
Separation Results between Fixed-Kernel and Feature-Learning Probability Metrics

Carles Domingo-Enrich

Courant Institute of Mathematical Sciences (NYU)
cd2754@nyu.edu

Youssef Mroueh

IBM Research AI
mroueh@us.ibm.com

Abstract

Several works in implicit and explicit generative modeling empirically observed that feature-learning discriminators outperform fixed-kernel discriminators in terms of the sample quality of the models. We provide separation results between probability metrics with fixed-kernel and feature-learning discriminators using the function classes \mathcal{F}_2 and \mathcal{F}_1 respectively, which were developed to study overparametrized two-layer neural networks. In particular, we construct pairs of distributions over hyper-spheres that can not be discriminated by fixed kernel (\mathcal{F}_2) integral probability metric (IPM) and Stein discrepancy (SD) in high dimensions, but that can be discriminated by their feature learning (\mathcal{F}_1) counterparts. To further study the separation we provide links between the \mathcal{F}_1 and \mathcal{F}_2 IPMs with sliced Wasserstein distances. Our work suggests that fixed-kernel discriminators perform worse than their feature learning counterparts because their corresponding metrics are weaker.

1 Introduction

The field of generative modeling, whose aim is to generate artificial samples of some target distribution given some true samples of it, is broadly divided into two types of models (Mohamed and Lakshminarayanan, 2017): explicit generative models, which involve learning an estimate of the log-density of the target distribution which is then sampled (e.g. energy-based models), and implicit generative models, where samples are generated directly by transforming some latent variable (e.g. generative adversarial networks (Goodfellow et al., 2014), normalizing flows (Rezende and Mohamed, 2015)).

Several works have observed experimentally that both in implicit and in explicit generative models, using ‘adaptive’ or ‘feature-learning’ function classes as discriminators yields better generative performance than ‘lazy’ or ‘kernel’ function classes. Within implicit models, Li et al. (2017) show that generative moment matching networks (GMMN) generate significantly better samples for the CIFAR-10 and MNIST datasets when using maximum mean discrepancy (MMD) with learned instead of fixed features. For a related method and in a similar spirit, Santos et al. (2019) show that for image generation, fixed-feature discriminators are only successful when we take an amount of features exponential in the intrinsic dimension of the dataset. Genevay et al. (2018) study implicit generative models with the Sinkhorn divergence as discriminator, and they also show that other than for simple datasets like MNIST, learning the Wasserstein cost is crucial for good performance.

As to explicit models, Grathwohl et al. (2020) train energy-based models with a Stein discrepancy based on neural networks and show improved performance with respect to kernel classes. Chang et al. (2020) show that Stein variational gradient descent (SVGD) fails in high dimensions, and that learning the kernel helps. Given the abundant experimental evidence, the aim of this work is to provide some theoretical results that showcase the advantages of feature-learning over kernel discriminators. For

the sake of simplicity, we compare the discriminative behavior of two function classes \mathcal{F}_1 and \mathcal{F}_2 , arising from infinite-width two-layer neural networks with different norms penalties on its weights (Bach, 2017). \mathcal{F}_1 displays an adaptive behavior, while \mathcal{F}_2 is an RKHS which consequently has a lazy behavior. Namely, our main contributions are:

- (i) We construct a sequence of pairs of distributions over hyperspheres \mathbb{S}^{d-1} of increasing dimensions, such that the \mathcal{F}_2 integral probability metric (IPM) between the pair decreases exponentially in the dimension, while the \mathcal{F}_1 IPM remains high.
- (ii) We construct a sequence of pairs of distributions over \mathbb{S}^{d-1} such that the \mathcal{F}_2 Stein discrepancy (SD) between the pair decreases exponentially in the dimension, while the \mathcal{F}_1 SD remains high.
- (iii) We prove polynomial upper and lower bounds between the \mathcal{F}_1 IPM and the max-sliced Wasserstein distance for distributions over Euclidean balls. For a class $\tilde{\mathcal{F}}_2$ related to \mathcal{F}_2 , we prove similar upper and lower bounds between the $\tilde{\mathcal{F}}_2$ IPM and the sliced Wasserstein distance for distributions over Euclidean balls.

Our findings reinforce the idea that generative models with kernel discriminators have worse performance because their corresponding metrics are weaker and thus unable to distinguish between different distributions, especially in high dimensions.

2 Related work

A recent line of research has studied the question of how neural networks compare to kernel methods, with a focus on supervised learning problems. Bach (2017) shows the approximation benefits of the \mathcal{F}_1 space for adapting to low-dimensional structures compared to the (kernel) space \mathcal{F}_2 ; an analysis that we leverage. The function space \mathcal{F}_1 was also studied by Ongie et al. (2019); Savarese et al. (2019); Williams et al. (2019), which focus on the ReLU activation function. More recently, several works showed that wide neural networks trained with gradient methods may behave like kernel methods in certain regimes (see, e.g., Jacot et al., 2018). Examples of works that compare ‘active/feature-learning’ and ‘kernel/lazy’ regimes for supervised learning include Chizat and Bach (2020); Ghorbani et al. (2019); Wei et al. (2020); Woodworth et al. (2020), and Domingo-Enrich et al. (2021) for energy-based models. We are not aware of any works that study how feature-learning function classes and kernel classes differ as discriminators for IPMs or Stein discrepancies.

It turns out that the \mathcal{F}_2 integral probability metric that we study is in fact MMD for certain kernels that often admit a closed form (Roux and Bengio, 2007; Cho and Saul, 2009; Bach, 2017). MMDs are probability metrics that were first introduced by Gretton et al. (2007, 2012) for kernel two-sample tests, and that have enjoyed ample success with the advent of deep-learning-based generative modeling as discriminating metrics: Li et al. (2015) and Dziugaite et al. (2015) introduced GMMN, which differ from GANs in that the discriminator network is replaced by a fixed-kernel MMD. Li et al. (2017) introduces an improvement on GMMN by using the MMD loss on learned features. From this viewpoint, our separation results in Sec. 5 can be interpreted as instances in which a given fixed-kernel MMD provably has less discriminative power than adaptive discriminators.

