{x}_i^T \mathbf{v}_{0[j]}) (\mathbf{x}_k^T \mathbf{v}_{0[j]}) < 0g} \leq \frac{3p\theta}{2\pi} \right\} \\
 & \text{(by the union bound)} \\
 & 2 \exp\left(-\frac{p}{24}\right) + 2 \exp\left(-\frac{p\theta}{12}\right) \text{ (by Eq. (80) and Eq. (81)).}
 \end{aligned}$$

The result of Lemma 38 thus follows. \square

G.2. Proof of Lemma 39

We first prove a useful lemma. Recall the definition of C_d in Eq. (60).

Lemma 42. *We have*

$$\frac{2^{\frac{D}{2}} \sqrt{d-1}}{nC_d} \geq \left[\frac{d-1}{n}, \frac{\pi(d-1)}{n} \right].$$

Proof. By Lemma 32 and Eq. (60), we have

$$C_d \geq \left[\frac{2^{\frac{D}{2}}}{\pi}, 2^{\frac{D}{2}} \sqrt{2d} \right].$$

Thus, we have

$$\frac{2^{\frac{D}{2}} \sqrt{d-1}}{nC_d} \geq \left[\frac{d-1}{n}, \frac{\pi(d-1)}{n} \right].$$

\square

Now we are ready to proof Lemma 39.

Proof. Recall the definition of θ in Lemma 39. Draw n caps on S^{d-1} centered at $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ with the colatitude angle φ where

$$\varphi = \frac{\theta}{2} = \frac{\pi}{2} \left(\frac{2^{\frac{D}{2}} \sqrt{d-1}}{nC_d} \right)^{\frac{1}{d-1}} \text{ (by Eq. (60)).} \quad (82)$$

By Lemma 42 and $n \geq \pi(d-1)$, we have $\varphi \geq [0, \pi/2]$. Thus, by Lemma 41, we have

$$\sin \varphi \geq \frac{2\varphi}{\pi} = \left(\frac{2^{\frac{D}{2}} \sqrt{d-1}}{nC_d} \right)^{\frac{1}{d-1}}. \quad (83)$$

By Lemma 8, the area of each cap is

$$A = \frac{1}{2} \lambda_{d-1}(S^{d-1}) I_{\sin^2 \varphi} \left(\frac{d-1}{2}, \frac{1}{2} \right).$$

Applying Lemma 35 and Eq. (83), we thus have

$$A \frac{1}{2} \lambda_{d-1}(S^{d-1}) \frac{C_d}{2(d-1)} (\sin^2 \varphi)^{\frac{d-1}{2}} = \frac{1}{n} \lambda_{d-1}(S^{d-1}).$$

In other words, we have

$$\frac{\lambda_{d-1}(S^{d-1})}{A} = n.$$

By the pigeonhole principle, we know there exist at least two different caps that overlap, i.e., the angle between them is at most 2φ . The result of this lemma thus follows by Eq. (82). \square

H. Proof of Proposition 5

We follow the sketch of proof in Section 5. Recall the definition of the pseudo ground-truth function $f_{\mathbf{V}_0}^g$ in Definition 2, and the corresponding $\mathbf{V} \in \mathbb{R}^{dp}$ that

$$\mathbf{V}[j] = \int_{S^{d-1}} \mathbf{1}_{\tilde{r}_{\mathbf{z}^T \mathbf{V}_0[j]} > 0g} \mathbf{z} \frac{g(\mathbf{z})}{p} d\mu(\mathbf{z}), \text{ for all } j \in \{1, 2, \dots, p\}. \quad (84)$$

We first show that the pseudo ground-truth can be written in a linear form.

Lemma 43. $\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V} = f_{\mathbf{V}_0}^g(\mathbf{x})$ for all $\mathbf{x} \in S^{d-1}$.

Proof. For all $\mathbf{x} \in S^{d-1}$, we have

$$\begin{aligned} \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V} &= \sum_{j=1}^p \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}[j] \mathbf{V}[j] \\ &= \sum_{j=1}^p \mathbf{1}_{\tilde{r}_{\mathbf{x}^T \mathbf{V}_0[j]} > 0g} \mathbf{x}^T \int_{S^{d-1}} \mathbf{1}_{\tilde{r}_{\mathbf{z}^T \mathbf{V}_0[j]} > 0g} \mathbf{z} \frac{g(\mathbf{z})}{p} d\mu(\mathbf{z}) \quad (\text{by Eq. (1) and Eq. (84)}) \\ &= \int_{S^{d-1}} \sum_{j=1}^p \mathbf{1}_{\tilde{r}_{\mathbf{x}^T \mathbf{V}_0[j]} > 0g} \mathbf{x}^T \mathbf{1}_{\tilde{r}_{\mathbf{z}^T \mathbf{V}_0[j]} > 0g} \mathbf{z} \frac{g(\mathbf{z})}{p} d\mu(\mathbf{z}) \\ &= \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{jC_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0}}{p} g(\mathbf{z}) d\mu(\mathbf{z}) \quad (\text{by Eq. (6)}) \\ &= f_{\mathbf{V}_0}^g(\mathbf{x}) \quad (\text{by Definition 2}). \end{aligned}$$

\square

Let $\mathbf{P} := \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}\mathbf{H}$. Since $\mathbf{P}^2 = \mathbf{P}$ and $\mathbf{P} = \mathbf{P}^T$, we know that \mathbf{P} is an orthogonal projection to the row-space of \mathbf{H} . Next, we give an expression for the test error. Note that even though Proposition 4 assumes no noise, below we state a more general version below with noise (which will be useful later).

Lemma 44. *If the ground-truth is $f(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V}$ for all \mathbf{x} , then we have*

$$\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} - \mathbf{I}) \mathbf{V} + \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1} \boldsymbol{\epsilon}, \text{ for all } \mathbf{x}.$$

Proof. Because $f(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V}$, we have $\mathbf{y} = \mathbf{H} \mathbf{V} + \boldsymbol{\epsilon}$. Thus, we have

$$\begin{aligned} \mathbf{V}^{\ell_2} &= \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y} \quad (\text{by Eq. (3)}) \\ &= \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1} (\mathbf{H} \mathbf{V} + \boldsymbol{\epsilon}). \end{aligned}$$

Further, we have

$$\begin{aligned} \mathbf{V}^{\ell_2} \quad \mathbf{V} &= (\mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}\mathbf{H} \quad \mathbf{I}) \mathbf{V} + \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}\boldsymbol{\epsilon} \\ &= (\mathbf{P} \quad \mathbf{I}) \mathbf{V} + \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}\boldsymbol{\epsilon}. \end{aligned}$$

Finally, using Eq. (4), we have

$$\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V}^{\ell_2} - \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V} = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} \quad \mathbf{I}) \mathbf{V} + \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}\mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}\boldsymbol{\epsilon}.$$

□

When there is no noise, Lemma 44 reduces to $\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} \quad \mathbf{I}) \mathbf{V}$. As we described in Section 5, $(\mathbf{P} \quad \mathbf{I}) \mathbf{V}$ has the interpretation of the distance from \mathbf{V} to the row-space of \mathbf{H} . We then have the following.

Lemma 45. For all $\mathbf{a} \succeq \mathbb{R}^n$, we have

$$j\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} \quad \mathbf{I}) \mathbf{V} j \stackrel{\rho_{\bar{p}}}{\leq} k \mathbf{V} \quad \mathbf{H}\mathbf{a} k_2.$$

Proof. Recall that $\mathbf{P} = \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}\mathbf{H}$. Thus, we have

$$\mathbf{P}\mathbf{H}^T = \mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}\mathbf{H}\mathbf{H}^T = \mathbf{H}^T. \quad (85)$$

We then have

$$\begin{aligned} k(\mathbf{P} \quad \mathbf{I}) \mathbf{V} k_2 &= k\mathbf{P} \mathbf{V} \quad \mathbf{V} k_2 \\ &= k\mathbf{P}(\mathbf{H}^T\mathbf{a} + \mathbf{V} \quad \mathbf{H}^T\mathbf{a}) \quad \mathbf{V} k_2 \\ &= k\mathbf{P}\mathbf{H}^T\mathbf{a} + \mathbf{P}(\mathbf{V} \quad \mathbf{H}^T\mathbf{a}) \quad \mathbf{V} k_2 \\ &= k\mathbf{H}^T\mathbf{a} + \mathbf{P}(\mathbf{V} \quad \mathbf{H}^T\mathbf{a}) \quad \mathbf{V} k_2 \text{ (by Eq. (85))} \\ &= k(\mathbf{P} \quad \mathbf{I})(\mathbf{V} \quad \mathbf{H}^T\mathbf{a}) k_2 \\ &= k \mathbf{V} \quad \mathbf{H}^T\mathbf{a} k_2 \text{ (because } \mathbf{P} \text{ is an orthogonal projection)}. \end{aligned}$$

Therefore, we have

$$\begin{aligned} j\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} \quad \mathbf{I}) \mathbf{V} j &= k\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} \quad \mathbf{I}) \mathbf{V} k_2 \\ &= k\mathbf{h}_{\mathbf{V}_0, \mathbf{x}} k_2 \quad k(\mathbf{P} \quad \mathbf{I}) \mathbf{V} k_2 \text{ (by Lemma 12)} \\ &= \stackrel{\rho_{\bar{p}}}{\leq} k \mathbf{V} \quad \mathbf{H}\mathbf{a} k_2 \text{ (by Lemma 11)}. \end{aligned}$$

□

Now we are ready to prove Proposition 5.

Proof. Because there is no noise, we have $\boldsymbol{\epsilon} = \mathbf{0}$. Thus, by Lemma 44, we have

$$\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) = \mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} \quad \mathbf{I}) \mathbf{V}. \quad (86)$$

We then have, for all $\mathbf{a} \succeq \mathbb{R}^n$,

$$\begin{aligned} &\Pr_{\mathbf{X}} \left\{ \left| \hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) \right| \geq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \\ &= \Pr_{\mathbf{X}} \left\{ j\mathbf{h}_{\mathbf{V}_0, \mathbf{x}}(\mathbf{P} \quad \mathbf{I}) \mathbf{V} j \geq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \\ &= \Pr_{\mathbf{X}} \left\{ \stackrel{\rho_{\bar{p}}}{\leq} k \mathbf{V} \quad \mathbf{H}\mathbf{a} k_2 \geq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \text{ (by Lemma 45)}. \end{aligned} \quad (87)$$

It only remains to find the vector \mathbf{a} . Define $\mathbf{K}_i \in \mathbb{R}^{dp}$ for $i = 1, 2, \dots, n$ as

$$\mathbf{K}_i[j] := \mathbf{1}_{\langle \mathbf{X}_i^T \mathbf{v}_0[j] \rangle > 0} g(\mathbf{X}_i) \frac{g(\mathbf{X}_i)}{p}, \quad j = 1, 2, \dots, p.$$

By Eq. (84), for all $j = 1, 2, \dots, p$, we have

$$\mathbb{E}_{\mathbf{X}_i} [\mathbf{K}_i[j]] = \mathbf{V}[j]. \quad (88)$$

Because $k_{\mathbf{X}_i k_2} = 1$, we have

$$k_{\mathbf{K}_i[j] k_2} = \frac{kgk_1}{p}.$$

Thus, we have

$$k_{\mathbf{K}_i k_2} = \sqrt{\sum_{j=1}^p k_{\mathbf{K}_i[j] k_2}^2} = \frac{kgk_1}{p},$$

i.e.,

$$\rho_{\bar{p}} k_{\mathbf{K}_i k_2} = kgk_1. \quad (89)$$

We now construct the vector \mathbf{a} . Define $\mathbf{a} \in \mathbb{R}^n$ whose i -th element is $\mathbf{a}_i = \frac{g(\mathbf{X}_i)}{np}$, $i = 1, 2, \dots, n$. Notice that \mathbf{a} is well-defined because $kgk_1 < 1$. Then, for all $j \in \{1, 2, \dots, p\}$, we have

$$\begin{aligned} (\mathbf{H}^T \mathbf{a})[j] &= \sum_{i=1}^n \mathbf{H}_i^T[j] \mathbf{a}_i \\ &= \sum_{i=1}^n \mathbf{1}_{\langle \mathbf{X}_i^T \mathbf{v}_0[j] \rangle > 0} g(\mathbf{X}_i) \frac{g(\mathbf{X}_i)}{np} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{K}_i[j], \end{aligned}$$

i.e.,

$$\mathbf{H}^T \mathbf{a} = \frac{1}{n} \sum_{i=1}^n \mathbf{K}_i. \quad (90)$$

Thus, by Eq. (89) and Lemma 16 (with $X_i = \rho_{\bar{p}} \mathbf{K}_i$, $U = kgk_1$, and $k = n$), we have

$$\Pr_{\mathbf{X}} \left\{ \rho_{\bar{p}} \left\| \left(\frac{1}{n} \sum_{i=1}^n \mathbf{K}_i \right) - \mathbb{E}_{\mathbf{X}} \mathbf{K}_1 \right\|_2 \leq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \geq 2e^2 \exp \left(-\frac{\rho_{\bar{p}}^2 n}{8kgk_1^2} \right).$$

Further, by Eq. (90) and Eq. (88), we have

$$\Pr_{\mathbf{X}} \left\{ \rho_{\bar{p}} k_{\mathbf{H}^T \mathbf{a}} \leq \mathbf{V}[k_2] \leq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \geq 2e^2 \exp \left(-\frac{\rho_{\bar{p}}^2 n}{8kgk_1^2} \right). \quad (91)$$

Plugging Eq. (91) into Eq. (87), we thus have

$$\Pr_{\mathbf{X}} \left\{ \left| \hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) \right| \leq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \geq 2e^2 \exp \left(-\frac{\rho_{\bar{p}}^2 n}{8kgk_1^2} \right).$$

□

I. Proof of Theorem 1

We first prove a useful lemma.

