

Learning with invariances in random features and kernel models

Song Mei

Department of Statistics, University of California, Berkeley

SONGMEI@BERKELEY.EDU

Theodor Misiakiewicz

Department of Statistics, Stanford University

MISIAKIE@STANFORD.EDU

Andrea Montanari

Department of Statistics and Department of Electrical Engineering, Stanford University

MONTANARI@STANFORD.EDU

Editors: Mikhail Belkin and Samory Kpotufe

Abstract

A number of machine learning tasks entail a high degree of invariance: the data distribution does not change if we act on the data with a certain group of transformations. For instance, labels of images are invariant under translations of the images. Certain neural network architectures—for instance, convolutional networks—are believed to owe their success to the fact that they exploit such invariance properties. With the objective of quantifying the gain achieved by invariant architectures, we introduce two classes of models: invariant random features and invariant kernel methods. The latter includes, as a special case, the neural tangent kernel for convolutional networks with global average pooling. We consider uniform covariates distributions on the sphere and hypercube and a general invariant target function. We characterize the test error of invariant methods in a high-dimensional regime in which the sample size and number of hidden units scale as polynomials in the dimension, for a class of groups that we call ‘degeneracy α ’, with $\alpha \leq 1$. We show that exploiting invariance in the architecture saves a d^α factor (d stands for the dimension) in sample size and number of hidden units to achieve the same test error as for unstructured architectures. Finally, we show that output symmetrization of an unstructured kernel estimator does not give a significant statistical improvement; on the other hand, data augmentation with an unstructured kernel estimator is equivalent to an invariant kernel estimator and enjoys the same improvement in statistical efficiency.

Keywords: Invariant function estimation, Random features, Kernel methods, convolutional neural tangent kernel, high dimensional limit

1. Introduction

Consider the following image classification problem. We are given data $\{(\mathbf{x}_i, y_i)\}_{i \leq n}$ where $\mathbf{x}_i \in \mathbb{R}^d$ is an image, and $y_i \in \mathbb{R}$ is its label. We would like to learn a function $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ to predict labels of new unseen images. Throughout this paper we will measure prediction error in terms of the square loss $R(\hat{f}) := \mathbb{E}\{(y_{\text{new}} - \hat{f}(\mathbf{x}_{\text{new}}))^2\}$.

We can think of $\mathbf{x} \in \mathbb{R}^d$ as a pixel representation of an image. For instance if this is a grayscale (one channel) two-dimensional image, \mathbf{x} can represent the pixel values on a $d_1 \times d_2$ grid with $d = d_1 d_2$. For mathematical convenience, we here work with the cartoon example of one-dimensional ‘images’ (or ‘signals’) with d pixels arranged on a line. Most of our results cover two-dimensional images as well.

We assume a model whereby the labels are $y_i = f_*(\mathbf{x}_i) + \varepsilon_i$, with noise ε_i independent of \mathbf{x}_i with $\mathbb{E}(\varepsilon_i) = 0$ and $\mathbb{E}(\varepsilon_i^2) = \sigma_\varepsilon^2$. In many applications, the target function f_* is invariant under translations of the image: if \mathbf{x}' is obtained by translating image \mathbf{x} , then $f_*(\mathbf{x}') = f_*(\mathbf{x})$. We

will consider here periodic shifts (in the case of one-dimensional images): for $\mathbf{x} \in \mathbb{R}^d$, $g_\ell \cdot \mathbf{x} := (x_{\ell+1}, \dots, x_d, x_1, \dots, x_\ell)$ denotes its ℓ -shift. Invariance implies $f_*(\mathbf{x}) = f_*(g_\ell \cdot \mathbf{x})$ for all ℓ and \mathbf{x} .

Convolutional neural networks are the state-of-the-art architecture for image classification and related computer vision tasks, and they are believed to exploit the translation invariance in a crucial way (Krizhevsky et al., 2012). Consider the simple example of two-layer convolutional networks with global average pooling. The network computes a nonlinear convolution of N filters $\mathbf{w}_1, \dots, \mathbf{w}_N$ with the image \mathbf{x} . The results are then combined linearly with coefficients a_1, \dots, a_N :

$$f_{\text{CNN}}(\mathbf{x}) = \frac{1}{d} \sum_{i=1}^N a_i \sum_{\ell=1}^d \sigma(\langle \mathbf{w}_i, g_\ell \cdot \mathbf{x} \rangle). \quad (1)$$

This simple convolutional network can be compared with a standard fully-connected two-layer network with the same number of parameters: $f_{\text{NN}}(\mathbf{x}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle)$. It is clear that —when the target function f_* is translation invariant— the convolutional model $f_{\text{CNN}}(\mathbf{x})$ is at least as powerful as $f_{\text{NN}}(\mathbf{x})$ in terms of approximation, since it is invariant by construction (see Appendix A.1 for a simple formal argument).

The main objective of this paper is to quantify the advantage of architectures —such as convolutional ones— that enforce invariance. We are interested in characterizing the gain both in approximation error and in generalization error. We consider a general type of invariance, defined by a group \mathcal{G}_d that is represented as a subgroup of $\mathcal{O}(d)$, the orthogonal group in d dimensions. This means that each element $g \in \mathcal{G}_d$ is identified with an orthogonal matrix (which we will also denote by g), and group composition corresponds to matrix multiplication. The group element $g \in \mathcal{G}_d$ acts on \mathbb{R}^d via $\mathbf{x} \mapsto g \cdot \mathbf{x}$. We will consider two simple distributions for the signals \mathbf{x} : $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ (the uniform distribution over the sphere in d dimensions with radius \sqrt{d}) and $\mathbf{x} \sim \text{Unif}(\mathcal{Q}^d)$ (with $\mathcal{Q}^d = \{+1, -1\}^d$ the discrete hypercube in d dimensions). We will write $(\mathcal{A}_d, \tau_d) \in \{(\mathbb{S}^{d-1}(\sqrt{d}), \text{Unif}), (\mathcal{Q}^d, \text{Unif})\}$ for either of these two probability spaces. In the case of $\mathcal{A}_d = \mathcal{Q}^d$, we will further require the action of \mathcal{G}_d to preserve \mathcal{Q}^d .

In order to gain some insights on the behavior of actual neural networks, we consider two classes of linear ‘overparametrized’ models: invariant random features models and invariant kernel machines. We next describe these two approaches.

Invariant random feature models. Given an activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ and a group \mathcal{G}_d endowed with invariant (Haar) measure π_d , we define the invariant random features (RF) function class

$$\mathcal{F}_{\text{RF,inv}}^N(\mathbf{W}, \mathcal{G}_d) = \left\{ f(\mathbf{x}) = \sum_{i=1}^N a_i \int_{\mathcal{G}_d} \sigma(\langle \mathbf{w}_i, g \cdot \mathbf{x} \rangle) \pi_d(dg) : a_i \in \mathbb{R}, i \in [N] \right\}. \quad (2)$$

Here $\mathbf{W} := (\mathbf{w}_1, \dots, \mathbf{w}_N)$ is the set of first layer weights which are fixed and not optimized over. We draw them randomly with $(\sqrt{d} \cdot \mathbf{w}_i)_{i \leq N} \sim_{\text{iid}} \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ or $\text{Unif}(\mathcal{Q}^d)$ depending on whether the feature vectors are $\mathbf{x}_i \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ or $\text{Unif}(\mathcal{Q}^d)$. If we let \mathcal{G}_d be the cyclic group $\text{Cyc}_d := \{g_0, g_1, \dots, g_{d-1}\}$ (here g_ℓ is the shift by ℓ positions), we obtain a random features version of the convolutional network of Eq. (1). Other examples will be presented in Section 2.

Given data $\{(\mathbf{x}_i, y_i)\}_{i \leq n}$, we consider to fit the second-layer coefficients $(a_i)_{i \leq N}$ in Eq. (2) using the random features ridge regression (RFRR). Notice that the estimated function \hat{f} is invariant by construction, $\hat{f}(\mathbf{x}) = \hat{f}(g \cdot \mathbf{x})$. We will denote the space of square integrable \mathcal{G}_d -invariant functions on $\mathcal{A}_d \in \{\mathbb{S}^{d-1}(\sqrt{d}), \mathcal{Q}^d\}$ by $L^2(\mathcal{A}_d, \mathcal{G}_d)$.

Invariant kernel machines. We then consider kernel ridge regression (KRR) in the reproducing kernel Hilbert space (RKHS) defined by a \mathcal{G}_d -invariant kernel. By this we mean a kernel $H \in L^2(\mathcal{A}_d \times \mathcal{A}_d)$ such that, for all $g, g' \in \mathcal{G}_d$, the following holds for every $\mathbf{x}_1, \mathbf{x}_2$:

$$H(\mathbf{x}_1, \mathbf{x}_2) = H(g \cdot \mathbf{x}_1, g' \cdot \mathbf{x}_2). \quad (3)$$

Note that, as a consequence of this property, any function that is not in $L^2(\mathcal{A}_d, \mathcal{G}_d)$ (i.e. any function that is not invariant) has infinite RKHS norm: indeed this provides an alternate characterization of invariant kernel methods. Among \mathcal{G}_d -invariant kernels, we focus on the subclass that is obtained by averaging an inner product kernel over the group \mathcal{G}_d

$$H_{\text{inv}}(\mathbf{x}_1, \mathbf{x}_2) = \int_{\mathcal{G}_d} h(\langle \mathbf{x}_1, g \cdot \mathbf{x}_2 \rangle / d) \pi_d(dg). \quad (4)$$

Invariant kernel machines can be regarded as large-width ($N \rightarrow \infty$) limits of invariant random features methods. Vice versa, the latter can be regarded as randomized approximations of invariant kernel methods. Moreover, invariant kernel methods also capture the large-width limits of other models, for instance, neural tangent models associated to convolutional networks (c.f. Section A.3).

We focus on a type of groups \mathcal{G}_d that we call *groups of degeneracy α* .

Definition 1 (Groups of degeneracy α) *Let $V_{d,k}$ be the subspace of degree- k polynomials that are orthogonal to polynomials of degree at most $(k-1)$ in $L^2(\mathcal{A}_d)$, and denote by $V_{d,k}(\mathcal{G}_d)$ the subspace of $V_{d,k}$ formed by polynomials that are \mathcal{G}_d -invariant. We say that \mathcal{G}_d has degeneracy α if for any integer $k \geq \alpha$ we have $\dim(V_{d,k}) / \dim(V_{d,k}(\mathcal{G}_d)) \asymp d^\alpha$ (i.e., there exists $0 < c_k \leq C_k < \infty$ such that $c_k \leq \dim(V_{d,k}) / \dim(V_{d,k}(\mathcal{G}_d)) / d^\alpha \leq C_k$ for any $d \geq 2$).*

This definition includes as special cases the cyclic group for one and two-dimensional signals (see Section 2), which have both degeneracy 1. Note that we can define an equivalence relation between degree- k polynomials: for $p_k, p'_k \in V_{d,k}$, we have $p_k \sim p'_k$ if and only if there exists $g \in \mathcal{G}_d$ such that $p_k(g \cdot \mathbf{x}) = p'_k(\mathbf{x})$. The dimension of the quotient space $\dim(V_{d,k} / V_{d,k}(\mathcal{G}_d))$ is then exactly equal to the ratio $\dim(V_{d,k}) / \dim(V_{d,k}(\mathcal{G}_d))$. For a group \mathcal{G}_d with degeneracy α , we can think about d^α as the ‘effective dimension’ of the group seen through its action on polynomials. The effective dimension of the group is not necessary equal to the size of the group (e.g., see Example 3 which is an infinite group with degeneracy 1). We will see below that this effective dimension is exactly equal to the factor that we save in sample size and number of hidden units by using invariant architectures.

We compare invariant methods to standard (non-invariant) random features models with inner product activation, defined as

$$\mathcal{F}_{\text{RF}}^N(\mathbf{W}) = \left\{ f(\mathbf{x}) = \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) : a_i \in \mathbb{R}, i \in [N] \right\}, \quad (5)$$

and standard inner product kernels $H(\mathbf{x}_1, \mathbf{x}_2) = h_d(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d)$. For groups with degeneracy $\alpha \leq 1$, we obtain a fairly complete characterization of the gain achieved by using invariant models, when the target function is an arbitrary invariant function $f_* \in L^2(\mathcal{A}_d; \mathcal{G}_d)$.

Invariance gain: underparametrized case. Consider the invariant RF class (2) in the underparametrized regime $N \ll n$. We prove that the test error is dominated by the approximation error. Namely,

if $d^{\ell-\alpha} \ll N \ll d^{\ell+1-\alpha}$, then the test error (c.f. Eq. (7)) gives $R(f_*; \lambda) \approx \|\bar{P}_{>\ell} f_*\|_{L^2}^2$, where $\bar{P}_{>\ell}$ is the projection orthogonal to the subspace of degree ℓ polynomials. In order to achieve the same risk, standard (non-invariant) RF models would require $d^\ell \ll N \ll d^{\ell+1}$: invariance saves a d^α factor in the network width to achieve the same risk.

Invariance gain: overparametrized case. Consider next the overparametrized regime $n \ll N$. In this case the test error is dominated by the statistical error. Namely, if $d^{\ell-\alpha} \ll n \ll d^{\ell+1-\alpha}$, then the test error gives $R(f_*; \lambda) \approx \|\bar{P}_{>\ell} f_*\|_{L^2}^2$. In order to achieve the same risk, standard (non-invariant) RF models would require $d^\ell \ll n \ll d^{\ell+1}$: invariance saves a d^α factor in the sample size to achieve the same risk.

These results are precisely presented in Theorem 2 and summarized in Table 1. We establish the same gain for invariant kernel methods in Theorem 5. While we focused in this paper on groups with degeneracy $\alpha \leq 1$ (which include our primary motivating examples, cyclic group in one or two dimensions), we expect similar results to hold for groups with $\alpha > 1$ (indeed our current proof techniques can handle the case $\alpha > 1$ at the price of adding the condition $N, n \geq d^{O(\alpha)}$ in our theorems). We defer this to future work.

Output symmetrization and data augmentation. Output symmetrization and data augmentation are two alternative approaches to incorporate invariances in machine learning models. We show that the performance of output symmetrization of standard KRR does not improve over standard KRR, and hence is sub-optimal compared to invariant KRR. On the other hand, it was shown that (c.f. Li et al. (2019)) data augmentation is mathematically equivalent to invariant KRR for discrete groups. As a consequence, our theoretical results characterize the statistical gain by performing data augmentation.

It is important to mention that our treatment omits an important characteristic of convolutional architectures: the fact that the filters w_i of Eq. (1) have a short window size $q \ll d$. Namely, they have only q non-zero entries, for instance the first q entries. Using short-window filters has some interesting consequences, which can be investigated using the same approach developed here. We will report on these in a forthcoming article, and instead focus here on the impact of invariance.

Our analysis is enabled by a simple yet important observation, which might generalize to other settings. The subspaces $V_{d,k}$ of degree- k polynomials (see Definition 1) are eigenspaces for inner product kernels. At the same time, they are preserved under the symmetry group \mathcal{G}_d . Namely, define $f^{(g)}(\mathbf{x}) = f(g \cdot \mathbf{x})$, we have $f^{(g)} \in V_{d,k}$ for any $f \in V_{d,k}$, $g \in \mathcal{G}_d$. This observation is crucial in determining the eigendecomposition of the relevant kernels.

Let us finally emphasize, that the factor- d gain in sample size for degeneracy-one groups is not correctly predicted by a naive ‘data augmentation heuristics’. The latter would suggest a gain of the order of $|\mathcal{G}_d|$ or of the size of orbits of \mathcal{G}_d . As shown by the example of band limited functions (see below) $|\mathcal{G}_d|$ can be ∞ but the degeneracy can still be one (and hence the gain is d).

1.1. Related literature

Invariant function estimation

A number of mathematical works emphasized the role of invariance in neural network architectures. Among others, Mallat (2012); Bruna and Mallat (2013); Mallat (2016) propose architectures

To fit a degree ℓ polynomial	Inner product random features	Invariant random features
Underparameterized regime ($N \ll n$)	$N \gg d^\ell$	$N \gg d^{\ell-\alpha}$
Overparameterized regime ($n \ll N$)	$n \gg d^\ell$	$n \gg d^{\ell-\alpha}$

Table 1: Sample size n and number of features N required to fit a \mathcal{G}_d -invariant polynomial of degree ℓ using ridge regression with the standard random features model (Eq. (5)) and the invariant random features model (Eq. (2)), for group \mathcal{G}_d of degeneracy $\alpha \leq 1$.

(‘deep scattering networks’) that explicitly achieve invariance to a rich group of transformations. However, these papers do not characterize the statistical error of these approaches.

The recent paper Li et al. (2020) constructs a simple data distribution on which a gap is proven between the sample complexity for convolutional architectures, and the one for standard (fully connected) architectures. This result differs from ours in several aspects. Most importantly, we study the risk for estimating general invariant functions using invariant kernels and random features, while Li et al. (2020) obtain results for a specific distribution using CNNs. Also, the weight sharing structure in Li et al. (2020) is different from the one in Eq. (1).

Another work Chen et al. (2020) studied the statistical benefits of data augmentation in the parametric setting via a group theory framework. Our result is different in the sense that we consider the non-parametric setting to estimate an invariant function using kernel methods.

To the best of our knowledge, our paper is the first that characterizes the precise statistical benefit of using invariant random features and kernel models.

Convolutional neural networks and convolutional kernels

A recent line of work (Jacot et al., 2018; Li and Liang, 2018; Du et al., 2019b,a; Allen-Zhu et al., 2019b,a; Arora et al., 2019a; Zou et al., 2020; Oymak and Soltanolkotabi, 2020) studied the training dynamics of overparametrized neural networks under certain random initialization, and showed that it converges to a kernel estimator, which corresponds to the “neural tangent kernel”. The convolutional neural tangent kernel, which corresponds to the tangent kernel of convolutional neural networks, was studied in Arora et al. (2019b); Li et al. (2019); Bietti and Mairal (2019). The connection between convolutional kernel ridge regression and data augmentation was pointed out in Li et al. (2019).

The network in Eq. (1) corresponds to a two-layer convolutional neural network with global average pooling, which is a special case of the convolutional network that was defined as in Arora et al. (2019b).

Random features and kernel methods

A number of authors have studied the generalization error of kernel machines (Caponnetto and De Vito, 2007; Jacot et al., 2020; Liang et al., 2020b,a) (Wainwright, 2019, Theorem 13.17) and random features models (Rahimi and Recht, 2009; Rudi and Rosasco, 2017; Ma et al., 2020; Bach, 2017). However, these results are not fine-grained enough to characterize the separation between invariant kernels (or random feature models) and standard inner product kernels, for several reasons. First, some of these results concern restricted target functions with bounded RKHS norm. Second, we establish a gap that holds pointwise, i.e. for any given target function f_* , while most of earlier

work only obtain minimax lower bounds. Finally, we need the upper and lower bounds match up to a $1 + o_d(1)$ factor, while earlier results only match up to unspecified constants.

The recent paper [Jacot et al. \(2020\)](#) provides sharp predictions for kernel machines, but it assumes that a certain random kernel matrix behaves like a random matrix with Gaussian components: proving an equivalence of this type is the central mathematical challenge we face here.

Our analysis builds on the general results of [Ghorbani et al. \(2021\)](#); [Mei et al. \(2021\)](#). In particular, [Mei et al. \(2021\)](#) provides general conditions under which the risk of random features and kernel methods can be characterized precisely. Checking these conditions for invariant methods requires to prove certain concentration properties for the entries of the relevant kernels. We achieve this goal for the cyclic group with general activations, and for degeneracy- α groups (for $\alpha \leq 1$) with polynomial activations. Generalizing these results to other groups, data distributions, and activations is a promising direction.

2. Examples

In this section, we provide three examples of our general setting. We show in Appendix D that all these groups have degeneracy 1 and therefore satisfy the assumptions of our general theorems.

Example 1 (One-dimensional images) *The cyclic group has elements $\text{Cyc}_d = \{g_0, g_1, \dots, g_{d-1}\}$ where g_i is a shift by i pixels. For any $\mathbf{x} = (x_1, \dots, x_d)^\top \in \mathcal{A}_d$, the action of group element g_i on \mathbf{x} is defined by $g_i \cdot \mathbf{x} = (x_{i+1}, x_{i+2}, \dots, x_d, x_1, x_2, \dots, x_i)^\top \in \mathcal{A}_d$. (In particular, g_i is identified with an orthogonal transformation in \mathbb{R}^d .) The measure π_d is the uniform probability measure on Cyc_d , i.e.,*

$$\int_{\text{Cyc}_d} f(g) \pi_d(dg) = \frac{1}{d} \sum_{i=0}^{d-1} f(g_i).$$

We will refer to the invariant functions $L^2(\mathcal{A}_d, \text{Cyc}_d)$ as the ‘cyclic functions’.

Example 2 (Two-dimensional images) *Let $d = d_1 \times d_2$. We identify $\mathcal{X}_{d_1 \times d_2} = \{\mathbf{X} \in \mathbb{R}^{d_1 \times d_2} : \|\mathbf{X}\|_F^2 = d\}$ with $\mathbb{S}^{d-1}(\sqrt{d})$ (simply by ‘vectorizing’ the matrix). The two-direction cyclic group has elements $\text{Cyc2D}_{d_1, d_2} = \{g_{ij} : 0 \leq i < d_1, 0 \leq j < d_2\}$. For any $\mathbf{X} = (X_{ij})_{i \in [d_1], j \in [d_2]} \in \mathcal{X}_{d_1 \times d_2}$, the action of group element $g_{ij} \in \text{Cyc2D}_{d_1, d_2}$ on \mathbf{X} is defined by*

$$g_{ij} \cdot \mathbf{X} = \begin{bmatrix} X_{i+1, j+1} & \dots & X_{i+1, d_2} & X_{i+1, 1} & \dots & X_{i+1, j} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ X_{d_1, j+1} & \dots & X_{d_1, d_2} & X_{d_1, 1} & \dots & X_{d_1, j} \\ X_{1, j+1} & \dots & X_{1, d_2} & X_{1, 1} & \dots & X_{1, j} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ X_{i, j+1} & \dots & X_{i, d_2} & X_{i, 1} & \dots & X_{i, j} \end{bmatrix}.$$

Again, this is an orthogonal transformation in $\mathcal{X}_{d_1 \times d_2} \cong \mathbb{S}^{d-1}(\sqrt{d})$, and $\text{Cyc2D}_{d_1, d_2}$ is isomorphic to a subgroup of $\mathcal{O}(d)$. The measure π_d is the uniform probability measure on $\text{Cyc2D}_{d_1, d_2}$. We will refer to the invariant functions $L^2(\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc2D}_d)$ as the ‘two-direction cyclic functions’.

Example 3 (The translation invariant function class on band-limited signals) *Suppose we have one-dimensional signals with very high resolution, but the signals are band-limited: their Fourier*

transforms have only d non-zero coefficients. We assume that the labels of the band-limited signals are invariant under translations. The following model captures this setting.

Let $\{\varphi_j\}_{j \in [d]} \subseteq \mathcal{F}([0, 1])$ be the real Fourier basis functions in $L^2([0, 1], \text{Unif})$. That is, we define $\varphi_1(t) = 1$, and for $p = 1, 2, \dots, \lfloor d/2 \rfloor$ (we assume d is odd), $\varphi_{2p}(t) = \sqrt{2} \cos(2\pi pt)$, $\varphi_{2p+1}(t) = \sqrt{2} \sin(2\pi pt)$. We define the band-limited covariate subspace $\mathbb{W}_d \subseteq L^2([0, 1], \text{Unif})$ to be (\mathbb{W} stands for waves)

$$\mathbb{W}_d = \left\{ x \in L^2([0, 1]) : x(t) = \sum_{j=1}^d \hat{x}_j \varphi_j(t), \hat{\mathbf{x}} = (\hat{x}_1, \dots, \hat{x}_d) \in \mathbb{S}^{d-1}(\sqrt{d}) \right\}.$$

Then the space \mathbb{W}_d can be identified with the space $\mathbb{S}^{d-1}(\sqrt{d})$.

Let $\text{Sft}_d = \{g_u, u \in [0, 1]\} \simeq \text{SO}(2)$ be the translation group that can act on \mathbb{W}_d . For any $x \in \mathbb{W}_d$, the action of group element $g_u \in \text{Sft}_d$ on x is defined by

$$[g_u \cdot x](t) = x(t - u).$$

Equivalently, the action of group element $g_u \in \text{Sft}_d$ on $\hat{\mathbf{x}} \in \mathbb{S}^{d-1}(\sqrt{d})$ is defined by

$$g_u \cdot \hat{\mathbf{x}} = (\hat{x}_1, \cos(2\pi u)\hat{x}_2 + \sin(2\pi u)\hat{x}_3, -\sin(2\pi u)\hat{x}_2 + \cos(2\pi u)\hat{x}_3, \dots).$$

That means, Sft_d can be interpreted as a subgroup of $\mathcal{O}(d)$. The measure π_d is the uniform distribution on Sft_d , i.e.,

$$\int_{\text{Sft}_d} f(g) \pi_d(dg) = \int_{[0, 1]} f(g_s) ds.$$

The function class $L^2(\mathbb{W}_d, \text{Sft}_d)$, or equivalently $L^2(\mathbb{S}^{d-1}(\sqrt{d}), \text{SO}(2))$, can be regarded as the translation invariant function class on band-limited signals.

3. Invariant random feature models

Let \mathcal{G}_d be a group of degeneracy α with $\alpha \leq 1$ as defined in Definition 1 and f_d be a function that is invariant under the action of \mathcal{G}_d , i.e., $f_d \in L^2(\mathcal{A}_d, \mathcal{G}_d)$. We consider fitting the data with the invariant random features model defined in Eq. (2) using ridge regression, which we call invariant RFRR. Namely, we learn a function $\hat{f}_{N, \lambda}^{\text{inv}}(\mathbf{x}; \hat{\mathbf{a}}(\lambda)) = \sum_{1 \leq j \leq N} \hat{a}_j \int_{\mathcal{G}_d} \sigma(\langle \mathbf{w}_j, g \cdot \mathbf{x} \rangle) \pi_d(dg)$ with

$$\hat{\mathbf{a}}(\lambda) = \arg \min_{\mathbf{a}} \left\{ \sum_{i=1}^n (y_i - \hat{f}_{N, \lambda}^{\text{inv}}(\mathbf{x}_i; \mathbf{a}))^2 + \frac{N\lambda}{d^\alpha} \|\mathbf{a}\|_2^2 \right\}, \quad (6)$$

where the regularization parameter λ can depend on the dimension d . (The factor d^α in the ridge penalty is introduced to compensate for the effect of averaging the random features over \mathcal{G}_d .) We further denote the test error of invariant RFRR by

$$R_{\text{RF,inv}}(f_d, \mathbf{X}, \mathbf{W}, \varepsilon, \lambda) := \mathbb{E}_{\mathbf{x}} \left[\left(f_d(\mathbf{x}) - \hat{f}_{N, \lambda}^{\text{inv}}(\mathbf{x}; \hat{\mathbf{a}}(\lambda)) \right)^2 \right]. \quad (7)$$

We will make the following assumption on σ .

Assumption 1 (Conditions on σ , n , N , and $(\mathcal{A}_d, \mathcal{G}_d)$ at level $(s, S) \in \mathbb{N}^2$) For $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, we assume the following conditions hold.

(a) For $(\mathcal{A}_d, \mathcal{G}_d) = (\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$, we assume σ to be $(\min(s, S) + 1) \vee 3$ differentiable and there exists constants $c_0 > 0$ and $c_1 < 1$ such that $|\sigma^{(k)}(u)| \leq c_0 e^{c_1 u^2/2}$ for any $2 \leq k \leq (\min(s, S) + 1) \vee 3$. Moreover, there exists an integer $p > 1/\delta$ such that $n \leq N^{1-\delta}$ or $N \leq n^{1-\delta}$ and $|\sigma(x)|, |\sigma'(x)| \leq c_0 \exp(c_1 x^2/(8p))$.

For general $(\mathcal{A}_d, \mathcal{G}_d)$, we assume that σ is a (finite degree) polynomial function.

(b) The Hermite coefficients $\mu_k(\sigma) \equiv \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)\text{He}_k(G)]$ verify $\mu_k(\sigma) \neq 0$ for any $0 \leq k \leq \min(s, S)$ (see Appendix H for definitions).

(c) We assume that σ is not a polynomial with degree less or equal to $\max(s, S)$.

For $k \in \mathbb{N}$, we denote by $\bar{\mathbb{P}}_{\leq k} : L^2(\mathcal{A}_d) \rightarrow L^2(\mathcal{A}_d)$ the orthogonal projection operator onto the subspace of polynomials of degree at most k , and $\bar{\mathbb{P}}_{> k} = \mathbf{I} - \bar{\mathbb{P}}_{\leq k}$ (see Appendix H for details). We denote $f(d) = o_{d, \mathbb{P}}(g(d))$ if $f(d)/g(d)$ converges to 0 in probability as $d \rightarrow \infty$.

Theorem 2 (Test error of invariant RFRR) Let \mathcal{G}_d be a group of degeneracy $\alpha \leq 1$ and let $\{f_d \in L^2(\mathcal{A}_d, \mathcal{G}_d)\}_{d \geq 1}$ be a sequence of \mathcal{G}_d -invariant functions. Assume $d^{s-\alpha+\delta} \leq n \leq d^{s+1-\alpha-\delta}$ and $d^{S-\alpha+\delta} \leq N \leq d^{S+1-\alpha-\delta}$ for fixed integers s, S and some $\delta > 0$. Let σ be an activation function that satisfies Assumption 1 at level (s, S) . Then the following hold for the test error of invariant RFRR (see Eq. (7)):

(a) (Overparametrized regime) Assume $N \geq nd^\delta$ for some $\delta > 0$. Then for any regularization parameter $\lambda = O_d(1)$ (including $\lambda = 0$) and any $\eta > 0$, we have

$$R_{\text{RF,inv}}(f_d, \mathbf{X}, \mathbf{W}, \varepsilon, \lambda) = \|\bar{\mathbb{P}}_{> s} f_d\|_{L^2}^2 + o_{d, \mathbb{P}}(1) \cdot (\|f_d\|_{L^{2+\eta}}^2 + \sigma_\varepsilon^2). \quad (8)$$

(b) (Underparametrized regime) Assume $n \geq Nd^\delta$ for some $\delta > 0$. Then for any regularization parameter $\lambda = O_d(n/N)$ (including $\lambda = 0$) and any $\eta > 0$, we have,

$$R_{\text{RF,inv}}(f_d, \mathbf{X}, \mathbf{W}, \varepsilon, \lambda) = \|\bar{\mathbb{P}}_{> S} f_d\|_{L^2}^2 + o_{d, \mathbb{P}}(1) \cdot (\|f_d\|_{L^{2+\eta}}^2 + \sigma_\varepsilon^2). \quad (9)$$

In particular, this theorem applies to the one-dimensional and two-dimensional cyclic groups, and band-limited functions listed in Section 2. We refer readers to Appendix A.2 for an informal intuition and Appendix B.2 for the proof of this result.

We can compare these bounds with ridge regression on the standard random features model of Eq. (5). Theorem 2 in Mei et al. (2021) (with Assumption 1) shows that the same test error holds as in Theorem 2 but with $d^{s+\delta} \leq n \leq d^{s+1-\delta}$ and $d^{S+\delta} \leq N \leq d^{S+1-\delta}$. We thus gain a factor d^α in the sample and feature complexity by using invariant features compared to non invariant ones.

Remark 3 Assumption 1 requires the activation function to be polynomial, except for the cyclic group, for which only differentiability conditions are assumed. These conditions are sufficient for the general assumptions in Mei et al. (2021) to hold: for the sake of length, we only verify them for non-polynomial activation functions in the case of the cyclic group in one dimension. However we believe that the differentiability condition (and indeed weaker conditions) should be sufficient

for general groups. For example, the current proofs can be modified to apply to more general subgroups of the permutation group on d elements (e.g., cyclic group in higher dimension). We defer these improvements to future work.