Other related work includes the Stein discrepancy literature. Stein’s method (Stein, 1972) dates to the 1970s. Gorham and Mackey (2015) introduced a computational approach to compute the Stein discrepancy in order to assess sample quality. Later, Chwialkowski et al. (2016), Liu et al. (2016) and Gorham and Mackey (2017) introduced the more practical kernelized Stein discrepancy (KSD) for goodness-of-fit tests. Liu and Wang (2016) introduced SVGD, the first method to use the KSD to obtain samples from a distribution. Barp et al. (2019) employed KSD to train parametric generative models, and Grathwohl et al. (2020) trained models replacing KSD by a neural-network-based SD.

Our work also touches on sliced and spiked Wasserstein distances. Sliced Wasserstein distances were introduced first by Kolouri et al. (2016); Kolouri et al. (2019). Spiked Wasserstein distances, which are a generalization, were studied later by Paty and Cuturi (2019), and they also appear in Niles-Weed and Rigollet (2019) as a good statistical estimator. Nadjahi et al. (2020) and Lin et al. (2021) have studied statistical properties of sliced and spiked Wasserstein distances, respectively.

3 Framework

3.1 Notation

If V is a normed vector space, we use $\mathcal{B}_V(\beta)$ to denote the closed ball of V of radius β , and $\mathcal{B}_V := \mathcal{B}_V(1)$ for the unit ball. If K denotes a subset of the Euclidean space, $\mathcal{P}(K)$ is the set of Borel probability measures, $\mathcal{M}(K)$ is the space of finite signed Radon measures and $\mathcal{M}^+(K)$ is the set of finite positive Radon measures. If γ is a signed Radon measure over K , then $\|\gamma\|_{\text{TV}}$ is the total variation (TV) norm of γ . Throughout the paper, and unless otherwise specified, $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ denotes a generic non-linear activation function. We use $(\cdot)_+ : \mathbb{R} \rightarrow \mathbb{R}$ to denote the ReLU activation, defined as $(x)_+ = \max\{x, 0\}$. τ denotes the uniform probability measure over a space that depends on the context. We use \mathbb{S}^d for the d -dimensional hypersphere and \log for the natural logarithm.

3.2 Overparametrized two-layer neural network spaces

Feature-learning regime. We define \mathcal{F}_1 as the Banach space of functions $f : K \rightarrow \mathbb{R}$ such that for some $\gamma \in \mathcal{M}(\mathbb{S}^d)$, for all $x \in K$ we have $f(x) = \int_{\mathbb{S}^d} \sigma(\langle \theta, x \rangle) d\gamma(\theta)$ for some signed Radon measure γ (Bach, 2017). The norm of \mathcal{F}_1 is defined as $\|f\|_{\mathcal{F}_1} = \inf \left\{ \|\gamma\|_{\text{TV}} \mid f(\cdot) = \int_{\mathbb{S}^d} \sigma(\langle \theta, \cdot \rangle) d\gamma(\theta) \right\}$.

Kernel regime. We define \mathcal{F}_2 as the (reproducing kernel) Hilbert space of functions $f : K \rightarrow \mathbb{R}$ such that for some absolutely continuous $\rho \in \mathcal{M}(\mathbb{S}^d)$ with $\frac{d\rho}{d\tau} \in \mathcal{L}^2(\mathbb{S}^d)$ (where τ is the uniform probability measure over \mathbb{S}^d), we have that for all $x \in K$, $f(x) = \int_{\mathbb{S}^d} \sigma(\langle \theta, x \rangle) d\rho(\theta)$. The norm of \mathcal{F}_2 is defined as $\|f\|_{\mathcal{F}_2}^2 = \inf \left\{ \int_{\mathbb{S}^d} h(\theta)^2 d\tau(\theta) \mid f(\cdot) = \int_{\mathbb{S}^d} \sigma(\langle \theta, \cdot \rangle) h(\theta) d\tau(\theta) \right\}$. As an RKHS, the kernel of \mathcal{F}_2 is

$$k(x, y) = \int_{\mathbb{S}^d} \sigma(\langle x, \theta \rangle) \sigma(\langle y, \theta \rangle) d\tau(\theta). \quad (1)$$

Such kernels admit closed form expressions for different choices of activation functions, among which ReLU (Roux and Bengio, 2007; Cho and Saul, 2009; Bach, 2017).

Remark that since $\int |h(\theta)| d\tau(\theta) \leq (\int h(\theta)^2 d\tau(\theta))^{1/2}$ by the Cauchy-Schwarz inequality, we have $\mathcal{F}_2 \subset \mathcal{F}_1$. In particular, when σ is the ReLU unit, Bach (2017) shows that two-layer networks with a single neuron belong to \mathcal{F}_1 but not to \mathcal{F}_2 , and their L^2 approximations in \mathcal{F}_2 have exponentially high norm in the dimension. Informally, one should understand \mathcal{F}_1 as the space of two-layer networks where both the input layer and output layer parameters are trained, in the limit of an infinite number of neurons. On the other hand, \mathcal{F}_2 is the space of infinite-width two-layer networks where only the output layer parameters are trained while the input layer parameters are sampled uniformly on the sphere and kept fixed.

4 \mathcal{F}_1 and \mathcal{F}_2 Integral Probability Metrics

Let K be a subset of \mathbb{R}^{d+1} . Integral probability metrics (IPM) are pseudometrics on $\mathcal{P}(K)$ of the form

$$d_{\mathcal{F}}(\mu, \nu) = \sup_{f \in \mathcal{F}} \mathbb{E}_{x \sim \mu} f(x) - \mathbb{E}_{x \sim \nu} f(x),$$

where \mathcal{F} is a class of functions from K to \mathbb{R} .