Lemma 46. *If $kgk_1 < 1$, then for any \mathbf{x} , we must have*

$$\int_{S^{d-1}} \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v} > 0}, \mathbf{x}^T \mathbf{v} > 0} g(\mathbf{z}) d\mu(\mathbf{z}) d\lambda(\mathbf{v}) = \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi - \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} g(\mathbf{z}) d\mu(\mathbf{z}).$$

Proof. This follows from Fubini's Theorem and by a change of order of the integral. Specifically, because $kgk_1 < 1$, we have

$$\int_{S^{d-1}} jg(\mathbf{z}) j d\mu(\mathbf{z}) < 1.$$

Thus, we have

$$\int_{S^{d-1}} \int_{S^{d-1}} jg(\mathbf{z}) j d\mu(\mathbf{z}) \lambda(\mathbf{v}) < 1.$$

Because $|\mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v} > 0}, \mathbf{x}^T \mathbf{v} > 0}| \leq 1$ when $\mathbf{x} \in S^{d-1}$ and $\mathbf{z} \in S^{d-1}$, we have

$$\int_{S^{d-1}} \int_{S^{d-1}} |\mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v} > 0}, \mathbf{x}^T \mathbf{v} > 0} g(\mathbf{z})| d\mu(\mathbf{z}) \lambda(\mathbf{v}) \leq \int_{S^{d-1}} \int_{S^{d-1}} jg(\mathbf{z}) j d\mu(\mathbf{z}) \lambda(\mathbf{v}) < 1.$$

Thus, by Fubini's theorem, we can exchange the sequence of integral, i.e., we have

$$\begin{aligned} & \int_{S^{d-1}} \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v} > 0}, \mathbf{x}^T \mathbf{v} > 0} g(\mathbf{z}) d\mu(\mathbf{z}) d\lambda(\mathbf{v}) \\ &= \int_{S^{d-1}} \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v} > 0}, \mathbf{x}^T \mathbf{v} > 0} g(\mathbf{z}) d\lambda(\mathbf{v}) d\mu(\mathbf{z}) \\ &= \int_{S^{d-1}} \left(\int_{S^{d-1}} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v} > 0}, \mathbf{x}^T \mathbf{v} > 0} d\lambda(\mathbf{v}) \right) \mathbf{x}^T \mathbf{z} g(\mathbf{z}) d\mu(\mathbf{z}) \\ &= \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi - \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} g(\mathbf{z}) d\mu(\mathbf{z}) \text{ (by Lemma 17)}. \end{aligned}$$

□

The following proposition characterizes generalization performance when $\epsilon = \mathbf{0}$, i.e., the bias term in Eq. (18).

Proposition 47. *Assume no noise ($\epsilon = \mathbf{0}$), a ground truth $f = f_g \in F^{\ell_2}$ where $kgk_1 < 1$, $n \geq 2$, $m \geq \left[1, \frac{\ln n}{\ln \frac{n}{2}}\right]$, $d \leq n^4$, and $p \leq 6J_m(n, d) \ln \left(4n^{1+\frac{1}{m}}\right)$. Then, for any $q \geq [1, 1)$ and for almost every $\mathbf{x} \in S^{d-1}$, we must have*

$$\begin{aligned} & \Pr_{\mathbf{v}_0, \mathbf{X}} \left\{ |j\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x})j| \leq n^{-\frac{1}{2}(1-\frac{1}{q})} \right. \\ & \left. + \left(1 + \sqrt{J_m(n, d)n}\right) p^{-\frac{1}{2}(1-\frac{1}{q})} \right\} \\ & \leq 2e^2 \left(\exp\left(\frac{p_q \bar{n}}{8kgk_1^2}\right) + \exp\left(\frac{p_q \bar{p}}{8kgk_1^2}\right) \right. \\ & \left. + \exp\left(\frac{p_q \bar{p}}{8nkgk_1^2}\right) \right) + \frac{2}{n}. \end{aligned}$$

Proof. We split the whole proof into 5 steps as follows.

Step 1: use pseudo ground-truth as a ‘‘intermediary’’

Recall Definition 2 where we define the pseudo ground-truth $f_{\mathbf{V}_0}^g$. We then define the output of the pseudo ground-truth for training input as

$$\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) := [f_{\mathbf{V}_0}^g(\mathbf{X}_1) \ f_{\mathbf{V}_0}^g(\mathbf{X}_2) \ \dots \ f_{\mathbf{V}_0}^g(\mathbf{X}_n)]^T.$$

The rest of the proof will use the pseudo ground-truth as a ‘‘intermediary’’ to connect the ground-truth f and the model output \hat{f}^{ℓ_2} . Specifically, we have

$$\begin{aligned} \hat{f}^{\ell_2}(\mathbf{x}) &= \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{V}^{\ell_2} \\ &= \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{F}(\mathbf{X}) \text{ (by Eq. (17) and } \epsilon = \mathbf{0}) \\ &= \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) + \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} (\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})). \end{aligned} \quad (92)$$

Thus, we have

$$\begin{aligned} & | \hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) | \\ &= \left| \hat{f}^{\ell_2}(\mathbf{x}) - f_{\mathbf{V}_0}^g(\mathbf{x}) + f_{\mathbf{V}_0}^g(\mathbf{x}) - f(\mathbf{x}) \right| \\ &= \left| \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - f_{\mathbf{V}_0}^g(\mathbf{x}) + \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} (\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})) \right. \\ &\quad \left. + f_{\mathbf{V}_0}^g(\mathbf{x}) - f(\mathbf{x}) \right| \text{ (by Eq. (92))} \\ &= \underbrace{\left| \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} \mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - f_{\mathbf{V}_0}^g(\mathbf{x}) \right|}_{\text{term A}} + \underbrace{\left| \mathbf{h}_{\mathbf{V}_0, \mathbf{x}} \mathbf{H}^T (\mathbf{H} \mathbf{H}^T)^{-1} (\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})) \right|}_{\text{term B}} \\ &\quad + \underbrace{\left| f_{\mathbf{V}_0}^g(\mathbf{x}) - f(\mathbf{x}) \right|}_{\text{term C}}. \end{aligned} \quad (93)$$

In Eq. (93), we can see that the term A corresponds to the test error of the pseudo ground-truth, the term B corresponds to the impact of the difference between the pseudo ground-truth and the real ground-truth in the training data, and the term C corresponds to the impact of the difference between pseudo ground-truth and real ground-truth in the test data. Using the terminology of bias-variance decomposition, we refer to term A as the ‘‘pseudo bias’’ and term B as the ‘‘pseudo variance’’.

Step 2: estimate term A

We have

$$\begin{aligned} \Pr_{\mathbf{x}, \mathbf{V}_0} \left\{ \text{term A} \geq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} &= \int_{\mathbf{V}_0 \in \mathbb{R}^{d_p}} \Pr_{\mathbf{X}} \left\{ \text{term A} \geq n^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \mid \mathbf{V}_0 \right\} d\lambda(\mathbf{V}_0) \\ &= \int_{\mathbf{V}_0 \in \mathbb{R}^{d_p}} 2e^2 \exp\left(-\frac{P_q \bar{n}}{8kgk_\gamma^2}\right) d\lambda(\mathbf{V}_0) \text{ (by Proposition 5)} \\ &= 2e^2 \exp\left(-\frac{P_q \bar{n}}{8kgk_\gamma^2}\right). \end{aligned} \quad (94)$$

Step 3: estimate term C

For all $j = 1, 2, \dots, p$, define

$$K_j^{\mathbf{x}} := \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{V}_0[j]} > 0, \mathbf{x}^T \mathbf{V}_0[j] > 0} g(\mathbf{z}) d\mu(\mathbf{z}).$$

We now show that $K_j^{\mathbf{x}}$ is bounded and with mean equal to f_g , where $f_g = \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi - \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} g(\mathbf{z}) d\mu(\mathbf{z})$ defined by Definition 1. Specifically, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{V}_0} K_j^{\mathbf{x}} &= \mathbb{E}_{\mathbf{v}} \left(\int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v}} > 0, \mathbf{x}^T \mathbf{v} > 0} g(\mathbf{z}) d\mu(\mathbf{z}) \right) \\ &= \int_{S^{d-1}} \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{f_{\mathbf{z}^T \mathbf{v}} > 0, \mathbf{x}^T \mathbf{v} > 0} g(\mathbf{z}) d\mu(\mathbf{z}) d\lambda(\mathbf{v}) \\ &= \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi - \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} g(\mathbf{z}) d\mu(\mathbf{z}) \text{ (by Lemma 46)} \\ &= f_g(\mathbf{x}) \text{ (by Definition 1)}. \end{aligned} \quad (95)$$

From Definition 2, we have

$$\begin{aligned}
 f_{\mathbf{V}_0}^g(\mathbf{x}) &= \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{j C_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0}}{p} g(\mathbf{z}) d\mu(\mathbf{z}) \text{ (by Definition 2)} \\
 &= \frac{1}{p} \sum_{j=1}^p \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{\mathbf{f}_{\mathbf{z}^T \mathbf{V}_0[j]} > 0, \mathbf{x}^T \mathbf{V}_0[j] > 0} g(\mathbf{z}) d\mu(\mathbf{z}) \text{ (by Eq. (6))} \\
 &= \frac{1}{p} \sum_{j=1}^p K_j^{\mathbf{x}}.
 \end{aligned} \tag{96}$$

Because $\mathbf{V}_0[j]$'s are *i.i.d.*, $K_j^{\mathbf{x}}$'s are also *i.i.d.*. Thus, we have

$$\mathbb{E}_{\mathbf{V}_0} f_{\mathbf{V}_0}^g(\mathbf{x}) = f_g(\mathbf{x}). \tag{97}$$

Further, for any $j \geq 1, 2, \dots, p$, we have

$$\begin{aligned}
 k K_j^{\mathbf{x}} k_2 &= j K_j^{\mathbf{x}} j \text{ (because } K_j^{\mathbf{x}} \text{ is a scalar)} \\
 &= \left| \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \mathbf{1}_{\mathbf{f}_{\mathbf{z}^T \mathbf{V}_0[j]} > 0, \mathbf{x}^T \mathbf{V}_0[j] > 0} g(\mathbf{z}) d\mu(\mathbf{z}) \right| \\
 &\quad \int_{S^{d-1}} \left| \mathbf{x}^T \mathbf{z} \mathbf{1}_{\mathbf{f}_{\mathbf{z}^T \mathbf{V}_0[j]} > 0, \mathbf{x}^T \mathbf{V}_0[j] > 0} g(\mathbf{z}) \right| d\mu(\mathbf{z}) \\
 &\quad \int_{S^{d-1}} \left| \mathbf{x}^T \mathbf{z} \mathbf{1}_{\mathbf{f}_{\mathbf{z}^T \mathbf{V}_0[j]} > 0, \mathbf{x}^T \mathbf{V}_0[j] > 0} \right| j g(\mathbf{z}) j d\mu(\mathbf{z}) \\
 &\quad \int_{S^{d-1}} j g(\mathbf{z}) j d\mu(\mathbf{z}) \\
 &= k g k_1.
 \end{aligned} \tag{98}$$

Thus, by Lemma 16, we have

$$\Pr_{\mathbf{V}_0} \left\{ \left\| \left(\frac{1}{p} \sum_{j=1}^p K_j^{\mathbf{x}} \right) - \mathbb{E}_{\mathbf{V}_0} K_1 \right\|_2 \geq p^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \leq 2e^2 \exp\left(-\frac{p}{8kgk_1^2}\right).$$

Further, by Eq. (96) and Eq. (95), we have

$$\Pr_{\mathbf{V}_0} \left\{ \left| f_{\mathbf{V}_0}^g(\mathbf{x}) - f_g(\mathbf{x}) \right| \geq p^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \leq 2e^2 \exp\left(-\frac{p}{8kgk_1^2}\right).$$