For the cyclic group, the current assumptions already include the interesting examples of the sigmoid $\sigma(x) = 1/(1 + \exp(x - c))$ and smoothed ReLU $\sigma(x) = \mathbb{E}_{G \sim \mathcal{N}(0, \varepsilon^2)}[(x - c + G)_+]$.

Note that Assumption 1.b) is necessary for the RKHS associated to the feature map σ to include all polynomials of degree less or equal to $\min(\mathfrak{s}, S)$.

Remark 4 Consider two-dimensional images with $d = D \times D$ (Example 2) and functions f_d that are invariant with respect to the group of cyclic translations along the horizontal direction only. It can be shown that this group has degeneracy $\alpha = 1/2$, and in fact $\dim(V_{d,k})/\dim(V_{d,k}(\mathcal{G}_d)) \asymp D = d^{1/2}$. Our theory also applies to this group.

4. Invariant kernel machines

Note that any invariant kernel of the form (4) can be written as a kernel of the form:

$$H_{d,\text{inv}}(\mathbf{x}_1, \mathbf{x}_2) = \int_{\mathcal{G}_d} \mathbb{E}_{\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})} [\sigma(\langle \mathbf{x}_1, \mathbf{w} \rangle) \sigma(\langle \mathbf{x}_2, g \cdot \mathbf{w} \rangle)] \pi_d(dg). \quad (10)$$

To see this, note that any inner product kernel h can be decomposed as

$$h(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d) = \mathbb{E}_{\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})} [\sigma(\langle \mathbf{x}_1, \mathbf{w} \rangle) \sigma(\langle \mathbf{x}_2, \mathbf{w} \rangle)]$$

for some activation function σ , which amounts to taking the square root of the positive semidefinite operator associated to h . Substituting in Eq. (4), we get the desired representation.

Consider Kernel ridge regression with regularization parameter λ associated to $H_{d,\text{inv}}$, that we call invariant KRR. Namely, we learn a function $\hat{f}_\lambda^{\text{inv}}(\mathbf{x}; \hat{\mathbf{u}}(\lambda)) = \sum_{i \in [n]} \hat{u}_i H_{d,\text{inv}}(\mathbf{x}_i, \mathbf{x})$ where

$$\hat{\mathbf{u}}(\lambda) = \arg \min_{\mathbf{u}} \left\{ \sum_{i=1}^n (y_i - \hat{f}_\lambda^{\text{inv}}(\mathbf{x}_i; \mathbf{u}))^2 + \frac{\lambda}{d^\alpha} \|\hat{f}_\lambda^{\text{inv}}(\cdot; \mathbf{u})\|_{\mathcal{H}}^2 \right\}. \quad (11)$$

with $\|\cdot\|_{\mathcal{H}}$ the RKHS norm associated to $H_{d,\text{inv}}$. We further denote the test error of invariant KRR by

$$R_{\text{KR,inv}}(f_d, \mathbf{X}, \varepsilon, \lambda) := \mathbb{E}_{\mathbf{x}} \left[\left(f_d(\mathbf{x}) - \hat{f}_\lambda^{\text{inv}}(\mathbf{x}; \hat{\mathbf{u}}) \right)^2 \right]. \quad (12)$$

Theorem 5 (Test error of invariant KRR) Let \mathcal{G}_d be a group of degeneracy $\alpha \leq 1$ and $\{f_d \in L^2(\mathcal{A}_d, \mathcal{G}_d)\}_{d \geq 1}$ be a sequence of \mathcal{G}_d -invariant functions. Assume $d^{\mathfrak{s}-\alpha+\delta} \leq n \leq d^{\mathfrak{s}+1-\alpha-\delta}$ for some fixed integer $\mathfrak{s} \geq 1$ and some $\delta > 0$. Let σ be an activation function that satisfies Assumption 1 at level $(\mathfrak{s}, \mathfrak{s})$ (and $N = \infty$) and let $H_{d,\text{inv}}$ be the associated invariant kernel as defined in Eq. (10). Then, the following holds for the test error of invariant KRR (c.f. Eq. (12)): for any $\lambda = O_d(1)$ (including $\lambda = 0$ identically) any $\eta > 0$, we have

$$R_{\text{KR,inv}}(f_d, \mathbf{X}, \varepsilon, \lambda) = \|\bar{\mathbf{P}}_{>\mathfrak{s}} f_d\|_{L^2}^2 + o_{d,\mathbb{P}}(1) \cdot (\|f_d\|_{L^{2+\eta}}^2 + \sigma_\varepsilon^2). \quad (13)$$

We can compare the performance of this kernel against a standard (inner product) kernel $H_d(\mathbf{x}, \mathbf{y}) = h_d(\langle \mathbf{x}, \mathbf{y} \rangle / d)$. Then Theorem 4 in [Ghorbani et al. \(2021\)](#) shows that the above theorem holds but with $d^{s+\delta} \leq n \leq d^{s+1-\delta}$. We gain a factor d^α in sample complexity by using an invariant kernel.

Remark 6 Recall that the neural tangent kernel (NTK) associated to a function $f(\mathbf{x}; \Theta)$ with random initialization Θ_0 is defined as

$$H_{\text{NT}}(\mathbf{x}, \mathbf{y}) := \mathbb{E}_{\Theta_0} \left[\langle \nabla_{\Theta} f(\mathbf{x}; \Theta_0), \nabla_{\Theta} f(\mathbf{y}; \Theta_0) \rangle \right].$$

The neural tangent kernel associated to a multi-layers fully connected network is an inner-product kernel (as long as the weights are initialized to be isotropic Gaussian.) In contrast, the NTK associated to the CNN of Eq. (1) is an example of invariant kernel, and is covered by Theorem 5 (see Appendix A.3 for more details).

5. Comparison with alternative approaches

To provide further context, it is useful to compare invariant random features and kernel models with other approaches. Here we consider two alternatives: (i) *output symmetrization*, which uses a non-invariant method for training and then symmetrizes the estimated function over the group \mathcal{G}_d to obtain an invariant function; (ii) *data augmentation*, which trains the model on a dataset augmented by samples obtained by applying group transformations to the original data. As shown in [Li et al. \(2019\)](#), data augmentation is mathematically equivalent to invariant kernel methods, so that it is superior to standard kernel methods (with inner-product kernels). On the other hand, we show that output symmetrization of standard kernel estimators does not significantly improve over the standard kernel estimator, and is fundamentally sub-optimal comparing to invariant kernel methods.

5.1. Output symmetrization

Given an estimator \hat{f} , the symmetrization operator $\mathcal{S}\hat{f}$ computes the average of \hat{f} over the group:

$$(\mathcal{S}\hat{f})(\mathbf{x}) \equiv \int_{\mathcal{G}_d} \hat{f}(g \cdot \mathbf{x}) \pi_d(dg). \quad (14)$$

When the target function f_d is \mathcal{G}_d -invariant, one might naively think that the symmetrization operation will significantly improve the performance of standard kernel estimators (standard RFRR and KRR). Indeed, when $f_d \in L^2(\mathcal{A}_d, \mathcal{G}_d)$, Jensen's inequality gives $\|f_d - \mathcal{S}\hat{f}\|_{L^2}^2 = \|\mathcal{S}(f_d - \hat{f})\|_{L^2}^2 \leq \|f_d - \hat{f}\|_{L^2}^2$. However, the proposition below (which is proved in Section A.4) shows that $\mathcal{S}\hat{f}$ is not significantly better when \hat{f} is a standard kernel estimator.

Proposition 7 Let $f_d \in L^2(\mathcal{A}_d, \mathcal{G}_d)$ be a sequence of target functions. For any sequence of estimators \hat{f}_d satisfying $\|\hat{f}_d - \bar{\mathbb{P}}_{\leq \ell} f_d\|_{L^2}^2 \leq \varepsilon$, we have

$$\begin{aligned} \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2}^2 - 2\varepsilon \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2} &\leq \|f_d - \mathcal{S}\hat{f}_d\|_{L^2}^2 \\ &\leq \|f_d - \hat{f}_d\|_{L^2}^2 \leq \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2}^2 + 2\varepsilon \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2} + \varepsilon^2. \end{aligned} \quad (15)$$

Now consider —to be definite— a setting in which $N \geq nd^\delta$ and $d^{\ell+\delta} \leq n \leq d^{\ell+1-\delta}$, and \mathcal{G}_d is a group with degeneracy 1. For any $f_d \in L^2(\mathcal{A}_d, \mathcal{G}_d)$ with $\|f_d\|_{L^{2+\eta}}^2 = O_d(1)$, the results of [Mei et al. \(2021\)](#) imply that standard RFRR (c.f. Eq. (5)) with sufficiently small regularization returns a function \hat{f}_{RF} with $\|\bar{\mathbb{P}}_{\leq \ell} f_d - \hat{f}_{\text{RF}}\|_{L^2}^2 = o_{d,\mathbb{P}}(1)$. Consequently, Proposition 7 implies that we have

$$\|f_d - \mathcal{S}\hat{f}_{\text{RF}}\|_{L^2}^2 = \|f_d - \hat{f}_{\text{RF}}\|_{L^2}^2 + o_{d,\mathbb{P}}(1) = \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2}^2 + o_{d,\mathbb{P}}(1),$$

while Theorem 2 implies that invariant RFRR $\hat{f}_{\text{RF}}^{\text{inv}}$ with sufficiently small regularization achieves a substantially smaller risk:

$$\|f_d - \hat{f}_{\text{RF}}^{\text{inv}}\|_{L^2}^2 = \|\bar{\mathbb{P}}_{> \ell+1} f_d\|_{L^2}^2 + o_{d,\mathbb{P}}(1).$$

5.2. Data augmentation

We consider full data augmentation whereby we replace each sample (y_i, \mathbf{x}_i) in the dataset by $|\mathcal{G}_d|$ samples $\{(y_i, g \cdot \mathbf{x}_i) : g \in \mathcal{G}_d\}$ (for simplicity we consider here the case of a finite group \mathcal{G}_d), and perform standard KRR on the augmented dataset. One might naively think that this is not as effective as enforcing invariance in the kernel structure. After all, we are only requiring invariance to hold at the sampled points. However, [Li et al. \(2019\)](#) showed that these two approaches are in fact equivalent.

We compare KRR using the kernel $H(\mathbf{x}, \mathbf{y}) = h(\langle \mathbf{x}, \mathbf{y} \rangle / d)$ on the augmented dataset, with invariant KRR on the original dataset using the symmetrized kernel $H_{\text{inv}}(\mathbf{x}, \mathbf{y}) = \int_{\mathcal{G}_d} h(\langle \mathbf{x}, g \cdot \mathbf{y} \rangle / d) \pi_d(dg)$. Denote by $\hat{f}_\lambda^{\text{data}}$ and $\hat{f}_\lambda^{\text{inv}}$ the KRR estimates with the standard kernel H and full data augmentation, and with the invariant kernel H_{inv} respectively.

Proposition 8 ([Li et al. \(2019\)](#)) *Let \mathcal{G} be a finite group, and H, H_{inv} as defined above. Then we have $\hat{f}_\lambda^{\text{data}} = \hat{f}_\lambda^{\text{inv}}$.*

A couple of remarks are in order. First, this equivalence is general (holds for any dataset $\{(y_i, \mathbf{x}_i)\}_{i \leq n}$), and is in fact a consequence of the algebraic structure of ridge regressions. Second, while this result establishes that the two approaches are mathematically equivalent, there are computational advantages for invariant KRR. Indeed, full data augmentation increases the size of the kernel matrix from n to $n|\mathcal{G}_d|$ which is computationally more expensive. Finally, this equivalence shows that data augmentation with standard KRR is superior to output symmetrization of standard KRR.

6. Numerical illustration

To check our predictions, we first consider the setting of $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ with $d = 30$, and three cyclic invariant polynomials $f_{d,\text{lin}}, f_{d,\text{quad}}, f_{d,\text{cube}} \in L^2(\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$ defined as

$$f_{d,\text{lin}} = \frac{1}{\sqrt{d}} \sum_{i=1}^d x_i, \quad f_{d,\text{quad}} = \frac{1}{\sqrt{d}} \sum_{i=1}^d x_i x_{i+1}, \quad f_{d,\text{cube}} = \frac{1}{\sqrt{d}} \sum_{i=1}^d x_i x_{i+1} x_{i+2}, \quad (16)$$

where the sub-index i in x_i should be understood in the modulo d sense ($d+1 = 1 \pmod{d}$). We compare the performance between two kernels: a standard (inner product) kernel $H_d(\mathbf{x}, \mathbf{y}) := h_d(\langle \mathbf{x}, \mathbf{y} \rangle / d)$ that we take to be the neural tangent kernel associated to a depth-5 neural network

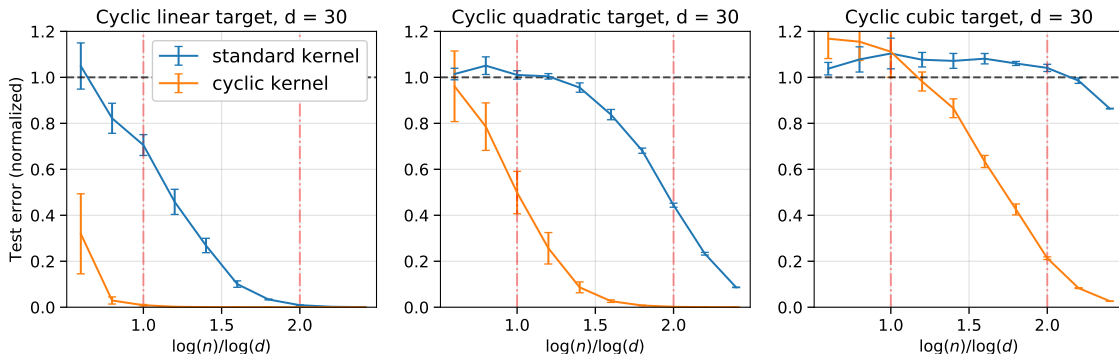


Figure 1: Learning cyclic polynomials (cf. Eq. (16)) over the d -dimensional sphere, $d = 30$, using KKR with a standard (inner-product) kernel and a cyclic invariant kernel, and regularization parameter $\lambda = 0^+$. We report the average and the standard deviation of the test error over 10 realizations, against the sample size n .

with fully connected layers and ReLU activations $\sigma(x) = \max(x, 0)$. We compare this with its cyclically invariant counterpart $H_{d,\text{Cyc}}(\mathbf{x}, \mathbf{y}) = d^{-1} \sum_{0 \leq i < d} h_d(\langle \mathbf{x}, g_i \cdot \mathbf{y} \rangle / d)$, where $g_i \in \text{Cyc}_d$ is the shift by i positions as defined in Example 1. Note that the precise number of layers L is not important. As long as L is fixed in the large N, n limit, our predictions remain unchanged, and the simulations appear to confirm this.

In Figure 1, we report the test errors of fitting each cyclic polynomials with KRR with the two kernels, and regularization parameter $\lambda = 0^+$ (min-norm interpolation). We consider $\sigma_\varepsilon = 0$ and we report the risk averaged over 10 instances against the number of samples n . We observe that the risk in fitting $f_{d,\text{lin}}$, $f_{d,\text{quad}}$ and $f_{d,\text{cube}}$, using KRR with the cyclic kernel $H_{d,\text{Cyc}}$ drops when $n = \Theta_d(1)$, $n = \Theta_d(d)$ and $n = \Theta_d(d^2)$ respectively. In contrast, the risk of KRR with the standard kernel drops when $n = \Theta_d(d)$, $n = \Theta_d(d^2)$ and $n = \Theta_d(d^3)$ respectively. This matches well the predictions of Theorem 5.

We next investigate the relevance of our results for real data. We consider the MNIST dataset ($d = 28 \times 28 = 784$, $n_{\text{train}} = 60000$, $n_{\text{test}} = 10000$ and 10 classes). We encoded class labels by $y_i \in \{-4.5, -3.5, \dots, 3.5, 4.5\}$. We make these data invariant under cyclic translations in two dimensions (Example 2): for each samples in the training and test sets, we replace the image by a uniformly generated 2 dimensional (cyclic) translation of the image (see Fig. 5 in Appendix A.5.2). In this cyclic invariant MNIST data set, the labels are therefore invariant under the action of $\text{Cyc}2\text{D}_{28,28}$.

Images are highly anisotropic in pixel space \mathbb{R}^{784} . In particular, directions corresponding to low-frequency components of the Fourier transform of \mathbf{x} have significantly larger variance than directions corresponding to high-frequency components. Nevertheless, Ghorbani et al. (2020), showed that the analysis of random features and kernel models of Ghorbani et al. (2021); Mei et al. (2021) extends to certain anisotropic models provided the ambient dimension d is replaced by a suitably defined effective dimension d_{eff} .

In order to explore the role of data anisotropy, we pre-process images as follows. We compute the discrete Fourier transform components of the images in the training set and select the $T \in \{20, 70, 120, 200, 400, 784\}$ components with the highest average absolute value. For each T , we

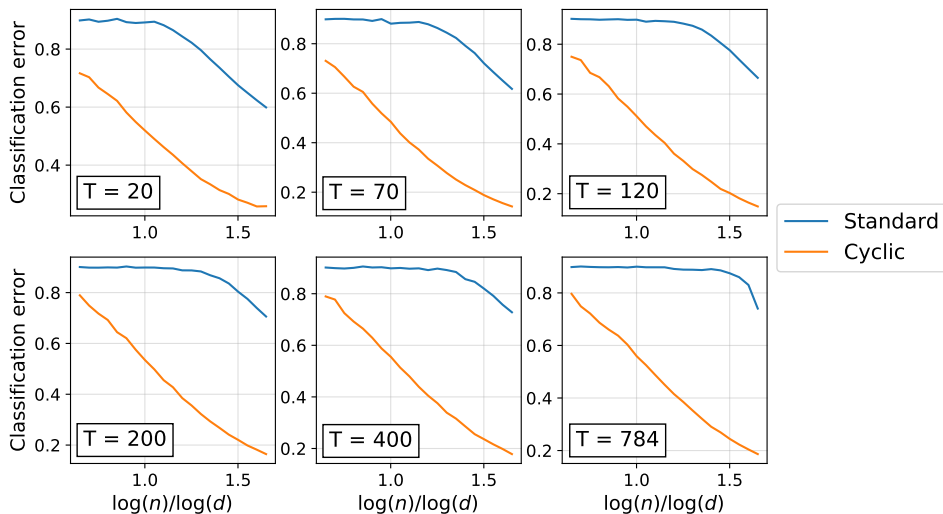


Figure 2: Classification error for the cyclic invariant MNIST dataset. For each frequencies content T , we plot the classification error averaged over 5 instances against the number of samples $\log(n)/\log(d)$, for KRR using a standard (inner-product) kernel and a cyclic invariant kernel and regularization parameter $\lambda = 0^+$.

then construct training and test sets in which we project each image onto the top T frequencies (see Fig. 3 in Appendix A.5.2). When T is small, we expect all the non-zero frequencies to have comparable variance and therefore $d_{\text{eff}} \approx T$. For larger T , we include frequencies of progressively small variance, and therefore d_{eff} should saturate.

For each frequency content T , we compare the performance of two kernels: a standard inner-product kernel $H_d(\mathbf{x}, \mathbf{y}) := h_d(\langle \mathbf{x}, \mathbf{y} \rangle / d)$ and its cyclic counterpart given by $H_{d, \text{Cyc}}(\mathbf{x}, \mathbf{y}) = 1/(28^2) \sum_{0 \leq i, j < 28} h_d(\langle \mathbf{x}, g_{ij} \cdot \mathbf{y} \rangle / d)$, where $g_{ij} \in \text{Cyc}2D_{28, 28}$. We choose H_d to be the neural tangent kernel associated to a two-layers neural network, and hence $H_{d, \text{Cyc}}$ is the one associated to a CNN analogous to (1) (but in two dimensions). We compute the KRR estimates with regularization parameter $\lambda = 0^+$. In Fig. 2, we report the classification error averaged over 5 instances against the number of samples $\log(n)/\log(d)$.

We observe that the cyclic invariant kernel vastly outperform the inner product kernel: the same test error is achieved at a significantly smaller sample size, in qualitative agreement with our general theory. In order to quantify this gap, for each T we fit two curves to the test error of the two kernels, which differ uniquely in an horizontal shift (see Appendix A.5.2). We estimate the sample complexity gain by the difference between these shifts, and denote this estimate by d_{eff} .

It is visually clear that d_{eff} increases with T , as expected. We plot d_{eff} as a function of T in Fig. 6 in Appendix A.5.2. We observe that the behavior of d_{eff} roughly matches our expectations: it grows linearly at small T and eventually saturates.

Acknowledgments

This work was supported by NSF through award DMS-2031883 and from the Simons Foundation through Award 814639 for the Collaboration on the Theoretical Foundations of Deep Learning. We also acknowledge NSF grants CCF-2006489, IIS-1741162 and the ONR grant N00014-18-1-2729.

References

- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in Neural Information Processing Systems*, pages 6155–6166, 2019a.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 6676–6688. 2019b.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 322–332. PMLR, 09–15 Jun 2019a. URL <http://proceedings.mlr.press/v97/arora19a.html>.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pages 8139–8148, 2019b.
- Francis Bach. On the equivalence between kernel quadrature rules and random feature expansions. *The Journal of Machine Learning Research*, 18(1):714–751, 2017.
- William Beckner. Inequalities in Fourier analysis. *Annals of Mathematics*, pages 159–182, 1975.
- William Beckner. Sobolev inequalities, the Poisson semigroup, and analysis on the sphere S^n . *Proceedings of the National Academy of Sciences*, 89(11):4816–4819, 1992.
- Alberto Bietti and Francis Bach. Deep equals shallow for relu networks in kernel regimes. *arXiv preprint arXiv:2009.14397*, 2020.
- Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. *arXiv preprint arXiv:1905.12173*, 2019.
- Aline Bonami. Etude des coefficients de Fourier des fonctions de $L^p(G)$. In *Annales de l’institut Fourier*, volume 20, pages 335–402, 1970.
- Joan Bruna and Stéphane Mallat. Invariant scattering convolution networks. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1872–1886, 2013.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- Shuxiao Chen, Edgar Dobriban, and Jane H Lee. A group-theoretic framework for data augmentation. *Journal of Machine Learning Research*, 21(245):1–71, 2020.

- Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685. PMLR, 2019a.
- Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019b. URL <https://openreview.net/forum?id=S1eK3i09YQ>.
- Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural networks outperform kernel methods? *Advances in Neural Information Processing Systems*, 33, 2020.
- Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers neural networks in high dimension. *The Annals of Statistics*, 49(2):1029–1054, 2021.
- Leonard Gross. Logarithmic sobolev inequalities. *American Journal of Mathematics*, 97(4):1061–1083, 1975.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, pages 8580–8589, 2018.
- Arthur Jacot, Berfin Şimşek, Francesco Spadaro, Clément Hongler, and Franck Gabriel. Kernel alignment risk estimator: Risk prediction from training data. *arXiv preprint arXiv:2006.09796*, 2020.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, pages 8168–8177, 2018.
- Zhiyuan Li, Ruosong Wang, Dingli Yu, Simon S Du, Wei Hu, Ruslan Salakhutdinov, and Sanjeev Arora. Enhanced convolutional neural tangent kernels. *arXiv preprint arXiv:1911.00809*, 2019.
- Zhiyuan Li, Yi Zhang, and Sanjeev Arora. Why are convolutional nets more sample-efficient than fully-connected nets? *arXiv preprint arXiv:2010.08515*, 2020.
- Tengyuan Liang, Alexander Rakhlin, and Xiyu Zhai. On the multiple descent of minimum-norm interpolants and restricted lower isometry of kernels. In *Conference on Learning Theory*, pages 2683–2711. PMLR, 2020a.
- Tengyuan Liang, Alexander Rakhlin, et al. Just interpolate: Kernel “ridgeless” regression can generalize. *Annals of Statistics*, 48(3):1329–1347, 2020b.
- Chao Ma, Stephan Wojtowytsch, Lei Wu, et al. Towards a mathematical understanding of neural network-based machine learning: what we know and what we don’t. *arXiv preprint arXiv:2009.10713*, 2020.

- Stéphane Mallat. Group invariant scattering. *Communications on Pure and Applied Mathematics*, 65(10):1331–1398, 2012.
- Stéphane Mallat. Understanding deep convolutional networks. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2065):20150203, 2016.
- Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Generalization error of random features and kernel methods: hypercontractivity and kernel matrix concentration. *arXiv preprint arXiv:2101.10588*, 2021.
- Ryan O’Donnell. *Analysis of boolean functions*. Cambridge University Press, 2014.
- Samet Oymak and Mahdi Soltanolkotabi. Toward moderate overparameterization: Global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 1(1):84–105, 2020.
- Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Advances in neural information processing systems*, pages 1313–1320, 2009.
- Mark Rudelson, Roman Vershynin, et al. Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18, 2013.
- Alessandro Rudi and Lorenzo Rosasco. Generalization properties of learning with random features. In *Advances in Neural Information Processing Systems*, pages 3215–3225, 2017.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep relu networks. *Machine Learning*, 109(3):467–492, 2020.

Contents

1	Introduction	1
1.1	Related literature	4
2	Examples	6
3	Invariant random feature models	7
4	Invariant kernel machines	9
5	Comparison with alternative approaches	10
5.1	Output symmetrization	10
5.2	Data augmentation	11
6	Numerical illustration	11
A	Some details in the main text	19
A.1	Approximation power of invariant networks	19
A.2	Intuition for the proofs of Theorems 2 and 5	19
A.3	Convolutional neural tangent kernel	20
A.4	Proof of Proposition 7	21
A.5	Details of numerical simulations	22
A.5.1	Synthetic data	22
A.5.2	Cyclic invariant MNIST data set	22
B	Proof of the main theorems	25
B.1	Notations	25
B.2	Proof of Theorem 2	25
B.3	Proof of Theorem 5	29
C	Decomposition of invariant functions	31
C.1	The invariant function class and the symmetrization operator	31
C.2	Orthogonal polynomials on invariant function class	31
C.3	A representation lemma	32
C.4	Gegenbauer decomposition of invariant features and kernels	32
D	Counting the degeneracy	34
D.1	Counting the degeneracy of Cyc_d and $\text{Cyc}2D_{d_1, d_2}$ (Example 1 and 2)	34
D.1.1	Proof of Proposition 12	34
D.1.2	Auxiliary lemmas	36
D.2	Counting the degeneracy of band-limited function class (Example 3)	37
E	Concentration for invariant groups with degeneracy $\alpha \leq 1$	40
E.1	Main proposition	40
E.2	Auxiliary Lemmas	42

F	Kernel concentration for the cyclic group and general σ	47
F.1	Main propositions	47
F.2	Auxiliary lemmas	50
G	Hypercontractivity of general activation σ for $(\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$	55
G.1	Proof of Proposition 29	55
G.2	Proof in the Gaussian case	58
G.3	Technical lemmas	61
H	Technical background of function spaces	64
H.1	Functions on the sphere	64
	H.1.1 Functional spaces over the sphere	64
	H.1.2 Gegenbauer polynomials	64
	H.1.3 Hermite polynomials	66
H.2	Functions on the hypercube	66
	H.2.1 Fourier basis	66
	H.2.2 Hypercubic Gegenbauer	67
H.3	Hypercontractivity of Gaussian measure and uniform distributions on the sphere and the hypercube	68

Appendix A. Some details in the main text

A.1. Approximation power of invariant networks

In the proposition below, we show that the approximation power of two-layers \mathcal{G}_d -invariant neural networks are always no worse than two-layers fully-connected neural networks when the target function is \mathcal{G}_d -invariant.

Proposition 9 *Let $\sigma \in C(\mathbb{R})$ be an activation function. Let $\mathcal{A}_d \in \{\mathbb{S}^{d-1}(\sqrt{d}), \mathcal{Q}^d\}$. Let \mathcal{G}_d be a subgroup of $\mathcal{O}(d)$ that preserves \mathcal{A}_d . Let π_d be the Haar measure of \mathcal{G}_d . Let $f_* \in L^2(\mathcal{A}_d; \mathcal{G}_d)$ be a \mathcal{G}_d -invariant function. Define the function classes of two-layers invariant neural networks and two-layers fully-connected neural networks by*

$$\mathcal{F}_{\text{NN}, \mathcal{G}_d, N} = \left\{ f(\mathbf{x}) = \sum_{i=1}^N a_i \int_{\mathcal{G}_d} \sigma(\langle \boldsymbol{\theta}_i, g \cdot \mathbf{x} \rangle / \sqrt{d}) \pi_d(dg) : \boldsymbol{\theta}_i \in \mathcal{A}_d, a_i \in \mathbb{R} \right\}, \quad (17)$$

$$\mathcal{F}_{\text{NN}, N} = \left\{ f(\mathbf{x}) = \sum_{i=1}^N a_i \sigma(\langle \boldsymbol{\theta}_i, \mathbf{x} \rangle / \sqrt{d}) : \boldsymbol{\theta}_i \in \mathcal{A}_d, a_i \in \mathbb{R} \right\}. \quad (18)$$

Then we have

$$\inf_{f \in \mathcal{F}_{\text{NN}, \mathcal{G}_d, N}} \|f_* - f\|_{L^2}^2 \leq \inf_{f \in \mathcal{F}_{\text{NN}, N}} \|f_* - f\|_{L^2}^2.$$

Proof [Proof of Proposition 9]

We define the symmetrization operator $\mathcal{S} : L^2(\mathcal{A}_d) \rightarrow L^2(\mathcal{A}_d; \mathcal{G}_d)$ by

$$(\mathcal{S}f)(\mathbf{x}) = \int_{\mathcal{G}_d} f(g \cdot \mathbf{x}) \pi_d(dg).$$

Since $f_* \in L^2(\mathcal{A}_d; \mathcal{G}_d)$, by Jensen's inequality, for any $f \in L^2(\mathcal{A}_d)$, we have

$$\|f_* - \mathcal{S}f\|_{L^2}^2 = \|\mathcal{S}(f_* - f)\|_{L^2}^2 \leq \|f_* - f\|_{L^2}^2.$$

Moreover, for any $f \in \mathcal{F}_{\text{NN}, N}$, we have $\mathcal{S}f \in \mathcal{F}_{\text{NN}, \mathcal{G}_d, N}$. This gives

$$\inf_{f \in \mathcal{F}_{\text{NN}, \mathcal{G}_d, N}} \|f_* - f\|_{L^2}^2 \leq \inf_{f \in \mathcal{F}_{\text{NN}, N}} \|f_* - \mathcal{S}f\|_{L^2}^2 \leq \inf_{f \in \mathcal{F}_{\text{NN}, N}} \|f_* - f\|_{L^2}^2.$$

This concludes the proof. ■

A.2. Intuition for the proofs of Theorems 2 and 5

Theorem 2 and 5 are consequences of general theorems proved in Mei et al. (2021). The d^α improvement between invariant and non-invariant models can be understood as follows: consider an inner-product activation $\sigma(\langle \mathbf{x}, \boldsymbol{\theta} \rangle / \sqrt{d})$ with $\mathbf{x}, \boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)$ (where we denoted $\boldsymbol{\theta} = \sqrt{d} \cdot \mathbf{w}$), then we have the following eigendecomposition

$$\sigma(\langle \mathbf{x}, \boldsymbol{\theta} \rangle / \sqrt{d}) = \sum_{k=0}^{\infty} \xi_{d,k}^2 \sum_{l=1}^{B(\mathcal{A}_d; k)} Y_{kl}^{(d)}(\mathbf{x}) Y_{kl}^{(d)}(\boldsymbol{\theta}),$$

where $\{Y_{kl}^{(d)}\}_{l \in [B(\mathcal{A}_d; k)]}$ form an orthonormal basis of $V_{d,k}$, the subspace of degree- k polynomials on \mathcal{A}_d (see Section H for background on functional spaces on the sphere and hypercube). The eigenvalues of σ are given by $\{\xi_{d,k}\}_{k \geq 0}$ with each having degeneracy $B(\mathcal{A}_d; k)$.