\mathcal{F}_2 IPM or \mathcal{F}_2 MMD. One possible choice for \mathcal{F} is the unit ball $\mathcal{B}_{\mathcal{F}_2}$ of \mathcal{F}_2 . Since \mathcal{F}_2 is an RKHS with kernel k , the corresponding IPM is in fact a maximum mean discrepancy (MMD) (Gretton et al., 2007) and it can be shown (Lemma 1 in App. A) to take the form

$$d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) = \int_{\mathbb{S}^d} \left(\int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right)^2 d\tau(\theta).$$

Notice that for any feature $\theta \in \mathbb{S}^d$, $\mathbb{E}_{x \sim p} \sigma(\langle x, \theta \rangle)$ can be seen as a generalized moment of p . $d_{\mathcal{B}_{\mathcal{F}_2}}$ can be seen as the L^2 distance between generalized moments of μ and ν as functions of $\theta \in \mathbb{S}^d$.

\mathcal{F}_1 IPM. An alternative choice for \mathcal{F} is the unit ball $\mathcal{B}_{\mathcal{F}_1}$ of \mathcal{F}_1 . The IPM for the unit ball of \mathcal{F}_1 can be developed (Lemma 2 in App. A) into

$$d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) = \sup_{\theta \in \mathbb{S}^d} \left| \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right|. \quad (2)$$

Observe that $d_{\mathcal{B}_{\mathcal{F}_1}}$ is the L^∞ distance between generalized moments of μ and μ as functions of $\theta \in \mathbb{S}^d$. That is, instead of averaging over features, all the weight is allocated to the feature at which the generalized moment difference is larger.

We will provide separate results for two interesting choices for K : (i) for $K = \mathbb{S}^d$, we obtain neural network discriminators without bias term which are amenable to analysis using the theory of spherical harmonics; and (ii) for $K = \mathbb{R}^d \times \{1\}$, we obtain neural networks discriminators with a bias term which is encoded by the last component (notice that probability measures over \mathbb{R}^d can be mapped trivially to probability measures over $\mathbb{R}^d \times \{1\}$). We will write $\mathcal{F}_1(K)$ or $\mathcal{F}_2(K)$ for specific K when it is not clear by the context.

A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is α -positive homogeneous function if for all $r \geq 0, x \in \mathbb{R}$, $f(rx) = r^\alpha f(x)$. One-dimensional α -positive homogeneous functions can be written in a general form as

$$f(x) = a(x)_+^\alpha + b(-x)_+^\alpha. \quad (3)$$

where $a, b \in \mathbb{R}$ are arbitrary. When the activation function σ is α -positive homogeneous, Theorem 1 shows that the \mathcal{F}_1 and \mathcal{F}_2 IPMs are distances when $K = \mathbb{R}^d \times \{1\}$ if a, b fulfill a certain condition which is satisfied by the ReLU activation, but they are *not* distances when $K = \mathbb{S}^d$. See Theorem 6 and Theorem 7 in App. B for the proof.

Theorem 1. *For any non-negative integer α , let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an α -positive homogeneous activation function of the form (3). If $(-1)^\alpha a - b \neq 0$ and $K = \mathbb{R}^d \times \{1\}$, both the \mathcal{F}_1 and \mathcal{F}_2 IPMs are distances on $\mathcal{P}(K)$. If $K = \mathbb{S}^d$, both the \mathcal{F}_1 and \mathcal{F}_2 IPMs are not distances on $\mathcal{P}(K)$, as there exist pairs of different measures for which the IPMs evaluate to zero.*

In other words, Theorem 1 states that certain fixed-kernel and feature-learning infinite neural networks with RELU or leaky RELU non-linearity, yield distances when we include a bias term, but not when the inputs lie in a hypersphere. This result sheds light on when the “neural net distance” introduced by Arora et al. (2017) is indeed a distance.

5 Separation between the \mathcal{F}_1 and \mathcal{F}_2 IPMs

In this section for each dimension $d \geq 2$, we construct a pair of probability measures μ_d, ν_d over $\mathcal{P}(\mathbb{S}^{d-1})$ such that the \mathcal{F}_1 IPM between μ_d and ν_d stays constant along the dimension, while the \mathcal{F}_2 IPM decreases exponentially.

Legendre harmonics and Legendre polynomials. Let $e_d \in \mathbb{R}^d$ be the d -th vector of the canonical basis. There is a unique homogeneous harmonic polynomial $L_{k,d}$ of degree k over \mathbb{R}^d such that: (i) $L_{k,d}(Ax) = L_{k,d}(x)$ for all orthogonal matrices that leave e_d invariant, and (ii) $L_{k,d}(e_d) = 1$. This polynomial receives the name of *Legendre harmonic*, and its restriction to \mathbb{S}^{d-1} is indeed a spherical harmonic of order k . If we express an arbitrary $\xi_{(d)} \in \mathbb{S}^{d-1}$ as $\xi_{(d)} = te_d + (1 - t^2)^{1/2}\xi_{(d-1)}$, where $\xi_{(d-1)} \perp e_d$, we can define the *Legendre polynomial* of degree k in dimension d as $P_{k,d}(t) := L_{k,d}(\xi_{(d)})$ by the invariance of $L_{k,d}$ (it is not straightforward that $P_{k,d}(t)$ is a polynomial on t , see Sec. 2.1.2 of Atkinson and Han (2012)). Conversely, $L_{k,d}(x) = P_{k,d}(\langle e_d, x \rangle)$ for any $x \in \mathbb{S}^{d-1}$, and by homogeneity, $L_{k,d}(x) = \|x\|^k P_{k,d}(\langle e_d, x \rangle / \|x\|)$ for any $x \in \mathbb{R}^d$. Legendre polynomials can also be characterized as the orthogonal sequence of polynomials on $[-1, 1]$ such that $P_{k,d}(1) = 1$ and $\int_{-1}^1 P_{k,d}(t)P_{l,d}(t)(1 - t^2)^{\frac{d-3}{2}} dt = 0$, for $k \neq l$.