Because $f \stackrel{\text{a.e.}}{=} f_g$, we have

$$\Pr_{\mathbf{V}_0} \left\{ \left| f_{\mathbf{V}_0}^g(\mathbf{x}) - f(\mathbf{x}) \right| \geq p^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \leq 2e^2 \exp\left(-\frac{p}{8kgk_1^2}\right).$$

Because $f_{\mathbf{V}_0}^g$ does not change with \mathbf{X} , we thus have

$$\Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term C} \geq p^{-\frac{1}{2}} \left(1 - \frac{1}{q}\right) \right\} \leq 2e^2 \exp\left(-\frac{p}{8kgk_1^2}\right). \tag{99}$$

Step 4: estimate term B

Our idea is to treat $\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})$ as a special form of noise, and then apply Proposition 4. We first bound the magnitude of this special noise. For $j = 1, 2, \dots, p$, we define

$$\mathbf{K}_j := [K_j^{\mathbf{X}_1} \ K_j^{\mathbf{X}_2} \ \dots \ K_j^{\mathbf{X}_n}]^T.$$

Then, we have

$$k\mathbf{K}_j k_2 = \sqrt{\sum_{i=1}^n kK_j^{\mathbf{X}_i} k_2^2} \stackrel{\rho_q}{\approx} \sqrt{nk}gk_1 \text{ (by Eq. (98)).}$$

Similar to how we get Eq. (99) in Step 3, we have

$$\Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \|\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_2 \leq p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \geq 2e^2 \exp\left(-\frac{\rho_q}{8nkgk_1^2}\right). \quad (100)$$

Thus, we have

$$\begin{aligned} & \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term B} \leq \sqrt{J_m(n, d)np} p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \\ = & \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term B} \leq \sqrt{J_m(n, d)np} p^{\frac{1}{2}(1 - \frac{1}{q})}, \|\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_2 \leq p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \\ & + \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term B} \leq \sqrt{J_m(n, d)np} p^{\frac{1}{2}(1 - \frac{1}{q})}, \|\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_2 < p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \\ & \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \|\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_2 \leq p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \\ & + \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term B} \leq \sqrt{J_m(n, d)n} \|\mathbf{F}_{\mathbf{V}_0}^g(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_2 \right\} \\ & 2e^2 \exp\left(-\frac{\rho_q}{8nkgk_1^2}\right) + \frac{2}{m} \stackrel{\rho_q}{\approx} \frac{2}{n} \text{ (by Eq. (100) and Proposition 4).} \end{aligned} \quad (101)$$

Step 5: estimate $\int j\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x})j$

We have

$$\begin{aligned} & \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \int j\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x})j \leq n^{\frac{1}{2}(1 - \frac{1}{q})} + \frac{1 + \sqrt{J_m(n, d)n}}{\rho_q} \right\} \\ & \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term A} + \text{term B} + \text{term C} \leq n^{\frac{1}{2}(1 - \frac{1}{q})} + \frac{1 + \sqrt{J_m(n, d)n}}{\rho_q} \right\} \text{ (by Eq. (93))} \\ & \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ \left\{ \text{term A} \leq n^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \cap \left\{ \text{term B} \leq \sqrt{J_m(n, d)np} p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \right. \\ & \left. \cap \left\{ \text{term C} \leq p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \right\} \\ & \Pr_{\mathbf{X}, \mathbf{V}_0} \left\{ \text{term A} \leq n^{\frac{1}{2}(1 - \frac{1}{q})} \right\} + \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term B} \leq \sqrt{J_m(n, d)np} p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \\ & + \Pr_{\mathbf{V}_0, \mathbf{X}} \left\{ \text{term C} \leq p^{\frac{1}{2}(1 - \frac{1}{q})} \right\} \text{ (by the union bound)} \\ & 2e^2 \left(\exp\left(-\frac{\rho_q}{8kgk_7^2}\right) + \exp\left(-\frac{\rho_q}{8kgk_1^2}\right) + \exp\left(-\frac{\rho_q}{8nkgk_1^2}\right) \right) + \frac{2}{m} \stackrel{\rho_q}{\approx} \frac{2}{n} \\ & \text{(by Eqs. (94)(99)(101)).} \end{aligned}$$

The last step exactly gives the conclusion of this proposition. □

Theorem 1 thus follows by Proposition 4, Proposition 47, Eq. (18), and the union bound.

J. Proof of Proposition 2 (lower bound for ground-truth functions outside $\overline{F^{\ell_2}}$)

We first show what $\hat{f}_7^{\ell_2}$ looks like. Define $\mathbf{H}^7 \in \mathbb{R}^{n \times n}$ where its (i, j) -th element is

$$\mathbf{H}_{i,j}^7 = \mathbf{X}_i^T \mathbf{X}_j \frac{\pi - \arccos(\mathbf{X}_i^T \mathbf{X}_j)}{2\pi}.$$

Notice that

$$\left(\frac{\mathbf{H}\mathbf{H}^T}{p}\right)_{i,j} = \frac{1}{p} \sum_{k=1}^p \mathbf{X}_i^T \mathbf{X}_j \mathbf{1}_{\mathbf{x}_i^T \mathbf{v}_{o[k]} > 0, \mathbf{x}_j^T \mathbf{v}_{o[k]} > 0} = \mathbf{X}_i^T \mathbf{X}_j \frac{jC_{\mathbf{x}_i, \mathbf{x}_j}^{\mathbf{v}_o}}{p}.$$

By Lemma 21, we have that $\left(\frac{\mathbf{H}\mathbf{H}^T}{p}\right)_{i,j}$ converges in probability to $(\mathbf{H}^T)_{i,j}$ as $p \rightarrow \infty$ uniformly in i, j . In other words,

$$\max_{i,j} \left| \left(\frac{\mathbf{H}\mathbf{H}^T}{p}\right)_{i,j} - (\mathbf{H}^T)_{i,j} \right| \xrightarrow{p \rightarrow \infty} 0, \quad (102)$$

Let $\mathbf{e}_i, i = 1, \dots, n$ denote the standard basis in \mathbb{R}^n . For $i = 1, 2, \dots, n$, define

$$g_{i,p} := np \mathbf{e}_i^T (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y}, \quad (103)$$

which is a number. Further, define

$$[g_{1,p} \ g_{2,p} \ \dots \ g_{n,p}]^T = np (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y}.$$

Further, define the number

$$g_{i,1} := n \mathbf{e}_i^T (\mathbf{H}^T)^{-1} \mathbf{y},$$

and

$$[g_{1,1} \ g_{2,1} \ \dots \ g_{n,1}]^T = n (\mathbf{H}^T)^{-1} \mathbf{y}.$$

Notice that $(\mathbf{H}^T)^{-1}$ exists because of Eq. (102) and Lemma 7.

By Eq. (102), we have

$$\max_{i \in \{1,2,\dots,n\}} |g_{i,p} - g_{i,1}| \xrightarrow{p \rightarrow \infty} 0. \quad (104)$$

For any given \mathbf{X} , we define $\hat{f}_1^{\ell_2}(\cdot) : S^{d-1} \rightarrow \mathbb{R}$ as

$$\hat{f}_1^{\ell_2}(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n \mathbf{x}^T \mathbf{X}_i \frac{\pi - \arccos(\mathbf{x}^T \mathbf{X}_i)}{2\pi} g_{i,1}. \quad (105)$$

By the definition of the Dirac delta function $\delta_a(\cdot)$ with peak position at a , we can write $\hat{f}_1^{\ell_2}(\mathbf{x})$ as an integral

$$\hat{f}_1^{\ell_2}(\mathbf{x}) = \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi - \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} \frac{1}{n} \sum_{i=1}^n g_{i,1} \delta_{\mathbf{X}_i}(\mathbf{z}) d\mu(\mathbf{z}).$$

Notice that $g_{i,1}$ only depends on the training data and does not change with p (and thus is finite). Therefore, we have $\hat{f}_1^{\ell_2} \in F^{\ell_2}$. It remains to show why $\hat{f}_1^{\ell_2}$ converges to $f_1^{\ell_2}$ in probability. The following lemma shows what $f_1^{\ell_2}$ looks like.

Lemma 48. $f_1^{\ell_2}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^T \mathbf{X}_i \frac{jC_{\mathbf{x}_i, \mathbf{x}}^{\mathbf{v}_o}}{p} g_{i,p} = \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{jC_{\mathbf{x}, \mathbf{z}}^{\mathbf{v}_o}}{p} \frac{1}{n} \sum_{i=1}^n g_{i,p} \delta_{\mathbf{X}_i}(\mathbf{z}) d\mu(\mathbf{z}).$

Proof. For any $\mathbf{x} \in S^{d-1}$, we have

$$\begin{aligned} \hat{f}_1^{\ell_2}(\mathbf{x}) &= \mathbf{h}_{\mathbf{v}_o, \mathbf{x}}^T \mathbf{V}^{\ell_2} \\ &= \mathbf{h}_{\mathbf{v}_o, \mathbf{x}}^T \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y} \quad (\text{by Eq. (3)}) \\ &= \mathbf{h}_{\mathbf{v}_o, \mathbf{x}}^T \sum_{i=1}^n \mathbf{H}_i^T \mathbf{e}_i^T (\mathbf{H}\mathbf{H}^T)^{-1} \mathbf{y} \\ &= \frac{1}{np} \sum_{i=1}^n \mathbf{h}_{\mathbf{v}_o, \mathbf{x}}^T \mathbf{H}_i^T g_{i,p} \quad (\text{by Eq. (103)}) \\ &= \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \mathbf{x}^T \mathbf{X}_i \mathbf{1}_{\mathbf{x}_i^T \mathbf{v}_{o[j]} > 0, \mathbf{x}^T \mathbf{v}_{o[j]} > 0} g_{i,p}. \end{aligned}$$

By Eq. (6), we thus have

$$\hat{f}^{\ell_2}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^T \mathbf{X}_i \frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} g_{i,p}. \quad (106)$$

By the definition of the Dirac delta function, we have

$$\hat{f}^{\ell_2}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^T \mathbf{X}_i \frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} g_{i,p} = \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{jC_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0}}{p} \frac{1}{n} \sum_{i=1}^n g_{i,p} \delta_{\mathbf{X}_i}(\mathbf{z}) d\mu(\mathbf{z}).$$

□

Now we are ready to prove the statement of Proposition 2, i.e., uniformly over all $\mathbf{x} \in S^{d-1}$, $\hat{f}^{\ell_2}(\mathbf{x}) \xrightarrow{P} \hat{f}_7^{\ell_2}(\mathbf{x})$ as $p \rightarrow \infty$ (notice that we have already shown that $\hat{f}_7^{\ell_2} \in F^{\ell_2}$). To be more specific, we restate that uniform convergence as the following lemma.

Lemma 49. For any given \mathbf{X} , $\sup_{\mathbf{x} \in S^{d-1}} | \hat{f}^{\ell_2}(\mathbf{x}) - \hat{f}_7^{\ell_2}(\mathbf{x}) | \xrightarrow{P} 0$ as $p \rightarrow \infty$.

Proof. For any $\zeta > 0$, define two events:

$$\mathcal{J}_1 := \left\{ \sup_{\mathbf{x}, \mathbf{z} \in S^{d-1}} \left| \frac{jC_{\mathbf{z}, \mathbf{x}}^{\mathbf{V}_0}}{p} - \frac{\pi \arccos(\mathbf{x}^T \mathbf{z})}{2\pi} \right| < \zeta \right\},$$

$$\mathcal{J}_2 := \left\{ \max_{i \in \{1, 2, \dots, n\}} |jg_{i,p} - g_{i,1}| < \zeta \right\}.$$

By Lemma 21, there exists a threshold p_0 such that for any $p > p_0$,

$$\Pr[\mathcal{J}_1] > 1 - \zeta.$$

By Eq. (104), there exists a threshold p_1 such that for any $p > p_1$,

$$\Pr[\mathcal{J}_2] > 1 - \zeta.$$

Thus, by the union bound, when $p > \max\{p_0, p_1\}$, we have

$$\Pr[\mathcal{J}_1 \cap \mathcal{J}_2] > 1 - 2\zeta. \quad (107)$$

When $\mathcal{J}_1 \cap \mathcal{J}_2$ happens, we have

$$\begin{aligned} & \sup_{\mathbf{x} \in S^{d-1}} | \hat{f}^{\ell_2}(\mathbf{x}) - \hat{f}_7^{\ell_2}(\mathbf{x}) | \\ &= \sup_{\mathbf{x} \in S^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{x}^T \mathbf{X}_i \left(\frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} g_{i,p} - \frac{\pi \arccos(\mathbf{x}^T \mathbf{X}_i)}{2\pi} g_{i,1} \right) \right| \\ & \quad (\text{by Lemma 48 and Eq. (105)}) \\ &= \sup_{\mathbf{x} \in S^{d-1}, i \in \{1, 2, \dots, n\}} \left| \left(\frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} g_{i,p} - \frac{\pi \arccos(\mathbf{x}^T \mathbf{X}_i)}{2\pi} g_{i,1} \right) \right| \quad (\text{because } |j\mathbf{x}^T \mathbf{X}_i| \leq 1) \\ &= \sup_{\mathbf{x} \in S^{d-1}, i \in \{1, 2, \dots, n\}} \left| \left(\frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} - \frac{\pi \arccos(\mathbf{x}^T \mathbf{X}_i)}{2\pi} \right) g_{i,1} + (g_{i,p} - g_{i,1}) \frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} \right| \\ &= \sup_{\mathbf{x} \in S^{d-1}, i \in \{1, 2, \dots, n\}} \left| \left(\frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} - \frac{\pi \arccos(\mathbf{x}^T \mathbf{X}_i)}{2\pi} \right) g_{i,1} \right| + \left| (g_{i,p} - g_{i,1}) \frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} \right| \\ & \leq \zeta \left(\max_i |jg_{i,1}| + 1 \right) \quad (\text{because } \mathcal{J}_1 \cap \mathcal{J}_2 \text{ happens, } \frac{jC_{\mathbf{X}_i, \mathbf{x}}^{\mathbf{V}_0}}{p} \in [0, 1], \text{ and } \frac{\pi \arccos(\mathbf{x}^T \mathbf{X}_i)}{2\pi} \in [0, 0.5]). \end{aligned}$$