As mentioned in the introduction, the symmetry group \mathcal{G}_d preserves $V_{d,k}$ (see Section C.2) and the invariant activation function has the following eigendecomposition

$$\bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}) := \int_{\mathcal{G}_d} \sigma(\langle \mathbf{x}, g \cdot \boldsymbol{\theta} \rangle / \sqrt{d}) \pi_d(dg) = \sum_{k=0}^{\infty} \xi_{d,k}^2 \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\mathbf{x}) \bar{Y}_{kl}^{(d)}(\boldsymbol{\theta}),$$

where the $\{\bar{Y}_{kl}^{(d)}\}_{l \in [B(\mathcal{A}_d; k)]}$ form an orthonormal basis of $V_{d,k}(\mathcal{G}_d)$, the subspace of degree- k invariant polynomials on \mathcal{A}_d . The eigenvalues of $\bar{\sigma}$ are given by $\{\xi_{d,k}\}_{k \geq 0}$ with each having degeneracy $D(\mathcal{A}_d; k)$.

Hence $\bar{\sigma}$ has the same eigenvalues $\xi_{d,k}$ as σ , but with degeneracy smaller by a factor

$$\frac{B(\mathcal{A}_d; k)}{D(\mathcal{A}_d; k)} = \Theta_d(d^\alpha).$$

In other words, in order to fit degree ℓ polynomials using invariant methods, one needs to fit a factor d^α less eigendirections, which translates to a factor d^α improvement in the sample and features complexity.

This intuition is verified rigorously in the proof of these theorems in Appendix B.

A.3. Convolutional neural tangent kernel

Proposition 10 *Let $\sigma \in C^1(\mathbb{R})$ be an activation function. Let \mathcal{G}_d be a discrete subgroup of $\mathcal{O}(d)$ with Haar measure π_d . Let f_N be an invariant neural network*

$$f_N(\mathbf{x}; \boldsymbol{\Theta}) = \sum_{i=1}^N a_i \int_{\mathcal{G}_d} \sigma(\langle \mathbf{w}_i, g \cdot \mathbf{x} \rangle) \pi_d(dg).$$

Let $a_i^0 \sim_{i.i.d.} \mathcal{N}(0, 1)$ and $\mathbf{w}_i^0 \sim_{i.i.d.} \text{Unif}(\mathbb{S}^{d-1})$ independently, and $\boldsymbol{\Theta}^0 = (a_1^0, \dots, a_N^0, \mathbf{w}_1^0, \dots, \mathbf{w}_N^0)$. Then there exists $h_d : [-1, 1] \rightarrow \mathbb{R}$, such that for any $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$, we have almost surely

$$\lim_{N \rightarrow \infty} \langle \nabla_{\boldsymbol{\Theta}} f_N(\mathbf{x}; \boldsymbol{\Theta}^0), \nabla_{\boldsymbol{\Theta}} f_N(\mathbf{y}; \boldsymbol{\Theta}^0) \rangle / N = \int_{\mathcal{G}_d} h_d(\langle \mathbf{x}, g \cdot \mathbf{y} \rangle / d) \pi_d(dg).$$

Proof [Proof of Proposition 10] For $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$, define

$$\begin{aligned} h_d^{(1)}(\langle \mathbf{x}, \mathbf{y} \rangle / d) &= \mathbb{E}_{\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})} [\sigma(\langle \mathbf{w}, \mathbf{x} \rangle) \sigma(\langle \mathbf{w}, \mathbf{y} \rangle)], \\ h_d^{(2)}(\langle \mathbf{x}, \mathbf{y} \rangle / d) &= \mathbb{E}_{\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})} [\sigma'(\langle \mathbf{w}, \mathbf{x} \rangle) \sigma'(\langle \mathbf{w}, \mathbf{y} \rangle) \langle \mathbf{x}, \mathbf{y} \rangle]. \end{aligned}$$

By the technical backgrounds in Section H, we can see that $h_d^{(1)}$ and $h_d^{(2)}$ can be well-defined.

Calculating the derivative of the neural network with respect to $\mathbf{a} = (a_1, \dots, a_N)$, we have

$$\frac{1}{N} \langle \nabla_{\mathbf{a}} f(\mathbf{x}; \boldsymbol{\Theta}^0), \nabla_{\mathbf{a}} f(\mathbf{y}; \boldsymbol{\Theta}^0) \rangle = \int_{\mathcal{G}_d \times \mathcal{G}_d} \frac{1}{N} \sum_{i=1}^N \left[\sigma(\langle \mathbf{w}_i, g \cdot \mathbf{x} \rangle) \sigma(\langle \mathbf{w}_i, g' \cdot \mathbf{y} \rangle) \right] \pi_d(dg) \pi_d(dg').$$

Since \mathcal{G}_d is a discrete group, by law of large numbers, we have

$$\begin{aligned}
 & \lim_{N \rightarrow \infty} \frac{1}{N} \langle \nabla_{\mathbf{a}} f(\mathbf{x}; \Theta^0), \nabla_{\mathbf{a}} f(\mathbf{y}; \Theta^0) \rangle \\
 &= \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E}_{\mathbf{w}} [\sigma(\langle \mathbf{w}, g \cdot \mathbf{x} \rangle) \sigma(\langle \mathbf{w}, g' \cdot \mathbf{y} \rangle)] \pi_d(\mathrm{d}g) \pi_d(\mathrm{d}g') \\
 &= \int_{\mathcal{G}_d \times \mathcal{G}_d} h_d^{(1)}(\langle g \cdot \mathbf{x}, g' \cdot \mathbf{y} \rangle / d) \pi_d(\mathrm{d}g) \pi_d(\mathrm{d}g') \\
 &= \int_{\mathcal{G}_d} h_d^{(1)}(\langle \mathbf{x}, g \cdot \mathbf{y} \rangle / d) \pi_d(\mathrm{d}g).
 \end{aligned}$$

Moreover, calculating the derivative of the neural network with respect to $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_N)$, we have

$$\begin{aligned}
 & \frac{1}{N} \langle \nabla_{\mathbf{W}} f(\mathbf{x}; \Theta^0), \nabla_{\mathbf{W}} f(\mathbf{y}; \Theta^0) \rangle \\
 &= \int_{\mathcal{G}_d \times \mathcal{G}_d} \frac{1}{N} \sum_{i=1}^N \left[(a_i^0)^2 \sigma'(\langle \mathbf{w}_i, g \cdot \mathbf{x} \rangle) \sigma'(\langle \mathbf{w}_i, g' \cdot \mathbf{y} \rangle) \langle g \cdot \mathbf{x}, g' \cdot \mathbf{y} \rangle \right] \pi_d(\mathrm{d}g) \pi_d(\mathrm{d}g')
 \end{aligned}$$

Since \mathcal{G}_d is a discrete group, by law of large numbers, we have

$$\begin{aligned}
 & \lim_{N \rightarrow \infty} \frac{1}{N} \langle \nabla_{\mathbf{W}} f(\mathbf{x}; \Theta^0), \nabla_{\mathbf{W}} f(\mathbf{y}; \Theta^0) \rangle \\
 &= \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E}_{\mathbf{w}} [\sigma'(\langle \mathbf{w}, g \cdot \mathbf{x} \rangle) \sigma'(\langle \mathbf{w}, g' \cdot \mathbf{y} \rangle) \langle g \cdot \mathbf{x}, g' \cdot \mathbf{y} \rangle] \pi_d(\mathrm{d}g) \pi_d(\mathrm{d}g') \\
 &= \int_{\mathcal{G}_d \times \mathcal{G}_d} h_d^{(2)}(\langle g \cdot \mathbf{x}, g' \cdot \mathbf{y} \rangle / d) \pi_d(\mathrm{d}g) \pi_d(\mathrm{d}g') \\
 &= \int_{\mathcal{G}_d} h_d^{(2)}(\langle \mathbf{x}, g \cdot \mathbf{y} \rangle / d) \pi_d(\mathrm{d}g).
 \end{aligned}$$

Taking $h_d = h_d^{(1)} + h_d^{(2)}$ concludes the proof. ■

A.4. Proof of Proposition 7

Let \hat{f}_d be an estimator satisfying

$$\varepsilon^2 := \|\bar{\mathbb{P}}_{\leq \ell} f_d - \hat{f}_d\|_{L^2}^2 = \|\bar{\mathbb{P}}_{\leq \ell} f_d - \bar{\mathbb{P}}_{\leq \ell} \hat{f}_d\|_{L^2}^2 + \|\bar{\mathbb{P}}_{> \ell} \hat{f}_d\|_{L^2}^2. \quad (19)$$

By Jensen's inequality and by the equation above, we have

$$\|\mathcal{S} \bar{\mathbb{P}}_{> \ell} \hat{f}_d\|_{L^2}^2 \leq \|\bar{\mathbb{P}}_{> \ell} \hat{f}_d\|_{L^2}^2 \leq \varepsilon^2. \quad (20)$$

As a consequence, we have

$$\begin{aligned}
 & \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2}^2 - 2\varepsilon \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2} \stackrel{(1)}{\leq} \|\bar{\mathbb{P}}_{> \ell} f_d - \mathcal{S} \bar{\mathbb{P}}_{> \ell} \hat{f}_d\|_{L^2}^2 \stackrel{(2)}{=} \|\bar{\mathbb{P}}_{> \ell} f_d - \bar{\mathbb{P}}_{> \ell} \mathcal{S} \hat{f}_d\|_{L^2}^2 \\
 & \stackrel{(3)}{\leq} \|f_d - \mathcal{S} \hat{f}_d\|_{L^2}^2 \stackrel{(4)}{=} \|\mathcal{S}(f_d - \hat{f}_d)\|_{L^2}^2 \stackrel{(5)}{\leq} \|f_d - \hat{f}_d\|_{L^2}^2 \\
 & \stackrel{(6)}{=} \|\bar{\mathbb{P}}_{\leq \ell} f_d - \bar{\mathbb{P}}_{\leq \ell} \hat{f}_d\|_{L^2}^2 + \|\bar{\mathbb{P}}_{> \ell} f_d - \bar{\mathbb{P}}_{> \ell} \hat{f}_d\|_{L^2}^2 \stackrel{(7)}{=} \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2}^2 + \varepsilon^2 + 2\varepsilon \|\bar{\mathbb{P}}_{> \ell} f_d\|_{L^2}.
 \end{aligned} \quad (21)$$

Here, (1) is by Eq. (20); (2) is by the fact that \mathcal{S} is exchangeable with $\bar{\mathbb{P}}_{>\ell}$ (c.f. Section C); (3) is by the fact that $\bar{\mathbb{P}}_{>\ell}$ is a projection operator; (4) is by the fact that f_d is \mathcal{G}_d -invariant; (5) is by Jensen’s inequality; (6) is by orthogonal decomposition; (7) is by Eq. (19). This concludes the proof.

A.5. Details of numerical simulations

A.5.1. SYNTHETIC DATA

We consider the standard (inner-product) kernel $H_d(\mathbf{x}, \mathbf{y}) = h_{\text{NTK}}(\langle \mathbf{x}, \mathbf{y} \rangle / d)$ to be the neural tangent kernel associated to a depth-5 neural network with fully connected layers and ReLU activation $\sigma(x) = \max(x, 0)$. This can be obtained iteratively as follow (see Jacot et al. (2018) and Bietti and Bach (2020)): define for $u \in [-1, 1]$,

$$h_0(u) = \frac{1}{\pi}(\pi - \arccos(u)), \quad h_1(u) = u \cdot h_0(u) + \frac{1}{\pi}\sqrt{1 - u^2},$$

and $h_{\text{NTK}}(u) = h_{\text{NTK}}^5(u)$ with $h_{\text{NTK}}^1(u) = h^1(u) = u$ and for $k = 2, \dots, 5$,

$$\begin{aligned} h^k(u) &= h_1(h^{k-1}(u)), \\ h_{\text{NTK}}^k(u) &= h_{\text{NTK}}^{k-1}(u)h_0(h^{k-1}(u)) + h^k(u). \end{aligned}$$

We compute the cyclic invariant kernel by summing over all cyclic translations $g \in \text{Cyc}_d$:

$$H_{d,\text{inv}}(\mathbf{x}, \mathbf{y}) = \frac{1}{d} \sum_{g \in \text{Cyc}_d} h_{\text{NTK}}(\langle \mathbf{x}, g \cdot \mathbf{y} \rangle / d).$$

A.5.2. CYCLIC INVARIANT MNIST DATA SET

We consider the MNIST data set of 28×28 grayscale images ($d = 784$) of handwritten digits, which contains 60000 training images and 10000 testing images. We pre-process the images in three steps:

- (a) We compute the discrete Fourier transform of the images in the training set and compute the average absolute value of the frequency components (see left frame of Fig. 3). For each $T \in \{20, 70, 120, 200, 400, 784\}$, we select $\Omega_T \subset [28] \times [28]$ to be the set of the top T frequencies (i.e., the T frequencies with highest absolute value averaged on the training set).
- (b) For each T , we construct a train and test sets in which we project each image onto Ω_T (i.e., we set all the frequency components not in Ω_T to 0). We displayed in Fig. 4 two digits and their projection on the top T frequencies Ω_T for different T .
- (c) For each image in the training and test sets, we replace the image by a uniformly generated 2 dimensional (cyclic) translation of the image. We display some examples in Fig. 5.

We further normalize the images so that $\|\mathbf{x}\|_2 = 1$ and center the labels $y_i \in \mathcal{Y}$ where $\mathcal{Y} = \{-4.5, -3.5, \dots, 3.5, 4.5\}$. In order to compute the classification error, we round the prediction value to the nearest label in \mathcal{Y} .

We use the inner-product kernel $H_d(\mathbf{x}, \mathbf{y}) = h_{\text{NTK}}(\langle \mathbf{x}, \mathbf{y} \rangle / d)$ where h_{NTK} is the neural tangent kernel associated to a 2-layers neural network with fully connected layers and ReLU activation $\sigma(x) = \max(x, 0)$, which given by

$$h_{\text{NTK}}(u) = u \cdot \left(\pi - \frac{\arccos(u)}{\pi} \right) + \frac{1}{\pi}\sqrt{1 - u^2}.$$

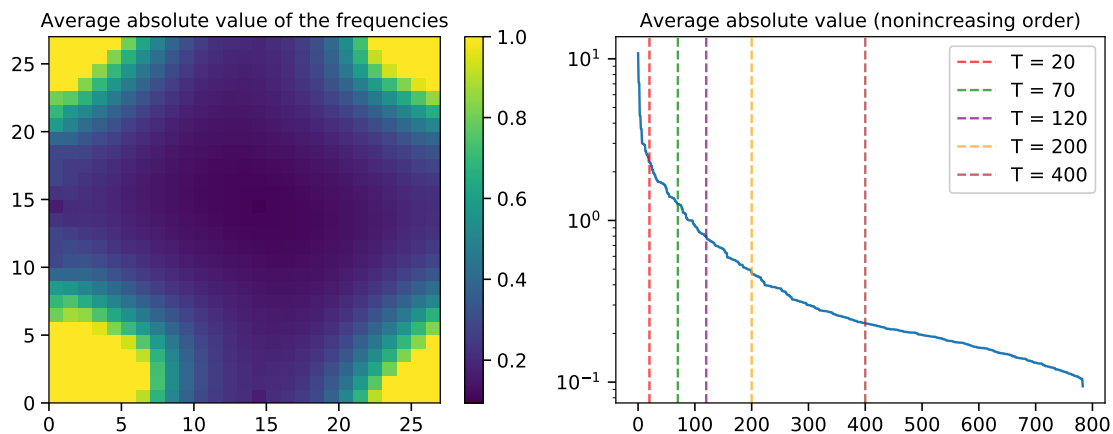


Figure 3: Left frame: the absolute value of the frequency components of MNIST images averaged over the training set (threshold at 1 in the figure). Coordinates on the bottom left-hand side correspond to lower frequency components while coordinates closer to the top right-hand side represent the high frequency directions. Right frame: average absolute value of the frequencies in nonincreasing order. The vertical lines correspond to the different T chosen ($T = 784$ corresponds to keeping all the frequencies).

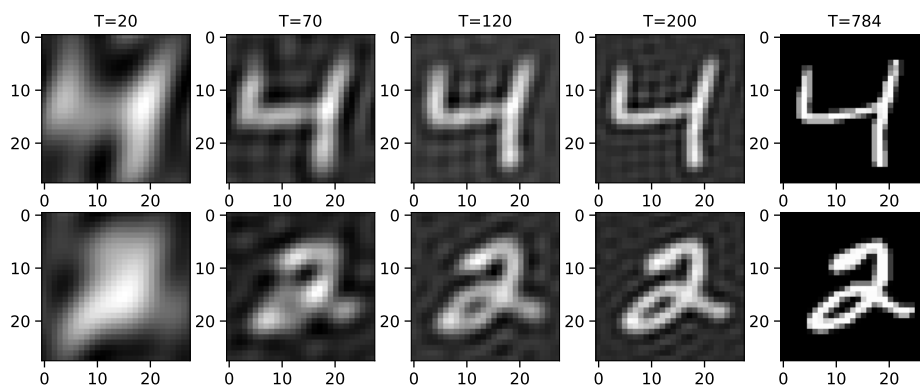


Figure 4: Examples of two images projected on the top T frequencies.

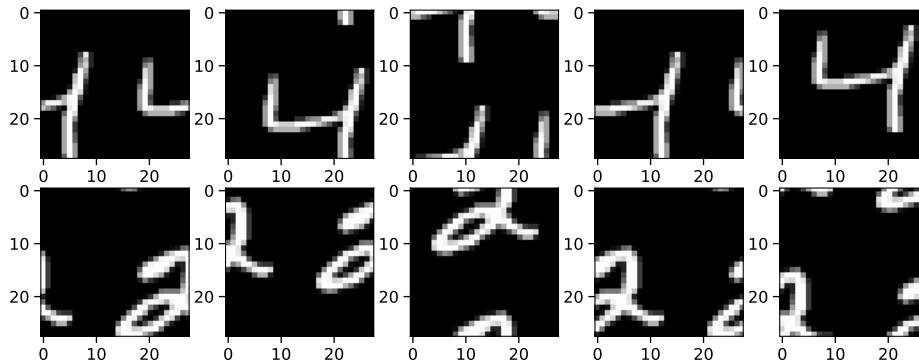


Figure 5: Examples of random 2-dimensional cyclic translations of the images (for $T = 784$).

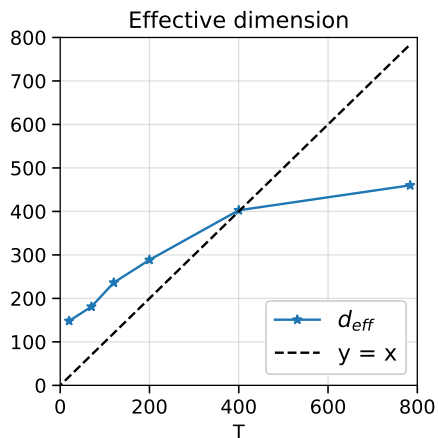


Figure 6: Estimated sample size gap between standard and invariant kernel methods, for the translationally invariant MNIST dataset, as a function of the frequency content T .

The cyclic invariant kernel is computed by summing over all two-dimensional cyclic translations $g_{ij} \in \text{Cyc2D}_{28,28}$:

$$H_{d,\text{inv}}(\mathbf{x}, \mathbf{y}) = \frac{1}{28^2} \sum_{i,j=0}^{27} h_{\text{NTK}}(\langle \mathbf{x}, g_{ij} \cdot \mathbf{y} \rangle / d).$$

For each T , we estimate the effective dimension d_{eff} by fitting two parallel lines through the classification error points of the standard and cyclic kernels at the same time (keeping only the points where the curves decrease). The estimated (log) effective dimension is then given by the difference of the offsets. We report these estimates for different T in Fig. 6.

Appendix B. Proof of the main theorems

In this section, we present the proofs of Theorem 2 and 5 stated in the main text. The rest of the appendices are organized as follow:

- Appendix C presents key properties of the decomposition of invariant functions, while Appendix H reviews some technical background on the functional spaces on the sphere and the hypercube.
- Appendix D proves that the examples of symmetry group listed in Section 2 (one and two-dimensional cyclic groups and band-limited functions) have degeneracy 1.
- Appendix E presents a key concentration result on the diagonal elements of polynomial invariant kernels. In particular, the results of Appendix E are the only ones required in the proofs of Theorems 2 and 5 in the case of polynomial activations for general symmetry group \mathcal{G}_d of degeneracy $\alpha \leq 1$.
- Appendices F and G provides necessary results to extend the proofs to non-polynomial activations in the case of $(\mathcal{A}_d, \mathcal{G}_d) = (\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$.

B.1. Notations

For a positive integer, we denote by $[n]$ the set $\{1, 2, \dots, n\}$. For vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, we denote $\langle \mathbf{u}, \mathbf{v} \rangle = u_1 v_1 + \dots + u_d v_d$ their scalar product, and $\|\mathbf{u}\|_2 = \langle \mathbf{u}, \mathbf{u} \rangle^{1/2}$ the ℓ_2 norm. Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, we denote $\|\mathbf{A}\|_{\text{op}} = \max_{\|\mathbf{u}\|_2=1} \|\mathbf{A}\mathbf{u}\|_2$ its operator norm and by $\|\mathbf{A}\|_F = (\sum_{i,j} A_{ij}^2)^{1/2}$ its Frobenius norm. If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a square matrix, the trace of \mathbf{A} is denoted by $\text{Tr}(\mathbf{A}) = \sum_{i \in [n]} A_{ii}$.

We use $O_d(\cdot)$ (resp. $o_d(\cdot)$) for the standard big-O (resp. little-o) relations, where the subscript d emphasizes the asymptotic variable. Furthermore, we write $f = \Omega_d(g)$ if $g(d) = O_d(f(d))$, and $f = \omega_d(g)$ if $g(d) = o_d(f(d))$. Finally, $f = \Theta_d(g)$ if we have both $f = O_d(g)$ and $f = \Omega_d(g)$.

We use $O_{d,\mathbb{P}}(\cdot)$ (resp. $o_{d,\mathbb{P}}(\cdot)$) the big-O (resp. little-o) in probability relations. Namely, for $h_1(d)$ and $h_2(d)$ two sequences of random variables, $h_1(d) = O_{d,\mathbb{P}}(h_2(d))$ if for any $\varepsilon > 0$, there exists $C_\varepsilon > 0$ and $d_\varepsilon \in \mathbb{Z}_{>0}$, such that

$$\mathbb{P}(|h_1(d)/h_2(d)| > C_\varepsilon) \leq \varepsilon, \quad \forall d \geq d_\varepsilon,$$

and respectively: $h_1(d) = o_{d,\mathbb{P}}(h_2(d))$, if $h_1(d)/h_2(d)$ converges to 0 in probability. Similarly, we will denote $h_1(d) = \Omega_{d,\mathbb{P}}(h_2(d))$ if $h_2(d) = O_{d,\mathbb{P}}(h_1(d))$, and $h_1(d) = \omega_{d,\mathbb{P}}(h_2(d))$ if $h_2(d) = o_{d,\mathbb{P}}(h_1(d))$. Finally, $h_1(d) = \Theta_{d,\mathbb{P}}(h_2(d))$ if we have both $h_1(d) = O_{d,\mathbb{P}}(h_2(d))$ and $h_1(d) = \Omega_{d,\mathbb{P}}(h_2(d))$.

B.2. Proof of Theorem 2

Let \mathcal{G}_d be a group of degeneracy $\alpha \leq 1$. Consider $\mathbf{x}, \boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)$, $d^{s-\alpha+\delta_0} \leq n \leq d^{s-\alpha+1-\delta_0}$, $d^{s-\alpha+\delta_0} \leq N \leq d^{s-\alpha+1-\delta_0}$ and an activation function σ that satisfies Assumption 1 at level (s, S) . Denote

$$\bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}) = \int_{\mathcal{G}_d} \sigma(\langle \boldsymbol{\theta}, g \cdot \mathbf{x} \rangle / \sqrt{d}) \pi_d(dg).$$

Theorem 2 is a consequence of Theorem 1 in Mei et al. (2021) where we take $\mathcal{X}_d = \Omega_d = \mathcal{A}_d$, $\nu_d = \tau_d = \text{Unif}(\mathcal{A}_d)$ and $\mathcal{D}_d = \mathcal{V}_d = L^2(\mathcal{A}_d, \mathcal{G}_d) \subset L^2(\mathcal{A}_d)$. The proof amounts to checking that $\bar{\sigma}$ indeed verifies the feature map concentration and spectral gap assumptions (see Section 2.2 in Mei et al. (2021)). We borrow some of the notations introduced in Mei et al. (2021) and refer the reader to their Section 2.1.

Proof [Proof of Theorem 2] For the sake of simplicity, we consider the overparametrized case $N(d) \geq n(d)d^\delta$ for some $\delta > 0$, and therefore $S \geq s$. The underparametrized case $d^\delta N(d) \leq n(d)$ is treated analogously.

Step 1. Diagonalization of the activation function $\bar{\sigma}$ and choosing $m = m(d)$, $M = M(d)$.

We can decompose the inner product activation σ in the basis of Gegenbauer polynomials (see Section H for definitions):

$$\sigma(\langle \mathbf{x}, \boldsymbol{\theta} \rangle / \sqrt{d}) = \sum_{k=0}^{\infty} \xi_{d,k} B(\mathcal{A}_d; k) Q_k^{(d)}(\langle \mathbf{x}, \boldsymbol{\theta} \rangle),$$

where (with $\mathbf{e} \in \mathcal{A}_d$ arbitrary)

$$\xi_{d,k}(\sigma) = \mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)} [\sigma(\langle \mathbf{e}, \boldsymbol{\theta} \rangle / \sqrt{d}) Q_k^{(d)}(\langle \mathbf{e}, \boldsymbol{\theta} \rangle)].$$

From Assumption 1.(a) that $|\sigma(x)| \leq c_0 \exp(c_1 x^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$ (which is trivially verified for a polynomial activation function), there exists a constant $C > 0$ such that (see for example Lemma 5 in Ghorbani et al. (2021))

$$\|\sigma(\langle \mathbf{e}, \cdot \rangle / \sqrt{d})\|_{L^2(\mathcal{A}_d)} = \sum_{k=1}^{\infty} \xi_{d,k}^2 B(\mathcal{A}_d; k) \leq C. \quad (22)$$

We have for fixed k , $B(\mathcal{A}_d; k) = \Theta(d^k)$. Furthermore, for non-polynomial activation functions in the case of $(\mathcal{A}_d, \mathcal{G}_d) = (\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$, we use $\sup_{k>s} B(\mathbb{S}^{d-1}; k)^{-1} = O_d(d^{-s-1})$ (Lemma 1 in Ghorbani et al. (2021)). We deduce that

$$\sup_{k>s} \xi_{d,k}^2 = O_d(d^{-s-1}), \quad (23)$$

$$\sup_{k>S} \xi_{d,k}^2 = O_d(d^{-S-1}). \quad (24)$$

From the correspondence between Gegenbauer and Hermite polynomials when $d \rightarrow \infty$ (see Eq. (111) in Section H.1.3), Assumption 1.(b) implies that $\xi_{d,k}^2 = \Theta_d(d^{-k})$ for $k = 0, \dots, s$.

Let us diagonalize $\bar{\sigma}$ in the basis of \mathcal{G}_d -invariant polynomials $\{\bar{Y}_{kl}\}_{k \geq 0, \ell \in [D(\mathcal{A}_d; k)]}$ (see Section C for definitions). From Lemma 11 stated in Section C.3, we have

$$\begin{aligned} \bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}) &= \sum_{k=0}^{\infty} \xi_{d,k} B(\mathcal{A}_d; k) \int_{\mathcal{G}_d} Q_k^{(d)}(\langle \mathbf{x}, g \cdot \boldsymbol{\theta} \rangle) \pi_d(dg) \\ &= \sum_{k=0}^{\infty} \xi_{d,k} \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\mathbf{x}) \bar{Y}_{kl}^{(d)}(\boldsymbol{\theta}). \end{aligned} \quad (25)$$

Denote $(\lambda_{d,j})_{j \geq 1}$ the eigenvalues of $\bar{\sigma}$ in non increasing order of their absolute value (namely, the $\xi_{d,k}$'s which have degeneracies $D(\mathcal{A}_d; k)$). Set m and M to be the number of eigenvalues $\lambda_{d,j}^2$

that are bigger than $d^{-s-1+\delta}$ and $d^{-S-1+\delta}$ respectively, for a constant $\delta > 0$ that will be set sufficiently small (see Step 4). From the above discussion, $(\lambda_{d,j})_{j \leq m}$ corresponds exactly to all the eigenvalues associated to invariant polynomials of degree less or equal to s , while $(\lambda_{d,j})_{j \leq M}$ does not contain any eigenvalues associated to invariant polynomials of degree bigger or equal to $S + 1$. Hence,

$$m = \sum_{k=0}^s D(\mathcal{A}_d; k) = \Theta_d(d^{s-\alpha}), \quad M \leq \sum_{k=0}^S D(\mathcal{A}_d; k) = O_d(d^{S-\alpha}), \quad (26)$$

where we used that \mathcal{G}_d has degeneracy α so that $D(\mathcal{A}_d; k) = \Theta_d(d^{-\alpha}) \cdot B(\mathcal{A}_d; k)$.

Step 2. Diagonal elements of the truncated kernel.

We introduce the kernel associated to activation $\bar{\sigma}$:

$$H_d(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}_{\boldsymbol{\theta}}[\bar{\sigma}(\mathbf{x}_1; \boldsymbol{\theta})\bar{\sigma}(\mathbf{x}_2; \boldsymbol{\theta})] = \sum_{k=0}^{\infty} \xi_{d,k}^2 D(\mathcal{A}_d; k) \Upsilon_k^{(d)}(\mathbf{x}_1, \mathbf{x}_2),$$

where we denote

$$\Upsilon_k^{(d)}(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{D(\mathcal{A}_d; k)} \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\mathbf{x}_1) \bar{Y}_{kl}^{(d)}(\mathbf{y}_1).$$

Similarly, we introduce a kernel in the feature space

$$U_d(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbb{E}_{\mathbf{x}}[\bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}_1)\bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}_2)] = \sum_{k=0}^{\infty} \xi_{d,k}^2 D(\mathcal{A}_d; k) \Upsilon_k^{(d)}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2).$$

We denote $\mathbb{H}_d, \mathbb{U}_d : L^2(\mathcal{A}_d, \mathcal{G}_d) \rightarrow L^2(\mathcal{A}_d, \mathcal{G}_d)$ the kernel operators with kernel representation H_d and U_d , and denote $\mathbb{H}_{d, > m}$ and $\mathbb{U}_{d, > M}$ the kernel operators where the biggest m and M eigenvalues respectively are set to 0. Recalling the discussion on the choice of m and M , denote $E = \{k : \xi_{d,k}^2 \leq d^{-S-1+\delta}\}$: E contain all integers bigger or equal to $S + 1$ and none smaller or equal to s .

The diagonal elements of the truncated kernels are then given by

$$\begin{aligned} H_{d, > m}(\mathbf{x}, \mathbf{x}) &= \sum_{k=s+1}^{\infty} \xi_{d,k}^2 D(\mathcal{A}_d; k) \Upsilon_k^{(d)}(\mathbf{x}, \mathbf{x}), \\ U_{d, > M}(\boldsymbol{\theta}, \boldsymbol{\theta}) &= \sum_{k \in E} \xi_{d,k}^2 D(\mathcal{A}_d; k) \Upsilon_k^{(d)}(\boldsymbol{\theta}, \boldsymbol{\theta}), \end{aligned} \quad (27)$$

and

$$\begin{aligned} \text{Tr}(\mathbb{H}_{d, > m}) &= \mathbb{E}_{\mathbf{x}}[H_{d, > m}(\mathbf{x}, \mathbf{x})] = \sum_{k=s+1}^{\infty} \xi_{d,k}^2 D(\mathcal{A}_d; k), \\ \text{Tr}(\mathbb{U}_{d, > M}) &= \mathbb{E}_{\boldsymbol{\theta}}[U_{d, > M}(\boldsymbol{\theta}, \boldsymbol{\theta})] = \sum_{k \in E} \xi_{d,k}^2 D(\mathcal{A}_d; k). \end{aligned}$$

From Assumption 1.(c), σ is not a polynomial of degree less or equal to S . Hence, there exists $\ell > S$ such that $\mu_{\ell}(\sigma) \neq 0$ and therefore $\xi_{d,\ell}^2 D(\mathcal{A}_d; \ell) = \Theta(d^{-\alpha})$. Furthermore, from Eq. (22) and the assumption that \mathcal{G}_d is of degeneracy α (for polynomial activation functions, see Proposition 22 in Section G for general σ in the case of $(\mathcal{A}_d, \mathcal{G}_d) = (\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$), we have

$$\text{Tr}(\mathbb{H}_{d, > m}) = \Theta(d^{-\alpha}), \quad \text{Tr}(\mathbb{U}_{d, > M}) = \Theta(d^{-\alpha}). \quad (28)$$

Step 3. Checking the feature map concentration property at level $\{N(d), M(d), n(d), m(d)\}_{d \geq 1}$.