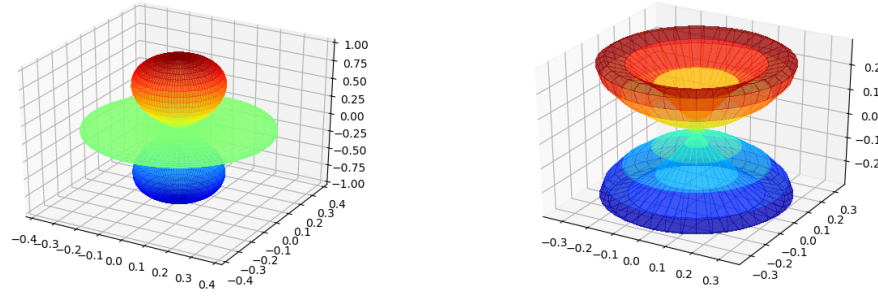


Figure 1: 3D polar plot representing the densities of the measures μ_d (left) and ν_d (right), for the choices $d = 3, k = 4$. In each direction, the distance from the origin to the surface is proportional to the density of the measure.

The pair μ_d and ν_d . We define μ_d and ν_d as the probability measures over \mathbb{S}^{d-1} with densities

$$\begin{aligned} \frac{d\mu_d}{d\lambda} &= \begin{cases} \frac{\gamma_{k,d} L_{k,d}(x)}{|\mathbb{S}^{d-1}|} & \text{if } L_{k,d}(x) > 0 \\ 0 & \text{if } L_{k,d}(x) \leq 0 \end{cases}, \\ \frac{d\nu_d}{d\lambda} &= \begin{cases} 0 & \text{if } L_{k,d}(x) > 0 \\ \frac{-\gamma_{k,d} L_{k,d}(x)}{|\mathbb{S}^{d-1}|} & \text{if } L_{k,d}(x) \leq 0 \end{cases}. \end{aligned} \quad (4)$$

for some $k \geq 2$ and some $\gamma_{k,d} \geq 0$, where λ is the Hausdorff measure over \mathbb{S}^{d-1} . Namely,

Proposition 1. *If we choose $\gamma_{k,d} = 2 \left(\int_{\mathbb{S}^{d-1}} |L_{k,d}(x)| d\tau(x) \right)^{-1}$, then μ_d and ν_d are probability measures.*

Figure 1 shows a representation of the measures μ_d, ν_d for $d = 3, k = 4$, where one can see that they allocate mass in different regions of the sphere. We are now ready to state our separation result, which is proved in [App. D](#).

Theorem 2. *Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an activation function that is bounded in $[-1, 1]$. For any $d \geq 2$ and $k \geq 1$, if we set $\gamma_{k,d}$ as in [Proposition 1](#) we have that*

$$d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) = \frac{2 \left| \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right|}{\int_{-1}^1 |P_{k,d}(t)| (1-t^2)^{\frac{d-3}{2}} dt}, \quad (5)$$

and

$$\frac{d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)}{d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)} = \sqrt{N_{k,d}} = \sqrt{\frac{(2k+d-2)(k+d-3)!}{k!(d-2)!}}, \quad (6)$$

where $N_{k,d}$ is the dimension of the space of spherical harmonics of order k over \mathbb{S}^{d-1} . That is,

$$\log \left(\frac{d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)}{d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)} \right) = \frac{1}{2} \left(k \log \left(\frac{k+d-3}{k} \right) + (d-2) \log \left(\frac{k+d-3}{d-2} \right) \right) + O(\log(k+d)). \quad (7)$$

From (7) we see that choosing the parameter k of the same order as d , $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)$ is exponentially larger than $d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)$ in the dimension d . Equation (5) holds regardless of the choice of the activation function σ , and decreases very slowly in d for the ReLu activation, as shown in [Figure 2](#). This result suggests that in high dimensions there exist high frequency densities that can be distinguished by feature-learning IPM discriminators but not by their fixed-kernel counterpart, and that may explain the differences in generative modeling performance for GMMN and Sinkhorn divergence ([Sec. 1](#)).

The key idea for the proof of [Theorem 2](#) is that the Legendre harmonics $L_{k,d}$ have constant L^∞ norm equal to 1 (see equation (26) in [App. C](#)), but their L^2 norm decreases as $1/N_{k,d}$ (see equation (32) in [App. D](#)). The proof boils down to relating $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)$ to the L^∞ norm of $L_{k,d}$, and $d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)$ to its L^2 norm.

6 Separation between \mathcal{F}_1 and \mathcal{F}_2 Stein discrepancies

The arguments to derive the separation result in [Sec. 5](#) can be leveraged to obtain a similar separation for the Stein discrepancy, which helps explain why for Stein discrepancy energy-based models (EBMs) and SVGD feature learning yields improved performance.

6.1 Stein operator and Stein discrepancy

As shown by [Domingo-Enrich et al. \(2021\)](#), for a probability measure ν on the sphere \mathbb{S}^{d-1} with a continuous and almost everywhere differentiable density $\frac{d\nu}{d\tau}$, the *Stein operator* $\mathcal{A}_\nu : \mathbb{S}^{d-1} \rightarrow \mathbb{R}^{d \times d}$ is defined as

$$(\mathcal{A}_\nu h)(x) = \left(\nabla \log \left(\frac{d\nu}{d\tau}(x) \right) - (d-1)x \right) h(x)^\top + \nabla h(x), \quad (8)$$

for any $h : \mathbb{S}^{d-1} \rightarrow \mathbb{R}^d$ that is continuous and almost everywhere differentiable, where ∇ denotes the Riemannian gradient. That is, for any $h : \mathbb{S}^{d-1} \rightarrow \mathbb{R}^d$ that is continuous and almost everywhere differentiable, the Stein identity holds: $\mathbb{E}_\nu[(\mathcal{A}_\nu h)(x)] = 0$.