Because $\max_i |jg_{i,1}|$ is fixed when \mathbf{X} is given, $\zeta \left(\max_i |jg_{i,1}| + 1 \right)$ can be arbitrarily small as long as ζ is small enough. The conclusion of this lemma thus follows by Eq. (107). □

If the ground-truth function $f \not\in \overline{F^{\ell_2}}$ (or equivalently, $D(f, F^{\ell_2}) > 0$), then the MSE of $\hat{f}_\tau^{\ell_2}$ (with respect to the ground-truth function f) is at least $D(f, F^{\ell_2})$ (because $\hat{f}_\tau^{\ell_2} \geq F^{\ell_2}$). Therefore, we have proved Proposition 2. Below we state an even stronger result than part (ii) of Proposition 2, i.e., it captures not only the MSE of $\hat{f}_\tau^{\ell_2}$, but also that of \hat{f}^{ℓ_2} for sufficiently large p .

Lemma 50. *For any given \mathbf{X} and $\zeta > 0$, there exists a threshold p_0 such that for all $p > p_0$, $\Pr \left\{ \sup_{\mathbf{x} \in S^{d-1}} |D(f, \hat{f}^{\ell_2}) - D(f, F^{\ell_2})| < \zeta \right\} > 1 - \zeta$.*

Proof. By Lemma 49, for any $\zeta > 0$, there must exist a threshold p_0 such that for all $p > p_0$,

$$\Pr \left\{ \sup_{\mathbf{x} \in S^{d-1}} |j\hat{f}^{\ell_2}(\mathbf{x}) - \hat{f}_\tau^{\ell_2}(\mathbf{x})| < \zeta \right\} > 1 - \zeta.$$

When $\sup_{\mathbf{x} \in S^{d-1}} |j\hat{f}^{\ell_2}(\mathbf{x}) - \hat{f}_\tau^{\ell_2}(\mathbf{x})| < \zeta$, we have

$$D(\hat{f}^{\ell_2}, \hat{f}_\tau^{\ell_2}) = \sqrt{\int_{S^{d-1}} \left(\hat{f}^{\ell_2}(\mathbf{x}) - \hat{f}_\tau^{\ell_2}(\mathbf{x}) \right)^2 d\mu(\mathbf{x})} < \zeta.$$

Because $\hat{f}_\tau^{\ell_2} \geq F^{\ell_2}$, we have $D(\hat{f}_\tau^{\ell_2}, f) \geq D(f, F^{\ell_2})$. Thus, by the triangle inequality, we have $D(f, \hat{f}^{\ell_2}) \geq D(f, \hat{f}_\tau^{\ell_2}) - D(\hat{f}_\tau^{\ell_2}, \hat{f}^{\ell_2}) \geq D(f, F^{\ell_2}) - \zeta$. Putting these together, we have

$$\Pr \left\{ D(f, \hat{f}^{\ell_2}) \geq D(f, F^{\ell_2}) - \zeta \right\} > 1 - \zeta.$$

Notice that $\text{MSE} = (D(f, \hat{f}^{\ell_2}))^2$. The result of this lemma thus follows. \square

K. Details for Section 4 (hyper-spherical harmonics decomposition on S^{d-1})

K.1. Convolution on S^{d-1}

First, we introduce the definition of the convolution on S^{d-1} . In (Dokmanic & Petrinovic, 2009), the convolution on S^{d-1} is defined as follows.

$$f_1 \sim f_2(\mathbf{x}) := \int_{\text{SO}(d)} f_1(\mathbf{S}\mathbf{e})f_2(\mathbf{S}^{-1}\mathbf{x})d\mathbf{S},$$

where \mathbf{S} is a $d \times d$ orthogonal matrix that denotes a rotation in S^{d-1} , chosen from the set $\text{SO}(d)$ of all rotations. In the following, we will show Eq. (13). To that end, we have

$$g \sim h(\mathbf{x}) = \int_{\text{SO}(d)} g(\mathbf{S}\mathbf{e})h(\mathbf{S}^{-1}\mathbf{x})d\mathbf{S}. \quad (108)$$

Now, we replace $\mathbf{S}\mathbf{e}$ by \mathbf{z} . Thus, we have

$$\mathbf{S}\mathbf{e} = \mathbf{z} \Rightarrow \mathbf{e} = \mathbf{S}^{-1}\mathbf{z} \Rightarrow (\mathbf{S}^{-1}\mathbf{x})^T \mathbf{e} = (\mathbf{S}^{-1}\mathbf{x})^T \mathbf{S}^{-1}\mathbf{z} \Rightarrow (\mathbf{S}^{-1}\mathbf{x})^T \mathbf{e} = \mathbf{x}^T (\mathbf{S}^{-1})^T \mathbf{S}^{-1}\mathbf{z}.$$

Because \mathbf{S} is an orthonormal matrix, we have $\mathbf{S}^T = \mathbf{S}^{-1}$. Therefore, we have $(\mathbf{S}^{-1}\mathbf{x})^T \mathbf{e} = \mathbf{x}^T \mathbf{z}$. Thus, by Eq. (14), we have

$$h(\mathbf{S}^{-1}\mathbf{x}) = (\mathbf{S}^{-1}\mathbf{x})^T \mathbf{e} \frac{\pi}{2\pi} \arccos((\mathbf{S}^{-1}\mathbf{x})^T \mathbf{e}) = \mathbf{x}^T \mathbf{z} \frac{\pi}{2\pi} \arccos(\mathbf{x}^T \mathbf{z}). \quad (109)$$

By plugging Eq. (109) into Eq. (108), we have

$$g \sim h(\mathbf{x}) = \int_{S^{d-1}} g(\mathbf{z}) \mathbf{x}^T \mathbf{z} \frac{\pi}{2\pi} \arccos(\mathbf{x}^T \mathbf{z}) d\mu(\mathbf{z}).$$

Eq. (13) thus follows.

The following lemma shows the intrinsic symmetry of such a convolution.

Lemma 51. Let $\mathbf{S} \in \mathbb{R}^{d \times d}$ denotes any rotation in \mathbb{R}^d . If $f(\mathbf{x}) \in F^{\ell_2}$, then $f(\mathbf{S}\mathbf{x}) \in F^{\ell_2}$.

Proof. Because $f(\mathbf{x}) \in F^{\ell_2}$, we can find g such that

$$f(\mathbf{x}) = \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi}{2\pi} \arccos(\mathbf{x}^T \mathbf{z}) g(\mathbf{z}) d\mu(\mathbf{z}).$$

Thus, we have

$$\begin{aligned} f(\mathbf{S}\mathbf{x}) &= \int_{S^{d-1}} (\mathbf{S}\mathbf{x})^T \mathbf{z} \frac{\pi}{2\pi} \arccos((\mathbf{S}\mathbf{x})^T \mathbf{z}) g(\mathbf{z}) d\mu(\mathbf{z}) \\ &= \int_{S^{d-1}} \mathbf{x}^T (\mathbf{S}^T \mathbf{z}) \frac{\pi}{2\pi} \arccos(\mathbf{x}^T (\mathbf{S}^T \mathbf{z})) g(\mathbf{z}) d\mu(\mathbf{z}) \\ &= \int_{S^{d-1}} \mathbf{x}^T (\mathbf{S}^T \mathbf{z}) \frac{\pi}{2\pi} \arccos(\mathbf{x}^T (\mathbf{S}^T \mathbf{z})) g(\mathbf{S}\mathbf{S}^T \mathbf{z}) d\mu(\mathbf{z}) \\ &\quad (\text{because } \mathbf{S} \text{ is a rotation, we have } \mathbf{S}\mathbf{S}^T = \mathbf{I}) \\ &= \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi}{2\pi} \arccos(\mathbf{x}^T \mathbf{z}) g(\mathbf{S}\mathbf{z}) d\mu(\mathbf{S}\mathbf{z}) \quad (\text{replace } \mathbf{S}^T \mathbf{z} \text{ by } \mathbf{z}) \\ &= \int_{S^{d-1}} \mathbf{x}^T \mathbf{z} \frac{\pi}{2\pi} \arccos(\mathbf{x}^T \mathbf{z}) g(\mathbf{S}\mathbf{z}) d\mu(\mathbf{z}) \quad (\text{by Assumption 1}) \end{aligned}$$

The result of this lemma thus follows. \square

K.2. Hyper-spherical harmonics

We follow the conventions of hyper-spherical harmonics in (Dokmanic & Petrinovic, 2009). We express $\mathbf{x} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_d] \in S^{d-1}$ in a set of hyper-spherical polar coordinates as follows.

$$\begin{aligned} \mathbf{x}_1 &= \sin \theta_{d-1} \sin \theta_{d-2} \dots \sin \theta_2 \sin \theta_1, \\ \mathbf{x}_2 &= \sin \theta_{d-1} \sin \theta_{d-2} \dots \sin \theta_2 \cos \theta_1, \\ \mathbf{x}_3 &= \sin \theta_{d-1} \sin \theta_{d-2} \dots \cos \theta_2, \\ &\vdots \\ \mathbf{x}_{d-1} &= \sin \theta_{d-1} \cos \theta_{d-2}, \\ \mathbf{x}_d &= \cos \theta_{d-1}. \end{aligned}$$

Notice that $\theta_1 \in [0, 2\pi)$ and $\theta_2, \theta_3, \dots, \theta_{d-1} \in [0, \pi)$. Let $\xi = [\theta_1 \ \theta_2 \ \dots \ \theta_{d-1}]$. In such coordinates, hyper-spherical harmonics are given by (Dokmanic & Petrinovic, 2009)

$$Y_{\mathbf{K}}^l(\xi) = A_{\mathbf{K}}^l \prod_{i=0}^{d-3} C_{k_i}^{\frac{d-i-2}{2} + k_{i+1}}(\cos \theta_{d-i-1}) \sin^{k_{i+1}} \theta_{d-i-1} e^{jk_d - 2\theta_1}, \quad (110)$$

where the normalization factor is

$$A_{\mathbf{K}}^l = \sqrt{\frac{1}{\binom{d}{2}} \prod_{i=0}^{d-3} 2^{2k_{i+1} + d - i - 4} \frac{(k_i - k_{i+1})! (d - i + 2k_i - 2)! 2^{\frac{d-i-2}{2} + k_{i+1}}}{\pi (k_i + k_{i+1} + d - i - 2)}},$$

and $C_d^\lambda(t)$ are the Gegenbauer polynomials of degree d . These Gegenbauer polynomials can be defined as the coefficients of α^n in the power-series expansion of the following function,

$$(1 - 2t\alpha + \alpha^2)^{-\lambda} = \sum_{i=0}^{\infty} C_i^\lambda(t) \alpha^i.$$

Further, the Gegenbauer polynomials can be computed by a three-term recursive relation,

$$(i+2)C_{i+2}^\lambda(t) = 2(\lambda+i+1)tC_{i+1}^\lambda(t) - (2\lambda+i)C_i^\lambda(t), \quad (111)$$

with $C_0^\lambda(t) = 1$ and $C_1^\lambda(t) = 2\lambda t$.

K.3. Calculate $C_{\mathbf{K}}^l(\xi)$ where $\mathbf{K} = \mathbf{0}$

Recall that $\mathbf{K} = (k_1, k_2, \dots, k_d)$ and $l = k_0$. By plugging $\mathbf{K} = \mathbf{0}$ into Eq. (110), we have

$$C_{\mathbf{0}}^l(\xi) = A_{\mathbf{0}}^l C_l^{\frac{d-2}{2}}(\cos \theta_d). \quad (112)$$

The following lemma gives an explicit form of Gegenbauer polynomials.