Let us first consider the case of a polynomial activation function σ . Denote D its degree and $u = u(d)$ the total (finite) number of nonzero eigenvalues of $\bar{\sigma}$ (which are associated to invariant polynomials of degree less or equal to D). Let us verify the feature map concentration property (Assumption 1 in Mei et al. (2021)) with sequence $u(d) \geq \max(m, M)$. Note that $u \geq \max(m, M)$, part (b) and (c) of the property are trivially verified in that case.

- (a) (*Hypercontractivity of finite eigenspaces on \mathcal{D}_d .*) The subspace of polynomials of degree less or equal to D on the hypercube and the sphere verifies the hypercontractivity property (see Lemmas 35 and 36 in Section H.3).
- (d) (*Concentration of diagonal elements.*) From Eq. (27) and Proposition 18 stated in Section E, we have

$$\begin{aligned} & \sup_{i \in [n]} \left| H_{d, > m}(\mathbf{x}_i, \mathbf{x}_i) - \mathbb{E}_{\mathbf{x}}[H_{d, > m}(\mathbf{x}, \mathbf{x})] \right| \\ & \leq \sum_{k=s+1}^D \xi_{d,k}^2 D(\mathcal{A}_d; k) \sup_{i \in [n]} \left| \Upsilon_k^{(d)}(\mathbf{x}_i, \mathbf{x}_i) - \mathbb{E}_{\mathbf{x}}[\Upsilon_k^{(d)}(\mathbf{x}, \mathbf{x})] \right| = o_{d, \mathbb{P}}(1) \cdot \mathbb{E}_{\mathbf{x}}[H_{d, > m}(\mathbf{x}, \mathbf{x})]. \end{aligned}$$

A similar computation shows the concentration of the diagonal elements of $U_{d, > M}$.

Let us now consider a non polynomial activation function σ in the case of $\mathcal{A}_d = \mathbb{S}^{d-1}(\sqrt{d})$ and $\mathcal{G}_d = \text{Cyc}_d$ (of degeneracy 1). Let us choose $\ell > 2S + 10$ such that $\mu_\ell(\sigma) \neq 0$ (it must exist otherwise σ would be a polynomial) and therefore $\xi_{d,\ell}^2 = \Theta_d(d^{-\ell})$. Consider $u = u(d)$ to be the number of eigenvalues such that $\lambda_{d,j}^2$ is strictly bigger than $\xi_{d,\ell}^2$. Then, $(\lambda_{d,j})_{j \leq u}$ do not contain any eigenvalues $\xi_{d,k}$ for $k \geq \ell$ and contain all $\xi_{d,k}$ for $k \leq s$. In particular, $u \geq \max(m, M)$. Denote $E = \{k : \xi_{d,k}^2 \leq \xi_{d,\ell}^2\}$: E contain all integers bigger or equal to ℓ .

Let us verify the feature map concentration property with the sequence $u(d)$ (part (a) is the same with D replaced by $\ell - 1$).

- (b) (*Properly decaying eigenvalues.*) We have

$$\begin{aligned} \text{Tr}(\mathbb{H}_{d, > u}) & \geq \xi_{d,\ell}^2 D(\mathbb{S}^{d-1}; \ell) = \Omega_d(d^{-\alpha}), \\ \text{Tr}(\mathbb{H}_{d, > u}^2) & = \sum_{k \in E} \xi_{d,k}^4 D(\mathbb{S}^{d-1}; k) \leq \xi_{d,\ell}^2 \text{Tr}(\mathbb{H}_{d, > u}). \end{aligned}$$

Hence,

$$\frac{\text{Tr}(\mathbb{H}_{d, > u})^2}{\text{Tr}(\mathbb{H}_{d, > u}^2)} \geq \xi_{d,\ell}^{-2} \cdot \text{Tr}(\mathbb{H}_{d, > u}) = \Omega_d(1) \cdot d^{2S+9} \geq \max(n, N)^{2+\delta}.$$

- (c) (*Hypercontractivity of the high degree part.*) Denote $\bar{\sigma}_{>u} = \bar{P}_E \bar{\sigma}$ the activation $\bar{\sigma}$ obtained by setting the first u eigenvalues to 0 (i.e., setting coefficients $k \notin E$ to zero in Eq. (25)). From Eq. (28), we need to show that for p as defined in Assumption 1.(a), we have

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}}[\bar{\sigma}_{>u}(\mathbf{x}; \boldsymbol{\theta})^{2p}]^{1/(2p)} = O_d(d^{-1/2+\delta}).$$

Denote $E_{\leq 4p} = E \cap \{0, \dots, 4p\}$ (recall that E contains all $k \geq \ell$) and decompose $\bar{\sigma}_{>u} = \bar{P}_{E_{\leq 4p}} \bar{\sigma} + \bar{P}_{>4p} \bar{\sigma}$. Then by triangle inequality we have,

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}} [\bar{\sigma}_{>u}(\mathbf{x}; \boldsymbol{\theta})^{2p}]^{1/(2p)} \leq \mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}} [\bar{P}_{E_{\leq 4p}} \bar{\sigma}(\mathbf{x}; \boldsymbol{\theta})^{2p}]^{1/(2p)} + \mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}} [\bar{P}_{>4p} \bar{\sigma}(\mathbf{x}; \boldsymbol{\theta})^{2p}]^{1/(2p)}.$$

Using hypercontractivity of polynomials of degree less or equal to $4p$, the first term is bounded by $O_d(d^{-1/2})$, while the second term is bounded in Proposition 29 in Section G.

(d) (*Concentration of diagonal elements.*) This is proved in Proposition 22 in Section F.

Step 4. Checking the spectral gap property at level $\{N(d), M(d), n(d), m(d)\}_{d \geq 1}$.

Let us now check the spectral gap property (Assumption 2 in Mei et al. (2021)).

(a) (*Number of samples.*) First by Eq. (26) and the assumption $d^{S-\alpha+\delta_0} \leq n \leq d^{S+1-\alpha-\delta_0}$, we have $m \leq n^{1-\delta}$ for $\delta > 0$ chosen sufficiently small. By the choice of m and recalling Eq. (28), we have

$$\begin{aligned} \lambda_{m+1}^{-2} \text{Tr}(\mathbb{H}_{d, > m}) &= \sup_{k \geq s+1} \{\xi_{d,k}^{-2}\} \cdot \text{Tr}(\mathbb{H}_{d, > m}) = \Omega_d(d^{S+1-\alpha}) \geq n^{1+\delta}, \\ \lambda_m^{-2} \text{Tr}(\mathbb{H}_{d, > m}) &= \xi_{d,s}^{-2} \text{Tr}(\mathbb{H}_{d, > m}) = O_d(1) \cdot d^{S-\alpha} \leq n^{1-\delta}, \end{aligned}$$

with $\delta > 0$ chosen sufficiently small.

(b) (*Number of features.*) By construction $M \geq m$. Furthermore, recalling Eq. (26) and the assumption $d^{S-\alpha+\delta_0} \leq N \leq d^{S+1-\alpha-\delta_0}$, we have $M \leq N^{1-\delta}$ for $\delta > 0$ chosen sufficiently small. By choice of M , $\lambda_{M+1}^2 \leq d^{-S-1+\delta}$. Hence,

$$\lambda_{M+1}^{-2} \text{Tr}(\mathbb{U}_{d, > M}) = \Omega_d(1) \cdot d^{S+1-\alpha-\delta} \geq N^{1+\delta},$$

for $\delta > 0$ chosen sufficiently small.

Finally notice that we used a different parametrization of λ in Eq. (6) and the condition in Mei et al. (2021) becomes $\lambda/d^\alpha = O_d(1) \cdot \text{Tr}(\mathbb{H}_{d, > m})$, i.e., $\lambda = O_d(1)$. This concludes the proof. \blacksquare

B.3. Proof of Theorem 5

We consider the same setting as in the previous section and consider

$$H_d(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}_{\boldsymbol{\theta}} [\bar{\sigma}(\mathbf{x}_1; \boldsymbol{\theta}) \bar{\sigma}(\mathbf{x}_2; \boldsymbol{\theta})].$$

Theorem 5 is a consequence of Theorem 4 in Mei et al. (2021) and the proof amounts to checking that H_d verifies the kernel concentration properties and eigenvalue condition (see Section 3.2 in Mei et al. (2021)). Note that some of the conditions were already covered in the proof of Theorem 2 and we will only mention the ones that still need to be verified. Furthermore, by the spectral gap property proven in Section B.2, the bound in Theorem 4 in Mei et al. (2021) (which is in term of a shrinkage operator) can indeed be rewritten as

$$R_{\text{KR,inv}}(f_d, \mathbf{X}, \lambda) = \|\bar{P}_{>s} f_d\|_{L^2}^2 + o_{d, \mathbb{P}}(1) \cdot (\|f_d\|_{L^{2+n}}^2 + \sigma_\varepsilon^2).$$

Proof [Proof of Theorem 5] We choose m as in the proof of Theorem 2.

Step 1. Checking the kernel concentration property at level $\{n(d), m(d)\}_{d \geq 1}$.

First notice that

$$\mathbb{E}_{\mathbf{x}}[H_{d, > m}(\mathbf{x}_i, \mathbf{x})] = \sum_{k=s+1}^{\infty} \xi_{d,k}^4 D(\mathcal{A}_d; k) \Upsilon_k^{(d)}(\mathbf{x}_i, \mathbf{x}),$$

and the concentration of the diagonal elements in the case of a polynomial activation function follows from the same argument as in Section B.2.

Hence, we only need to check this property in the case non polynomial activation function σ ($\mathcal{A}_d = \mathbb{S}^{d-1}(\sqrt{d})$ and $\mathcal{G}_d = \text{Cyc}_d$ of degeneracy 1). Let us choose u as in the proof of the feature map concentration property in Theorem 2.

- (Properly decaying eigenvalues.) We have

$$\begin{aligned} \text{Tr}(\mathbb{H}_{d, > u}^2) &\geq \xi_{d,\ell}^4 D(\mathbb{S}^{d-1}; \ell) = \Omega_d(1) \cdot d^{-\ell-1}, \\ \text{Tr}(\mathbb{H}_{d, > u}^4) &\leq \sup_{j \leq u} \{\lambda_{d,j}^6\} \text{Tr}(\mathbb{H}_{d, > u}) = O_d(1) \cdot d^{-3\ell}. \end{aligned}$$

Hence,

$$\frac{\text{Tr}(\mathbb{H}_{d, > u}^2)^2}{\text{Tr}(\mathbb{H}_{d, > u}^4)} = \Omega_d(1) \cdot d^{\ell-2} = \Omega_d(d^{2s}) \geq n^{2+\delta}.$$

- (Concentration of the diagonal elements of the kernel.) This is proven in Proposition 23 in Section F.

Step 2. Checking the eigenvalue condition at level $\{n(d), m(d)\}_{d \geq 1}$.

By the choice of m , we have

$$\lambda_{d, m+1}^{-4} \text{Tr}(\mathbb{H}_{d, > m}^2) = \frac{\sum_{k \geq s+1} \xi_{d,k}^4 D(\mathcal{A}_d; k)}{\sup_{k \geq s+1} \xi_{d,k}^4} \geq D(\mathcal{A}_d; s+1) = \Omega_d(d^{s+1-\alpha}) \geq n^{1+\delta}.$$

Again notice that we used a different parametrization of λ in Eq. (11) and the condition in Mei et al. (2021) becomes $\lambda/d^\alpha = O_d(1) \cdot \text{Tr}(\mathbb{H}_{d, > m})$, i.e., $\lambda = O_d(1)$. This concludes the proof. ■

Appendix C. Decomposition of invariant functions

In this section, we take $\mathcal{A}_d \in \{\mathbb{S}^{d-1}(\sqrt{d}), \mathcal{Q}^d\}$, and \mathcal{G}_d to be any group that is isomorphic to a subgroup of $\mathcal{O}(d)$ and that preserves \mathcal{A}_d . This section is mostly built on the technical background presented in Appendix H.

C.1. The invariant function class and the symmetrization operator

Let $L^2(\mathcal{A}_d)$ be the class of L^2 functions on \mathcal{A}_d equipped with uniform probability measure $\text{Unif}(\mathcal{A}_d)$. We define the invariant function class to be

$$L^2(\mathcal{A}_d, \mathcal{G}_d) = \left\{ f \in L^2(\mathcal{A}_d) : f(\mathbf{x}) = f(g \cdot \mathbf{x}), \forall \mathbf{x} \in \mathcal{A}_d, \forall g \in \mathcal{G}_d \right\}.$$

We define the symmetrization operator $\mathcal{S} : L^2(\mathcal{A}_d) \rightarrow L^2(\mathcal{A}_d, \mathcal{G}_d)$ to be

$$(\mathcal{S}f)(\mathbf{x}) = \int_{\mathcal{G}_d} f(g \cdot \mathbf{x}) \pi_d(dg).$$

C.2. Orthogonal polynomials on invariant function class

For either $\mathcal{A}_d \in \{\mathbb{S}^{d-1}(\sqrt{d}), \mathcal{Q}^d\}$, we define $V_{d, \leq k} \subseteq L^2(\mathcal{A}_d)$ to be the subspace spanned by all the degree ℓ polynomials, $V_{d, > k} \equiv V_{d, \leq k}^\perp \subseteq L^2(\mathcal{A}_d)$ to be the orthogonal complement of $V_{d, \leq k}$, and $V_{d, k} = V_{d, \leq k} \cap V_{d, \leq k-1}^\perp$. In words, $V_{d, k}$ contains all degree k polynomials that orthogonal to all polynomials of degree at most $k-1$. We further define $V_{d, < k} = V_{d, \leq k-1}$ and $V_{d, \geq k} = V_{d, > k-1}$.

Let $\bar{\mathbb{P}}_{\leq \ell}$ to be the projection operator on $L^2(\mathcal{A}_d, \text{Unif})$ that project a function onto $V_{d, \leq \ell}$, the space spanned by all the degree ℓ polynomials. Then it is easy to see that $\bar{\mathbb{P}}_{\leq \ell}$ and \mathcal{S} operator commute. This means, for any $f \in L^2(\mathcal{A}_d)$, we have

$$\bar{\mathbb{P}}_{\leq \ell}[\mathcal{S}(f)] = \mathcal{S}[\bar{\mathbb{P}}_{\leq \ell}(f)].$$

Similarly, we can define $\bar{\mathbb{P}}_\ell, \bar{\mathbb{P}}_{< \ell}, \bar{\mathbb{P}}_{> \ell}, \bar{\mathbb{P}}_{\geq \ell}$, which commute with \mathcal{S} . We denote $V_{d, \ell}(\mathcal{G}_d) \equiv \mathcal{P}_\ell(\mathcal{A}_d, \mathcal{G}_d)$ to be the space of polynomials in the images of $\bar{\mathbb{P}}_\ell \mathcal{S}$ (which is consistent with the definition of $V_{d, \ell}(\mathcal{G}_d)$ in Definition 1). Then we have

$$\mathcal{P}_\ell(\mathcal{A}_d, \mathcal{G}_d) = \bar{\mathbb{P}}_\ell(L^2(\mathcal{A}_d, \mathcal{G}_d)) = \mathcal{S}[\bar{\mathbb{P}}_\ell(L^2(\mathcal{A}_d))].$$

We denote $D(\mathcal{A}_d; k) = D(\mathcal{A}_d; \mathcal{G}_d; k) \equiv \dim(\mathcal{P}_k(\mathcal{A}_d, \mathcal{G}_d))$ to be the dimension of $\mathcal{P}_k(\mathcal{A}_d, \mathcal{G}_d)$. We denote $\{\bar{Y}_{kl}^{(d)}\}_{l \in [D(\mathcal{A}_d; k)]}$ to be a set of orthonormal polynomial basis in $\mathcal{P}_k(\mathcal{A}_d, \mathcal{G}_d)$. That means

$$\mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathcal{A}_d)} [\bar{Y}_{k_1 l_1}^{(d)}(\mathbf{x}) \bar{Y}_{k_2 l_2}^{(d)}(\mathbf{x})] = \mathbf{1}\{k_1 = k_2, l_1 = l_2\},$$

and

$$\bar{Y}_{kl}^{(d)}(\mathbf{x}) = \bar{Y}_{kl}^{(d)}(g \cdot \mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{A}_d, \quad \forall g \in \mathcal{G}_d.$$

C.3. A representation lemma

We have the following representation lemma. This lemma is important in the proofs of counting the degeneracy of groups (See Section D).

Lemma 11 (Convolution representation of projection operator) *Let $Q_k^{(d)}$ be the k -th Gegenbauer polynomial, or the k -th hypercubic Gegenbauer polynomial. For any fixed integer k , we have*

$$\frac{1}{D(\mathcal{A}_d; k)} \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\mathbf{x}) \bar{Y}_{kl}^{(d)}(\mathbf{y}) = \frac{B(\mathcal{A}_d; k)}{D(\mathcal{A}_d; k)} \int_{\mathcal{G}_d} Q_k^{(d)}(\langle \mathbf{x}, g \cdot \mathbf{y} \rangle) \pi_d(dg). \quad (29)$$

Proof [Proof of Lemma 11] Define

$$\Gamma_{1k}(\mathbf{x}, \mathbf{y}) = \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\mathbf{x}) \bar{Y}_{kl}^{(d)}(\mathbf{y}),$$

and

$$\Gamma_{2k}(\mathbf{x}, \mathbf{y}) = B(\mathcal{A}_d; k) \int_{\mathcal{G}_d} Q_k^{(d)}(\langle g \cdot \mathbf{x}, \mathbf{y} \rangle) \pi_d(dg).$$

Then Γ_{1k} and Γ_{2k} define two operators $\mathbb{T}_{1k}, \mathbb{T}_{2k} : L^2(\mathcal{A}_d) \rightarrow L^2(\mathcal{A}_d)$, i.e., for $j = 1, 2$,

$$\mathbb{T}_{jk}f(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim \text{Unif}(\mathcal{A}_d)}[\Gamma_{jk}(\mathbf{x}, \mathbf{y})f(\mathbf{y})].$$

Recall that $Q_k^{(d)}$ is a representation of the projector onto the subspace of degree- k spherical harmonics (see Eq. (104) in Section H.1.2). We deduce that

$$\mathbb{T}_{2k}f(\mathbf{x}) = \mathcal{S} \mathbb{E}_{\mathbf{y}}[B(\mathcal{A}_d; k) Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle) f(\mathbf{y})] = \mathcal{S} \bar{\mathbb{P}}_k f(\mathbf{x}),$$

and therefore $\mathbb{T}_{2k} = \mathcal{S} \bar{\mathbb{P}}_k$.

Furthermore, we have $\mathbb{T}_{1k} = \bar{\mathbb{P}}_k \mathcal{S}$. Indeed, the images of both \mathbb{T}_{1k} and $\bar{\mathbb{P}}_k \mathcal{S}$ are $\mathcal{P}_k(\mathcal{A}_d, \mathcal{G}_d)$, the space $\mathcal{P}_k(\mathcal{A}_d, \mathcal{G}_d)^\perp$ is the null space of both \mathbb{T}_{1k} and $\bar{\mathbb{P}}_k \mathcal{S}$, and $\mathbb{T}_{1k} \bar{Y}_{kp}^{(d)}(\mathbf{x}) = \bar{\mathbb{P}}_k \mathcal{S} \bar{Y}_{kp}^{(d)}(\mathbf{x}) = \bar{Y}_{kp}^{(d)}(\mathbf{x})$.

By the commutativity of $\bar{\mathbb{P}}_k$ and \mathcal{S} operator, we have $\mathbb{T}_{1k} = \bar{\mathbb{P}}_k \mathcal{S} = \mathcal{S} \bar{\mathbb{P}}_k = \mathbb{T}_{2k}$, and hence $\Gamma_{1k} = \Gamma_{2k}$. \blacksquare

C.4. Gegenbauer decomposition of invariant features and kernels

By Section H, for either $\mathcal{A}_d \in \{\mathbb{S}^{d-1}(\sqrt{d}), \mathcal{Q}^d\}$, for any activation function $\sigma \in L^2([-\sqrt{d}, \sqrt{d}], \tau_d^1)$ (where τ_d^1 is the distribution of $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle$ when $\mathbf{x}_1, \mathbf{x}_2 \sim_{iid} \text{Unif}(\mathcal{A}_d)$), we can define its coefficients $\xi_{d,k}(\sigma)$ defined by

$$\xi_{d,k}(\sigma) = \int_{[-\sqrt{d}, \sqrt{d}]} \sigma(x) Q_k^{(d)}(\sqrt{d}x) \tau_d^1(dx), \quad (30)$$

so that we have the following equation holds in $L^2([- \sqrt{d}, \sqrt{d}], \tau_d^1)$ sense

$$\sigma(x) = \sum_{k=0}^{\infty} \xi_{d,k}(\sigma) B(\mathcal{A}_d; k) Q_k^{(d)}(\sqrt{d}x).$$

For any group \mathcal{G}_d that is a subgroup of $\mathcal{O}(d)$, we define

$$\bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}) \equiv \int_{\mathcal{G}_d} \sigma(\langle \mathbf{x}, g \cdot \boldsymbol{\theta} \rangle / \sqrt{d}) \pi_d(dg).$$

Then, by the representation lemma (Lemma 11), we have

$$\begin{aligned} \bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}) &\equiv \sum_{k=0}^{\infty} \xi_{d,k}(\sigma) B(\mathcal{A}_d; k) \int_{\mathcal{G}_d} Q_k^{(d)}(\langle \mathbf{x}, g \cdot \boldsymbol{\theta} \rangle) \pi_d(dg) \\ &= \sum_{k=0}^{\infty} \xi_{d,k}(\sigma) \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\mathbf{x}) \bar{Y}_{kl}^{(d)}(\boldsymbol{\theta}). \end{aligned}$$

As a consequence, suppose we define

$$H_d(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)} [\bar{\sigma}(\mathbf{x}; \boldsymbol{\theta}) \bar{\sigma}(\mathbf{y}; \boldsymbol{\theta})].$$

Then we have

$$H_d(\mathbf{x}, \mathbf{y}) = \sum_{k=0}^{\infty} \xi_{d,k}(\sigma)^2 \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\mathbf{x}) \bar{Y}_{kl}^{(d)}(\mathbf{y}).$$

Appendix D. Counting the degeneracy

D.1. Counting the degeneracy of Cyc_d and $\text{Cyc}2\text{D}_{d_1, d_2}$ (Example 1 and 2)

Proposition 12 *Let $\mathcal{G}_d \in \{\text{Cyc}_d, \text{Cyc}2\text{D}_{d_1, d_2}\}$ with $d = d_1 \times d_2$. Let $\mathcal{A}_d \in \{\mathbb{S}^{d-1}(\sqrt{d}), \mathcal{Q}^d\}$. Then for any fixed $k \geq 1$, we have*

$$\dim(\mathcal{P}_k(\mathcal{A}_d, \mathcal{G}_d)) \equiv D(\mathcal{A}_d; k) = \Theta_d(d^{k-1}).$$

D.1.1. PROOF OF PROPOSITION 12

Here we state a key lemma that is used to prove Proposition 12.

Lemma 13 *Let $\mathcal{G}_d \in \{\text{Cyc}_d, \text{Cyc}2\text{D}_{d_1, d_2}\}$ with $d = d_1 \times d_2$. Denote*

$$F_k(\mathbf{z}) = \int_{\mathcal{G}_d} (\langle \mathbf{z}, g \cdot \mathbf{z} \rangle / d)^k \pi_d(dg).$$

Then for any fixed $k \geq 1$, we have

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})] = \Theta_d(d^{-1}), \quad (31)$$

$$\mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{Q}^d)}[F_k(\boldsymbol{\theta})] = \Theta_d(d^{-1}), \quad (32)$$

$$\mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{S}^{d-1}(\sqrt{d})}[F_k(\boldsymbol{\theta})] = \Theta_d(d^{-1}). \quad (33)$$

Proof [Proof of Lemma 13] We prove Eq. (31) and (33). The proof for Eq. (32) is similar to the proof of Eq. (31).

Let $\{L_\ell\}_{0 \leq \ell \leq d-1}$ be the matrix representation of group elements of Cyc_d or $\text{Cyc}2\text{D}_{d_1, d_2}$: when $\mathcal{G}_d = \text{Cyc}_d$, $g_\ell \in \text{Cyc}_d$ gives matrix representation L_ℓ for $0 \leq \ell \leq d-1$; when $\mathcal{G}_d = \text{Cyc}2\text{D}_{d_1, d_2}$, $g_{st} \in \text{Cyc}2\text{D}_{d_1, d_2}$ gives matrix representation $L_{s \times d_2 + t}$ for $0 \leq s \leq d_1 - 1$, $0 \leq t \leq d_2 - 1$. As a consequence, for either $\mathcal{G}_d \in \{\text{Cyc}_d, \text{Cyc}2\text{D}_{d_1, d_2}\}$, $L_0 = \mathbf{I}_d$ is the identity matrix. This gives

$$F_k(\mathbf{z}) = \|\mathbf{z}\|_2^{2k} / d^{k+1} + \sum_{l=1}^{d-1} \langle \mathbf{z}, L_l \mathbf{z} \rangle^k / d^{k+1}.$$

Step 1. The case $k = 1$. For either $\mathcal{G}_d \in \{\text{Cyc}_d, \text{Cyc}2\text{D}_{d_1, d_2}\}$, we have $\mathbb{E}[\langle \mathbf{z}, L_l \mathbf{z} \rangle] = 0$ for $1 \leq l \leq d-1$. As a consequence, we have

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_1(\mathbf{z})] = \mathbb{E}[\|\mathbf{z}\|_2^2 / d^2] + \sum_{l=1}^{d-1} \mathbb{E}[\langle \mathbf{z}, L_l \mathbf{z} \rangle / d^2] = \frac{1}{d}. \quad (34)$$

Step 2. The case $k = 2$. Note we have

$$\begin{aligned}
 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_2(\mathbf{z})] &= \mathbb{E}[\|\mathbf{z}\|_2^4/d^3] + \sum_{l=1}^{d-1} \mathbb{E}[\langle \mathbf{x}, L_l \mathbf{x} \rangle^2/d^3] \\
 &= \mathbb{E}\left[\left(\sum_{i=1}^d x_i^2\right)^2\right]/d^3 + \sum_{l=1}^{d-1} \mathbb{E}\left[\left(\sum_{i=1}^d x_i(L_l \mathbf{x})\right)^2\right]/d^3 \\
 &= \sum_{i,j=1}^d \mathbb{E}[x_i^2 x_j^2]/d^3 + \sum_{l=1}^{d-1} \sum_{i,j=1}^d \mathbb{E}[x_i(L_l \mathbf{x}) x_j(L_l \mathbf{x})]/d^3 \\
 &= \left(\frac{1}{d} + \frac{2}{d^2}\right) + \sum_{l=1}^{d-1} \sum_{i,j=1}^d \mathbb{E}[x_i(L_l \mathbf{x}) x_j(L_l \mathbf{x})]/d^3.
 \end{aligned}$$

Note that for either $\mathcal{G}_d \in \{\text{Cyc}_d, \text{Cyc}2\text{D}_{d_1, d_2}\}$, for any $i \in [d]$ and $1 \leq l \leq d-1$, the random variable $(L_l \mathbf{x})_i$ is independent from x_i . This gives

$$0 \leq \sum_{l=1}^{d-1} \sum_{i,j=1}^d \mathbb{E}[x_i(L_l \mathbf{x}) x_j(L_l \mathbf{x})]/d^3 \leq \frac{2(d-1)d}{d^3} = \Theta_d(d^{-1}).$$

As a consequence, we have

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_2(\mathbf{z})] = \Theta_d(d^{-1}). \quad (35)$$

Step 3. The case $k \geq 3$. By the moment formula of the χ^2 distribution, we have

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[(\|\mathbf{z}\|_2^2/d)^k] = 1 + o_d(1).$$

Moreover, for either $\mathcal{G}_d \in \{\text{Cyc}_d, \text{Cyc}2\text{D}_{d_1, d_2}\}$, for any $l \neq 0$, we have

$$\mathbb{E}[\langle \mathbf{z}, L_l \mathbf{z} \rangle]/d = 0.$$

As a consequence, by the Hanson-Wright inequality as in Lemma 14, for any fixed $k \geq 3$ and $\varepsilon > 0$, we have

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} \left[\sup_{1 \leq l \leq d-1} (\langle \mathbf{z}, L_l \mathbf{z} \rangle/d)^k \right] = O_d(d^{-k/2+\varepsilon}).$$

Therefore, for $k \geq 3$, we have

$$\left| \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})] - \frac{1}{d} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[(\|\mathbf{z}\|_2^2/d)^k] \right| \leq \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} \left[\sup_{1 \leq l \leq d-1} (\langle \mathbf{z}, L_l \mathbf{z} \rangle/d)^k \right] = o_d(1/d),$$

so that

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})] = 1/d + o_d(1/d). \quad (36)$$

Combining Eq. (34), (35), and (36) proves Eq. (31).