If \mathcal{H} is a class of functions from \mathbb{S}^{d-1} to \mathbb{R}^d , the Stein discrepancy ([Gorham and Mackey, 2015](#); [Liu et al., 2016](#)) for \mathcal{H} is a non-symmetric functional defined on pairs of probability measures over K as

$$\text{SD}_{\mathcal{H}}(\nu_1, \nu_2) = \sup_{h \in \mathcal{H}} \mathbb{E}_{\nu_1}[\text{Tr}(\mathcal{A}_{\nu_2} h(x))]. \quad (9)$$

When $\mathcal{H} = \mathcal{B}_{\mathcal{H}_0^d} = \{(h_i)_{i=1}^d \in \mathcal{H}_0^d \mid \sum_{i=1}^d \|h_i\|_{\mathcal{H}_0}^2 \leq 1\}$ for some reproducing kernel Hilbert space (RKHS) \mathcal{H}_0 with kernel k with continuous second order partial derivatives, there exists a closed form for the problem (9) and the corresponding object is known as kernelized Stein discrepancy (KSD) [Liu et al. \(2016\)](#); [Gorham and Mackey \(2017\)](#). When the domain is \mathbb{S}^{d-1} , the KSD takes the following form (Lemma 5, [Domingo-Enrich et al. \(2021\)](#)):

$$\text{KSD}(\nu_1, \nu_2) = \text{SD}_{\mathcal{B}_{\mathcal{H}_0^d}}^2(\nu_1, \nu_2) = \mathbb{E}_{x, x' \sim \nu_1} [u_{\nu_2}(x, x')],$$

where we have $u_\nu(x, x') = (s_\nu(x) - (d-1)x)^\top (s_\nu(x') - (d-1)x')k(x, x') + (s_\nu(x) - (d-1)x)^\top \nabla_{x'} k(x, x') + (s_\nu(x') - (d-1)x')^\top \nabla_x k(x, x') + \text{Tr}(\nabla_{x, x'} k(x, x'))$, and we use $\tilde{u}_\nu(x, x')$ to denote the sum of the first three terms (remark that the fourth term does not depend on ν). Here we have used the notation $s_\nu(x) = \nabla \log(\frac{d\nu}{d\tau}(x))$, which is known as the score function.

6.2 Separation result

We show a separation result between the two cases:

- \mathcal{F}_1 Stein discrepancy: $\mathcal{H} = \mathcal{B}_{\mathcal{F}_1^d} = \{(h_i)_{i=1}^d \in \mathcal{F}_1^d \mid \sum_{i=1}^d \|h_i\|_{\mathcal{F}_1}^2 \leq 1\}$. This discriminator set initially appeared as a particular configuration in the framework of [Huggins and Mackey \(2018\)](#), and its statistical properties for energy based model training were later studied by [Domingo-Enrich et al. \(2021\)](#).
- \mathcal{F}_2 Stein discrepancy: $\mathcal{H} = \mathcal{B}_{\mathcal{F}_2^d} = \{(h_i)_{i=1}^d \in \mathcal{F}_2^d \mid \sum_{i=1}^d \|h_i\|_{\mathcal{F}_2}^2 \leq 1\}$. Since \mathcal{F}_2 is an RKHS, this corresponds to a KSD with the kernel k . However, particular care must be taken in checking that the kernel k has continuous second order partial derivatives, which might not always be the case (i.e. with $\alpha = 1$).

The pair μ_d and ν_d . For $d \geq 2$, we set μ_d to be the uniform Borel probability measure over \mathbb{S}^{d-1} . We define ν_d as the probability measure over \mathbb{S}^{d-1} with density

$$\frac{d\nu_d}{d\lambda}(x) = \frac{\exp(\gamma_{k,d} L_{k,d}(x))}{\int_{\mathbb{S}^{d-1}} \exp(\gamma_{k,d} L_{k,d}(x)) d\lambda(x)} \quad (10)$$

for some $\gamma_{k,d} \in \mathbb{R}$ that we will specify later on and some $k \geq 2$.

Theorem 3. *Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an α -positive homogeneous activation function of the form (3) such that $a + (-1)^{k+1}b \neq 0$. For all $k \geq 1$, $d \geq 2$ we can choose $\gamma_{k,d} \in [-1, 1]$ such that $SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d) = 1$ and*

$$\frac{SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d)}{SD_{\mathcal{B}_{\mathcal{F}_2^d}}(\mu_d, \nu_d)} \geq \frac{\frac{k(d+k-3)}{\alpha+1}}{\sqrt{\frac{2}{N_{k,d}} \left(k(k+d-2) \left(\frac{d+\alpha-2}{\alpha+1} \right)^2 + \left(\frac{k(d+k-3)}{\alpha+1} \right)^2 \right)}} \quad (11)$$

That is,

$$\log \left(\frac{SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d)}{SD_{\mathcal{B}_{\mathcal{F}_2^d}}(\mu_d, \nu_d)} \right) \geq \frac{1}{2} \left(k \log \left(\frac{k+d-3}{k} \right) + (d-2) \log \left(\frac{k+d-3}{d-2} \right) \right) + O(\log(k+d)) \quad (12)$$

As in Theorem 2, from (12) we see that choosing the parameter k of the same order as d , $SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d)$ is exponentially larger than $SD_{\mathcal{B}_{\mathcal{F}_2^d}}(\mu_d, \nu_d)$ in the dimension d . This result suggests that in high dimensions there exist high frequency densities that can be distinguished by feature-learning Stein Discrepancy discriminators but not by their fixed-kernel counterpart, and that may explain the differences in generative modeling performance for Stein discrepancy EBM and SVGD (Sec. 1).

7 Bounds of \mathcal{F}_1 and \mathcal{F}_2 IPMs by sliced Wasserstein distances

\mathcal{F}_1 and \mathcal{F}_2 IPMs measure differences of densities by slicing the input space and then maximizing (resp. averaging) the appropriate quantities. Max-sliced and sliced Wasserstein distances work, which have been studied by several works, work in an analogous fashion; one projects the distributions onto one-dimensional subspaces, and then maximizes or averages over the subspaces. Unlike the Wasserstein distance, which has been used for generative models such as WGAN (Arjovsky et al., 2017) but whose estimation suffers from the curse of dimensionality, max-sliced and sliced Wasserstein enjoy parametric estimation rates which make them more suitable as discriminators.