Lemma 52.

$$C_i^\lambda(t) = \sum_{k=0}^{\lfloor \frac{i}{2} \rfloor} (-1)^k \frac{(i-k+\lambda)!}{(\lambda)k!(i-2k)!} (2t)^{i-2k}. \quad (113)$$

Proof. We use mathematical induction. We already know that $C_0^\lambda(t) = 1$ and $C_1^\lambda(t) = 2\lambda t$, which both satisfy Eq. (113). Suppose that $C_i^\lambda(t)$ and $C_{i+1}^\lambda(t)$ satisfy Eq. (113), i.e.,

$$\begin{aligned} C_i^\lambda(t) &= \sum_{k=0}^{\lfloor \frac{i}{2} \rfloor} (-1)^k \frac{(i-k+\lambda)!}{(\lambda)k!(i-2k)!} (2t)^{i-2k}, \\ C_{i+1}^\lambda(t) &= \sum_{k=0}^{\lfloor \frac{i+1}{2} \rfloor} (-1)^k \frac{(i-k+\lambda+1)!}{(\lambda)k!(i-2k+1)!} (2t)^{i-2k+1}. \end{aligned}$$

It remains to show that $C_{i+2}^\lambda(t)$ also satisfy Eq. (113). By Eq. (111), it suffices to show that

$$\begin{aligned} & (i+2) \sum_{k=0}^{\lfloor \frac{i+2}{2} \rfloor} (-1)^k \frac{(i-k+\lambda+2)!}{(\lambda)k!(i-2k+2)!} (2t)^{i-2k+2} \\ &= 2(\lambda+i+1)t \sum_{k=0}^{\lfloor \frac{i+1}{2} \rfloor} (-1)^k \frac{(i-k+\lambda+1)!}{(\lambda)k!(i-2k+1)!} (2t)^{i-2k+1} \\ &= (2\lambda+i) \sum_{k=0}^{\lfloor \frac{i}{2} \rfloor} (-1)^k \frac{(i-k+\lambda)!}{(\lambda)k!(i-2k)!} (2t)^{i-2k}. \end{aligned} \quad (114)$$

To that end, it suffices to show that the coefficients of $(2t)^{i-2k+2}$ are the same for both sides of Eq. (114), for $k = 0, 1, \dots, \lfloor \frac{i+2}{2} \rfloor$. For the first step, we verify the coefficients of $(2t)^{i-2k+2}$ for $k = 1, \dots, \lfloor \frac{i+1}{2} \rfloor$. We have

$$\begin{aligned} & \text{coefficients of } (2t)^{i-2k+2} \text{ on the right-hand-side of Eq. (114)} \\ &= (\lambda+i+1) (-1)^k \frac{(i-k+\lambda+1)!}{(\lambda)k!(i-2k+1)!} - (2\lambda+i) (-1)^{k-1} \frac{(i-k+\lambda+1)!}{(\lambda)(k-1)!(i-2k+2)!} \\ &= (-1)^k \frac{(i-k+\lambda+1)!}{(\lambda)k!(i-2k+2)!} ((\lambda+i+1)(i-2k+2) + (2\lambda+i)k) \\ &= (-1)^k \frac{(i-k+\lambda+1)!}{(\lambda)k!(i-2k+2)!} ((\lambda+i+1)(i+2) + (2\lambda+i)k - 2k(\lambda+i+1)) \\ &= (-1)^k \frac{(i-k+\lambda+1)!}{(\lambda)k!(i-2k+2)!} ((\lambda+i+1)(i+2) - k(i+2)) \\ &= (-1)^k \frac{(i-k+\lambda+1)!}{(\lambda)k!(i-2k+2)!} (\lambda - k + i + 1)(i+2) \\ &= (i+2) (-1)^k \frac{(i-k+\lambda+2)!}{(\lambda)k!(i-2k+2)!} \\ &= \text{coefficients of } (2t)^{i-2k+2} \text{ on the left-hand-side of Eq. (114)}. \end{aligned}$$

For the second step, we verify the coefficient of $(2t)^i$ for $k = 0$, i.e., the coefficient of $(2t)^{i+2}$. We have

$$\begin{aligned} & \text{coefficients of } (2t)^{i+2} \text{ on the right-hand-side of Eq. (114)} \\ &= (\lambda + i + 1) \frac{(i + \lambda + 1)}{(\lambda)(i + 1)!} \\ &= (i + 2) \frac{(i + 2 + \lambda)}{(\lambda)(i + 2)!} \\ &= \text{coefficients of } (2t)^{i+2} \text{ on the left-hand-side of Eq. (114).} \end{aligned}$$

For the third step, we verify the coefficient of $(2t)^i$ for $k = b\frac{i+2}{2}c = b\frac{i}{2}c + 1$. We consider two cases: 1) i is even, and 2) i is odd. When i is even, we have $b\frac{i}{2}c + 1 = \frac{i}{2} + 1$, i.e., $i - 2k + 2 = 0$. Thus, we have

$$\begin{aligned} & \text{coefficients of } (2t)^0 \text{ on the right-hand-side of Eq. (114)} \\ &= (2\lambda + i) \binom{i}{2} \frac{\left(\frac{i}{2} + \lambda\right)}{(\lambda) \left(\frac{i}{2}\right)!} \\ &= (i + 2) \binom{i}{2} \frac{\left(\frac{i}{2} + 1 + \lambda\right)}{(\lambda) \left(\frac{i}{2} + 1\right)!} \\ &= \text{coefficients of } (2t)^0 \text{ on the left-hand-side of Eq. (114).} \end{aligned}$$

When i is odd, we have $k = b\frac{i+1}{2}c + 1 = \frac{i+1}{2} = b\frac{i+1}{2}c$ and this case has already been verified in the first step.

In conclusion, the coefficients of $(2t)^i$ are the same for both sides of Eq. (114), for $k = 0, 1, \dots, b\frac{i+2}{2}c$. Thus, by mathematical induction, the result of this lemma thus follows. \square

Applying Lemma 52 in Eq. (112), we have

$$l_0(\xi) = A_0^l \sum_{k=0}^{b\frac{l}{2}c} \binom{l}{k} \frac{\binom{l-k+\frac{d-2}{2}}{k} (2 \cos \theta_{d-1})^{l-2k}}{\left(\frac{d-2}{2}\right)k!(l-2k)!}. \quad (115)$$

We give a few examples of $l_0(\xi)$ as follows.

$$\begin{aligned} 0_0(\xi) &= A_0^0, \\ 1_0(\xi) &= A_0^1 (d-2) \cos \theta_{d-1}, \\ 2_0(\xi) &= A_0^2 \frac{d-2}{2} (d \cos^2 \theta_{d-1} - 1), \\ 3_0(\xi) &= A_0^3 \frac{d-2}{2} d \left(\frac{d+2}{3} \cos^3 \theta_{d-1} - \cos \theta_{d-1} \right). \end{aligned}$$

K.4. Proof of Proposition 3

Recall that

$$h(\mathbf{x}) := \mathbf{x}^T \mathbf{e} \frac{\pi \arccos(\mathbf{x}^T \mathbf{e})}{2\pi}, \quad \mathbf{e} := [0 \ 0 \ \dots \ 0 \ 1]^T \in \mathbb{R}^d.$$

Notice that $\mathbf{x}^T \mathbf{e} = \cos \theta_{d-1}$. Thus, we have

$$h(\mathbf{x}) = \cos \theta_{d-1} \frac{\pi \arccos(\cos \theta_{d-1})}{2\pi}.$$

The arccos function has a Taylor Series Expansion:

$$\arccos(a) = \frac{\pi}{2} \sum_{i=0}^{\infty} \frac{(2i)!}{2^{2i} (i!)^2} \frac{a^{2i+1}}{2i+1},$$

which converges when $d \geq a + 1$. Thus, we have

$$h(\mathbf{x}) = \frac{1}{4} \cos \theta_{d-1} + \frac{1}{2\pi} \sum_{i=0}^7 \frac{(2i)!}{2^{2i} (i!)^2} \frac{\cos^{2i+2} \theta_{d-1}}{2i+1}. \quad (116)$$

By comparing terms of even and odd power of $\cos \theta_{d-1}$ in Eq. (115) and Eq. (116), we immediately see that $h(\mathbf{x}) \in \mathcal{B}_0^l(\mathbf{x})$ when $l = 1$, and $h(\mathbf{x}) \notin \mathcal{B}_0^l(\mathbf{x})$ when $l = 3, 5, 7, \dots$. It remains to examine whether $h(\mathbf{x}) \in \mathcal{B}_0^l(\mathbf{x})$ or $h(\mathbf{x}) \in \mathcal{B}_0^l(\mathbf{x})$ for $l \geq 0, 1, 2, 4, 6, \dots, g$. We first introduce the following lemma.

Lemma 53. *Let a and b be two non-negative integers. Define the function*

$$Q(a, b) := \int_{S^{d-1}} \cos^a(\theta_{d-1}) \mathbf{b}_0^b(\xi) d\mu(\mathbf{x}).$$

We must have

$$Q(2k, 2m) \begin{cases} > 0, & \text{if } m = k, \\ = 0, & \text{if } m > k. \end{cases} \quad (117)$$

Proof. We have

$$Q(2k, 0) = \int_{S^{d-1}} \cos^{2k}(\theta_{d-1}) \mathbf{0}_0^0(\xi) d\mu(\mathbf{x}) = A_0^0 \int_{S^{d-1}} \cos^{2k}(\theta_{d-1}) d\mu(\mathbf{x}) > 0.$$

Thus, to finish the proof, we only need to consider the case of $m = 1$ in Eq. (117). We then prove by mathematical induction on the first parameter of $Q(\cdot, \cdot)$, i.e., k in Eq. (117). When $m > 0$, we have

$$Q(0, 2m) = \int_{S^{d-1}} \mathbf{0}_0^{2m}(\xi) d\mu(\mathbf{x}) = \frac{1}{A_0^0} \int_{S^{d-1}} \mathbf{0}_0^0(\xi) \mathbf{0}_0^{2m}(\xi) d\mu(\mathbf{x}) = 0$$

(by the orthogonality of the basis).

Thus, Eq. (117) holds for all m when $k = 0$. Suppose that Eq. (117) holds when $k = i$. To complete the mathematical induction, it only remains to show that Eq. (117) also holds for all m when $k = i + 1$. By Eq. (111) and Eq. (112), for any l , we have

$$\cos(\theta_{d-1}) \mathbf{0}_0^{l+1}(\xi) = \frac{(l+2)A_0^{l+1}}{(d+2l)A_0^{l+2}} \mathbf{0}_0^{l+2}(\xi) + \frac{(d-2+l)A_0^{l+1}}{(d+2l)A_0^l} \mathbf{0}_0^l(\xi).$$

Thus, we have

$$Q(a+1, l+1) = q_{l,1} Q(a, l+2) + q_{l,2} Q(a, l), \quad (118)$$

where

$$q_{l,1} := \frac{(l+2)A_0^{l+1}}{(d+2l)A_0^{l+2}}, \quad q_{l,2} := \frac{(d-2+l)A_0^{l+1}}{(d+2l)A_0^l}.$$

It is obvious that $q_{l,1} > 0$ and $q_{l,2} > 0$. Applying Eq. (118) multiple times, we have

$$Q(2i+2, 2m) = q_{2m-1,1} Q(2i+1, 2m+1) + q_{2m-1,2} Q(2i+1, 2m-1), \quad (119)$$

$$Q(2i+1, 2m+1) = q_{2m,1} Q(2i, 2m+2) + q_{2m,2} Q(2i, 2m), \quad (120)$$

$$Q(2i+1, 2m-1) = q_{2m-2,1} Q(2i, 2m) + q_{2m-2,2} Q(2i, 2m-2). \quad (121)$$

(Notice that we have already let $m = 1$, so all $q_{\cdot,1}, q_{\cdot,2}, Q(\cdot, \cdot)$ in those equations are well-defined.) By plugging Eq. (120) and Eq. (121) into Eq. (119), we have

$$Q(2i+2, 2m) = q_{2m,1} q_{2m-1,1} Q(2i, 2m+2) + (q_{2m-1,1} q_{2m,2} + q_{2m-1,2} q_{2m-2,1}) Q(2i, 2m) + q_{2m-1,2} q_{2m-2,2} Q(2i, 2m-2). \quad (122)$$

To prove that Eq. (117) holds when $k = i + 1$ for all m , we consider two cases, Case 1: $m = i + 1$, and Case 2: $m > i + 1$. Notice that by the induction hypothesis, we already know that Eq. (117) holds when $k = i$ for all m .

Case 1. When $m = i + 1$, we have $m - 1 = i$. Thus, by the induction hypothesis for $k = i$, we have $Q(2i, 2m - 2) > 0$ (by $m - 1 = i$), which implies that the third term of the right-hand-side of Eq. (122) is positive. Further, by the induction hypothesis for $k = i$, we also know that $Q(2i, 2m + 2) = 0$ and $Q(2i, 2m) = 0$ (regardless of the value of m), which means that the first and the second term of Eq. (122) is non-negative. Thus, by considering all three terms in Eq. (122) together, we have $Q(2i + 2, 2m) > 0$ when $m = i + 1$.