Step 4. From Gaussian to spherical. Note that when $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, we have $\|\mathbf{z}\|_2^2 \sim \chi^2(d)$ which is independent of $\sqrt{d} \cdot \mathbf{z}/\|\mathbf{z}\|_2 \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$. Hence, we have

$$\begin{aligned}
 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})] &= \mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{S}^{d-1}(\sqrt{d}), \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\boldsymbol{\theta})(\|\mathbf{z}\|_2^{2k}/d^k)] \\
 &= \mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{S}^{d-1}(\sqrt{d})}[F_k(\boldsymbol{\theta})] \cdot \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\|\mathbf{z}\|_2^{2k}/d^k].
 \end{aligned}$$

Note that for fixed $k \geq 1$, the moment formula for χ^2 distribution gives

$$\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\|\mathbf{z}\|_2^{2k}/d^k] = 1 + o_d(1).$$

Combining with Eq. (31), we have

$$\mathbb{E}_{\boldsymbol{\theta} \sim \mathbb{S}^{d-1}(\sqrt{d})}[F_k(\boldsymbol{\theta})] = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})]/\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\|\mathbf{z}\|_2^{2k}/d^k] = \Theta_d(d^{-1}).$$

This proves Eq. (33). ■

Proof [Proof of Proposition 12] Denote

$$P_k(\boldsymbol{\theta}) \equiv \frac{1}{B(\mathcal{A}_d; k)} \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\boldsymbol{\theta})^2 = \int_{\mathcal{G}_d} Q_k^{(d)}(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle) \pi_d(dg).$$

By Lemma 11, for any fixed $k \geq 1$, we have

$$\mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)}[P_k(\boldsymbol{\theta})] = \frac{D(\mathcal{A}_d; k)}{B(\mathcal{A}_d; k)}. \quad (37)$$

By Lemma 15, we have

$$P_k(\boldsymbol{\theta}) = \sum_{m=0}^k a_{d,k,m} F_m(\boldsymbol{\theta}),$$

where $|a_{d,k,m}| \leq C_{k,m}/d^{(k-m)/2}$. As a result, we have

$$\mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)}[P_k(\boldsymbol{\theta})] = \sum_{m=0}^k a_{d,k,m} \mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)}[F_m(\boldsymbol{\theta})] = \Theta(d^{-1}).$$

Combining with Eq. (37) shows that $D(\mathcal{A}_d; k) = \Theta(d^{-1}B(\mathcal{A}_d; k)) = \Theta(d^{k-1})$. This concludes the proof. ■

D.1.2. AUXILIARY LEMMAS

Lemma 14 (Hanson-Wright inequality) *There exists a universal constant $c > 0$, such that for any $t > 0$ and $d \in \mathbb{N}$, and any permutation matrix $L \in \mathbb{R}^{d \times d}$ be any permutation matrix, when $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ or $\mathbf{x} \sim \text{Unif}(\mathcal{Q}^d)$, we have*

$$\mathbb{P}\left(\left|\langle \mathbf{x}, L \cdot \mathbf{x} \rangle - \mathbb{E}[\langle \mathbf{x}, L \cdot \mathbf{x} \rangle]\right|/d \geq t\right) \leq 2 \cdot \exp\{-cd \cdot \min(t^2, t)\}.$$

Proof [Proof of Lemma 14] Note that for any permutation matrix L , we have $\|L\|_F \leq \sqrt{d}$, and $\|L\|_{\text{op}} \leq 1$. By the Hanson-Wright inequality of vectors with independent sub-Gaussian entries (for example, see Theorem 1.1 of Rudelson et al. (2013)), we have

$$\mathbb{P}\left(\left|\langle \mathbf{x}, L\mathbf{x} \rangle - \mathbb{E}[\langle \mathbf{x}, L\mathbf{x} \rangle]\right|/d > t\right) \leq 2 \exp\{-cd \cdot \min(t^2, t)\}.$$

This concludes the proof ■

Lemma 15 Let $Q_k^{(d)}$ be either the k 'th Gegenbauer polynomial or the k 'th hypercubic Gegenbauer polynomial (as defined in Section H). Let coefficients of monomials in $Q_k^{(d)}(d \cdot x)$ to be $\{a_{d,k,m}\}_{0 \leq m \leq k}$. That is, we have

$$Q_k^{(d)}(x) = \sum_{m=0}^k a_{d,k,m} (x/d)^m.$$

Then, for any fixed k , there exists constant $C(k)$, such that

$$|a_{d,k,m}| \leq C(k)/d^{(k-m)/2}.$$

Moreover, we have

$$\lim_{d \rightarrow \infty} a_{d,k,k} = 1.$$

Finally, for k and m in different parity, we have

$$a_{d,k,m} = 0.$$

Proof [Proof of Lemma 15] The proof holds by the following equation

$$\lim_{d \rightarrow \infty} \text{Coeff} \left\{ B(\mathcal{A}_d; k)^{1/2} Q_k^{(d)}(\sqrt{d} \cdot x) \right\} = \text{Coeff} \left\{ \frac{1}{\sqrt{k!}} \text{He}_k(x) \right\}.$$

when $Q_k^{(d)}$ is either Gegenbauer polynomial or Hypercubic Gegenbauer polynomial (See Eq. (110) and Eq. (112)). ■

D.2. Counting the degeneracy of band-limited function class (Example 3)

Proposition 16 Follow the notations of Example 3. Then for any fixed $k \geq 1$, we have

$$D(\mathbb{S}^{(d-1)}; k) = \Theta_d(d^{k-1}).$$

Here we state Lemma 17 that is used to prove Proposition 16. Given Lemma 17, the proof of Proposition 16 is the same as the proof of Proposition 12.

Lemma 17 Follow the notations of Example 3. Denote

$$F_k(z) = \int_{\text{Sft}_d} (\langle z, g \cdot z \rangle / d)^k \pi_d(dg).$$

Then for any fixed $k \geq 1$, we have

$$\mathbb{E}_{z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} [F_k(z)] = \Theta_d(d^{-1}).$$

Proof [Proof of Lemma 17]

We prove the lemma for the case when d is odd. We denote $u_1 = z_1^2$, and $u_i = z_{2i}^2 + z_{2i+1}^2$ for $i = 2, \dots, (d-1)/2$. Then we have

$$\begin{aligned}
 \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})] &= d^{-k} \cdot \mathbb{E}_{\mathbf{z}} \left\{ \int_{[0,1]} \left(\sum_{j=0}^{(d-1)/2} u_j \cos(2\pi j t) \right)^k dt \right\} \\
 &= d^{-k} \cdot \mathbb{E}_{\mathbf{z}} \left\{ \int_{[0,1]} \sum_{j_1, \dots, j_k=0}^{(d-1)/2} \left(\prod_{s \in [k]} u_{j_s} \cos(2\pi j_s t) \right) dt \right\} \\
 &= d^{-k} \cdot \sum_{j_1, \dots, j_k=0}^{(d-1)/2} \mathbb{E}_{\mathbf{z}} \left\{ \prod_{s \in [k]} u_{j_s} \right\} \left(\int_{[0,1]} \prod_{s \in [k]} \cos(2\pi j_s t) dt \right).
 \end{aligned} \tag{38}$$

Step 1. Bound Z function. First, we denote

$$Z(j_1, \dots, j_k) = \mathbb{E}_{\mathbf{z}} \left\{ \prod_{s \in [k]} u_{j_s} \right\}.$$

We have

$$\begin{aligned}
 \sup_{j_1, \dots, j_k} Z(j_1, \dots, j_k) &\leq \sup_{j_1, \dots, j_k} \prod_{s \in [k]} \mathbb{E}[u_{j_s}^{2k}]^{1/(2k)} \\
 &\leq \sup_{j \in \{0, 1, \dots, (d-1)/2\}} \mathbb{E}[u_j^{2k}]^{1/2} \leq 2^k \cdot \mathbb{E}_{G \sim \mathcal{N}(0,1)}[G^{2k}]^{1/2} \equiv M_k.
 \end{aligned} \tag{39}$$

Moreover, we have

$$\inf_{j_1, \dots, j_k} Z(j_1, \dots, j_k) \geq \mathbb{E}_{G \sim \mathcal{N}(0,1)}[G^2] = 1. \tag{40}$$

Step 2. Bound $|\mathcal{I}|$. Further, we denote

$$\mathcal{I} = \left\{ (j_1, \dots, j_k) \in \{0, \dots, (d-1)/2\}^k : \exists (\varepsilon_i)_{i \in [k]} \in \{\pm 1\}^k, \sum_{i=1}^k \varepsilon_i j_i = 0 \right\},$$

Then it is easy to see that

$$[(d+1)/2]^{k-1} \leq |\mathcal{I}| \leq 2 \cdot (d+1)^{k-1}. \tag{41}$$

Step 3. Bound E function. Next, we denote

$$E(j_1, \dots, j_k) = \int_{[0,1]} \prod_{s \in [k]} \cos(2\pi j_s t) dt.$$

It is easy to see that

$$\sup_{j_1, \dots, j_k} |E(j_1, \dots, j_k)| \leq 1. \tag{42}$$

Moreover, for any $(j_1, \dots, j_k) \notin \mathcal{I}$, we have $E(j_1, \dots, j_k) = 0$. For any $(j_1, \dots, j_k) \in \mathcal{I}$, we have

$$E(j_1, \dots, j_k) = \frac{1}{2^k} \int_{[0,1]} \prod_{s \in [k]} [\exp(i2\pi j_s t) + \exp(-i2\pi j_s t)] dt \geq 1/2^k. \tag{43}$$

The last inequality used the fact that $(j_1, \dots, j_k) \in \mathcal{I}$.

Step 4. Concludes the proof. Therefore, combining Eq. (38) (39) (41) (42), we have

$$\begin{aligned} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})] &\leq d^{-k} \cdot M_k \cdot \sum_{j_1, \dots, j_k=0}^{(d-1)/2} |E(j_1, \dots, j_k)| \\ &\leq d^{-k} \cdot M_k \cdot |\mathcal{I}| = O_d(d^{-1}). \end{aligned}$$

Combining Eq. (38) (40) (41) (43), we have

$$\begin{aligned} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{z})] &\geq d^{-k} \cdot \sum_{j_1, \dots, j_k=0}^{(d-1)/2} |E(j_1, \dots, j_k)| \\ &\geq d^{-k} \cdot |\mathcal{I}|/2^k = \Omega_d(d^{-1}). \end{aligned}$$

This concludes the proof. ■

Appendix E. Concentration for invariant groups with degeneracy $\alpha \leq 1$

Let $Q_k^{(d)}$ be the k 'th Gegenbauer polynomial on $\mathcal{A}_d \in \{\mathbb{S}^{d-1}(\sqrt{d}), \mathcal{Q}^d\}$ (see Section H for definitions). Let \mathcal{G}_d be an invariant group with degeneracy α . That means, for any fixed $k \geq \alpha$, we have $B(\mathcal{A}_d; k)/[D(\mathcal{A}_d; k)d^\alpha] = \Theta_d(1)$. For $k \in \mathbb{N}_{\geq 0}$, we denote

$$\Upsilon_k(\boldsymbol{\theta}) = \frac{1}{D(\mathcal{A}_d; k)} \sum_{l=1}^{D(\mathcal{A}_d; k)} \overline{Y}_{kl}^{(d)}(\boldsymbol{\theta}) \overline{Y}_{kl}^{(d)}(\boldsymbol{\theta}) = \frac{B(\mathcal{A}_d; k)}{D(\mathcal{A}_d; k)} \int_{\mathcal{G}_d} Q_k^{(d)}(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle) \pi_d(dg). \quad (44)$$

Then we have

$$\mathbb{E}[\Upsilon_k(\boldsymbol{\theta})] = 1.$$

In this section, we show that Υ_k concentration around its mean, for any fixed $k \geq 2$ and $\alpha \leq 1$.

E.1. Main proposition

Proposition 18 *Let \mathcal{G}_d be an invariant group with degeneracy $\alpha \leq 1$. Let $(\boldsymbol{\theta}_i)_{i \in [N]} \sim \text{Unif}(\mathcal{A}_d)$ where $N = O_d(d^p)$ for some fixed integer p . Let Υ_k be as defined in Eq. (44). Then for any fixed $k \geq 2$, we have*

$$\sup_{i \in [N]} \left| \Upsilon_k(\boldsymbol{\theta}_i) - 1 \right| = o_{d, \mathbb{P}}(1).$$

Proof [Proof of Proposition 18] Let us first focus on the sphere case $\mathcal{A}_d = \mathbb{S}^{d-1}(\sqrt{d})$. Let $(\mathbf{x}_i)_{i \in [N]} \sim_{iid} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Without loss of generality, we assume \mathbf{x}_i and $\boldsymbol{\theta}_i$ are coupled such that $\boldsymbol{\theta}_i = \sqrt{d} \cdot \mathbf{x}_i / \|\mathbf{x}_i\|_2$. Denote

$$F_k(\mathbf{z}) = \int_{\mathcal{G}_d} (\langle \mathbf{z}, g \cdot \mathbf{z} \rangle / d)^k \pi_d(dg).$$

Let $\{a_{d,k,m}\}_{0 \leq m \leq k}$ be the coefficients of monomials in $Q_k^{(d)}(d \cdot x)$. That is, we have

$$Q_k^{(d)}(x) = \sum_{m=0}^k a_{d,k,m} (x/d)^m.$$

Then

$$\int_{\mathcal{G}_d} Q_k^{(d)}(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle) \pi_d(dg) = \sum_{m=0}^k a_{d,k,m} F_m(\boldsymbol{\theta}).$$

Moreover, by Lemma 15, we have $|a_{d,k,m}| \leq C_{k,m}/d^{(k-m)/2}$, $\lim_{d \rightarrow \infty} a_{d,k,k} = 1$, and $a_{d,k,m} = 0$ for k and m of different parity.

Then we have

$$\begin{aligned}
 & \sup_{i \in [N]} |\Upsilon_k(\boldsymbol{\theta}_i) - \mathbb{E}[\Upsilon_k(\boldsymbol{\theta}_i)]| \\
 &= \frac{B(\mathbb{S}^{d-1}; k)}{D(\mathbb{S}^{d-1}; k)} \sup_{i \in [N]} \left| \int_{\mathcal{G}_d} Q_k(\langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle) \pi_d(dg) - \mathbb{E} \left[\int_{\mathcal{G}_d} Q_k(\langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle) \pi_d(dg) \right] \right| \\
 &\leq C \times d^\alpha \times \sum_{m=1}^k a_{d,k,m} \times \sup_{i \in [N]} \left| F_m(\boldsymbol{\theta}_i) - \mathbb{E}[F_m(\boldsymbol{\theta}_i)] \right| \\
 &\leq C \times d^\alpha \times \sum_{m=1}^k a_{d,k,m} \times \sup_{i \in [N]} \left| F_m(\mathbf{x}_i) - \mathbb{E}[F_m(\mathbf{x}_i)] \right| \cdot [d^m / \|\mathbf{x}_i\|_2^{2m}].
 \end{aligned}$$

By the concentration of χ^2 -distribution, for any $\varepsilon > 0$, the following event happens with high probability

$$\mathcal{E}_1 \equiv \left\{ \sup_{i \in [N]} \left| \|\mathbf{x}_i\|_2^2 / d - 1 \right| \leq 1/d^{1/2-\varepsilon} \right\}.$$

Moreover, combining Lemma 19 with Lemma 20, for any fixed $m \geq 2$, we have

$$\mathbb{E}[(F_m(\mathbf{x}) - \mathbb{E}[F_m(\mathbf{x})])^2] \leq C_m d^{-1-3\alpha/2}.$$

By the hypercontractivity property of Gaussian distribution as per Lemma 37, for any $\varepsilon > 0$, taking q sufficiently large, we have

$$\begin{aligned}
 & \mathbb{E} \left[\sup_{i \in [N]} \left| F_m(\mathbf{x}_i) - \mathbb{E}[F_m(\mathbf{x}_i)] \right| \right] \leq \mathbb{E} \left[\sum_{i=1}^N \left(F_m(\mathbf{x}_i) - \mathbb{E}[F_m(\mathbf{x}_i)] \right)^{2q} \right]^{1/(2q)} \\
 &\leq C(q) \cdot d^{p/(2q)} \cdot \mathbb{E}[(F_m(\mathbf{x}) - \mathbb{E}[F_m(\mathbf{x})])^2]^{1/2} \leq C d^{-1-3\alpha/2+\varepsilon}
 \end{aligned}$$

By Markov's inequality, we deduce that the following event happens with high probability

$$\mathcal{E}_2 \equiv \left\{ \forall 2 \leq m \leq k, \sup_{i \in [N]} \left| F_m(\mathbf{x}_i) - \mathbb{E}[F_m(\mathbf{x}_i)] \right| \leq C d^{-1-3\alpha/2+\varepsilon} \right\}.$$

Finally, by Lemma 19, we have

$$\mathbb{E}[F_1(\mathbf{x})^2] \leq C d^{-2\alpha},$$

and by the hypercontractivity property of low degree polynomials with Gaussian measure (Lemma 37), for any $\varepsilon > 0$, taking q sufficiently large, we have

$$\begin{aligned}
 & \mathbb{E} \left[\sup_{i \in [N]} \left| F_1(\mathbf{x}_i) - \mathbb{E}[F_1(\mathbf{x}_i)] \right| \right] \leq c \mathbb{E} \left[\sum_{i=1}^N F_1(\mathbf{x}_i)^{2q} \right]^{1/(2q)} + \mathbb{E}[F_1(\mathbf{x})^2]^{1/2} \\
 &\leq C(q) \cdot d^{p/(2q)} \cdot \mathbb{E}[F_1(\mathbf{x})^2]^{1/2} + \mathbb{E}[F_1(\mathbf{x})^2]^{1/2} \leq C d^{-\alpha+\varepsilon}.
 \end{aligned}$$

As a result, the following event happens with high probability

$$\mathcal{E}_3 \equiv \left\{ \sup_{i \in [N]} \left| F_1(\mathbf{x}_i) - \mathbb{E}[F_1(\mathbf{x}_i)] \right| \leq C d^{-\alpha+\varepsilon} \right\}.$$

When all the events \mathcal{E}_1 , \mathcal{E}_2 , and \mathcal{E}_3 happen, for any $k \geq 2$, we have

$$\begin{aligned} & \sup_{i \in [N]} |\Upsilon_k(\boldsymbol{\theta}_i) - \mathbb{E}[\Upsilon_k(\boldsymbol{\theta}_i)]| \\ & \leq C \times d^\alpha \times \sum_{m=1}^k a_{d,k,m} \times \sup_{i \in [N]} \left| F_m(\mathbf{x}_i) - \mathbb{E}[F_m(\mathbf{x}_i)] \right| \cdot [d^m / \|\mathbf{x}_i\|_2^{2m}] \\ & \leq C \times d^\alpha \times \left[d^{-(k-1)/2} d^{-\alpha+\varepsilon} + \sum_{m=2}^k d^{-(k-m)/2} \times d^{-1-3\alpha/2+\varepsilon} \right] = o_d(1). \end{aligned}$$

The case of the hypercube $\mathcal{A}_d \sim \mathcal{Q}^d$ follows similarly without introducing the gaussian measure and using Lemma 21 instead of Lemma 20. \blacksquare

E.2. Auxiliary Lemmas

Lemma 19 *Let \mathcal{G}_d be an invariant group with degeneracy $\alpha \leq 1$. Denote*

$$F_k(\mathbf{z}) = \int_{\mathcal{G}_d} (\langle \mathbf{z}, g \cdot \mathbf{z} \rangle / d)^k \pi_d(dg).$$

Then for any fixed $s \in [1, \infty)$ and integer $k \geq 1$, we have

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} [F_k(\mathbf{x})^s]^{1/s} = O_d(d^{-\alpha}), \quad (45)$$

$$\mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)} [F_k(\boldsymbol{\theta})^s]^{1/s} = O_d(d^{-\alpha}). \quad (46)$$

Proof [Proof of Lemma 19]

For $\boldsymbol{\theta} \in \mathcal{A}_d$, denote

$$P_k(\boldsymbol{\theta}) \equiv \frac{1}{B(\mathcal{A}_d; k)} \sum_{l=1}^{D(\mathcal{A}_d; k)} \bar{Y}_{kl}^{(d)}(\boldsymbol{\theta})^2 = \int_{\mathcal{G}_d} Q_k^{(d)}(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle) \pi_d(dg).$$

By Lemma 11 and by the assumption that \mathcal{G}_d is an invariant group with degeneracy $\alpha \leq 1$, i.e., $B(\mathcal{A}_k; k) / [D(\mathcal{A}_k; k) d^\alpha] = \Theta_d(1)$, for any fixed $k \geq 1$, we have

$$\mathbb{E}[P_k(\boldsymbol{\theta})] = \frac{D(\mathcal{A}_d; k)}{B(\mathcal{A}_d; k)} = O_d(d^{-\alpha}).$$

Throughout the proof, we will denote $L^s = L^s(\mathcal{A}_d)$ to be the L^s space with respect to distribution $\boldsymbol{\theta} \sim \text{Unif}(\mathcal{A}_d)$.

By the hypercontractivity of low degree polynomials on the sphere and the hypercube, as per Lemmas 35 and 36, for any $s \geq 1$, we have

$$\begin{aligned}
 \|P_k\|_{L^s} &= \frac{D(\mathcal{A}_d; k)}{B(\mathcal{A}_d; k)} \left\| \frac{B(\mathcal{A}_d; k)}{D(\mathcal{A}_d; k)} P_k \right\|_{L^s} \leq \frac{C_{k,s}}{d^\alpha} \left\| \frac{B(\mathcal{A}_d; k)}{D(\mathcal{A}_d; k)} P_k \right\|_{L^2} \\
 &= \frac{C_{k,s}}{d^\alpha} \left[\frac{1}{D(\mathcal{A}_d; k)^2} \sum_{l_1, l_2=1}^{D(\mathcal{A}_d; k)} \mathbb{E} [\overline{Y}_{kl_1}^{(d)}(\boldsymbol{\theta})^2 \overline{Y}_{kl_2}^{(d)}(\boldsymbol{\theta})^2] \right]^{1/2} \\
 &\leq \frac{C_{k,s}}{d^\alpha} \left[\frac{1}{D(\mathcal{A}_d; k)^2} \sum_{l_1, l_2=1}^{D(\mathcal{A}_d; k)} \mathbb{E} [\overline{Y}_{kl_1}^{(d)}(\boldsymbol{\theta})^4]^{1/2} \mathbb{E} [\overline{Y}_{kl_2}^{(d)}(\boldsymbol{\theta})^4]^{1/2} \right]^{1/2} \\
 &\leq \frac{C_{k,s}}{d^\alpha} \left[\frac{1}{D(\mathcal{A}_d; k)^2} \sum_{l_1, l_2=1}^{D(\mathcal{A}_d; k)} \mathbb{E} [\overline{Y}_{kl_1}^{(d)}(\boldsymbol{\theta})^2] \mathbb{E} [\overline{Y}_{kl_2}^{(d)}(\boldsymbol{\theta})^2] \right]^{1/2} = \frac{C_{k,s}}{d^\alpha}.
 \end{aligned} \tag{47}$$

Let $\{a_{d,k,m}\}_{0 \leq m \leq k}$ be the coefficients of monomials in $Q_k^{(d)}(d \cdot x)$. That is, we have

$$Q_k^{(d)}(x) = \sum_{m=0}^k a_{d,k,m} (x/d)^m.$$

Then

$$P_k = \sum_{m=0}^k a_{d,k,m} F_m. \tag{48}$$

Moreover, by Lemma 15, we have $|a_{d,k,m}| \leq C_{k,m}/d^{(k-m)/2}$, $\lim_{d \rightarrow \infty} a_{d,k,k} = 1$, and $a_{d,k,m} = 0$ for k and m have different parity.

We conclude the proof by induction over k . Note we have $F_0(\boldsymbol{\theta}) \equiv 1$. Moreover, for any $s \geq 1$, by Eq. (47) and (48) (and note that $a_{d,1,0} = 0$ and $\lim_{d \rightarrow \infty} a_{d,1,1} \rightarrow 1$), we have

$$\|F_1\|_{L^s} = \frac{1}{a_{d,1,1}} \|P_1\|_{L^s} \leq C_s/d^\alpha.$$

Fix a $k \geq 2$. Assume that, for any $1 \leq u \leq k-1$, we have $\|F_u\|_{L^s} \leq C_{u,s}/d^\alpha$ for $s \geq 1$, by Eq. (48) and (47), and the fact that $|a_{d,k,m}| \leq C_{k,m}/d^{(k-m)/2}$ and $\lim_{d \rightarrow \infty} a_{d,k,k} = 1$, we have

$$\begin{aligned}
 \|F_k\|_{L^s} &= \left\| \frac{1}{a_{d,k,k}} P_k - \sum_{m=1}^{k-1} a_{d,k,m} F_m - a_{d,k,0} \right\|_{L^s} \\
 &\leq \left\| \frac{1}{a_{d,k,k}} P_k \right\|_{L^s} + \sum_{m=1}^{k-1} |a_{d,k,m}| \|F_m\|_{L^s} + |a_{d,k,0}| \\
 &\leq C/d^\alpha + \left[\sum_{m=1}^{k-1} C/d^{(k-m)/2} \right] \cdot C/d^\alpha + C/d^{k/2} \leq C_{k,s}/d^\alpha,
 \end{aligned}$$

where we recall that we assume $\alpha \leq 1$.

Finally, for the case of $\mathcal{A}_d = \mathbb{S}^{d-1}(\sqrt{d})$, recalling that we can write

$$F_k(\mathbf{x}) = (\|\mathbf{x}\|_2^2/d)^k F_k(\boldsymbol{\theta}),$$

where $\boldsymbol{\theta} = \mathbf{x}/\|\mathbf{x}\|_2 \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ is independent of $\|\mathbf{x}\|_2$ in the case of $\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$. Hence, we get by Cauchy-Schwarz inequality

$$\mathbb{E}_{\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{x})^s]^{1/s} = \|F_k\|_{L^{2s}} \cdot \mathbb{E}_{\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)} \left[(\|\mathbf{x}\|_2^2/d)^{2sk} \right]^{1/(2s)} \leq C_{k,s}/d^\alpha,$$

by hypercontractivity of low degree polynomials for Gaussian measure (Lemma 37). \blacksquare

Lemma 20 *Let $\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$. Let \mathcal{G}_d be a general invariant group. Let $F_k(\mathbf{z})$ be defined as in Lemma 19. Then for any fixed $k \geq 1$, there exists a constant C_k , such that*

- If k is odd, then

$$\text{Var}_{\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{x})] \leq \frac{C_k}{d} \mathbb{E}_{\mathbf{x}}[F_{k-1}(\mathbf{x})^2].$$

- If k is even, then

$$\text{Var}_{\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)}[F_k(\mathbf{x})] \leq \frac{C_k}{d} \left(\mathbb{E}_{\mathbf{x}}[F_{k-2}(\mathbf{x})^2] \wedge \mathbb{E}_{\mathbf{x}}[F_k(\mathbf{x})^2]^{(2k-1)/(2k)} \right).$$

Proof [Proof of Lemma 20]

By the Gaussian Poincaré inequality, we have

$$\mathbb{E}_{\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)}[(F_k(\mathbf{x}) - \mathbb{E}[F_k(\mathbf{x})])^2] \leq \mathbb{E}[\|\nabla F_k(\mathbf{x})\|_2^2].$$

We have

$$\nabla F_k(\mathbf{x}) = k \int_{\mathcal{G}_d} (\langle \mathbf{x}, g \cdot \mathbf{x} \rangle / d)^{k-1} [(g \cdot \mathbf{x} + g^{-1} \cdot \mathbf{x}) / d] \pi_d(dg),$$

which gives

$$\begin{aligned} & \mathbb{E}[\|\nabla F_k(\mathbf{x})\|_2^2] \\ &= \frac{4k^2}{d} \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E} \left[(\langle \mathbf{x}, g_1 \cdot \mathbf{x} \rangle / d)^{k-1} (\langle \mathbf{x}, g_2 \cdot \mathbf{x} \rangle / d)^{k-1} \langle \mathbf{x}, g_1 g_2 \cdot \mathbf{x} \rangle / d \right] \pi_d(dg_1) \pi_d(dg_2). \end{aligned} \quad (49)$$

Case 1: Odd k . When k is odd, we have

$$\begin{aligned} & \mathbb{E}[\|\nabla F_k(\mathbf{x})\|_2^2] \\ & \leq \frac{4k^2}{d} \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E} \left[(\langle \mathbf{x}, g_1 \cdot \mathbf{x} \rangle / d)^{k-1} (\langle \mathbf{x}, g_2 \cdot \mathbf{x} \rangle / d)^{k-1} \|\mathbf{x}\|_2^2 / d \right] \pi_d(dg_1) \pi_d(dg_2) \\ & = \frac{4k^2}{d} \mathbb{E}[F_{k-1}(\mathbf{x})^2 (\|\mathbf{x}\|_2^2 / d)] \leq \frac{4k^2}{d} \mathbb{E}[F_{k-1}(\mathbf{x})^4]^{1/2} \mathbb{E}[(\|\mathbf{x}\|_2^2 / d)^2]^{1/2} \leq \frac{C_k}{d} \mathbb{E}[F_{k-1}(\mathbf{x})^2], \end{aligned}$$

where we used in the second line Cauchy-Schwarz inequality and that the matrix representations of g are orthogonal matrices, and in the last inequality the hypercontractivity of low degree polynomials for Gaussian measures (Lemma 37).

Case 2: Even k . Bound 1. When k is even, we have the following first bound

$$\begin{aligned} & \mathbb{E}[\|\nabla F_k(\mathbf{x})\|_2^2] \\ & \leq \frac{4k^2}{d} \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E} \left[(\langle \mathbf{x}, g_1 \cdot \mathbf{x} \rangle / d)^{k-2} (\langle \mathbf{x}, g_2 \cdot \mathbf{x} \rangle / d)^{k-2} \|\mathbf{x}\|_2^6 / d^3 \right] \pi_d(dg_1) \pi_d(dg_2) \\ & = \frac{4k^2}{d} \mathbb{E}[F_{k-2}(\mathbf{x})^2 (\|\mathbf{x}\|_2^6 / d^3)] \leq \frac{4k^2}{d} \mathbb{E}[F_{k-2}(\mathbf{x})^4]^{1/2} \mathbb{E}[(\|\mathbf{x}\|_2^6 / d^3)^2]^{1/2} \leq \frac{C_k}{d} \mathbb{E}[F_{k-2}(\mathbf{x})^2]. \end{aligned}$$

Case 3: Even k . Bound 2. When k is even, we have the following second bound, which follows by Hölder's inequality $\mathbb{E}[XY] \leq \mathbb{E}[|X|^{k/(k-1)}]^{(k-1)/k} \cdot \mathbb{E}[|Y|^k]^{1/k}$,

$$\begin{aligned} & \mathbb{E}[\|\nabla F_k(\mathbf{x})\|_2^2] \\ & = \frac{4k^2}{d} \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E} \left[(\langle \mathbf{x}, g_1 \cdot \mathbf{x} \rangle / d)^{k-1} (\langle \mathbf{x}, g_2 \cdot \mathbf{x} \rangle / d)^{k-1} \langle \mathbf{x}, g_1 g_2 \cdot \mathbf{x} \rangle / d \right] \pi_d(dg_1) \pi_d(dg_2) \\ & \leq \frac{4k^2}{d} \mathbb{E} \left[\int_{\mathcal{G}_d \times \mathcal{G}_d} (\langle \mathbf{x}, g_1 \cdot \mathbf{x} \rangle / d)^k (\langle \mathbf{x}, g_2 \cdot \mathbf{x} \rangle / d)^k \pi_d(dg_1) \pi_d(dg_2) \right]^{(k-1)/k} \\ & \quad \times \mathbb{E} \left[\int_{\mathcal{G}_d} (\langle \mathbf{x}, g \cdot \mathbf{x} \rangle / d)^k \pi_d(dg) \right]^{1/k} \\ & = \frac{4k^2}{d} \mathbb{E}[F_k(\mathbf{x})^2]^{(k-1)/k} \times \mathbb{E}[F_k(\mathbf{x})]^{1/k} \leq \frac{4k^2}{d} \mathbb{E}[F_k(\mathbf{x})^2]^{(2k-1)/(2k)}. \end{aligned}$$

Combining these two bounds yields the result for k even. ■

Lemma 21 *Let $\boldsymbol{\theta} \sim \text{Unif}(\mathcal{Q}^d)$. Let \mathcal{G}_d be a general invariant group that preserves \mathcal{Q}^d . Let $F_k(\mathbf{z})$ be defined as in Lemma 19. Then for any fixed $k \geq 1$, there exists a constant C_k , such that*

- If k is odd, then

$$\text{Var}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{Q}^d)}[F_k(\boldsymbol{\theta})] \leq \frac{C_k}{d} \mathbb{E}_{\boldsymbol{\theta}}[F_{k-1}(\boldsymbol{\theta})^2] + \frac{C_k}{d^3}.$$

- If k is even, then

$$\text{Var}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{Q}^d)}[F_k(\boldsymbol{\theta})] \leq \frac{C_k}{d} \left(\mathbb{E}_{\boldsymbol{\theta}}[F_{k-2}(\boldsymbol{\theta})^2] \wedge \mathbb{E}_{\boldsymbol{\theta}}[F_k(\boldsymbol{\theta})^2]^{(2k-1)/(2k)} \right) + \frac{C_k}{d^3}.$$

Proof [Proof of Lemma 21]

The proof is similar to the proof of Lemma 20. By the discrete Poincaré inequality, we have

$$\mathbb{E}_{\boldsymbol{\theta} \sim \text{Unif}(\mathcal{Q}_d)} \left[(F_k(\boldsymbol{\theta}) - \mathbb{E}[F_k(\boldsymbol{\theta})])^2 \right] \leq \mathbb{E}_{\boldsymbol{\theta}} \left[\sum_{i=1}^d D_i F_k(\boldsymbol{\theta})^2 \right],$$

where D_i denote the discrete derivative defined as

$$D_i f(\boldsymbol{\theta}) = \frac{f(\boldsymbol{\theta}) - f(\boldsymbol{\theta}_{-i})}{2},$$

with $\boldsymbol{\theta}_{-i} = (\theta_1, \dots, \theta_{i-1}, -\theta_i, \theta_{i+1}, \dots, \theta_d)$. Let $\varphi_g(\boldsymbol{\theta}) = (\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle / d)^k$, then

$$\begin{aligned} D_i \varphi_g(\boldsymbol{\theta}) &= \frac{(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle / d)^k - (\langle \boldsymbol{\theta}_{-i}, g \cdot \boldsymbol{\theta} \rangle / d)^k}{2} + \frac{(\langle \boldsymbol{\theta}_{-i}, g \cdot \boldsymbol{\theta} \rangle / d)^k - (\langle \boldsymbol{\theta}_{-i}, g \cdot \boldsymbol{\theta}_{-i} \rangle / d)^k}{2} \\ &=: D_{i,1} \varphi_g(\boldsymbol{\theta}) + D_{i,2} \varphi_g(\boldsymbol{\theta}). \end{aligned}$$

We have $\langle \boldsymbol{\theta}_{-i}, g \cdot \boldsymbol{\theta} \rangle = \langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle - 2\theta_i(g \cdot \boldsymbol{\theta})_i$. By Taylor expansion, the first term verifies (recall that $g \cdot \boldsymbol{\theta} \in \mathcal{Q}^d$ and $\theta_i^2(g \cdot \boldsymbol{\theta})_i^2 = 1$)

$$D_{i,1} \varphi_g(\boldsymbol{\theta}) = k(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle / d)^{k-1}(\theta_i(g \cdot \boldsymbol{\theta})_i / d) - k(k-1)X_{i,1}(\boldsymbol{\theta}, g)^{k-2} / d^2,$$

where $X_{i,1}(\boldsymbol{\theta}, g)$ is on the line segment between $\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle / d$ and $\langle \boldsymbol{\theta}_{-i}, g \cdot \boldsymbol{\theta} \rangle / d$. Similarly, Taylor expansion on the second term yields

$$D_{i,2} \varphi_g(\boldsymbol{\theta}) = k(\langle \boldsymbol{\theta}_{-i}, g \cdot \boldsymbol{\theta}_{-i} \rangle / d)^{k-1}(\theta_i(g^{-1} \cdot \boldsymbol{\theta}_{-i})_i / d) + k(k-1)X_{i,2}(\boldsymbol{\theta}, g)^{k-2} / d^2,$$

where $X_{i,2}(\boldsymbol{\theta}, g)$ is on the line segment between $\langle \boldsymbol{\theta}_{-i}, g \cdot \boldsymbol{\theta}_{-i} \rangle / d$ and $\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta}_{-i} \rangle / d$.