The goal of this section is to show that \mathcal{F}_1 IPMs are equivalent to max-sliced Wasserstein distances up to a constant power, while sliced Wasserstein distances are similarly equivalent to a fixed-kernel IPM with a kernel that is slightly different from the \mathcal{F}_2 kernel. These bounds are helpful to get a quantitative understanding of how strong feature-learning and fixed-kernel IPMs are, and provide a novel bridge between sliced optimal transport and generative modeling discriminators.

7.1 Spiked and sliced Wasserstein distances

Throughout this section k denotes an integer such that $1 \leq k \leq d$. The Stiefel manifold \mathcal{V}_k is the set of matrices $U \in \mathbb{R}^{k \times d}$ such that $UU^\top = I_{k \times k}$ (i.e. the rows of U are orthonormal). We define the k -dimensional projection robust p -Wasserstein distance between $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ as

$$\overline{\mathcal{W}}_{p,k}(\mu, \nu)^p = \max_{U \in \mathcal{V}_k} \min_{\pi \in \Gamma(\mu, \nu)} \int \|Ux - Uy\|^p d\pi(x, y), \quad (13)$$

where $\Gamma(\mu, \nu)$ denotes the set of couplings between μ, ν , i.e. of measures $\mathcal{P}(K \times K)$ with projections μ and ν . This is the distance studied by Niles-Weed and Rigollet (2019) as a good estimator for the Wasserstein distance for a certain class of target densities with low dimensional structure.

The integral k -dimensional projection robust p -Wasserstein distance between $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ is defined as

$$\underline{\mathcal{W}}_{p,k}(\mu, \nu)^p = \int_{\mathcal{V}_k} \left(\min_{\pi \in \Gamma(\mu, \nu)} \int \|Ux - Uy\|^p d\pi(x, y) \right) d\tau(U), \quad (14)$$

where τ is the uniform measure over \mathcal{V}_k . [Nadjahi et al. \(2020\)](#) studied statistical aspects of this distance in the case in which $k = 1$, while [Lin et al. \(2021\)](#) considers the case with general k . Notice trivially that $\overline{\mathcal{W}}_{p,k}(\mu, \nu) \geq \underline{\mathcal{W}}_{p,k}(\mu, \nu)$.

Sliced Wasserstein distances are spiked Wasserstein distances with $k = 1$, but they were studied first chronologically ([Bonneel et al., 2014](#); [Kolouri et al., 2016](#); [Kolouri et al., 2019](#)). Namely, the *sliced Wasserstein distance* is the integral 1-dimensional projection robust Wasserstein distance $\underline{\mathcal{W}}_{p,k}$, and the *max-sliced Wasserstein distance* is the 1-dimensional projection robust Wasserstein distance $\overline{\mathcal{W}}_{p,k}$. Some arguments are easier for the case $k = 1$ because the Stiefel manifold is the sphere \mathbb{S}^{d-1} .

7.2 Results

We prove in [Theorem 4](#) that for $K = \{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\} \times \{1\}$, for which the \mathcal{F}_1 space corresponds to overparametrized two-layer neural networks with bias, the \mathcal{F}_1 IPM can be upper and lower-bounded by the projection robust Wasserstein distance $\overline{\mathcal{W}}_{1,k}(\mu, \nu)$ up to a constant power (not depending on the dimension).

Theorem 4. *Let $\delta > 0$ be larger than a certain constant depending on k and α . Let $\sigma(x) = (x)_+^\alpha$ be the α -th power of the ReLu activation function, where α is a non-negative integer. Let μ, ν be Borel probability measures with support included in $\{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\} \times \{1\}$. Let $d_{\mathcal{B}_{\mathcal{F}_1}}$ be as defined in (2) and $\overline{\mathcal{W}}_{1,k}$ as defined in (13). Then,*

$$\delta \overline{\mathcal{W}}_{1,k}(\mu, \nu) \geq \delta d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) \geq \overline{\mathcal{W}}_{1,k}(\mu, \nu) - 2C(k, \alpha) \delta^{-\frac{1}{\alpha + (k-1)/2}} \log(\delta), \quad (15)$$

where $C(k, \alpha)$ is a constant that depends only on k and α . If we optimize the lower bound in (15) with respect to δ , we obtain $\overline{\mathcal{W}}_{1,k}(\mu, \nu) \geq d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) \geq \tilde{\Omega}(\overline{\mathcal{W}}_{1,k}(\mu, \nu)^{\alpha + \frac{k+1}{2}})$ where $\tilde{\Omega}$ hides log factors.

While for the \mathcal{F}_2 IPM the link with the sliced Wasserstein distance is not straightforward, it can be established when we switch from uniform τ to an alternative feature measure $\tilde{\tau}$. We define the class $\tilde{\mathcal{F}}_2$ of functions $\mathbb{R}^d \rightarrow \mathbb{R}$ as the RKHS associated with the following kernel

$$\begin{aligned} \tilde{k}(x, y) &= \int_{\mathbb{S}^d} \sigma(\langle (x, 1), \theta \rangle) \sigma(\langle (y, 1), \theta \rangle) d\tilde{\tau}(\theta) = \\ &= \frac{1}{\pi} \int_{\mathbb{S}^{d-1}} \int_{-1}^1 \sigma(\langle (x, 1), (\sqrt{1-t^2}\xi, t) \rangle) \sigma(\langle (y, 1), (\sqrt{1-t^2}\xi, t) \rangle) (1-t^2)^{-1/2} dt d\tau_{(d-1)}(\xi). \end{aligned}$$

Proposition 2. ($\tilde{\tau}$ as a rescaling of uniform measure) *The measure $d\tilde{\tau}(\sqrt{1-t^2}\xi, t) = \frac{1}{\pi}(1-t^2)^{-1/2} dt d\tau_{(d-1)}(\xi)$ is a probability measure. For comparison, the uniform measure over \mathbb{S}^d can be written as $d\tau(\sqrt{1-t^2}\xi, t) = \frac{\Gamma((d+1)/2)}{\sqrt{\pi}\Gamma(d/2)}(1-t^2)^{\frac{d-1}{2}} dt d\tau_{(d-1)}(\xi)$.*