Case 2. When $m > i + 1$, we have $m + 1 > i$, $m > i$, and $m - 1 > i$. Thus, by the induction hypothesis for $k = i$, we have $Q(2i, 2m + 2) = Q(2i, 2m) = Q(2i, 2m - 2) = 0$. Therefore, by Eq. (122), we have $Q(2i + 2, 2m) = 0$.

In summary, Eq. (117) holds when $k = i + 1$ for all m . The mathematical induction is completed and the result of this lemma follows. \square

By Lemma 53, for all $k \geq 0$, we have

$$\begin{aligned} & \int_{S^{d-1}} \frac{1}{2\pi} \sum_{i=0}^k \frac{(2i)!}{2^{2i}(i!)^2} \frac{\cos^{2i+2} \theta_d}{2i+1} \mathbf{0}^{2k}(\xi) d\mu(\mathbf{x}) \\ &= \frac{1}{2\pi} \sum_{i=0}^k \frac{(2i)!}{2^{2i}(i!)^2} \frac{1}{2i+1} \int_{S^{d-1}} \cos^{2i+2} \theta_d \mathbf{0}^{2k}(\xi) d\mu(\mathbf{x}) \\ &> 0. \end{aligned}$$

Thus, by Eq. (116), we know that $h(\mathbf{x}) \in \mathcal{B}_0^l(\mathbf{x})$ for all $l \geq 0, 2, 4, \dots, g$.

K.5. A special case: when $d = 2$

When $d = 2$, S^{d-1} denotes a unit circle. Therefore, every \mathbf{x} corresponds to an angle $\varphi \in [\pi, \pi]$ such that $\mathbf{x} = [\cos \varphi \ \sin \varphi]^T$. In this situation, the hyper-spherical harmonics are the well-known Fourier series, i.e., $1, \cos(\theta), \sin(\theta), \cos(2\theta), \sin(2\theta), \dots$. Thus, we can explicitly calculate all Fourier coefficients of h more easily.

Similarly to Appendix K.1, we first write down the convolution for $d = 2$, which is also in a simpler form. For any function $f_g \in F^{\ell_2}$, we have

$$\begin{aligned} f_g(\varphi) &= \frac{1}{2\pi} \int_{\varphi-\pi}^{\varphi+\pi} \frac{\pi}{2\pi} \frac{j\theta}{2\pi} \frac{\varphi^j}{2\pi} \cos(\theta - \varphi) g(\theta) d\theta \\ &= \frac{1}{2\pi} \int_{\pi}^{\pi} \frac{\pi}{2\pi} \frac{j\theta^j}{2\pi} \cos \theta g(\theta + \varphi) d\theta \text{ (replace } \theta \text{ by } \theta - \varphi) \\ &= \frac{1}{2\pi} \int_{\pi}^{\pi} \frac{\pi}{2\pi} \frac{j\theta^j}{2\pi} \cos \theta g(\varphi - \theta) d\theta \text{ (replace } \theta \text{ by } -\theta). \end{aligned}$$

Define $h(\theta) := \frac{\pi}{2\pi} \frac{j\theta^j}{2\pi} \cos \theta$. We then have

$$f_g(\varphi) = \frac{1}{2\pi} h(\varphi) \sim g(\varphi),$$

where \sim denotes (continuous) circular convolution. Let $c_{f_g}(k), c_h(k)$ and $c_g(k)$ (where $k = \dots, -1, 0, 1, \dots$) denote the (complex) Fourier series coefficients for $f_g(\varphi), h(\varphi)$, and $g(\varphi)$, correspondingly. Specifically, we have

$$f_g(\varphi) = \sum_{k=-\infty}^{\infty} c_{f_g}(k) e^{ik\varphi}, \quad h(\varphi) = \sum_{k=-\infty}^{\infty} c_h(k) e^{ik\varphi}, \quad g(\varphi) = \sum_{k=-\infty}^{\infty} c_g(k) e^{ik\varphi}.$$

Thus, we have

$$c_{f_g}(k) = c_h(k) c_g(k). \quad (123)$$

Now we calculate $c_h(k)$, i.e., the Fourier decomposition of $h(\cdot)$. We have

$$\begin{aligned} c_h(k) &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\pi}{2\pi} j\theta j \cos \theta e^{ik\theta} d\theta \\ &= \frac{1}{4\pi} \int_{-\pi}^{\pi} \left(1 - \frac{j\theta j}{\pi}\right) \frac{e^{i(k+1)\theta} + e^{i(k-1)\theta}}{2} d\theta \\ &= \frac{1}{8\pi^2} \int_{-\pi}^{\pi} j\theta j \left(e^{i(k+1)\theta} + e^{i(k-1)\theta}\right) d\theta + \frac{1}{8\pi} \int_{-\pi}^{\pi} \left(e^{i(k+1)\theta} + e^{i(k-1)\theta}\right) d\theta. \end{aligned}$$

It is easy to verify that

$$\int x e^{cx} dx = e^{cx} \left(\frac{cx-1}{c^2}\right), \quad c \neq 0.$$

Thus, we have

$$\begin{aligned} c_h(1) &= \frac{1}{8\pi^2} \int_{-\pi}^{\pi} j\theta j (e^{i2\theta} + 1) d\theta + \frac{1}{4} \\ &= \frac{1}{8\pi^2} \left(\pi^2 \int_{-\pi}^0 \theta e^{i2\theta} d\theta + \int_0^{\pi} \theta e^{i2\theta} d\theta \right) + \frac{1}{4} \\ &= \frac{1}{8\pi^2} \left(\pi^2 + \frac{i2\pi}{4} + \frac{i2\pi}{4} \right) + \frac{1}{4} \\ &= \frac{1}{8} + \frac{1}{4} \\ &= \frac{1}{8}. \end{aligned}$$

Similarly, we have

$$c_h(-1) = \frac{1}{8}.$$

Now we consider the situation of $n \neq 1$. We have

$$\begin{aligned} \int_{-\pi}^0 j\theta j e^{i(k+1)\theta} d\theta &= e^{i(k+1)\theta} \left. \frac{i(k+1)\theta}{(k+1)^2} - \frac{1}{(k+1)^2} \right|_{-\pi}^0 = \frac{1}{(k+1)^2} + \frac{1}{(k+1)^2} e^{i(k+1)\pi}, \\ \int_0^{\pi} j\theta j e^{i(k+1)\theta} d\theta &= e^{i(k+1)\theta} \left. \frac{i(k+1)\theta}{(k+1)^2} - \frac{1}{(k+1)^2} \right|_0^{\pi} = \frac{1}{(k+1)^2} + \frac{1+i(k+1)\pi}{(k+1)^2} e^{i(k+1)\pi}. \end{aligned}$$

Notice that $e^{-i(k+1)\pi} = e^{i(k+1)2\pi} e^{i(k+1)\pi} = e^{i(k+1)\pi}$. Therefore, we have

$$\int_{-\pi}^{\pi} j\theta j e^{i(k+1)\theta} d\theta = \frac{2}{(k+1)^2} \left(e^{i(k+1)\pi} - 1 \right).$$

Similarly, we have

$$\int_{-\pi}^{\pi} j\theta j e^{i(k-1)\theta} d\theta = \frac{2}{(k-1)^2} \left(e^{i(k-1)\pi} - 1 \right).$$

In summary, we have

$$\begin{aligned} c_h(k) &= \begin{cases} \frac{1}{8}, & k = 1 \\ \frac{1}{4\pi^2} \left(\frac{1}{(k+1)^2} + \frac{1}{(k-1)^2} \right) \left(e^{i(k+1)\pi} - 1 \right), & \text{otherwise} \end{cases} \\ &= \begin{cases} \frac{1}{8}, & k = 1 \\ \frac{1}{2\pi^2} \left(\frac{1}{(k+1)^2} + \frac{1}{(k-1)^2} \right), & k = 0, 2, 4, \dots \\ 0, & k = 3, 5, \dots \end{cases} \end{aligned}$$

By Eq. (123), we thus have

$$c_{f_g}(k) = \begin{cases} \frac{1}{8}c_g(k), & k = 1 \\ \frac{1}{2\pi^2} \left(\frac{1}{(k+1)^2} + \frac{1}{(k-1)^2} \right) c_g(k), & k = 0, 2, 4, \dots \\ 0, & k = 3, 5, \dots \end{cases}$$

In other words, when $d = 2$, functions in F^{ℓ_2} can only contain frequencies $0, \theta, 2\theta, 4\theta, 6\theta, \dots$, and cannot contain other frequencies $3\theta, 5\theta, 7\theta, \dots$.

K.6. Details of Remark 2

As we discussed in Remark 2, a ReLU activation function with bias that operates on $\mathbf{x} \in \mathbb{R}^{d-1}$, $kgk_2^2 = \frac{d-1}{\beta}$ can be equivalently viewed as one without bias that operates on $\mathbf{x} \in S^{d-1}$, but with the last element of \mathbf{x} fixed at $1/\beta$. Note that by fixing the last element of $\mathbf{x} \in S^{d-1}$ at a constant $1/\beta$, we essentially consider ground-truth functions with a much smaller domain $D := \left\{ \mathbf{x} = \begin{bmatrix} \tilde{\mathbf{x}} \\ 1/\beta \end{bmatrix} \mid \mathbf{x} \in \mathbb{R}^{d-1}, kgk_2^2 = \frac{d-1}{\beta} \right\} \subset S^{d-1}$. Correspondingly, define a vector $\mathbf{a} \in \mathbb{R}^{d-1}$ and $a_0 \in \mathbb{R}$ such that $\mathbf{a} = \begin{bmatrix} \tilde{\mathbf{a}} \\ a_0 \end{bmatrix} \in \mathbb{R}^d$. We claim that for any $\mathbf{a} \in \mathbb{R}^d$ and for all non-negative integer l , a ground-truth function $f(\mathbf{x}) = (\mathbf{x}^T \mathbf{a})^l, \mathbf{x} \in D$ must be learnable. In other words, all polynomials can be learned in the constrained domain D . Towards this end, recall that we have already shown that polynomials (of $\mathbf{x} \in S^{d-1}$) to the power of $l = 0, 1, 2, 4, 6, \dots$ are learnable. Thus, it suffices to prove that polynomials of $\mathbf{x} \in D$ to the power of $l = 3, 5, 7, \dots$ can be represented by a finite sum of those to the power of $l = 0, 1, 2, 4, 6, \dots$. The idea is to utilize the fact that the binomial expansion of $(\mathbf{x}^T \mathbf{a} + \frac{a_0}{\beta})^l$ contains $(\mathbf{x}^T \mathbf{a})^k$ for all $k = 0, 1, 2, 3, \dots, l$. Here we give an example for writing $(\mathbf{x}^T \mathbf{a})^3$ as a linear combination of learnable components. Other values of $l = 5, 7, 9, \dots$ can be proved in a similar way. Notice that

$$\begin{aligned} (\mathbf{x}^T \mathbf{a})^3 &= \frac{1}{4} \left((\mathbf{x}^T \mathbf{a} + 1)^4 - (\mathbf{x}^T \mathbf{a})^4 - 6(\mathbf{x}^T \mathbf{a})^2 - 4(\mathbf{x}^T \mathbf{a})^2 - 1 \right) \text{ (by the binomial expansion of } (\mathbf{x}^T \mathbf{a} + 1)^4 \text{)} \\ &= \frac{1}{4} \left(\left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ \beta \end{bmatrix} \right)^4 - \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right)^4 - 6 \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right)^2 - 4 \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right) - 1 \right). \end{aligned} \quad (124)$$

Thus, for all $\mathbf{x} = \begin{bmatrix} \tilde{\mathbf{x}} \\ 1/\beta \end{bmatrix}$ and $\mathbf{a} = \begin{bmatrix} \tilde{\mathbf{a}} \\ a_0 \end{bmatrix}$, we have

$$\begin{aligned} (\mathbf{x}^T \mathbf{a})^3 &= \left(\mathbf{x}^T \mathbf{a} + \frac{a_0}{\beta} \right)^3 \\ &= (\mathbf{x}^T \mathbf{a})^3 + 3 \left(\frac{a_0}{\beta} \right) (\mathbf{x}^T \mathbf{a})^2 + 3 \left(\frac{a_0}{\beta} \right)^2 (\mathbf{x}^T \mathbf{a}) + \left(\frac{a_0}{\beta} \right)^3 \\ &= (\mathbf{x}^T \mathbf{a})^3 + 3 \left(\frac{a_0}{\beta} \right) \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right)^2 + 3 \left(\frac{a_0}{\beta} \right)^2 \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right) + \left(\frac{a_0}{\beta} \right)^3 \\ &= \frac{1}{4} \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ \beta \end{bmatrix} \right)^4 - \frac{1}{4} \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right)^4 + \left(3 \left(\frac{a_0}{\beta} \right) - \frac{3}{2} \right) \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right)^2 \\ &\quad + \left(3 \left(\frac{a_0}{\beta} \right)^2 - 1 \right) \left(\mathbf{x}^T \begin{bmatrix} \tilde{\mathbf{a}} \\ 0 \end{bmatrix} \right) + \left(\left(\frac{a_0}{\beta} \right)^3 - \frac{1}{4} \right) \text{ (by Eq. (124)),} \end{aligned}$$

which is a sum of 5 learnable components (corresponding to the polynomials with power of 4, 4, 2, 1, and 0, respectively).