Using Jensen's inequality to separate each of the 4 terms in $D_i \varphi_g(\boldsymbol{\theta})$, using that $\boldsymbol{\theta}_{-i}$ and $\boldsymbol{\theta}$ have the same distribution, we get

$$\begin{aligned} &\mathbb{E}_{\boldsymbol{\theta}} \left[\sum_{i=1}^d D_i F_k(\boldsymbol{\theta})^2 \right] \\ &\leq \frac{32k^2}{d} \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E} \left[(\langle \boldsymbol{\theta}, g_1 \cdot \boldsymbol{\theta} \rangle / d)^{k-1} (\langle \boldsymbol{\theta}, g_2 \cdot \boldsymbol{\theta} \rangle / d)^{k-1} \langle \boldsymbol{\theta}, g_1 g_2 \cdot \boldsymbol{\theta} \rangle / d \right] \pi_d(dg_1) \pi_d(dg_2) \\ &\quad + \frac{16k^2(k-1)^2}{d^4} \sum_{s \in \{1,2\}} \sum_{i=1}^d \int_{\mathcal{G}_d \times \mathcal{G}_d} \mathbb{E} \left[X_{i,s}(\boldsymbol{\theta}, g_1)^{k-2} X_{i,s}(\boldsymbol{\theta}, g_2)^{k-2} \right] \pi_d(dg_1) \pi_d(dg_2). \end{aligned}$$

Noticing that $\sup_{i,s,\boldsymbol{\theta},g} |X_{i,s}(\boldsymbol{\theta}, g)| \leq 1$, the second term in the above equation can be bounded by C_k/d^3 . The first term in the above equation can be bounded using the same way as bounding the right hand side of Eq. (49) as in the proof of Lemma 20. This concludes the proof. \blacksquare

Appendix F. Kernel concentration for the cyclic group and general σ

Throughout this section, we will always take $\mathcal{G}_d = \text{Cyc}_d$ to be the cyclic group, and $\mathcal{A}_d = \mathbb{S}^{d-1}(\sqrt{d})$ to be the sphere. We will write in short $B(d, k) = B(\mathbb{S}^{d-1}(\sqrt{d}); k)$ and $D(d, k) = D(\mathbb{S}^{d-1}(\sqrt{d}); \text{Cyc}_d; k)$. We recall that the cyclic group has degeneracy 1, i.e., for each integers $k \geq 1$, $B(d, k)/D(d, k) = \Theta_d(d)$.

F.1. Main propositions

Let the Gegenbauer decomposition of σ be

$$\sigma(x) = \sum_{k=0}^{\infty} \xi_{d,k}(\sigma) B(d, k) Q_k^{(d)}(\sqrt{d}x).$$

For $S \subseteq \mathbb{N}$, we define

$$\sigma_{d,S}(x) = \sum_{k \in S} \xi_{d,k}(\sigma) B(d, k) Q_k^{(d)}(\sqrt{d}x). \quad (50)$$

For any $\|\boldsymbol{\theta}_1\|_2 = \|\boldsymbol{\theta}_2\|_2 = \sqrt{d}$ and any $S \subseteq \mathbb{Z}_{\geq 0}$, denote

$$h_{d,S}(\langle \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \rangle / d) \equiv \mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} [\sigma_{d,S}(\langle \boldsymbol{\theta}_1, \mathbf{x} \rangle / \sqrt{d}) \sigma_{d,S}(\langle \boldsymbol{\theta}_2, \mathbf{x} \rangle / \sqrt{d})].$$

Proposition 22 *Let $\ell \geq 2$ be a fixed integer. Assume that $\sigma \in C^{\ell \vee 3}(\mathbb{R})$ be a $\ell \vee 3$ 'th continuously differentiable function with derivatives satisfy $\sup_{0 \leq k \leq \ell \vee 3} \sigma^{(k)}(u) \leq c_0 \exp(c_1 u^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$.*

Define $H_{d,S} : \mathbb{S}^{d-1}(\sqrt{d}) \times \mathbb{S}^{d-1}(\sqrt{d}) \rightarrow \mathbb{R}$ via

$$H_{d,S}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \equiv \int_{\text{Cyc}_d} h_{d,S}(\langle \boldsymbol{\theta}_1, g \cdot \boldsymbol{\theta}_2 \rangle / d) \pi_d(dg). \quad (51)$$

Then, for $N = d^p$ for any fixed p , letting $(\boldsymbol{\theta}_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, we have

$$\sup_{i \in [N]} \left| H_{d, \geq \ell}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i) - \mathbb{E} H_{d, \geq \ell}(\boldsymbol{\theta}, \boldsymbol{\theta}) \right| = o_{d, \mathbb{P}}(1) \cdot \mathbb{E} H_{d, \geq \ell}(\boldsymbol{\theta}, \boldsymbol{\theta}). \quad (52)$$

Moreover, we have $\mathbb{E} H_{d, \geq \ell}(\boldsymbol{\theta}, \boldsymbol{\theta}) = O_d(d^{-1})$.

Proof [Proof of Proposition 22] We let $C, C_k, C_{k,\ell}$ be constants that depend on σ, k , and ℓ but independent of dimension d . The exact values of these constant can change from line to line.

Step 1. Finite subset $S \subseteq \{2, 3, \dots\}$. Note we have

$$H_{d,S}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \equiv \sum_{k \in S} \xi_{d,k}(\sigma)^2 B(d, k) \int_{\text{Cyc}_d} Q_k^{(d)}(\langle \boldsymbol{\theta}_1, g \cdot \boldsymbol{\theta}_2 \rangle) \pi_d(dg).$$

By Lemma 11 and Proposition 12, for any $S \subseteq \mathbb{N}$ with finite cardinality $|S| < \infty$, we have

$$\mathbb{E}[H_{d,S}(\boldsymbol{\theta}, \boldsymbol{\theta})] = \sum_{k \in S} \xi_{d,k}(\sigma)^2 D(d, k) = \Theta_d(d^{-1}).$$

Moreover, by Proposition 18, we have

$$\sup_{i \in [N]} \left| H_{d,S}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i) - \mathbb{E}[H_{d,S}(\boldsymbol{\theta}, \boldsymbol{\theta})] \right| \leq \sum_{k \in S} \xi_{d,k}(\sigma)^2 D(d, k) \sup_{i \in [N]} \left| \Upsilon_k(\boldsymbol{\theta}_i) - 1 \right| = o_d(1) \cdot \mathbb{E}[H_{d,S}(\boldsymbol{\theta}, \boldsymbol{\theta})].$$

Step 2. For general set $S = \{u : u \geq \ell\}$.

By Lemma 27, we have $\sup_{d \geq 1} \sup_{\gamma \in [-1, 1]} |h_{d, \geq \ell}^{(\ell)}(\gamma)| \leq C_\ell$. Therefore, for any $\gamma \in [-1, 1]$, we have

$$\left| h_{d, \geq \ell}(\gamma) - \sum_{k=0}^{\ell-1} \frac{1}{k!} h_{d, \geq \ell}^{(k)}(0) \gamma^k \right| \leq C_\ell \cdot |\gamma|^{\ell+1}. \quad (53)$$

By Lemma 28, for any $k \leq \ell - 1$, we have

$$\left| h_{d, \geq \ell}^{(k)}(0) \right| \leq C_{k, \ell} \cdot d^{-(\ell-k)/2}. \quad (54)$$

Moreover, by the Hanson-Wright inequality as in Lemma 14, since N is at most polynomial in d , then for any $\delta > 0$, we have

$$\sup_{1 \leq k \leq \ell+1} \sup_{g \in \text{Cyc}_d \setminus \mathbf{I}} \sup_{i \in [N]} \left| \langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle^k \right| \cdot d^{-k/2-\delta} = o_{d, \mathbb{P}}(1), \quad (55)$$

and

$$\sup_{1 \leq k \leq \ell+1} \sup_{g \in \text{Cyc}_d \setminus \mathbf{I}} \mathbb{E} \left[\left| \langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle^k \right| \right] \cdot d^{-k/2-\delta} = o_d(1). \quad (56)$$

Therefore, by Eq. (53), (54), (55) and (56), we have

$$\sup_{g \in \text{Cyc}_d \setminus \mathbf{I}} \sup_{i \in [N]} \left| h_{d, \geq \ell}(\langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle / d) \right| = O_{d, \mathbb{P}}(d^{-\ell/2+\delta}),$$

and

$$\sup_{g \in \text{Cyc}_d \setminus \mathbf{I}} \mathbb{E} \left[\left| h_{d, \geq \ell}(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle / d) \right| \right] = O_d(d^{-\ell/2+\delta}).$$

As a result, for any $\ell \geq 3$, we have

$$\begin{aligned} & \sup_{i \in [N]} \left| H_{d, \geq \ell}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i) - \mathbb{E} H_{d, \geq \ell}(\boldsymbol{\theta}, \boldsymbol{\theta}) \right| \\ &= \sup_{i \in [N]} \left| \int_{\text{Cyc}_d \setminus \mathbf{I}} h_{d, \geq \ell}(\langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle / d) \pi_d(dg) - \mathbb{E} \int_{\text{Cyc}_d \setminus \mathbf{I}} h_{d, \geq \ell}(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle / d) \pi_d(dg) \right| \\ &\leq \sup_{i \in [N]} \left| \int_{\text{Cyc}_d \setminus \mathbf{I}} h_{d, \geq \ell}(\langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle / d) \pi_d(dg) \right| + \left| \mathbb{E} \int_{\text{Cyc}_d \setminus \mathbf{I}} h_{d, \geq \ell}(\langle \boldsymbol{\theta}, g \cdot \boldsymbol{\theta} \rangle / d) \pi_d(dg) \right| \\ &\leq O_{d, \mathbb{P}}(d^{-\ell/2+\delta}) = o_{d, \mathbb{P}}(d^{-1}). \end{aligned}$$

By the arguments in Step 1, for any $\ell \geq 2$, we have

$$\sup_{i \in [N]} \left| H_{d, \geq \ell}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i) - \mathbb{E} H_{d, \geq \ell}(\boldsymbol{\theta}, \boldsymbol{\theta}) \right| \leq o_{d, \mathbb{P}}(d^{-1}).$$

Finally, for any $\ell \geq 2$, for any σ such that $\sigma_{d,\geq\ell}$ that is non-trivial (if $\sigma_{d,\geq\ell} = 0$, this proposition holds trivially), we have

$$\mathbb{E}H_{d,\geq\ell}(\boldsymbol{\theta}, \boldsymbol{\theta}) = \Theta_d(d^{-1}).$$

This proves the proposition. \blacksquare

Proposition 23 *Let $\ell \geq 2$ be a fixed integer. Assume that $\sigma \in C(\mathbb{R})$ be a continuous function with $|\sigma(u)| \leq c_0 \exp(c_1 u^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$.*

Define $H_{d,S} : \mathbb{S}^{d-1}(\sqrt{d}) \times \mathbb{S}^{d-1}(\sqrt{d}) \rightarrow \mathbb{R}$ via

$$H_{d,S}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \equiv \int_{\text{Cyc}_d} h_{d,S}(\langle \boldsymbol{\theta}_1, g \cdot \boldsymbol{\theta}_2 \rangle / d) \pi_d(dg). \quad (57)$$

Then, for $N = O(d^p)$ for any fixed p , letting $(\boldsymbol{\theta}_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, we have

$$\sup_{i \in [N]} \left| \mathbb{E}_{\boldsymbol{\theta}}[H_{d,\geq\ell}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] - \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d,\geq\ell}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] \right| = o_{d,\mathbb{P}}(1) \cdot \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d,\geq\ell}(\boldsymbol{\theta}', \boldsymbol{\theta})^2]. \quad (58)$$

Proof [Proof of Proposition 23]

Denoting $\mu_k(\sigma) = \mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)\text{He}_k(G)]$. Let $q = \min\{k \geq \ell : \mu_k(\sigma) \neq 0\}$ and let $u = q + 2$. We consider the case when $q < \infty$, since for $q = \infty$, the claim holds trivially. We have the expression

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}}[H_{d,\geq\ell}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] &= \sum_{k=\ell}^{\infty} \xi_{d,k}(\sigma)^4 B(d, k) \int_{\text{Cyc}_d} Q_k^{(d)}(\langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle) \pi_d(dg) \\ &= \mathbb{E}_{\boldsymbol{\theta}}[H_{d, [\ell, u]}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] + \mathbb{E}_{\boldsymbol{\theta}}[H_{d, \geq u}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2]. \end{aligned}$$

Step 1. Upper bounding $\mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, \geq u}(\boldsymbol{\theta}', \boldsymbol{\theta})^2]$ and $\mathbb{E}_{\boldsymbol{\theta}}[H_{d, \geq u}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2]$. We have

$$\begin{aligned} &\sup_{\boldsymbol{\theta}_i} \mathbb{E}_{\boldsymbol{\theta}}[H_{d, \geq u}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] \\ &= \sup_{\boldsymbol{\theta}_i} \sum_{k=u}^{\infty} \xi_{d,k}(\sigma)^4 \sum_{l \in [D(d, k)]} \bar{Y}_{kl}^{(d)}(\boldsymbol{\theta}_i)^2 \leq \sup_{\boldsymbol{\theta}_i} \sum_{k=u}^{\infty} \xi_{d,k}(\sigma)^4 \sum_{l \in [B(d, k)]} Y_{kl}^{(d)}(\boldsymbol{\theta}_i)^2 \\ &= \sum_{k=u}^{\infty} \xi_{d,k}(\sigma)^4 B(d, k) Q_k^{(d)}(d) = \sum_{k=u}^{\infty} \xi_{d,k}(\sigma)^4 B(d, k) \\ &\leq \left[\sup_{k \geq u} B(d, k)^{-1} \right] \cdot \left[\sum_{k=u}^{\infty} \xi_{d,k}(\sigma)^2 B(d, k) \right]^2 \\ &= B(d, u)^{-1} \mathbb{E}_{\mathbf{x}, \boldsymbol{\theta} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))}[\sigma(\langle \mathbf{x}, \boldsymbol{\theta} \rangle / \sqrt{d})^2] = \Theta_d(d^{-u}). \end{aligned} \quad (59)$$

This also gives

$$\mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, \geq u}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] = \Theta_d(d^{-u}). \quad (60)$$

Step 2. Upper bounding $\sup_{i \in [N]} |\mathbb{E}_{\boldsymbol{\theta}}[H_{d, [\ell, u]}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] - \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, [\ell, u]}(\boldsymbol{\theta}', \boldsymbol{\theta})^2]|$. By Proposition 18, we have

$$\begin{aligned} & \sup_{i \in [N]} \left| \mathbb{E}_{\boldsymbol{\theta}}[H_{d, [\ell, u]}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] - \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, [\ell, u]}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] \right| \\ & \leq \sum_{k=\ell}^{u-1} \xi_{d,k}(\sigma)^4 D(d, k) \sup_{i \in [N]} \left| \frac{B(d, k)}{D(d, k)} \int_{\text{Cyc}_d} Q_k^{(d)}(\langle \boldsymbol{\theta}_i, g \cdot \boldsymbol{\theta}_i \rangle) \pi_d(dg) - 1 \right| \\ & = o_{d, \mathbb{P}}(1) \cdot \left[\sum_{k=\ell}^{u-1} \xi_{d,k}(\sigma)^4 D(d, k) \right] = o_{d, \mathbb{P}}(1) \cdot \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, [\ell, u]}(\boldsymbol{\theta}', \boldsymbol{\theta})^2]. \end{aligned} \quad (61)$$

Step 3. Lower bounding $\mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, \geq \ell}(\boldsymbol{\theta}', \boldsymbol{\theta})^2]$. We have

$$\mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, \geq \ell}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] = \sum_{k=\ell}^{\infty} \xi_{d,k}(\sigma)^4 D(d, k) \geq \xi_{d,q}(\sigma)^4 D(d, q) = \Theta_d(d^{-q-1}). \quad (62)$$

The last equality is by Proposition 12, and the fact that

$$\lim_{d \rightarrow \infty} \xi_{d,q}(\sigma)^2 B(d, q) = \mu_q(\sigma)^2 / q! > 0.$$

Step 4. Complete the proof. By Eq. (59), (60), (61) and (62), we have

$$\begin{aligned} & \sup_{i \in [N]} \left| \mathbb{E}_{\boldsymbol{\theta}}[H_{d, \geq \ell}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] - \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, \geq \ell}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] \right| \\ & \leq \sup_{i \in [N]} \left| \mathbb{E}_{\boldsymbol{\theta}}[H_{d, [\ell, u]}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] - \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, [\ell, u]}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] \right| \\ & \quad + \sup_{i \in [N]} \left| \mathbb{E}_{\boldsymbol{\theta}}[H_{d, \geq \ell}(\boldsymbol{\theta}_i, \boldsymbol{\theta})^2] \right| + \left| \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, \geq \ell}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] \right| \\ & \leq o_{d, \mathbb{P}}(1) \cdot \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, [\ell, u]}(\boldsymbol{\theta}', \boldsymbol{\theta})^2] + \Theta_d(d^{-u}) = o_{d, \mathbb{P}}(1) \cdot \mathbb{E}_{\boldsymbol{\theta}, \boldsymbol{\theta}'}[H_{d, \geq \ell}(\boldsymbol{\theta}', \boldsymbol{\theta})^2]. \end{aligned}$$

This completes the proof. ■

F.2. Auxiliary lemmas

The following lemma is a reformulation of (Ghorbani et al., 2021, Lemma 5).

Lemma 24 *Assume $\sigma \in C(\mathbb{R})$ with $\sigma(u)^2 \leq c_0 \exp(c_1 u^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$. Then*

(a) $\mathbb{E}_{G \sim \mathcal{N}(0,1)}[\sigma(G)^2] < \infty$.

(b) *Let $\|\boldsymbol{x}\|_2 = \sqrt{d}$. Then there exists $d_0 = d_0(c_1)$ such that, for $\boldsymbol{w} \sim \text{Unif}(\mathbb{S}^{d-1})$,*

$$\sup_{d \geq d_0} \mathbb{E}_{\boldsymbol{w}}[\sigma(\langle \boldsymbol{w}, \boldsymbol{x} \rangle)^2] < \infty. \quad (63)$$

(c) Let $\|\mathbf{x}\|_2 = \sqrt{d}$, $\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})$ and $\tau \sim \chi(d)/\sqrt{d}$. Then we have

$$\lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}, \tau} \left[\left(\sigma(\tau \langle \mathbf{w}, \mathbf{x} \rangle) - \sigma(\langle \mathbf{w}, \mathbf{x} \rangle) \right)^2 \right] = 0. \quad (64)$$

Lemma 25 Assume that $\psi, \phi \in C(\mathbb{R})$ with $\psi(u)^2, \phi(u)^2 \leq c_0 \exp(c_1 u^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$. Denote

$$\begin{aligned} E_d[\psi, \phi](\gamma) &\equiv \mathbb{E}_{\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})} [\psi(\langle \mathbf{x}, \mathbf{w} \rangle) \phi(\langle \mathbf{x}', \mathbf{w} \rangle)], \\ E[\psi, \phi](\gamma) &\equiv \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)} [\psi(\langle \mathbf{x}, \mathbf{g} \rangle) \phi(\langle \mathbf{x}', \mathbf{g} \rangle)], \end{aligned}$$

where $\|\mathbf{x}\|_2 = \|\mathbf{x}'\|_2 = \sqrt{d}$ such that $\langle \mathbf{x}, \mathbf{x}' \rangle / d = \gamma$ (by an invariance argument, E_d and E do not depend on the choice of \mathbf{x} and \mathbf{x}'). Then we have

$$\lim_{d \rightarrow \infty} \sup_{\gamma \in [-1, 1]} \left| E_d[\psi, \phi](\gamma) - E[\psi, \phi](\gamma) \right| = 0, \quad (65)$$

and

$$\sup_{\gamma \in [-1, 1]} \left| E[\psi, \phi](\gamma) \right| < \infty. \quad (66)$$

Proof [Proof of Lemma 25] Let $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)$, $\mathbf{w} = \mathbf{g}/\|\mathbf{g}\|_2$ and $\tau = \|\mathbf{g}\|_2$. Then we have $\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})$, $\tau \sim \chi(d)/\sqrt{d}$ independently. We further denote

$$\bar{E}_d[\psi, \phi](\gamma) \equiv \mathbb{E}_{\mathbf{w}, \tau} [\psi(\tau \langle \mathbf{x}, \mathbf{w} \rangle) \phi(\langle \mathbf{x}', \mathbf{w} \rangle)]$$

where $\|\mathbf{x}\|_2 = \|\mathbf{x}'\|_2 = \sqrt{d}$ such that $\langle \mathbf{x}, \mathbf{x}' \rangle / d = \gamma$. Note we have

$$\begin{aligned} &\lim_{d \rightarrow \infty} \sup_{\gamma \in [-1, 1]} \left| E_d[\psi, \phi](\gamma) - \bar{E}_d[\psi, \phi](\gamma) \right| \\ &\leq \lim_{d \rightarrow \infty} \sup_{\gamma \in [-1, 1]} \left| \mathbb{E}_{\mathbf{w}, \tau} \left\{ \left[\psi(\tau \langle \mathbf{x}, \mathbf{w} \rangle) - \psi(\langle \mathbf{x}, \mathbf{w} \rangle) \right] \phi(\langle \mathbf{x}', \mathbf{w} \rangle) \right\} \right| \\ &\leq \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}, \tau} \left\{ \left[\psi(\tau \langle \mathbf{x}, \mathbf{w} \rangle) - \psi(\langle \mathbf{x}, \mathbf{w} \rangle) \right]^2 \right\}^{1/2} \mathbb{E}_{\mathbf{w}} [\phi(\langle \mathbf{x}', \mathbf{w} \rangle)^2]^{1/2} = 0, \end{aligned}$$

where the last equality is by (b) and (c) in Lemma 24. Moreover, we have

$$\begin{aligned} &\lim_{d \rightarrow \infty} \sup_{\gamma \in [-1, 1]} \left| \bar{E}_d[\psi, \phi](\gamma) - E[\psi, \phi](\gamma) \right| \\ &\leq \lim_{d \rightarrow \infty} \sup_{\gamma \in [-1, 1]} \left| \mathbb{E}_{\mathbf{w}, \tau} \left\{ \left[\phi(\tau \langle \mathbf{x}, \mathbf{w} \rangle) - \phi(\langle \mathbf{x}, \mathbf{w} \rangle) \right] \psi(\tau \langle \mathbf{x}', \mathbf{w} \rangle) \right\} \right| \\ &\leq \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}, \tau} \left\{ \left[\phi(\tau \langle \mathbf{x}, \mathbf{w} \rangle) - \phi(\langle \mathbf{x}, \mathbf{w} \rangle) \right]^2 \right\}^{1/2} \mathbb{E}_{G \sim \mathcal{N}(0, 1)} [\psi(G)^2]^{1/2} = 0, \end{aligned}$$

where the last equality is by (a) and (c) in Lemma 24. Combining the two equations above proves Eq. (65).

Finally, note that we have

$$\begin{aligned} \sup_{\gamma \in [-1, 1]} \left| E[\psi, \phi](\gamma) \right| &\leq \mathbb{E}_{\mathbf{g}} [\psi(\langle \mathbf{x}, \mathbf{g} \rangle)^2]^{1/2} \mathbb{E}_{\mathbf{g}} [\phi(\langle \mathbf{x}', \mathbf{g} \rangle)^2]^{1/2} \\ &= \mathbb{E}_{G \sim \mathcal{N}(0, 1)} [\psi(G)^2]^{1/2} \mathbb{E}_{G \sim \mathcal{N}(0, 1)} [\phi(G)^2]^{1/2} < \infty. \end{aligned}$$

This proves Eq. (66). ■

Lemma 26 Assume that $\sigma \in C^\ell(\mathbb{R})$ with derivatives satisfy $\sup_{0 \leq k \leq \ell} |\sigma^{(k)}(u)|^2 \leq c_0 \exp(c_1 u^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$. Denote

$$\begin{aligned} h_d(\gamma) &\equiv \mathbb{E}_{\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1})} [\sigma(\langle \mathbf{x}, \mathbf{w} \rangle) \sigma(\langle \mathbf{x}', \mathbf{w} \rangle)], \\ h(\gamma) &\equiv \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)} [\sigma(\langle \mathbf{x}, \mathbf{g} \rangle) \sigma(\langle \mathbf{x}', \mathbf{g} \rangle)], \end{aligned}$$

where $\|\mathbf{x}\|_2 = \|\mathbf{x}'\|_2 = \sqrt{d}$ such that $\langle \mathbf{x}, \mathbf{x}' \rangle / d = \gamma$ (by an invariance argument, h_d and h do not depend on the choice of \mathbf{x} and \mathbf{x}'). Then we have

$$\lim_{d \rightarrow \infty} \sup_{0 \leq k \leq \ell} \sup_{\gamma \in [-1, 1]} |h_d^{(k)}(\gamma) - h^{(k)}(\gamma)| = 0,$$

and

$$\sup_{0 \leq k \leq \ell} \sup_{\gamma \in [-1, 1]} |h^{(k)}(\gamma)| < \infty.$$

Proof [Proof of Lemma 26]

For $k = 0$, the result is implied by Lemma 25 by observing that $h'_d = E_d[\sigma, \sigma]$ and $h' = E[\sigma, \sigma]$.

For $k = 1$, the result is implied by Lemma 25 by the fact that $h'_d = E_d[u\sigma(u), \sigma'(u)]$ and $h' = E[u\sigma(u), \sigma'(u)]$, and there exist constants $c_0 > 0$ and $c_1 < 1$ such that $\sigma'(u), u\sigma(u) \leq c_0 e^{c_1 u^2/2}$. Indeed, for $\|\mathbf{x}\|_2 = \|\mathbf{x}'\|_2 = \sqrt{d}$ such that $\langle \mathbf{x}, \mathbf{x}' \rangle / d = \gamma$, we have (similarly for h')

$$\begin{aligned} h'_d(\gamma) &= \lim_{\delta \rightarrow 0} \delta^{-1} \left\{ \mathbb{E}_{\mathbf{w}} \left[\sigma(\langle \mathbf{x}, \mathbf{w} \rangle) \sigma(\langle (1 - \delta^2)^{1/2} \mathbf{x}' + \delta \mathbf{x}, \mathbf{w} \rangle) \right] - \mathbb{E}_{\mathbf{w}} \left[\sigma(\langle \mathbf{x}, \mathbf{w} \rangle) \sigma(\langle \mathbf{x}', \mathbf{w} \rangle) \right] \right\} \\ &= \mathbb{E}_{\mathbf{w}} \left[\sigma(\langle \mathbf{x}, \mathbf{w} \rangle) \sigma'(\langle \mathbf{x}', \mathbf{w} \rangle) \langle \mathbf{x}, \mathbf{w} \rangle \right] = E_d[u\sigma(u), \sigma'(u)](\gamma). \end{aligned}$$

By an induction argument, for any fixed k , $h_d^{(k)}$ can be identified by a fixed number of combinations of $E_d[\psi, \phi]$ with $\psi, \phi \in \Lambda_k \equiv \{u^s \sigma^{(t)}(u)\}_{0 \leq s, t \leq k}$. Further, for any fixed k , there exists $c_{0,k} > 0$ and $c_{1,k} < 1$ such that, for any $\psi \in \Lambda_k$, we have $\psi(u) \leq c_{0,k} e^{c_{1,k} u^2/2}$. Applying Lemma 25 proves the lemma. \blacksquare

Lemma 27 Assume that $\sigma \in C^k(\mathbb{R})$ with derivatives satisfy $\sup_{0 \leq s \leq k} |\sigma^{(s)}(u)|^2 \leq c_0 \exp(c_1 u^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$. For any $\|\boldsymbol{\theta}_1\|_2 = \|\boldsymbol{\theta}_2\|_2 = \sqrt{d}$ and any $\ell \geq 1$, denote

$$h_{d,S}(\langle \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \rangle / d) \equiv \mathbb{E}_{\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))} [\sigma_{d,S}(\langle \boldsymbol{\theta}_1, \mathbf{x} \rangle / \sqrt{d}) \sigma_{d,S}(\langle \boldsymbol{\theta}_2, \mathbf{x} \rangle / \sqrt{d})].$$

where $\sigma_{d,S}$ is given in Eq. (50). Then we have

$$\sup_{d \geq 1} \sup_{\gamma \in [-1, 1]} |h_{d, \geq \ell}^{(k)}(\gamma)| \leq C_{k, \ell}.$$

Proof [Proof of Lemma 27] Note we have

$$h_{d, \geq \ell}(\gamma) = h_d(\gamma) - h_{d, < \ell}(\gamma).$$

By Lemma 26, we have

$$\sup_{d \geq 1} \sup_{\gamma \in [-1, 1]} |h_d^{(k)}(\gamma)| \leq C_k.$$

Moreover, since $h_{d,<\ell}(\gamma)$ is a degree $\ell - 1$ polynomial on $[-1, 1]$ and its coefficients converge to the coefficients of $h_{<\ell}$ with $h_{<\ell}(\langle \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \rangle / d) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)}[\sigma_{d,<\ell}(\langle \boldsymbol{\theta}_1, \mathbf{x} \rangle / \sqrt{d}) \sigma_{d,<\ell}(\langle \boldsymbol{\theta}_2, \mathbf{x} \rangle / \sqrt{d})]$. Then, it is easy to see that

$$\sup_{d \geq 1} \sup_{\gamma \in [-1, 1]} \left| h_{d,<\ell}^{(k)}(\gamma) \right| \leq C_{k,\ell}.$$

This proves the lemma. \blacksquare

Lemma 28 *Assume that $\sigma \in C^\ell(\mathbb{R})$ with derivatives satisfy $\sup_{0 \leq s \leq \ell} |\sigma^{(s)}(u)|^2 \leq c_0 \exp(c_1 u^2/2)$ for some constants $c_0 > 0$ and $c_1 < 1$. Then there exists constant $C_{k,\ell}$, such that*

$$\left| h_{d,\geq\ell}^{(k)}(0) \right| \leq C_{k,\ell} \cdot d^{-(\ell-k)/2}.$$

Proof [Proof of Lemma 28]

We let $C, C_k, C_{k,\ell}$ be constants that depend on σ, k , and ℓ but independent of dimension d . The exact values of these constant can change from line to line.