That is, \mathcal{F}_2 and $\tilde{\mathcal{F}}_2$ are both fixed-kernel spaces with a similar kernel. They differ only in the weighing measure of the kernel; all the expressions which are valid in the \mathcal{F}_2 setting are also valid for $\tilde{\mathcal{F}}_2$ if we replace τ by $\tilde{\tau}$. In analogy with the \mathcal{F}_2 IPM, the $\tilde{\mathcal{F}}_2$ IPM is given below.

$$d_{\mathcal{B}_{\tilde{\mathcal{F}}_2}}^2(\mu, \nu) = \int_{\mathbb{S}^d} \left(\int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right)^2 d\tilde{\tau}(\theta). \quad (16)$$

Analogously to [Theorem 4](#), [Theorem 5](#) establishes that the $\tilde{\mathcal{F}}_2$ IPM is upper and lower-bounded by the sliced Wasserstein distance $\underline{\mathcal{W}}_{1,1}(\mu, \nu)$ up to a constant power (not depending on the dimension). The reason to introduce the space $\tilde{\mathcal{F}}_2$ is that in the proof, the argument that makes the connection with the sliced Wasserstein distance requires the base measure of the kernel to be $\tilde{\tau}$ and does not work for τ . However, we do not imply that a similar result for the \mathcal{F}_2 IPM is false.

Theorem 5. Let $\delta > 0$ be larger than a certain constant depending on k and α . Let $\sigma(x) = (x)_+^\alpha$ be the α -th power of the ReLU activation function, where α is a non-negative integer. Let μ, ν be Borel probability measures with support included in $\{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\} \times \{1\}$. Let $d_{\mathcal{B}_{\mathcal{F}_2}}$ be as defined in (16) and $\underline{\mathcal{W}}_{1,1}$ as defined in (14). Then,

$$\delta d_{\mathcal{F}_2}^{2/3}(\mu, \nu) \geq \left(\frac{5}{12\pi\alpha 2^{\alpha/2}} \right)^{1/3} \left(\underline{\mathcal{W}}_{1,1}(\mu, \nu) - 2C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta) \right). \quad (17)$$

and $\pi d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) \leq \underline{\mathcal{W}}_{1,1}(\mu, \nu)$. If we optimize the lower bound in (17) with respect to δ , we obtain $d_{\mathcal{F}_2}^{2/3}(\mu, \nu) \geq \tilde{\Omega}(\underline{\mathcal{W}}_{1,1}(\mu, \nu)^{1+\alpha})$.

8 Experiments

To validate and clarify our findings, we perform experiments of the settings studied [Sec. 5](#), [Sec. 6](#) and [Sec. 7](#). We use the ReLU activation function $\sigma(x) = (x)_+$, although remark that the results of [Sec. 5](#) hold for a generic activation function, and the results of [Sec. 6](#) and [Sec. 7](#) hold for non-negative integer powers of the ReLU activation. The empirical estimates in the plots are detailed in [App. G](#). They are averaged over 10 repetitions; the error bars show the maximum and minimum.

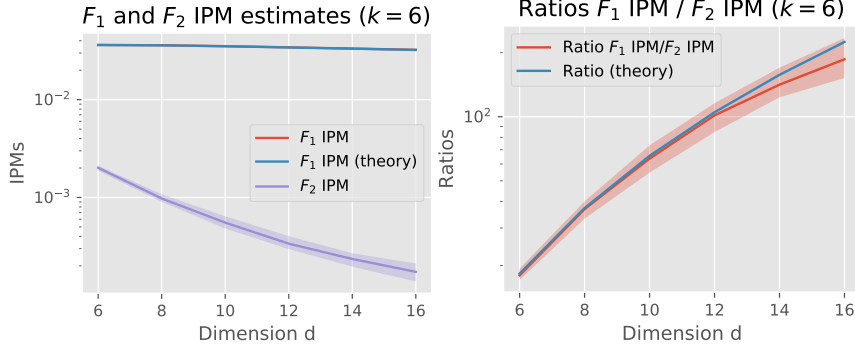


Figure 2: \mathcal{F}_1 and \mathcal{F}_2 IPM estimates for the pairs μ_d and ν_d defined in (4) for $k = 6$ and varying dimension d . (Left) The blue and red curves (superposed) show two different estimates of the \mathcal{F}_1 IPM. The purple curve shows estimates for the \mathcal{F}_2 IPM. (Right) The blue curve shows the theoretical ratio between the \mathcal{F}_1 and the \mathcal{F}_2 IPMs (see equation (6)). The red curve shows an empirical estimate of the ratio obtained by dividing the IPM estimates. 4400 million samples of μ_d and ν_d are used.

Separation between \mathcal{F}_1 and \mathcal{F}_2 IPMs. [Figure 2](#) shows \mathcal{F}_1 and \mathcal{F}_2 IPM estimates for the pairs μ_d and ν_d defined in (4) for the Legendre polynomial of degree $k = 6$ and varying dimension d , and its ratios. We observe that while the \mathcal{F}_1 IPM remains nearly constant in the dimension, the \mathcal{F}_2 IPM experiences a significant decrease. The ratios between IPMs closely track those predicted by our [Theorem 2](#), the mismatch being due to the overestimation of the \mathcal{F}_2 IPM caused by statistical errors. We were constrained in the values of k and d that we could choose; when the \mathcal{F}_2 IPM is small, which is the case when k and/or d are large, we need a high number of samples from the distributions μ_d, ν_d to make the statistical error smaller than $d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)$ and get a good estimate.