L. Discussion when g is a δ -function ($kgk_1 = 1$)

We now discuss what happens to the conclusion of Theorem 1 if g contains a δ -function, in which case $kgk_1 = 1$. In Eq. (10) of Theorem 1, only Term 1 and Term 4 (come from Proposition 5) will be affected when $kgk_1 = 1$. That is because only Proposition 5 requires $kgk_1 < 1$ during the proof of Theorem 1. To accommodate the situation when g contains a δ -function ($kgk_1 = 1$), we need a new version of Proposition 5. In other words, we need to know the performance of the overfitted NTK solution in learning the pseudo ground-truth when $kgk_1 = 1$.

Without loss of generality, we consider the situation that $g = \delta_{z_0}$. We have the following proposition.

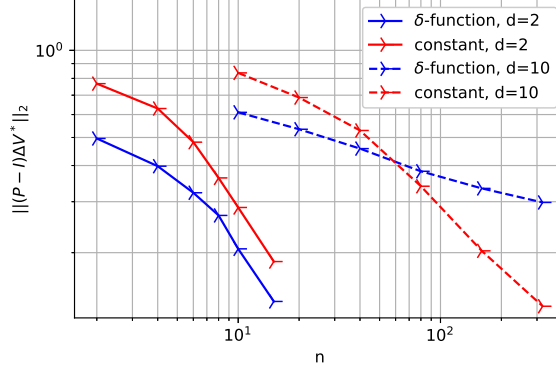


Figure 7. The curves of the model error $k(\mathbf{P} - \mathbf{I}) \mathbf{V} k_2$ for learning the pseudo ground-truth $f_{\mathbf{V}_0}^g$ with respect to n for different g and different d , where $p = 20000$, and $\epsilon = \mathbf{0}$. Every curve is the average of 10 random simulation runs.

Proposition 54. *If the ground-truth function is $f = f_{\mathbf{V}_0}^g$ in Definition 2 with $g = \delta_{\mathbf{z}_0}$ and $\epsilon = \mathbf{0}$, for any $\mathbf{x} \in S^{d-1}$ and $q \in (1, \infty)$, we have*

$$\Pr_{\mathbf{x}, \mathbf{V}_0} \left\{ \int f_{\hat{f}}^{\ell_2}(\mathbf{x}) - f(\mathbf{x}) \left(\sqrt{\frac{3}{4} + \frac{\pi^2}{2}} \right) \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right)^{\frac{1}{2(d-1)}} n^{-\frac{1}{2(d-1)(1-\frac{1}{q})}} \right\} \\ \leq \exp\left(-n^{\frac{1}{q}}\right) + 2 \exp\left(-\frac{p}{24} \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right)^{\frac{1}{d-1}} n^{-\frac{1}{d-1}(1-\frac{1}{q})}\right),$$

when

$$n \geq \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right)^{\frac{q}{q-1}}, \text{ i.e., } \left((d-1)B\left(\frac{d-1}{2}, \frac{1}{2}\right) \right) n^{-(1-\frac{1}{q})} \geq 1. \quad (125)$$

(Estimates of $B(\frac{d-1}{2}, \frac{1}{2})$ can be found in Lemma 32.)

Proposition 54 implies that when n is large and p is much larger than $n^{-\frac{1}{2(d-1)(1-\frac{1}{q})}}$, the test error between the pseudo ground-truth and learned result decreases with n at the speed $O(n^{-\frac{1}{2(d-1)(1-\frac{1}{q})}})$. Further, if we let q be large, then the decreasing speed with n is almost $O(n^{-\frac{1}{2(d-1)}})$. When $d = 3$, this speed is slower than $O(n^{-\frac{1}{2}})$ described in Proposition 5 (i.e., Term 1 in Eq. (10) of Theorem 1). When $d = 2$, the decreasing speed with respect to n is $O(n^{-\frac{1}{2}})$ for both Proposition 5 and Proposition 54. Nonetheless, Proposition 54 implies that the ground-truth functions $f_g \in F^{\ell_2}$ is still learnable even when g is a δ -function (i.e., $k_g k_{\gamma} = 1$), but the test error potentially suffers a slower convergence speed with respect to n when d is large.

In Fig. 7, we plot the curves of the model error $k(\mathbf{P} - \mathbf{I}) \mathbf{V} k_2$ for learning the pseudo ground-truth $f_{\mathbf{V}_0}^g$ with respect to n when $g = \delta_{\mathbf{z}_0}$ (two blue curves) and when g is constant (two red curves). We plot both the case when $d = 2$ (two solid curves) and the case when $d = 10$ (two dashed curves). By Lemma 44, the model error $k(\mathbf{P} - \mathbf{I}) \mathbf{V} k_2$ can represent the generalization performance for learning the pseudo ground-truth $f_{\mathbf{V}_0}^g$ when there is no noise. In Fig. 7, we can see that those two curves corresponding to $d = 10$ have different slopes and the other two curves corresponding to $d = 2$ have a similar slope, which confirms our prediction in the earlier paragraph (i.e., when $d = 2$ the test error will decay at the same speed regardless of whether g contains a δ -function or not, but when $d > 2$ the test error will decay more slowly when g contains a δ -function).

L.1. Proof of Proposition 54

We first show two useful lemmas.

Lemma 55. *For any $q \in (1, \infty)$, if $b \in [n^{-(1-\frac{1}{q})}, 1]$, then*

$$(1-b)^n \leq \exp\left(-n^{\frac{1}{q}}\right).$$

Proof. By Lemma 29, we have

$$\begin{aligned}
 & e^{-b} \leq 1 - b \\
 \Rightarrow & e^{-1} \leq (1 - b)^{\frac{1}{b}} \\
 \Rightarrow & \exp\left(-\frac{1}{n^q}\right) \leq (1 - b)^{n^{1/q}/b} \\
 \Rightarrow & \exp\left(-\frac{1}{n^q}\right) \leq (1 - b)^n \text{ because } b \geq [n^{-(1-1/q)}, 1].
 \end{aligned}$$

□

Lemma 56. Consider $\mathbf{x}_1 \in S^{d-1}$ where $\varphi = \arccos(\mathbf{x}_1^T \mathbf{z}_0)$. For any $\theta \in [\varphi, \pi]$, there must exist $\mathbf{x}_2 \in S^{d-1}$ such that $\arccos(\mathbf{x}_2^T \mathbf{z}_0) = \theta$ and

$$C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0} = C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}, \quad C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0} = C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}. \quad (126)$$

We will explain the intuition of Lemma 56 in Remark 8 right after we use the lemma. We put the proof of Lemma 56 in Section L.2.

Now we are ready to prove Proposition 54. Recall \mathbf{V} defined in Eq. (84). By Eq. (1) and $g = \delta_{\mathbf{z}_0}$, we have

$$\mathbf{V} = \frac{(\mathbf{h}_{\mathbf{V}_0, \mathbf{z}_0})^T}{p}.$$

Define

$$\begin{aligned}
 i &= \arg \min_{i \in \{1, 2, \dots, ng\}} \|\mathbf{X}_i - \mathbf{z}_0\|_2, \\
 \theta &= \arccos(\mathbf{X}_i^T \mathbf{z}_0).
 \end{aligned}$$

Thus, we have

$$\begin{aligned}
 \|\mathbf{X}_i - \mathbf{z}_0\|_2 &= \sqrt{2 - 2 \cos \theta} \quad (\text{by the law of cosines}) \\
 &= 2 \sin \frac{\theta}{2} \quad (\text{by the half angle identity}) \\
 &\leq \theta \quad (\text{by Lemma 41}).
 \end{aligned} \quad (127)$$

(Graphically, Eq. (127) means that a chord is not longer than the corresponding arc.)

As we discussed in the proof sketch of Proposition 5, we now construct the vector \mathbf{a} such that $\mathbf{H}^T \mathbf{a}$ is close to \mathbf{V} . Define $\mathbf{a} \in \mathbb{R}^n$ whose i -th element is

$$\mathbf{a}_i = \begin{cases} 1/p, & \text{if } i = i \\ 0, & \text{if } i \in \{1, 2, \dots, ng\} \setminus \{i\}. \end{cases}$$

Thus, we have $\mathbf{H}^T \mathbf{a} = (\mathbf{h}_{\mathbf{V}_0, \mathbf{X}_i})^T / p$. Therefore, we have

$$\begin{aligned}
 \|\mathbf{H}^T \mathbf{a} - \mathbf{V}\|_2^2 &= \sum_{j=1}^p k(\mathbf{H}^T \mathbf{a})[j] - \mathbf{V}[j] k_2^2 \\
 &= \frac{1}{p^2} \sum_{j=1}^p \left(\mathbf{1}_{\{\mathbf{X}_i^T \mathbf{v}_0[j] > 0, \mathbf{z}_0^T \mathbf{v}_0[j] > 0\}} k \|\mathbf{X}_i - \mathbf{z}_0\|_2^2 + \mathbf{1}_{\{\mathbf{X}_i^T \mathbf{v}_0[j] < 0, \mathbf{z}_0^T \mathbf{v}_0[j] < 0\}} \right) \\
 &= \frac{1}{p^2} \left(p \|\mathbf{X}_i - \mathbf{z}_0\|_2^2 + j C_{\mathbf{X}_i, \mathbf{z}_0}^{\mathbf{V}_0} + j C_{\mathbf{X}_i, \mathbf{z}_0}^{\mathbf{V}_0} \right) \quad (\text{by Eq. (6)}) \\
 &= \frac{1}{p^2} \left(p (\theta)^2 + j C_{\mathbf{X}_i, \mathbf{z}_0}^{\mathbf{V}_0} + j C_{\mathbf{X}_i, \mathbf{z}_0}^{\mathbf{V}_0} \right) \quad (\text{by Eq. (127)}).
 \end{aligned}$$

Thus, we have

$$\begin{aligned} \rho_{\overline{p}k\mathbf{H}\mathbf{a}} \quad \mathbf{V} \quad k_2 \quad & \sqrt{(\theta)^2 + \frac{jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} + jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0}}{p}} \\ & \sqrt{\pi\theta + \frac{jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} + jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0}}{p}} \quad (\text{because } \theta < \pi). \end{aligned} \quad (128)$$

Remark 7. We give a geometric interpretation of Eq. (128) when $d = 2$ by Fig. 4, where \widehat{OA} denotes z_0 , \widehat{OB} denotes \mathbf{X}_i . Then, $jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} + jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0}$ corresponds to the number of $\mathbf{V}_0[j]$'s whose direction is in the arc \widehat{CE} or the arc \widehat{FD} , and θ corresponds to the angle $\angle AOB$. Intuitively, when n increases, \mathbf{X}_i and z_0 get closer, so θ becomes smaller. At the same time, both the arc \widehat{CE} and the arc \widehat{FD} become shorter. Consequently, the value of Eq. (128) decreases as n increases. In the rest of the proof, we will quantitatively estimate the above relationship.

Recall C_d in Eq. (60). Define

$$\theta := \frac{\pi}{2} \left(\frac{2^{\rho} 2^{(d-1)}}{C_d} \right)^{\frac{1}{d-1}} n^{-\frac{1}{d-1}(1-\frac{1}{q})} \geq \left[0, \frac{\pi}{2} \right] \quad (\text{by Eq. (125)}). \quad (129)$$

For any $q \geq (1, 1)$, we define two events:

$$\begin{aligned} \mathcal{J}_1 &:= \left\{ \frac{jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} + jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0}}{p} \geq \frac{3\theta}{2\pi} \right\}, \\ \mathcal{J}_2 &:= f\theta < \theta g. \end{aligned}$$

If both \mathcal{J}_1 and \mathcal{J}_2 happen, by Eq. (128), we must then have

$$\begin{aligned} \rho_{\overline{p}k\mathbf{H}\mathbf{a}} \quad \mathbf{V} \quad k_2 \quad & \left(\sqrt{\frac{3}{2\pi} + \pi} \right)^{\rho_{\overline{p}k\mathbf{H}\mathbf{a}}} \\ & = \left(\sqrt{\frac{3}{4} + \frac{\pi^2}{2}} \right) \left(\frac{2^{\rho} 2^{(d-1)}}{C_d} \right)^{\frac{1}{2(d-1)}} n^{-\frac{1}{2(d-1)}(1-\frac{1}{q})}. \end{aligned}$$

Thus, by Lemma 44 and Lemma 45, if $f = f_{\mathbf{V}_0}^g$ and both \mathcal{J}_1 and \mathcal{J}_2 happen, then for any $\mathbf{x} \geq S^{d-1}$, we must have

$$j\hat{f}^{\ell_2}(\mathbf{x}) - f(\mathbf{x})j \leq \left(\sqrt{\frac{3}{4} + \frac{\pi^2}{2}} \right) \left(\frac{2^{\rho} 2^{(d-1)}}{C_d} \right)^{\frac{1}{2(d-1)}} n^{-\frac{1}{2(d-1)}(1-\frac{1}{q})}. \quad (130)$$

It then only remains to estimate the probability of $\mathcal{J}_1 \setminus \mathcal{J}_2$.