We let $\tilde{\tau}_d$ be the measure of $\langle \mathbf{e}_1, \mathbf{x} \rangle$ when $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ (hence converging weakly to a standard Gaussian), and $\tilde{Q}_k^{(d)}(x) = \sqrt{B(d,k)} Q_k^{(d)}(x/\sqrt{d})$ to be the rescaled Gegenbauer polynomials, forming an orthonormal system with respect to $\tilde{\tau}_d$. In particular $\tilde{Q}_k^{(d)}$ converges to the k -th Hermite polynomial. We let $\langle \cdot, \cdot \rangle$ denote the scalar product with respect to $\tilde{\tau}_d$.

By the definition of $h_{d,\geq\ell}$, we have

$$\langle h_{d,\geq\ell}(\cdot/\sqrt{d}), \tilde{Q}_k^{(d)} \rangle = 0, \quad \forall k \leq \ell - 1. \quad (67)$$

Let $\tilde{h}_{d,\geq\ell}(x)$ be obtained from $h_{d,\geq\ell}(x)$ by removing its Taylor expansion up to term $x^{\ell-1}$, i.e., we have

$$\tilde{h}_{d,\geq\ell}(x) = h_{d,\geq\ell}(x) - \sum_{k=0}^{\ell-1} \frac{h_{d,\geq\ell}^{(k)}(0)}{k!} x^k.$$

Then Eq. (67) gives

$$\sum_{j=0}^{\ell-1} \langle \tilde{Q}_k^{(d)}, x \rangle^j \left(\frac{h_{d,\geq\ell}^{(k)}(0)}{j! d^{j/2}} \right) = -\frac{\Delta_k(d)}{d^{\ell/2}}, \quad \forall k \leq \ell - 1, \quad (68)$$

$$\Delta_k(d) \equiv d^{\ell/2} \langle \tilde{h}_{d,\geq\ell}(\cdot/\sqrt{d}), \tilde{Q}_k \rangle.$$

We claim that $\sup_{d \geq 1} |\Delta_k(d)| \leq C_{k,\ell}$. Indeed, by Rodrigues formula, there exist non-negative constants $A_{d,k}, \tilde{A}_{d,k}$ with $\sup_{d \geq 1} A_{d,k} \vee \tilde{A}_{d,k} \leq C_k$, such that

$$\begin{aligned} \Delta_k(d) &= (-1)^k d^{\ell/2} A_{d,k} \int_{-1}^1 \tilde{h}_{d,\geq\ell}(x/\sqrt{d}) \frac{d^k}{dx^k} \left(1 - \frac{x^2}{d}\right)^{\frac{d-3}{2}+k} dx \\ &= A_{d,k} d^{(\ell-k)/2} \int_{-1}^1 \tilde{h}_{d,\geq\ell}^{(k)}(x/\sqrt{d}) \left(1 - \frac{x^2}{d}\right)^{\frac{d-3}{2}+k} dx \\ &= \tilde{A}_{d,k} d^{(\ell-k)/2} \cdot \mathbb{E}_{X_d \sim \tilde{\tau}_d} \left\{ \tilde{h}_{d,\geq\ell}^{(k)}(X_d/\sqrt{d}) \left(1 - \frac{X_d^2}{d}\right)^k \right\}. \end{aligned} \quad (69)$$

By the definition of $\tilde{h}_{d,\geq\ell}$, using the Taylor expansion in the integral form, we have

$$\tilde{h}_{d,\geq\ell}(\gamma) = \int_0^\gamma h_{d,\geq\ell}^{(\ell)}(u) \frac{(\gamma-u)^{\ell-1}}{(\ell-1)!} du,$$

and hence for any $k \leq \ell - 1$, we have

$$\tilde{h}_{d,\geq\ell}^{(k)}(\gamma) = \int_0^\gamma h_{d,\geq\ell}^{(\ell)}(u) \frac{(\gamma-u)^{\ell-1-k}}{(\ell-1-k)!} du,$$

so that for any $\gamma \in [-1, 1]$, we have

$$\sup_{d \geq 1} |\tilde{h}_{d,\geq\ell}^{(k)}(\gamma)| \leq C_{k,\ell} \cdot \sup_{d \geq 1} \sup_{u \in [-1,1]} |h_{d,\geq\ell}^{(\ell)}(u)| \cdot |\gamma|^{\ell-k} \leq C_{k,\ell} \cdot |\gamma|^{\ell-k}.$$

The last inequality is by Lemma 27 (here we used the assumption that $\sigma \in C^\ell(\mathbb{R})$). Therefore, by Eq. (69), we have (note X_d converges in distribution to a standard Gaussian random variable)

$$|\Delta_k(d)| \leq C_{k,\ell} \cdot \mathbb{E}_{X_d \sim \tilde{\tau}_d} \{|X_d|^{\ell-k}\} \leq C_{k,\ell}. \quad (70)$$

To conclude, we reconsider Eq. (68). Let $\mathbf{M}(d) = (M_{k,q}(d))_{0 \leq k, q \leq \ell-1} \in \mathbb{R}^{\ell \times \ell}$ be the matrix with entries $M_{k,q}(d) \equiv \langle \tilde{Q}_k^{(d)}, x^q \rangle$, $\boldsymbol{\xi}(d) = (\xi_q(d))_{0 \leq q \leq \ell-1} \in \mathbb{R}^\ell$ the vector with entries $\xi_q(d) \equiv h_{d,\geq\ell}^{(q)}(0)/(q!d^{q/2})$, and $\boldsymbol{\Delta}(d) = (\Delta_0(d), \dots, \Delta_{\ell-1}(d))^\top \in \mathbb{R}^\ell$. We can therefore rewrite this equation as

$$\mathbf{M}(d)\boldsymbol{\xi}(d) = \boldsymbol{\Delta}(d)/d^{\ell/2}. \quad (71)$$

As $d \rightarrow \infty$, $\mathbf{M}(d)$ converges entrywise to $\mathbf{M}(\infty) = (M_{k,q}(\infty))_{0 \leq k, q \leq \ell-1}$, whereby

$$M_{k,q}(\infty) \equiv \mathbb{E}_{G \sim \mathcal{N}(0,1)} [\text{He}_k(G)G^q] / \sqrt{k!}.$$

Since $\mathbf{M}(\infty)$ is non-singular (because the Hermite polynomials are a basis), it follows that $\sigma_{\min}(\mathbf{M}(d))$ is bounded away from zero for d large enough, and therefore $\sup_{d \geq 1} \sigma_{\max}(\mathbf{M}(d)^{-1}) < \infty$. Therefore combining with Eq. (70), we get

$$\|\boldsymbol{\xi}(d)\|_2 \leq C_\ell \cdot \|\boldsymbol{\Delta}(d)\|_2 \cdot d^{-\ell/2} \leq C_\ell \cdot d^{-\ell/2}. \quad (72)$$

Therefore, for any $0 \leq k \leq \ell - 1$, we have

$$\left| h_{d,\geq\ell}^{(k)}(0) \right| \leq k!d^{k/2} |\xi_k(d)| \leq C_{k,\ell} \cdot d^{-(\ell-k)/2}. \quad (73)$$

This proves the lemma. ■

Appendix G. Hypercontractivity of general activation σ for $(\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$

Let us consider an activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ and denote for $\mathbf{x} \in \mathbb{S}^{d-1}(\sqrt{d})$ and $\mathbf{w} \in \mathbb{S}^{d-1}(1)$,

$$\bar{\sigma}(\mathbf{x}; \mathbf{w}) = \int_{\text{Cyc}_d} \sigma(\langle \mathbf{x}, g \cdot \mathbf{w} \rangle) \pi_d(dg) = \frac{1}{d} \sum_{i=0}^{d-1} \sigma(\langle \mathbf{x}, \mathbf{L}^i \mathbf{w} \rangle),$$

where $\mathbf{L} \in \mathbb{R}^{d \times d}$ is the cyclic permutation matrix that shifts the coordinates by one (hence \mathbf{L}^i shifts the coordinates by i).

Denote $\bar{\sigma}_{>\ell} = \bar{\mathbb{P}}_{>\ell} \bar{\sigma}$ the projection of $\bar{\sigma}$ orthogonal to cyclic polynomials of degree less or equal to ℓ . From the discussion in Section C.2, we have

$$\bar{\mathbb{P}}_{>\ell} \bar{\sigma}(\cdot; \mathbf{w}) = \bar{\mathbb{P}}_{>\ell} \mathcal{S}[\sigma(\langle \cdot, \mathbf{w} \rangle / \sqrt{d})] = \mathcal{S} \bar{\mathbb{P}}_{>\ell} [\sigma(\langle \cdot, \mathbf{w} \rangle / \sqrt{d})],$$

where $\mathcal{S} : L^2(\mathbb{S}^{d-1}(\sqrt{d})) \rightarrow L^2(\mathbb{S}^{d-1}(\sqrt{d}), \text{Cyc}_d)$ is the symmetrization operator defined in Section C.1 and $\bar{\mathbb{P}}_{>\ell} : L^2(\mathbb{S}^{d-1}(\sqrt{d})) \rightarrow L^2(\mathbb{S}^{d-1}(\sqrt{d}))$ is the projection orthogonal to (general) polynomials of degree less or equal to ℓ in $L^2(\mathbb{S}^{d-1}(\sqrt{d}))$ (see Section H). Hence, denoting $\sigma_{>\ell} = \bar{\mathbb{P}}_{>\ell} \sigma$, we have

$$\bar{\sigma}_{>\ell}(\mathbf{x}; \mathbf{w}) = \frac{1}{d} \sum_{i=0}^{d-1} \sigma_{>\ell}(\langle \mathbf{x}, \mathbf{L}^i \mathbf{w} \rangle). \quad (74)$$

Proposition 29 *Consider fixed integers $m \geq 1$ and $\ell \geq 4m$. Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a differentiable activation function such that $|\sigma(x)|, |\sigma'(x)| \leq c_0 \exp(c_1 x^2 / (8m))$ for some constants $c_0 > 0$ and $c_1 < 1$. Let $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ and $\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(1))$, then for any $\varepsilon > 0$,*

$$\mathbb{E}_{\mathbf{x}, \mathbf{w}} [\bar{\sigma}_{>\ell}(\mathbf{x}; \mathbf{w})^{2m}]^{1/(2m)} = d^{\varepsilon-1/2} \cdot O_d(1). \quad (75)$$

G.1. Proof of Proposition 29

The goal of this proof is to replace $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ by $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)$ and using Proposition 30 (stated in Section G.2), which is the Gaussian equivalent of Proposition 29.

Recall that $\sigma_{>\ell}$ is defined as the projection of σ orthogonal to degree ℓ polynomials with respect to the distribution $\langle \mathbf{x}, \mathbf{e} \rangle$ with $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ and $\|\mathbf{e}\|_2 = 1$ arbitrary. We can write it explicitly in terms of Gegenbauer polynomials:

$$\sigma_{>\ell}(x) = \sigma(x) - \sum_{k=0}^{\ell} \xi_{d,k} B(\mathbb{S}^{d-1}; k) Q_k(\sqrt{d}x).$$

Let us introduce $\varphi_{>\ell}$ defined as the projection of σ orthogonal to degree ℓ polynomials with respect to the Gaussian measure. It is given explicitly by

$$\varphi_{>\ell}(x) = \sigma(x) - \sum_{k=0}^{\ell} \frac{\mu_k(\sigma)}{k!} \text{He}_k(x),$$

where He_k denote the k -th Hermite polynomial (see Section H.1.3 for definitions).

Consider the symmetrized activation functions

$$\begin{aligned}\bar{\sigma}_{>\ell}(\mathbf{x}; \mathbf{w}) &= \bar{\sigma}(\mathbf{x}; \mathbf{w}) - \sum_{k=0}^{\ell} \xi_{d,k} B(\mathbb{S}^{d-1}; k) \bar{Q}_k(\mathbf{x}; \mathbf{w}), \\ \bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w}) &= \bar{\sigma}(\mathbf{g}; \mathbf{w}) - \sum_{k=0}^{\ell} \frac{\mu_k(\sigma)}{k!} \bar{\text{He}}_k(\mathbf{g}; \mathbf{w}),\end{aligned}$$

where we denoted the symmetrized polynomials

$$\begin{aligned}\bar{Q}_k(\mathbf{x}; \mathbf{w}) &= \frac{1}{d} \sum_{i=0}^{d-1} Q_k(\sqrt{d} \langle \mathbf{x}, \mathbf{L}^i \mathbf{w} \rangle), \\ \bar{\text{He}}_k(\mathbf{g}; \mathbf{w}) &= \frac{1}{d} \sum_{i=0}^{d-1} \text{He}_k(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle).\end{aligned}$$

Consider $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$, $\mathbf{g} \sim \text{N}(0, \mathbf{I}_d)$ and $\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(1))$. Because $\langle \mathbf{x}, \mathbf{e} \rangle$ converges in distribution to a normal distribution, we expect the moments of $\bar{\sigma}_{>\ell}(\mathbf{x}; \mathbf{w})$ to converge to the moments of $\bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w})$. Let us show that this convergence occurs with rate $O_d(d^{\varepsilon-1/2})$. By triangle inequality, we have

$$\mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[(\bar{\sigma}_{>\ell}(\sqrt{d} \mathbf{g} / \|\mathbf{g}\|_2; \mathbf{w}) - \bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w}))^{2m} \right]^{1/(2m)} \leq R_1 + R_2 + R_3 + R_4,$$

with

$$\begin{aligned}R_1 &= \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[(\bar{\sigma}(\sqrt{d} \mathbf{g} / \|\mathbf{g}\|_2; \mathbf{w}) - \bar{\sigma}(\mathbf{g}; \mathbf{w}))^{2m} \right]^{1/(2m)}, \\ R_2 &= \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[(A_{\leq 2}(\sqrt{d} \mathbf{g} / \|\mathbf{g}\|_2; \mathbf{w}) - B_{\leq 2}(\mathbf{g}; \mathbf{w}))^{2m} \right]^{1/(2m)}, \\ R_3 &= \mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[A_{[3:\ell]}(\mathbf{x}; \mathbf{w})^{2m} \right]^{1/(2m)}, \\ R_4 &= \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[B_{[3:\ell]}(\mathbf{g}; \mathbf{w})^{2m} \right]^{1/(2m)},\end{aligned}$$

where we denoted $[3:\ell] = \{3, \dots, \ell\}$ and for any subset $S \subset \{0, \dots, \ell\}$,

$$\begin{aligned}A_S(\mathbf{x}; \mathbf{w}) &= \sum_{k \in S} \xi_{d,k} B(\mathbb{S}^{d-1}; k) \bar{Q}_k(\mathbf{x}; \mathbf{w}), \\ B_S(\mathbf{g}; \mathbf{w}) &= \sum_{k \in S} \frac{\mu_k(\sigma)}{k!} \bar{\text{He}}_k(\mathbf{g}; \mathbf{w}).\end{aligned}$$

Step 1. Bound on R_1 .

Denote $\tau = \|\mathbf{g}\|_2 / \sqrt{d}$ and $\mathbf{x} = \sqrt{d} \mathbf{g} / \|\mathbf{g}\|_2$, such that τ and \mathbf{x} are independent, and $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$. By the mean value theorem, there exists $\tilde{\tau}$ on the line segment between 1 and τ such that

$$\bar{\sigma}(\tau \cdot \mathbf{x}; \mathbf{w}) - \bar{\sigma}(\mathbf{x}; \mathbf{w}) = (\tau - 1) \langle \nabla_{\mathbf{x}} \bar{\sigma}(\tilde{\tau} \cdot \mathbf{x}; \mathbf{w}), \mathbf{x} \rangle.$$

By Cauchy-Schwarz inequality, we get

$$\begin{aligned} R_1 &= \mathbb{E}_{\tau, \mathbf{x}, \mathbf{w}} \left[\left(\bar{\sigma}(\tau \cdot \mathbf{x}; \mathbf{w}) - \bar{\sigma}(\mathbf{x}; \mathbf{w}) \right)^{2m} \right]^{1/(2m)} \\ &= \mathbb{E}_{\tau, \mathbf{x}, \mathbf{w}} \left[(\tau - 1)^{2m} \langle \nabla_{\mathbf{x}} \bar{\sigma}(\tilde{\tau} \cdot \mathbf{x}; \mathbf{w}), \mathbf{x} \rangle^{2m} \right]^{1/(2m)} \\ &\leq \mathbb{E}_{\tau} \left[(\tau - 1)^{4m} \right]^{1/(4m)} \cdot \mathbb{E}_{\tau, \mathbf{x}, \mathbf{w}} \left[\langle \nabla_{\mathbf{x}} \bar{\sigma}(\tilde{\tau} \cdot \mathbf{x}; \mathbf{w}), \mathbf{x} \rangle^{4m} \right]^{1/(4m)} \end{aligned}$$

Let us bound the first term:

$$\mathbb{E}_{\tau} \left[(\tau - 1)^{4m} \right]^{1/(4m)} \leq \mathbb{E}_{\tau} \left[(\tau^2 - 1)^{4m} \right]^{1/(4m)} \leq C_{4m} \mathbb{E}_{\tau} \left[(\tau^2 - 1)^2 \right]^{1/2} = C_{4m} \sqrt{\frac{2}{d}}, \quad (76)$$

where we used in the first inequality that $|\tau - 1| \leq |\tau^2 - 1|$ for $\tau \geq 0$; in the second inequality that $\tau^2 - 1$ is a degree 2 polynomial in $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)$ and verifies the hypercontractivity property of Lemma 32; last equality, that $d \cdot \tau^2 = \|\mathbf{g}\|_2^2$ follows a chisquared distribution of degree d .

For the second term, we have

$$\langle \nabla_{\mathbf{x}} \bar{\sigma}(\tilde{\tau} \cdot \mathbf{x}; \mathbf{w}), \mathbf{x} \rangle = \frac{1}{d} \sum_{i=0}^{d-1} \langle \mathbf{x}, \mathbf{L}^i \mathbf{w} \rangle \sigma^{(1)}(\langle \tilde{\tau} \cdot \mathbf{x}, \mathbf{L}^i \mathbf{w} \rangle).$$

Recall that $\tilde{\tau} \cdot \mathbf{x}$ is between $\tau \cdot \mathbf{x}$ and \mathbf{x} which have marginal distributions $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)$ and $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$ respectively. Denote x_1 the first coordinate of \mathbf{x} (therefore $\tau \cdot x_1 \sim \mathcal{N}(0, 1)$). By Jensen's inequality and using that by rotation $\langle \mathbf{x}, \mathbf{L}^i \mathbf{w} \rangle$ has the same distribution as x_1 , we get

$$\begin{aligned} \mathbb{E}_{\tau, \mathbf{x}, \mathbf{w}} \left[\langle \nabla_{\mathbf{x}} \bar{\sigma}(\tilde{\tau} \cdot \mathbf{x}; \mathbf{w}), \mathbf{x} \rangle^{4m} \right] &\leq \mathbb{E}_{\tau, x_1} \left[x_1^{4m} \sigma^{(1)}(\tilde{\tau} \cdot x_1)^{4m} \right] \\ &\leq C \cdot \mathbb{E}_{G \sim \mathcal{N}(0, 1)} \left[\max(G^{4m}, 1) \exp \left\{ c_1 \max(G^2, 1) / 2 \right\} \right] \\ &= O_d(1), \end{aligned} \quad (77)$$

where we used that $c_1 < 1$.

Combining Eqs. (76) and (77) yields

$$R_1 = d^{-1/2} \cdot O_d(1). \quad (78)$$

Step 2. Bound on R_3 .

We have

$$\begin{aligned} R_3 &= \mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[A_{[3:\ell]}(\mathbf{x}; \mathbf{w})^{2m} \right]^{1/(2m)} \leq C_{2m} \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathbf{w}} \left[A_{[3:\ell]}(\mathbf{x}; \mathbf{w})^2 \right]^m \right]^{1/(2m)} \\ &\leq C_m C_{2m} \mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[A_{[3:\ell]}(\mathbf{x}; \mathbf{w})^2 \right]^{1/2}, \end{aligned}$$

where in the first inequality we used hypercontractivity of low-degree polynomials on the sphere with respect to \mathbf{w} (Lemma 36), and in the second we used hypercontractivity of low-degree symmetric functions with respect to \mathbf{x} (Lemma 6 in Mei et al. (2021)). By Lemma 11, we have

$$\mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[A_{[3:\ell]}(\mathbf{x}; \mathbf{w})^2 \right] = \sum_{k=3}^{\ell} \xi_{d,k}^2 B(\mathbb{S}^{d-1}; k) \cdot \frac{D(\mathbb{S}^{d-1}; k)}{B(\mathbb{S}^{d-1}; k)} = O_d(d^{-1}).$$

We deduce that

$$R_3 = d^{-1/2} \cdot O_d(1). \quad (79)$$

Step 3. Bound on R_4 .

Similarly to R_3 , we have

$$\begin{aligned} R_4 &= \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[B_{[3:\ell]}(\mathbf{g}; \mathbf{w})^{2m} \right]^{1/(2m)} \leq C_{2m} \mathbb{E}_{\mathbf{w}} \left[\mathbb{E}_{\mathbf{g}} \left[B_{[3:\ell]}(\mathbf{g}; \mathbf{w})^2 \right]^m \right]^{1/(2m)} \\ &\leq C_m C_{2m} \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[B_{[3:\ell]}(\mathbf{g}; \mathbf{w})^2 \right]^{1/2}, \end{aligned}$$

where in the first inequality we used hypercontractivity of low-degree polynomials with respect to \mathbf{g} (Lemma 32), and in the second we used hypercontractivity of low-degree symmetric functions with respect to \mathbf{w} .

Following the proof of Proposition 30, by setting $m = 1$ and $\bar{\varphi}_{>2}(\mathbf{g}; \mathbf{w}) = B_{[3:\ell]}(\mathbf{g}; \mathbf{w})$, we have for any $\varepsilon > 0$,

$$R_4 = d^{\varepsilon-1/2} \cdot O_d(1). \quad (80)$$

Step 4. Conclude.

The bound on R_2 is more technical and we defer it to Section G.3. By Lemma 33, we have

$$R_2 = d^{-1/2} \cdot O_d(1). \quad (81)$$

Hence combining the bounds (78), (79), (80) and (81), we obtain for any $\varepsilon > 0$,

$$\mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[\bar{\sigma}_{>\ell}(\mathbf{x}; \mathbf{w})^{2m} \right]^{1/(2m)} \leq \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w})^{2m} \right]^{1/(2m)} + O_d(d^{\varepsilon-1/2}).$$

Using Proposition 30 concludes the proof.

G.2. Proof in the Gaussian case

Recall that we defined

$$\bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w}) = \frac{1}{d} \sum_{i=0}^{d-1} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle), \quad (82)$$

where

$$\varphi_{>\ell}(x) = \sigma(x) - \sum_{k=0}^{\ell} \frac{\mu_k(\sigma)}{k!} \text{He}_k(x). \quad (83)$$

Let us now state and prove the Gaussian version of Proposition 29.

Proposition 30 *Consider fixed integers $m \geq 1$ and $\ell \geq 4m$. Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an activation function such that $|\sigma(x)| \leq c_0 \exp(c_1 x^2 / (8m))$ for some constants $c_0 > 0$ and $c_1 < 1$. Let $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)$ and $\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(1))$, then*

$$\mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w})^{2m} \right]^{1/(2m)} = d^{-1/2} \cdot O_d(1). \quad (84)$$

Proof [Proof of Proposition 30] Let us expand $\bar{\varphi}_{>\ell}$ as in Eq. (82)

$$\mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[d^{2m} \cdot \bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w})^{2m} \right] = \sum_{0 \leq i_1, \dots, i_{2m} \leq d-1} \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\prod_{k \in [2m]} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^{i_k} \mathbf{w} \rangle) \right].$$

Let us consider the event

$$\mathcal{A}_\varepsilon := \left\{ \mathbf{w} \in \mathbb{S}^{d-1}(1); \sup_{k \in [d-1]} |\langle \mathbf{w}, \mathbf{L}^k \mathbf{w} \rangle| \leq C d^{\varepsilon-1/2} \right\},$$

and for each set of indices $\mathcal{I} = \{i_1, \dots, i_{2m}\}$, consider separately the expectation over \mathcal{A}_ε and $\mathcal{A}_\varepsilon^c$:

$$\mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\prod_{i \in \mathcal{I}} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle) \right] = A + B,$$

where

$$\begin{aligned} A &:= \mathbb{E}_{\mathbf{w}} \left[\mathbf{1}_{\mathcal{A}_\varepsilon} \mathbb{E}_{\mathbf{g}} \left[\prod_{i \in \mathcal{I}} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle) \right] \right], \\ B &:= \mathbb{E}_{\mathbf{w}} \left[\mathbf{1}_{\mathcal{A}_\varepsilon^c} \mathbb{E}_{\mathbf{g}} \left[\prod_{i \in \mathcal{I}} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle) \right] \right]. \end{aligned}$$

By Cauchy-Schwarz and Jensen's inequality, we have

$$B \leq \mathbb{P}(\mathcal{A}_\varepsilon^c)^{1/2} \cdot \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\prod_{i \in \mathcal{I}} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle)^2 \right]^{1/2},$$

with

$$\begin{aligned} \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\prod_{i \in \mathcal{I}} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle)^2 \right]^{1/2} &\leq \prod_{i \in \mathcal{I}} \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle)^{4m} \right]^{1/(4m)} \\ &= \mathbb{E}_G \left[\varphi_{>\ell}(G)^{4m} \right]^{1/2} = O_d(1), \end{aligned}$$

where we used Hölder's inequality and that $\varphi_{>\ell}$ is the sum of a degree ℓ polynomial and σ with $|\sigma(x)| \leq c_0 \exp(c_1 x^2 / (8m))$, with constants $c_0 > 0$ and $c_1 < 1$. Combining these bounds and Lemma 34, we deduce there exists a constant C independent of d and \mathcal{I} such that

$$B \leq C \exp(-cd^{2\varepsilon}). \quad (85)$$

Similarly, by Hölder's inequality, we have the following first bound on A :

$$A \leq \prod_{i \in \mathcal{I}} \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle)^{2m} \right]^{1/(2m)} = \mathbb{E}_G \left[\varphi_{>\ell}(G)^{2m} \right] \leq C. \quad (86)$$

Fix $\mathbf{w} \in \mathcal{A}_\varepsilon$. Denote \mathcal{I}_0 the set of distinct indices in \mathcal{I} and $p = |\mathcal{I}_0| \leq 2m$. Denote for each $i \in \mathcal{I}_0$, r_i the multiplicity of i in \mathcal{I} , and $g_i = \langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle$. We have $\sup_{i \neq j} |\mathbb{E}[g_i g_j]| \leq \sup_{k \in [d-1]} |\langle \mathbf{w}, \mathbf{L}^k \mathbf{w} \rangle| \leq C d^{\varepsilon-1/2}$ and $\mathbb{E}[g_i^2] = 1$. Hence, if there exists $i \in \mathcal{I}$ that appears only once, we have by taking $\psi(x) = \varphi_{>\ell}(x)$ and $q = 2m \leq \ell/2$ in Lemma 31 stated below

$$\left| \mathbb{E}_{\mathbf{g}} \left[\prod_{i \in \mathcal{I}} \varphi_{>\ell}(\langle \mathbf{g}, \mathbf{L}^i \mathbf{w} \rangle) \right] \right| \leq C' d^{(2m+1)(\varepsilon-1/2)},$$

where C' is independent of \mathbf{w} . We deduce that

$$A \leq C' d^{(2m+1)(\varepsilon-1/2)}. \quad (87)$$

There are at most $m^{2m} d^m$ sets of indices \mathcal{I} with no isolated index. Hence, combining the bounds (85), (86) and (87), we get

$$\mathbb{E}_{\mathbf{g}, \mathbf{w}} [d^{2m} \cdot \bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w})^{2m}] \leq C d^{2m} \exp(-cd^{2\varepsilon}) + C m^{2m} d^m + C' d^{2m} \cdot d^{(2m+1)(\varepsilon-1/2)}.$$

Taking $\varepsilon \leq 1/(4m+2)$, we get

$$\mathbb{E}_{\mathbf{g}, \mathbf{w}} [\bar{\varphi}_{>\ell}(\mathbf{g}; \mathbf{w})^{2m}]^{1/(2m)} = d^{-1/2} \cdot O_d(1),$$

which concludes the proof. \blacksquare

The proof of Proposition 30 relies on the following key lemma:

Lemma 31 *Let $q, p, m \geq 1$ be three integers such that $p \leq 2m$. Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be a function such that $|\psi(x)| \leq c_0 \exp(c_1 x^2 / (4m))$ for some constants $c_0 > 0$ and $c_1 < 1$. Furthermore, for all $k = 0, \dots, 2q$,*

$$\mu_k(\psi) = \mathbb{E}_G[\psi(G)\text{He}_k(G)] = 0,$$

where $G \sim \mathcal{N}(0, 1)$, i.e., ψ is orthogonal to all polynomials of degree less or equal to q with respect to the standard normal distribution. Let $\mathbf{g} = (g_1, \dots, g_p) \sim \mathcal{N}(0, \Sigma)$ with $\Sigma_{11} = \dots = \Sigma_{pp} = 1$ and $\sup_{i \neq j} |\Sigma_{ij}| \leq C d^{\varepsilon-1/2}$. Let (r_1, \dots, r_p) be p integers such that $r_1 + \dots + r_p = 2m$ and there exists k such that $r_k = 1$. Then there exists $C' > 0$ depending only on c_0, c_1, C, q, m such that

$$\left| \mathbb{E}_{\mathbf{g}} \left[\prod_{k \in [p]} \psi(g_k)^{r_k} \right] \right| \leq C' d^{(q+1)(\varepsilon-1/2)}. \quad (88)$$

Proof [Proof of Lemma 31] Without loss of generality, let us assume that $r_1 = 1$. Let us rewrite the expectation with respect to $\tilde{\mathbf{g}} \sim \mathcal{N}(0, \mathbf{I}_p)$:

$$\mathbb{E}_{\mathbf{g}} \left[\prod_{k \in [p]} \psi(g_k)^{r_k} \right] = \frac{1}{\sqrt{\det(\Sigma)}} \mathbb{E}_{\tilde{\mathbf{g}}} \left[\prod_{k \in [p]} \psi(\tilde{g}_k)^{r_k} \cdot \exp \left\{ \tilde{\mathbf{g}}^\top \mathbf{M} \tilde{\mathbf{g}} / 2 \right\} \right], \quad (89)$$

where we denoted $\mathbf{M} = \mathbf{I}_p - \Sigma^{-1}$.