Step 1: Estimate the probability of \mathcal{J}_1 conditional on \mathcal{J}_2 .

When \mathcal{J}_2 happens, we have $\theta < \theta$. By Lemma 56, we can find $\mathbf{x} \geq S^{d-1}$ such that the angle between \mathbf{x} and z_0 is exactly θ and

$$\frac{jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} + jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0}}{p} = \frac{jC_{\mathbf{x}, z_0}^{\mathbf{V}_0} + jC_{\mathbf{x}, z_0}^{\mathbf{V}_0}}{p}. \quad (131)$$

Remark 8. We give a geometric interpretation of Eq. (131) (i.e., Lemma 56) when $d = 2$ by Fig. 4. Recall in Remark 7 that, if we take \widehat{OA} as z_0 and \widehat{OB} as \mathbf{X}_i , then $jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} + jC_{\mathbf{X}_i, z_0}^{\mathbf{V}_0}$ corresponds to the number of $\mathbf{V}_0[j]$'s whose direction is in the arc \widehat{CE} or the arc \widehat{FD} . If we fix \widehat{OA} (i.e., z_0) and increase the angle $\angle AOB$ (corresponding to θ), then both the arc \widehat{CE} and the arc \widehat{FD} will become longer. In other words, if we replace \mathbf{X}_i by \mathbf{x} such that the angle θ (between z_0 and \mathbf{X}_i) increases to the angle θ (between z_0 and \mathbf{x}), then $C_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} = C_{\mathbf{x}, z_0}^{\mathbf{V}_0}$ and $C_{\mathbf{X}_i, z_0}^{\mathbf{V}_0} = C_{\mathbf{x}, z_0}^{\mathbf{V}_0}$, and thus Eq. (131) follows.

We next estimate the probability that the right-hand-side of Eq. (131) is greater than $\frac{3\theta}{2\pi}$. By Eq. (6), we have

$$\frac{jC_{\mathbf{x}, \mathbf{z}_0}^{\mathbf{V}_0} + jC_{\mathbf{x}, \mathbf{z}_0}^{\mathbf{V}_0}}{p} = \frac{1}{p} \sum_{j=1}^p \underbrace{\mathbf{1}_{\{ \mathbf{x}^T \mathbf{V}_0[j] > 0, \mathbf{z}_0^T \mathbf{V}_0[j] > 0 \text{ OR } \mathbf{x}^T \mathbf{V}_0[j] > 0, \mathbf{z}_0^T \mathbf{V}_0[j] < 0 \}}}_{\text{Term A}}. \quad (132)$$

Notice that the angle between \mathbf{x} and \mathbf{z}_0 is $\pi - \theta$, and the angle between \mathbf{x} and $-\mathbf{z}_0$ is also $\pi - \theta$. By Lemma 17 and Assumption 1, we know that the Term A in Eq. (132) follows Bernoulli distribution with the probability $2 \frac{\pi - (\pi - \theta)}{2\pi} = \frac{\theta}{\pi}$. By letting $\delta = 1/2$, $a = p$, $b = \frac{\theta}{\pi}$ in Lemma 14, we have

$$\Pr_{\mathbf{V}_0} \left\{ \left| \frac{jC_{\mathbf{x}, \mathbf{z}_0}^{\mathbf{V}_0} + jC_{\mathbf{x}, \mathbf{z}_0}^{\mathbf{V}_0}}{p} - \frac{p\theta}{\pi} \right| > \frac{p\theta}{2\pi} \right\} \leq 2 \exp \left(- \frac{p\theta}{12\pi} \right).$$

By Eq. (131), we then have

$$\Pr_{\mathbf{V}_0} [\mathcal{J}_1^c | \mathcal{J}_2] \leq \Pr_{\mathbf{V}_0} \left\{ \frac{jC_{\mathbf{x}, \mathbf{z}_0}^{\mathbf{V}_0} + jC_{\mathbf{x}, \mathbf{z}_0}^{\mathbf{V}_0}}{p} > \frac{3\theta}{2\pi} \right\} \leq 2 \exp \left(- \frac{p\theta}{12\pi} \right).$$

Step 2: Estimate the probability of \mathcal{J}_2 .

By Lemma 8 and Assumption 1, for any $i \in \{1, 2, \dots, n\}$ and because $\theta \in [0, \pi/2]$, we have

$$\begin{aligned} \Pr_{\mathbf{X}} \{ \arccos(\mathbf{X}_i^T \mathbf{z}_0) > \theta \} &= 1 - \frac{1}{2} I_{\sin^2 \theta} \left(\frac{d-1}{2}, \frac{1}{2} \right) \\ &= 1 - \frac{C_d}{2^{d-1} \Gamma(d/2)} \sin^{d-1} \theta \quad (\text{by Lemma 35}). \end{aligned}$$

Note that since $\Pr_{\mathbf{X}} \{ \arccos(\mathbf{X}_i^T \mathbf{z}_0) > \theta \} \geq 0$, we must have

$$\frac{C_d}{2^{d-1} \Gamma(d/2)} \sin^{d-1} \theta \leq 1. \quad (133)$$

Further, because all \mathbf{X}_i 's are *i.i.d.* for $i \in \{1, 2, \dots, n\}$, we have

$$\Pr_{\mathbf{X}} \{ \theta > \theta_g \} = \Pr_{\mathbf{X}} \left\{ \min_{i \in \{1, 2, \dots, n\}} \arccos(\mathbf{X}_i^T \mathbf{z}_0) > \theta \right\} = \left(1 - \frac{C_d}{2^{d-1} \Gamma(d/2)} \sin^{d-1} \theta \right)^n. \quad (134)$$

By Eq. (129) and Lemma 41, we then have

$$\sin \theta \leq \left(\frac{2^{D-1} \Gamma(D/2)}{C_d} \right)^{\frac{1}{D-1}} n^{-\frac{1}{D-1} (1 - \frac{1}{q})},$$

i.e.,

$$\frac{C_d}{2^{D-1} \Gamma(D/2)} \sin^{D-1} \theta \leq n^{-\frac{1}{D-1} (1 - \frac{1}{q})}.$$

Thus, by Eq. (133), Eq. (134), and Lemma 55, we have

$$\Pr_{\mathbf{X}} [\mathcal{J}_2^c] = \Pr_{\mathbf{X}} \{ \theta > \theta_g \} \leq \exp \left(- \frac{1}{n^{1/q}} \right).$$

Combining the results of Step 1 and Step 2, we thus have

$$\begin{aligned}
 \Pr_{\mathbf{x}, \mathbf{V}_0} [J_1 \setminus J_2] &= \Pr_{\mathbf{x}, \mathbf{V}_0} [J_1 | J_2] \Pr_{\mathbf{x}, \mathbf{V}_0} [J_2] \\
 &= \Pr_{\mathbf{V}_0} [J_1 | J_2] \Pr_{\mathbf{X}} [J_2] \text{ (because of } \mathbf{V}_0 \text{ and } \mathbf{X} \text{ are independent)} \\
 &\quad \left(1 - 2 \exp\left(-\frac{p\theta}{12\pi}\right) \right) \left(1 - \exp\left(-n^{\frac{1}{q}}\right) \right) \\
 &\quad \left(1 - \exp\left(-n^{\frac{1}{q}}\right) - 2 \exp\left(-\frac{p\theta}{12\pi}\right) \right) \\
 &= \left(1 - \exp\left(-n^{\frac{1}{q}}\right) - 2 \exp\left(-\frac{p}{24} \left(\frac{2^{\frac{D-2}{2}}(d-1)}{C_d}\right)^{\frac{1}{d-1}} n^{-\frac{1}{d-1}(1-\frac{1}{q})}\right) \right) \text{ (by Eq. (130)).}
 \end{aligned}$$

By Eq. (60), the conclusion of Proposition 54 thus follows.

L.2. Proof of Lemma 56

Proof. When $\mathbf{x}_1 = \mathbf{z}_0$, the conclusion of this lemma trivially holds because $C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0} = C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0} = ?$ (because $\mathbf{x}^T \mathbf{V}_0[j]$ and $\mathbf{z}_0^T \mathbf{V}_0[j]$ cannot be both positive or negative at the same time when $\mathbf{x}_1 = \mathbf{z}_0$). It remains to consider $\mathbf{x}_1 \neq \mathbf{z}_0$. Define

$$\mathbf{z}_{0,?} := \frac{\mathbf{x}_1 - (\mathbf{x}_1^T \mathbf{z}_0) \mathbf{z}_0}{\|\mathbf{x}_1 - (\mathbf{x}_1^T \mathbf{z}_0) \mathbf{z}_0\|}.$$

Thus, we have $\mathbf{z}_{0,?}^T \mathbf{z}_0 = 0$ and $\|\mathbf{z}_{0,?}\| = 1$, i.e., \mathbf{z}_0 and $\mathbf{z}_{0,?}$ are orthonormal basis vectors on the 2D plane L spanned by \mathbf{x}_1 and \mathbf{z}_0 . Thus, we can represent \mathbf{x}_1 as

$$\mathbf{x}_1 = \cos \varphi \mathbf{z}_0 + \sin \varphi \mathbf{z}_{0,?} \in L.$$

For any $\theta \in [\varphi, \pi]$, we construct \mathbf{x}_2 as

$$\mathbf{x}_2 := \cos \theta \mathbf{z}_0 + \sin \theta \mathbf{z}_{0,?} \in L.$$

In order to show $C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0} \subset C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$, we only need to prove any $j \in C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0}$ must in $C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$. For any $\mathbf{V}_0[j]$, $j = 1, 2, \dots, p$, define the angle $\theta_j \in [0, 2\pi]$ as the angle between \mathbf{z}_0 and $\mathbf{V}_0[j]$'s projected component \mathbf{v}_j on L ¹⁰, i.e.,

$$\mathbf{v}_j = \cos \theta_j \mathbf{z}_0 + \sin \theta_j \mathbf{z}_{0,?} \in L.$$

By the proof of Lemma 17, we know that $j \in C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0}$ if and only if $\theta_j \in (\frac{\pi}{2}, \frac{\pi}{2}) \setminus (\pi + \varphi - \frac{\pi}{2}, \pi + \varphi + \frac{\pi}{2}) \pmod{2\pi}$. Similarly, $j \in C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$ if and only if $\theta_j \in (\frac{\pi}{2}, \frac{\pi}{2}) \setminus (\pi + \theta - \frac{\pi}{2}, \pi + \theta + \frac{\pi}{2}) \pmod{2\pi}$. Because $\varphi \in [0, \pi]$ and $\theta \in [\varphi, \pi]$, we have

$$\left(\frac{\pi}{2}, \frac{\pi}{2} \right) \setminus \left(\pi + \varphi - \frac{\pi}{2}, \pi + \varphi + \frac{\pi}{2} \right) \subset \left(\frac{\pi}{2}, \frac{\pi}{2} \right) \setminus \left(\pi + \theta - \frac{\pi}{2}, \pi + \theta + \frac{\pi}{2} \right) \pmod{2\pi}.$$

Thus, whenever $j \in C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0}$, we must have $j \in C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$. Therefore, we conclude that $C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0} \subset C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$. Using a similar method, we can also show that $C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0} \subset C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$. The result of this lemma thus follows. \square

¹⁰Note that such an angle θ_j is well defined as long as $\mathbf{V}_0[j]$ is not perpendicular to L . The reason that we do not need to worry about those j 's such that $\mathbf{V}_0[j] \perp L$ is as follows. When $\mathbf{V}_0[j] \perp L$, we then have $\mathbf{x}_1^T \mathbf{V}_0[j] = \mathbf{x}_2^T \mathbf{V}_0[j] = \mathbf{z}_0^T \mathbf{V}_0[j] = 0$. Thus, those j 's do not belong to any set $C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0}$, $C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$, $C_{\mathbf{x}_1, \mathbf{z}_0}^{\mathbf{V}_0}$, or $C_{\mathbf{x}_2, \mathbf{z}_0}^{\mathbf{V}_0}$ in Eq. (126).