By Taylor expansion around 0 at order $q+1$, there exists $\zeta(\tilde{\mathbf{g}})$ between 0 and $\tilde{\mathbf{g}}^\top \mathbf{M} \tilde{\mathbf{g}} / 2$ such that

$$\exp \left\{ \tilde{\mathbf{g}}^\top \mathbf{M} \tilde{\mathbf{g}} / 2 \right\} = \sum_{s=0}^q \frac{1}{2^s s!} (\tilde{\mathbf{g}}^\top \mathbf{M} \tilde{\mathbf{g}})^s + \frac{1}{2^{q+1} (q+1)!} \exp\{\zeta(\tilde{\mathbf{g}})\} \cdot (\tilde{\mathbf{g}}^\top \mathbf{M} \tilde{\mathbf{g}})^{q+1}.$$

Notice that the terms $s = 0, \dots, q$ are polynomials of degree smaller or equal to $2q$ in $\tilde{\mathbf{g}}$. By the assumption of orthonormality of ψ to polynomials of degree less or equal to $2q$, we deduce

$$\begin{aligned} & \left| \mathbb{E}_{\mathbf{g}} \left[\prod_{k \in [p]} \psi(g_k)^{r_k} \right] \right| \\ &= \frac{1}{2^{q+1} (q+1)! \sqrt{\det(\Sigma)}} \left| \mathbb{E}_{\tilde{\mathbf{g}}} \left[\prod_{k \in [p]} \psi(\tilde{g}_k)^{r_k} \cdot \exp\{\zeta(\tilde{\mathbf{g}})\} \cdot (\tilde{\mathbf{g}}^\top \mathbf{M} \tilde{\mathbf{g}})^{q+1} \right] \right| \\ &\leq \frac{\|\mathbf{M}\|_{\text{op}}^{q+1}}{2^{q+1} (q+1)! \sqrt{\det(\Sigma)}} \mathbb{E}_{\tilde{\mathbf{g}}} \left[\prod_{k \in [p]} |\psi(\tilde{g}_k)|^{r_k} \cdot \exp\{\|\mathbf{M}\|_{\text{op}} \|\tilde{\mathbf{g}}\|_2^2\} \cdot \|\tilde{\mathbf{g}}\|_2^{2(q+1)} \right]. \end{aligned}$$

Furthermore, from the bound $|\psi(x)| \leq c_0 \exp(c_1 x^2 / (4m))$ and that $r_k \leq 2m$, we have

$$\left| \mathbb{E}_{\mathbf{g}} \left[\prod_{k \in [p]} \psi(g_k)^{r_k} \right] \right| \leq \frac{c_0^{2m} p^{2q} \|\mathbf{M}\|_{\text{op}}^{q+1}}{2^{q+1} (q+1)! \sqrt{\det(\boldsymbol{\Sigma})}} \mathbb{E}_G \left[G^{2q+2} \exp\{c_1 G^2 / 2 + \|\mathbf{M}\|_{\text{op}} G^2\} \right]^p. \quad (90)$$

From the assumptions on $\boldsymbol{\Sigma}$, we have $\|\boldsymbol{\Sigma} - \mathbf{I}_p\|_{\text{op}} \leq \|\boldsymbol{\Sigma} - \mathbf{I}_p\|_F \leq p \sup_{i \neq j} |\Sigma_{ij}| = O_d(d^{\varepsilon-1/2})$, and therefore $\|\mathbf{M}\|_{\text{op}} = O_d(d^{\varepsilon-1/2})$ and $\det(\boldsymbol{\Sigma})^{-1/2} = O_d(1)$.

From the assumption that $c_1 < 1$ and taking d sufficiently large such that $\|\mathbf{M}\|_{\text{op}} < (1 - c_1)/4$, the expectation on the right hand side of Eq. (90) is bounded by a constant. We deduce that

$$\mathbb{E}_{\mathbf{g}} \left[\prod_{k \in [p]} \psi(g_k)^{r_k} \right] = \|\mathbf{M}\|_{\text{op}}^{q+1} \cdot O_d(1) = d^{(q+1)(\varepsilon-1/2)} \cdot O_d(1),$$

which concludes the proof. \blacksquare

G.3. Technical lemmas

The first lemma is a straightforward consequence of the proof of Lemma 37 (we include a proof for completeness).

Lemma 32 *For any $\ell \in \mathbb{N}$ and $f \in L^2(\mathbb{R}^d, \gamma_d)$ to be a degree ℓ polynomial on \mathbb{R}^d , where $\gamma_d = \mathcal{N}(0, \mathbf{I}_d)$ is the isotropic Gaussian distribution. Then for any $q \geq 2$, we have*

$$\|f\|_{L^q(\mathbb{R}^d, \gamma_d)}^2 \leq (q-1)^\ell \cdot \|f\|_{L^2(\mathbb{R}^d, \gamma_d)}^2.$$

Proof [Proof of Lemma 32] Let $\boldsymbol{\varepsilon} = (\varepsilon_{i,j})_{i \in [d], j \in [D]} \sim \text{Unif}(\mathcal{Q}^{dD})$ and define for $i = 1, \dots, d$,

$$G_i = \frac{\varepsilon_{i,1} + \dots + \varepsilon_{i,D}}{\sqrt{D}}.$$

Consider f a degree ℓ polynomial on \mathbb{R}^d and define

$$\tilde{f}(\boldsymbol{\varepsilon}) = f(G_1, \dots, G_d).$$

From hypercontractivity of low degree polynomials on the hypercube (Lemma 35), we have

$$\|\tilde{f}\|_{L^q(\mathcal{Q}^{dD})}^2 \leq (q-1)^\ell \cdot \|\tilde{f}\|_{L^2(\mathcal{Q}^{dD})}^2. \quad (91)$$

Furthermore, by the multivariate central limit theorem, as $D \rightarrow \infty$ for d fixed, (G_1, \dots, G_d) converges in distribution to $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)$. By dominated convergence theorem, we have $\|\tilde{f}\|_{L^q(\mathcal{Q}^{dD})}^2 \rightarrow \|f\|_{L^q(\mathbb{R}^d, \gamma_d)}^2$ and $\|\tilde{f}\|_{L^2(\mathcal{Q}^{dD})}^2 \rightarrow \|f\|_{L^2(\mathbb{R}^d, \gamma_d)}^2$, and taking the limit in inequality (91) yields the result. \blacksquare

Lemma 33 *Follow the notations in Section G.1. We have*

$$\mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\left(A_{\leq 2}(\sqrt{d}\mathbf{g}/\|\mathbf{g}\|_2; \mathbf{w}) - B_{\leq 2}(\mathbf{g}; \mathbf{w}) \right)^{2m} \right]^{1/(2m)} = O_d(d^{-1/2}).$$

Proof [Proof of Lemma 33] Denote $\tau = \|\mathbf{g}\|_2/\sqrt{d}$ and $\mathbf{x} = \sqrt{d}\mathbf{g}/\|\mathbf{g}\|_2$. Recall that we defined

$$\begin{aligned} A_{\leq 2}(\mathbf{x}, \mathbf{w}) &= \xi_{d,0} + \xi_{d,1}B(\mathbb{S}^{d-1}; 1)\overline{Q}_1(\mathbf{x}; \mathbf{w}) + \xi_{d,2}B(\mathbb{S}^{d-1}; 2)\overline{Q}_2(\mathbf{x}; \mathbf{w}), \\ B_{\leq 2}(\tau \cdot \mathbf{x}, \mathbf{w}) &= \mu_0(\sigma) + \mu_1(\sigma)\overline{\text{He}}_1(\tau \cdot \mathbf{x}; \mathbf{w}) + \frac{\mu_2(\sigma)}{2}\overline{\text{He}}_2(\tau \cdot \mathbf{x}; \mathbf{w}), \end{aligned}$$

Let us bound the difference of each term separately.

Step 1. Bound 0th order term.

Following the same argument as in the bound of R_1 in Section G.1, we have

$$\begin{aligned} c_0 &:= |\mu_0(\sigma) - \xi_{d,0}| = \left| \mathbb{E}_{\tau, x_1} [\sigma(\tau \cdot x_1) - \sigma(x_1)] \right| \\ &\leq \mathbb{E}_{\tau} [(\tau - 1)^2]^{1/2} \mathbb{E}_{\tau, x_1} [x_1^2 \sigma'(\tilde{\tau} \cdot x_1)^2]^{1/2} = O_d(d^{-1/2}). \end{aligned} \quad (92)$$

Step 2. Bound 1st order term.

We have $\text{He}_1(x) = x$ and $B(\mathbb{S}^{d-1}; 1)^{1/2} \cdot Q_1(\sqrt{d}x) = x$. Hence,

$$\begin{aligned} c_1 &:= \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\left(\mu_1(\sigma)\overline{\text{He}}_1(\tau \cdot \mathbf{x}; \mathbf{w}) - \xi_{d,1}B(\mathbb{S}^{d-1}; 1)\overline{Q}_1(\mathbf{x}; \mathbf{w}) \right)^{2m} \right] \\ &= \mathbb{E}_{\tau} \left[\left(\tau \cdot \mu_1(\sigma) - \xi_{d,1}B(\mathbb{S}^{d-1}; 1)^{1/2} \right)^{2m} \right] \cdot \mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[B(\mathbb{S}^{d-1}; 1)^m \overline{Q}_1(\mathbf{x}; \mathbf{w})^{2m} \right]. \end{aligned}$$

Using the convergence of Gegenbauer coefficients to Hermite coefficients (see Eq. (110) in Section H.1.3), there exists a constant $C > 0$ such that

$$\mathbb{E}_{\tau} \left[\left(\tau \cdot \mu_1(\sigma) - \xi_{d,1}B(\mathbb{S}^{d-1}; 1)^{1/2} \right)^{2m} \right]^{1/(2m)} \leq C \left[\mathbb{E}_{\tau} [\tau^{2m}]^{1/(2m)} + 1 \right] = O_d(1), \quad (93)$$

where we used for example that low-degree polynomials of τ^2 are hypercontractive (see the bound on R_1 in Section G.1). From the same argument as in the bound of R_3 in Section G.1, we have

$$\mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[B(\mathbb{S}^{d-1}; 1)^m \overline{Q}_1(\mathbf{x}; \mathbf{w})^{2m} \right]^{1/(2m)} \leq C_{2m} \frac{D(\mathbb{S}^{d-1}; 1)^{1/2}}{B(\mathbb{S}^{d-1}; 1)^{1/2}} = O_d(d^{-1/2}). \quad (94)$$

Combining Eqs. (93) and (94) yields

$$c_1 = O_d(d^{-1/2}). \quad (95)$$

Step 3. Bound 2nd order term.

We have $\text{He}_2(x) = x^2 - 1$ and $B(\mathbb{S}^{d-1}; 2)^{1/2} \cdot Q_2(\sqrt{d}x) = a_{2,d} \cdot (x^2 - 1)$ with $a_{2,d} = \Theta_d(1)$. We can rewrite

$$\text{He}_2(\tau \cdot x_1) = \tau^2 B(\mathbb{S}^{d-1}; 2)^{1/2} \cdot Q_2(\sqrt{d}x)/a_{2,d} + \tau^2 - 1.$$

Hence, by triangle inequality,

$$\begin{aligned} c_1 &:= \mathbb{E}_{\mathbf{g}, \mathbf{w}} \left[\left(\mu_2(\sigma)\overline{\text{He}}_2(\tau \cdot \mathbf{x}; \mathbf{w})/2 - \xi_{d,2}B(\mathbb{S}^{d-1}; 2)\overline{Q}_2(\mathbf{x}; \mathbf{w}) \right)^{2m} \right]^{1/(2m)} \\ &\leq \mathbb{E}_{\tau} \left[\left(\tau \cdot \mu_2(\sigma)/(2a_{2,d}) - \xi_{d,2}B(\mathbb{S}^{d-1}; 2)^{1/2} \right)^{2m} \right]^{1/(2m)} \cdot \mathbb{E}_{\mathbf{x}, \mathbf{w}} \left[B(\mathbb{S}^{d-1}; 2)^m \overline{Q}_2(\mathbf{x}; \mathbf{w})^{2m} \right]^{1/(2m)} \\ &\quad + \frac{|\mu_2(\sigma)|}{2} \mathbb{E}_{\tau} [(\tau^2 - 1)^{2m}]^{1/(2m)}. \end{aligned}$$

The first term is bounded as the 1-st order term while the second term is bounded as the 0-th order term. Combining the two yields

$$c_2 = O_d(d^{-1/2}). \quad (96)$$

Step 4. Conclude.

Combining the bounds (92), (95) and (96), we get by triangle inequality

$$\mathbb{E}_{\tau, \mathbf{x}, \mathbf{w}} \left[(A_{\leq 2}(\mathbf{x}; \mathbf{w}) - B_{\leq 2}(\tau \cdot \mathbf{x}; \mathbf{w}))^{2m} \right]^{1/(2m)} \leq c_0 + c_1 + c_2 = O_d(d^{-1/2}),$$

which concludes the proof. ■

Lemma 34 *Let $\varepsilon > 0$ and $\mathbf{w} \sim \text{Unif}(\mathbb{S}^{d-1}(1))$. Then there exists $C, c > 0$ such that*

$$\mathbb{P}(\mathcal{A}_\varepsilon^c) \leq C \exp(-cd^{2\varepsilon}).$$

Proof [Proof of Lemma 34] Let us use the correspondence between uniform distribution and Gaussian distribution: $\mathbf{w} \sim \mathbf{z}/\|\mathbf{z}\|_2$, where $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_d)$. We have for $k = 1, \dots, d-1$,

$$\mathbb{P}(|\langle \mathbf{w}, \mathbf{L}^k \mathbf{w} \rangle| \geq t) = \mathbb{P}(|\langle \mathbf{z}, \mathbf{L}^k \mathbf{z} \rangle| / \|\mathbf{z}\|_2^2 \geq t) \leq \mathbb{P}(|\langle \mathbf{z}, \mathbf{L}^k \mathbf{z} \rangle| / d \geq t/2) + \mathbb{P}(\|\mathbf{z}\|_2^2 \leq d/2).$$

Note that for any $k \in [d-1]$, we have $\|\mathbf{L}^k\|_F \leq \sqrt{d}$ and $\|\mathbf{L}^k\|_{\text{op}} \leq 1$. By the Hanson-Wright inequality, for any $k \neq 0$, we have

$$\mathbb{P}\left(|\langle \mathbf{z}, \mathbf{L}^k \mathbf{z} \rangle| / d > t\right) \leq 2 \exp\{-cd \cdot \min(t^2, t)\}.$$

Furthermore, by standard concentration of the norm of Gaussian vectors, we have

$$\mathbb{P}(\|\mathbf{z}\|_2^2 \leq d/2) \leq C \exp(-cd).$$

Taking $t = Cd^{\varepsilon-1/2}$ and combining the above two bounds, we get

$$\mathbb{P}(|\langle \mathbf{w}, \mathbf{L}^k \mathbf{w} \rangle| \geq t) \leq 2 \exp(-cd^{2\varepsilon}) + C \exp(-cd).$$

Taking the union bounds over $k \in [d-1]$ concludes the proof. ■

Appendix H. Technical background of function spaces

H.1. Functions on the sphere

H.1.1. FUNCTIONAL SPACES OVER THE SPHERE

For $d \geq 3$, we let $\mathbb{S}^{d-1}(r) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = r\}$ denote the sphere with radius r in \mathbb{R}^d . We will mostly work with the sphere of radius \sqrt{d} , $\mathbb{S}^{d-1}(\sqrt{d})$ and will denote by τ_d the uniform probability measure on $\mathbb{S}^{d-1}(\sqrt{d})$. All functions in this section are assumed to be elements of $L^2(\mathbb{S}^{d-1}(\sqrt{d}), \tau_d)$, with scalar product and norm denoted as $\langle \cdot, \cdot \rangle_{L^2}$ and $\|\cdot\|_{L^2}$:

$$\langle f, g \rangle_{L^2} \equiv \int_{\mathbb{S}^{d-1}(\sqrt{d})} f(\mathbf{x}) g(\mathbf{x}) \tau_d(d\mathbf{x}). \quad (97)$$

For $\ell \in \mathbb{Z}_{\geq 0}$, let $\tilde{V}_{d,\ell}$ be the space of homogeneous harmonic polynomials of degree ℓ on \mathbb{R}^d (i.e. homogeneous polynomials $q(\mathbf{x})$ satisfying $\Delta q(\mathbf{x}) = 0$), and denote by $V_{d,\ell}$ the linear space of functions obtained by restricting the polynomials in $\tilde{V}_{d,\ell}$ to $\mathbb{S}^{d-1}(\sqrt{d})$. With these definitions, we have the following orthogonal decomposition

$$L^2(\mathbb{S}^{d-1}(\sqrt{d}), \tau_d) = \bigoplus_{\ell=0}^{\infty} V_{d,\ell}. \quad (98)$$

The dimension of each subspace is given by

$$\dim(V_{d,\ell}) = B(\mathbb{S}^{d-1}; \ell) = \frac{2\ell + d - 2}{d - 2} \binom{\ell + d - 3}{\ell}. \quad (99)$$

For each $\ell \in \mathbb{Z}_{\geq 0}$, the spherical harmonics $\{Y_{\ell,j}^{(d)}\}_{1 \leq j \leq B(\mathbb{S}^{d-1}; \ell)}$ form an orthonormal basis of $V_{d,\ell}$:

$$\langle Y_{ki}^{(d)}, Y_{sj}^{(d)} \rangle_{L^2} = \delta_{ij} \delta_{ks}.$$

Note that our convention is different from the more standard one, that defines the spherical harmonics as functions on $\mathbb{S}^{d-1}(1)$. It is immediate to pass from one convention to the other by a simple scaling. We will drop the superscript d and write $Y_{\ell,j} = Y_{\ell,j}^{(d)}$ whenever clear from the context.

We denote by \bar{P}_k the orthogonal projections to $V_{d,k}$ in $L^2(\mathbb{S}^{d-1}(\sqrt{d}), \tau_d)$. This can be written in terms of spherical harmonics as

$$\bar{P}_k f(\mathbf{x}) \equiv \sum_{l=1}^{B(\mathbb{S}^{d-1}; k)} \langle f, Y_{kl} \rangle_{L^2} Y_{kl}(\mathbf{x}). \quad (100)$$

We also define $\bar{P}_{\leq \ell} \equiv \sum_{k=0}^{\ell} \bar{P}_k$, $\bar{P}_{> \ell} \equiv \mathbf{I} - \bar{P}_{\leq \ell} = \sum_{k=\ell+1}^{\infty} \bar{P}_k$, and $\bar{P}_{< \ell} \equiv \bar{P}_{\leq \ell-1}$, $\bar{P}_{\geq \ell} \equiv \bar{P}_{> \ell-1}$.

H.1.2. GEGENBAUER POLYNOMIALS

The ℓ -th Gegenbauer polynomial $Q_{\ell}^{(d)}$ is a polynomial of degree ℓ . Consistently with our convention for spherical harmonics, we view $Q_{\ell}^{(d)}$ as a function $Q_{\ell}^{(d)} : [-d, d] \rightarrow \mathbb{R}$. The set $\{Q_{\ell}^{(d)}\}_{\ell \geq 0}$ forms

an orthogonal basis on $L^2([-d, d], \tilde{\tau}_d^1)$, where $\tilde{\tau}_d^1$ is the distribution of $\sqrt{d}\langle \mathbf{x}, \mathbf{e}_1 \rangle$ when $\mathbf{x} \sim \tau_d$, satisfying the normalization condition:

$$\langle Q_k^{(d)}(\sqrt{d}\langle \mathbf{e}_1, \cdot \rangle), Q_j^{(d)}(\sqrt{d}\langle \mathbf{e}_1, \cdot \rangle) \rangle_{L^2(\mathbb{S}^{d-1}(\sqrt{d}))} = \frac{1}{B(\mathbb{S}^{d-1}; k)} \delta_{jk}. \quad (101)$$

In particular, these polynomials are normalized so that $Q_\ell^{(d)}(d) = 1$. As above, we will omit the superscript (d) in $Q_\ell^{(d)}$ when clear from the context.

Gegenbauer polynomials are directly related to spherical harmonics as follows. Fix $\mathbf{v} \in \mathbb{S}^{d-1}(\sqrt{d})$ and consider the subspace of V_ℓ formed by all functions that are invariant under rotations in \mathbb{R}^d that keep \mathbf{v} unchanged. It is not hard to see that this subspace has dimension one, and coincides with the span of the function $Q_\ell^{(d)}(\langle \mathbf{v}, \cdot \rangle)$.

We will use the following properties of Gegenbauer polynomials

1. For $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$

$$\langle Q_j^{(d)}(\langle \mathbf{x}, \cdot \rangle), Q_k^{(d)}(\langle \mathbf{y}, \cdot \rangle) \rangle_{L^2} = \frac{1}{B(\mathbb{S}^{d-1}; k)} \delta_{jk} Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle). \quad (102)$$

2. For $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{d-1}(\sqrt{d})$

$$Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle) = \frac{1}{B(\mathbb{S}^{d-1}; k)} \sum_{i=1}^{B(\mathbb{S}^{d-1}; k)} Y_{ki}^{(d)}(\mathbf{x}) Y_{ki}^{(d)}(\mathbf{y}). \quad (103)$$

These properties imply that —up to a constant— $Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle)$ is a representation of the projector onto the subspace of degree $-k$ spherical harmonics

$$(\bar{P}_k f)(\mathbf{x}) = B(\mathbb{S}^{d-1}; k) \int_{\mathbb{S}^{d-1}(\sqrt{d})} Q_k^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle) f(\mathbf{y}) \tau_d(d\mathbf{y}). \quad (104)$$

For a function $\sigma \in L^2([- \sqrt{d}, \sqrt{d}], \tau_d^1)$ (where τ_d^1 is the distribution of $\langle \mathbf{e}_1, \mathbf{x} \rangle$ when $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d}))$), denoting its spherical harmonics coefficients $\xi_{d,k}(\sigma)$ to be

$$\xi_{d,k}(\sigma) = \int_{[- \sqrt{d}, \sqrt{d}]} \sigma(x) Q_k^{(d)}(\sqrt{d}x) \tau_d^1(dx), \quad (105)$$

then we have the following equation holds in $L^2([- \sqrt{d}, \sqrt{d}], \tau_d^1)$ sense

$$\sigma(x) = \sum_{k=0}^{\infty} \xi_{d,k}(\sigma) B(\mathbb{S}^{d-1}; k) Q_k^{(d)}(\sqrt{d}x).$$

For any rotationally invariant kernel $H_d(\mathbf{x}_1, \mathbf{x}_2) = h_d(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle / d)$, with $h_d(\sqrt{d} \cdot) \in L^2([- \sqrt{d}, \sqrt{d}], \tau_d^1)$, we can associate a self adjoint operator $\mathcal{H}_d : L^2(\mathbb{S}^{d-1}(\sqrt{d})) \rightarrow L^2(\mathbb{S}^{d-1}(\sqrt{d}))$ via

$$\mathcal{H}_d f(\mathbf{x}) \equiv \int_{\mathbb{S}^{d-1}(\sqrt{d})} h_d(\langle \mathbf{x}, \mathbf{x}_1 \rangle / d) f(\mathbf{x}_1) \tau_d(d\mathbf{x}_1). \quad (106)$$

By rotational invariance, the space V_k of homogeneous polynomials of degree k is an eigenspace of \mathcal{H}_d , and we will denote the corresponding eigenvalue by $\xi_{d,k}(h_d)$. In other words $\mathcal{H}_d f(x) \equiv \sum_{k=0}^{\infty} \xi_{d,k}(h_d) \overline{P}_k f$. The eigenvalues can be computed via

$$\xi_{d,k}(h_d) = \int_{[-\sqrt{d}, \sqrt{d}]} h_d(x/\sqrt{d}) Q_k^{(d)}(\sqrt{d}x) \tau_d^1(dx). \quad (107)$$

H.1.3. HERMITE POLYNOMIALS

The Hermite polynomials $\{\text{He}_k\}_{k \geq 0}$ form an orthogonal basis of $L^2(\mathbb{R}, \gamma)$, where $\gamma(dx) = e^{-x^2/2} dx / \sqrt{2\pi}$ is the standard Gaussian measure, and He_k has degree k . We will follow the classical normalization (here and below, expectation is with respect to $G \sim \text{N}(0, 1)$):

$$\mathbb{E}\{\text{He}_j(G) \text{He}_k(G)\} = k! \delta_{jk}. \quad (108)$$

As a consequence, for any function $g \in L^2(\mathbb{R}, \gamma)$, we have the decomposition

$$g(x) = \sum_{k=0}^{\infty} \frac{\mu_k(g)}{k!} \text{He}_k(x), \quad \mu_k(g) \equiv \mathbb{E}\{g(G) \text{He}_k(G)\}. \quad (109)$$

The Hermite polynomials can be obtained as high-dimensional limits of the Gegenbauer polynomials introduced in the previous section. Indeed, the Gegenbauer polynomials (up to a \sqrt{d} scaling in domain) are constructed by Gram-Schmidt orthogonalization of the monomials $\{x^k\}_{k \geq 0}$ with respect to the measure $\tilde{\tau}_d^1$, while Hermite polynomials are obtained by Gram-Schmidt orthogonalization with respect to γ . Since $\tilde{\tau}_d^1 \Rightarrow \gamma$ (here \Rightarrow denotes weak convergence), it is immediate to show that, for any fixed integer k ,

$$\lim_{d \rightarrow \infty} \text{Coeff}\{Q_k^{(d)}(\sqrt{d}x) B(\mathbb{S}^{d-1}; k)^{1/2}\} = \text{Coeff}\left\{\frac{1}{(k!)^{1/2}} \text{He}_k(x)\right\}. \quad (110)$$

Here and below, for P a polynomial, $\text{Coeff}\{P(x)\}$ is the vector of the coefficients of P . As a consequence, for any fixed integer k , we have

$$\mu_k(\sigma) = \lim_{d \rightarrow \infty} \xi_{d,k}(\sigma) (B(\mathbb{S}^{d-1}; k) k!)^{1/2}, \quad (111)$$

where $\mu_k(\sigma)$ and $\xi_{d,k}(\sigma)$ are given in Eq. (109) and (105).

H.2. Functions on the hypercube

Fourier analysis on the hypercube is a well studied subject [O'Donnell \(2014\)](#). The purpose of this section is to introduce some notations that make the correspondence with proofs on the sphere straightforward. For convenience, we will adopt the same notations as for their spherical case.

H.2.1. FOURIER BASIS

Denote $\mathcal{Q}^d = \{-1, +1\}^d$ the hypercube in d dimension. Let us denote τ_d to be the uniform probability measure on \mathcal{Q}^d . All the functions will be assumed to be elements of $L^2(\mathcal{Q}^d, \tau_d)$ (which

contains all the bounded functions $f : \mathcal{Q}^d \rightarrow \mathbb{R}$, with scalar product and norm denoted as $\langle \cdot, \cdot \rangle_{L^2}$ and $\| \cdot \|_{L^2}$:

$$\langle f, g \rangle_{L^2} \equiv \int_{\mathcal{Q}^d} f(\mathbf{x})g(\mathbf{x})\tau_d(d\mathbf{x}) = \frac{1}{2^n} \sum_{\mathbf{x} \in \mathcal{Q}^d} f(\mathbf{x})g(\mathbf{x}).$$

Notice that $L^2(\mathcal{Q}^d, \tau_d)$ is a 2^n dimensional linear space. By analogy with the spherical case we decompose $L^2(\mathcal{Q}^d, \tau_d)$ as a direct sum of $d + 1$ linear spaces obtained from polynomials of degree $\ell = 0, \dots, d$

$$L^2(\mathcal{Q}^d, \tau_d) = \bigoplus_{\ell=0}^d V_{d,\ell}.$$

For each $\ell \in \{0, \dots, d\}$, consider the Fourier basis $\{Y_{\ell,S}^{(d)}\}_{S \subseteq [d], |S|=\ell}$ of degree ℓ , where for a set $S \subseteq [d]$, the basis is given by

$$Y_{\ell,S}^{(d)}(\mathbf{x}) \equiv x^S \equiv \prod_{i \in S} x_i.$$

It is easy to verify that (notice that $x_i^k = x_i$ if k is odd and $x_i^k = 1$ if k is even)

$$\langle Y_{\ell,S}^{(d)}, Y_{k,S'}^{(d)} \rangle_{L^2} = \mathbb{E}[x^S \times x^{S'}] = \delta_{\ell,k} \delta_{S,S'}.$$

Hence $\{Y_{\ell,S}^{(d)}\}_{S \subseteq [d], |S|=\ell}$ form an orthonormal basis of $V_{d,\ell}$ and

$$\dim(V_{d,\ell}) = B(\mathcal{Q}^d; \ell) = \binom{d}{\ell}.$$

As above, we will omit the superscript (d) in $Y_{\ell,S}^{(d)}$ when clear from the context.

H.2.2. HYPERCUBIC GEGENBAUER

We consider the following family of polynomials $\{Q_\ell^{(d)}\}_{\ell=0,\dots,d}$ that we will call hypercubic Gegenbauer, defined as

$$Q_\ell^{(d)}(\langle \mathbf{x}, \mathbf{y} \rangle) = \frac{1}{B(\mathcal{Q}^d; \ell)} \sum_{S \subseteq [d], |S|=\ell} Y_{\ell,S}^{(d)}(\mathbf{x})Y_{\ell,S}^{(d)}(\mathbf{y}).$$

Notice that the right hand side only depends on $\langle \mathbf{x}, \mathbf{y} \rangle$ and therefore these polynomials are uniquely defined. In particular,

$$\langle Q_\ell^{(d)}(\langle \mathbf{1}, \cdot \rangle), Q_k^{(d)}(\langle \mathbf{1}, \cdot \rangle) \rangle_{L^2} = \frac{1}{B(\mathcal{Q}^d; k)} \delta_{\ell k}.$$

Hence $\{Q_\ell^{(d)}\}_{\ell=0,\dots,d}$ form an orthogonal basis of $L^2(\{-d, -d+2, \dots, d-2, d\}, \tilde{\tau}_d^1)$ where $\tilde{\tau}_d^1$ is the distribution of $\langle \mathbf{1}, \mathbf{x} \rangle$ when $\mathbf{x} \sim \tau_d$, i.e., $\tilde{\tau}_d^1 \sim 2\text{Bin}(d, 1/2) - d/2$.

We have

$$\langle Q_\ell^{(d)}(\langle \mathbf{x}, \cdot \rangle), Q_k^{(d)}(\langle \mathbf{y}, \cdot \rangle) \rangle_{L^2} = \frac{1}{B(\mathcal{Q}^d; k)} Q_k(\langle \mathbf{x}, \mathbf{y} \rangle) \delta_{\ell k}.$$

For a function $\sigma(\cdot/\sqrt{d}) \in L^2(\{-d, -d+2, \dots, d-2, d\}, \tilde{\tau}_d^1)$, denote its hypercubic Gegenbauer coefficients $\xi_{d,k}(\sigma)$ to be

$$\xi_{d,k}(\sigma) = \int_{\{-d, -d+2, \dots, d-2, d\}} \sigma(x/\sqrt{d}) Q_k^{(d)}(x) \tilde{\tau}_d^1(dx).$$

Notice that by weak convergence of $\langle \mathbf{1}, \mathbf{x} \rangle / \sqrt{d}$ to the normal distribution, we have also convergence of the (rescaled) hypercubic Gegenbauer polynomials to the Hermite polynomials, i.e., for any fixed k , we have

$$\lim_{d \rightarrow \infty} \text{Coeff}\{Q_k^{(d)}(\sqrt{d}x) B(\mathcal{Q}^d; k)^{1/2}\} = \text{Coeff}\left\{\frac{1}{(k!)^{1/2}} \text{He}_k(x)\right\}. \quad (112)$$

H.3. Hypercontractivity of Gaussian measure and uniform distributions on the sphere and the hypercube

By Holder's inequality, we have $\|f\|_{L^p} \leq \|f\|_{L^q}$ for any f and any $p \leq q$. The reverse inequality does not hold in general, even up to a constant. However, for some measures, the reverse inequality will hold for some sufficiently nice functions. These measures satisfy the celebrated hypercontractivity properties [Gross \(1975\)](#); [Bonami \(1970\)](#); [Beckner \(1975, 1992\)](#).

Lemma 35 (Hypercube hypercontractivity [Beckner \(1975\)](#)) *For any $\ell = \{0, \dots, d\}$ and $f_d \in L^2(\mathcal{Q}^d)$ to be a degree ℓ polynomial, then for any integer $q \geq 2$, we have*

$$\|f_d\|_{L^q(\mathcal{Q}^d)}^2 \leq (q-1)^\ell \cdot \|f_d\|_{L^2(\mathcal{Q}^d)}^2.$$

Lemma 36 (Spherical hypercontractivity [Beckner \(1992\)](#)) *For any $\ell \in \mathbb{N}$ and $f_d \in L^2(\mathbb{S}^{d-1})$ to be a degree ℓ polynomial, for any $q \geq 2$, we have*

$$\|f_d\|_{L^q(\mathbb{S}^{d-1})}^2 \leq (q-1)^\ell \cdot \|f_d\|_{L^2(\mathbb{S}^{d-1})}^2.$$

Lemma 37 (Gaussian hypercontractivity) *For any $\ell \in \mathbb{N}$ and $f \in L^2(\mathbb{R}, \gamma)$ to be a degree ℓ polynomial on \mathbb{R} , where γ is the standard Gaussian distribution. Then for any $q \geq 2$, we have*

$$\|f\|_{L^q(\mathbb{R}, \gamma)}^2 \leq (q-1)^\ell \cdot \|f\|_{L^2(\mathbb{R}, \gamma)}^2.$$

The Gaussian hypercontractivity is a direct consequence of hypercube hypercontractivity.