Separation between \mathcal{F}_1 and \mathcal{F}_2 SDs. [Figure 3](#) shows \mathcal{F}_1 and \mathcal{F}_2 SD estimates for the pairs μ_d and ν_d defined in equation [Subsec. 6.2](#) for the Legendre polynomial of degree $k = 5$ and varying dimension d , and its ratios. The fact that the empirical ratio is significantly above the theoretical lower bound indicates that our lower bound (although exponential) is not tight. This can be guessed by looking at the slackness in the inequalities of [Lemma 10](#) and [Lemma 11](#).

$\mathcal{F}_1, \mathcal{F}_2, \tilde{\mathcal{F}}_2$ IPMs versus max-sliced and sliced Wasserstein. [Figure 4](#) shows several metrics between a standard multivariate Gaussian and a Gaussian with unit variance in all directions except for one of smaller variance 0.1, in varying dimensions. We observe that while the \mathcal{F}_1 IPM and the max-sliced Wasserstein distance are constant, the $\mathcal{F}_2, \tilde{\mathcal{F}}_2$ IPMs and the sliced Wasserstein distance

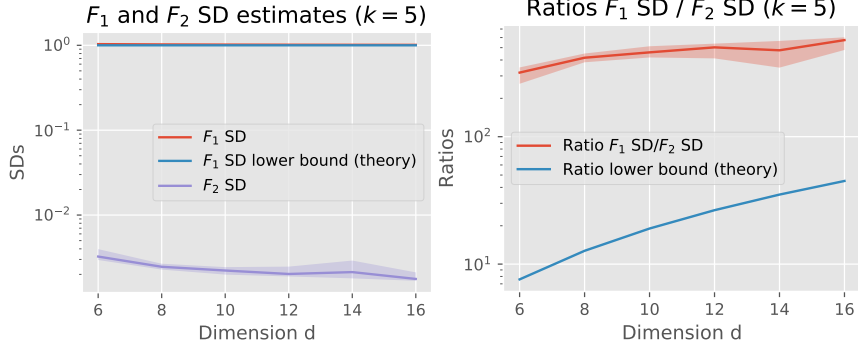


Figure 3: \mathcal{F}_1 and \mathcal{F}_2 SD estimates for the pairs μ_d and ν_d defined in Subsec. 6.2 for $k = 5$ and varying dimension d . (Left) The red curve shows an empirical estimate of the \mathcal{F}_1 SD, the blue curve shows a theoretical lower bound (Lemma 10 in App. E) on the \mathcal{F}_1 SD, the purple curve shows an estimate of the \mathcal{F}_2 SD. (Right) The blue curve represents the theoretical lower bound on the ratio between the \mathcal{F}_1 and the \mathcal{F}_2 SDs (see equation (11)), while the red curve shows an empirical estimate of the ratio obtained by dividing the SD estimates. 30 million samples are used.

decrease. For high dimensions they match the corresponding distances between two datasets of standard multivariate Gaussian, which means that the statistical noise precludes discrimination in these metrics.

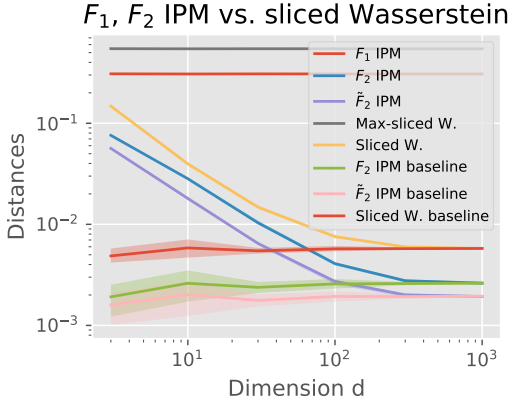


Figure 4: For varying dimension d , we plot \mathcal{F}_1 , \mathcal{F}_2 , $\tilde{\mathcal{F}}_2$ IPM, sliced and max-sliced Wasserstein estimates between a standard multivariate Gaussian and a Gaussian with unit variance in all directions except for one of smaller variance 0.1. The estimates are computed using 100000 samples of each distribution. For comparison, the same estimates are shown between a standard multivariate Gaussian and itself, using two different sets of 100000 samples.

9 Conclusions and discussion

We have shown pairs of distributions over hyperspheres for which the \mathcal{F}_1 IPM and SD are exponentially larger than the \mathcal{F}_2 IPM and SD. In parallel, we have also provided links between the \mathcal{F}_1 IPM and max-sliced Wasserstein distance, and between the $\tilde{\mathcal{F}}_2$ IPM and the sliced Wasserstein distance. The densities of the distributions constructed in Sections Sec. 5 and Sec. 6 are based on Legendre harmonics of increasing degree. Keeping in mind that spherical harmonics are the Fourier basis for $L^2(\mathbb{S}^{d-1})$ (in the sense that they constitute an orthonormal basis of eigenvalues of the Laplace-Beltrami operator), one can infer a simple overarching idea from our constructions: ‘ \mathcal{F}_1 discriminators are better than \mathcal{F}_2 discriminators at telling apart distributions whose densities have only high frequency differences. It would be interesting to develop this intuition into a more general theory. Another avenue of future work is to understand how deep discriminators perform versus shallow ones, in analogy with the work of Eldan and Shamir (2016) for regression.

Acknowledgements. We thank Joan Bruna for useful discussions. CD acknowledges partial support by “la Caixa” Foundation (ID 100010434), under agreement LCF/BQ/AA18/11680094.

References

- Arjovsky, M., Chintala, S., and Bottou, L. (2017). Wasserstein gan. *arXiv preprint arXiv:1701.07875*.
- Arora, S., Ge, R., Liang, Y., Ma, T., and Zhang, Y. (2017). Generalization and equilibrium in generative adversarial nets (GANs). In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 224–232. JMLR. org.
- Atkinson, K. and Han, W. (2012). *Spherical Harmonics and Approximations on the Unit Sphere: An Introduction*, volume 2044. Springer.
- Bach, F. (2017). Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53.
- Barp, A., Briol, F.-X., Duncan, A., Girolami, M., and Mackey, L. (2019). Minimum stein discrepancy estimators. In *Advances in Neural Information Processing Systems*, pages 12964–12976.
- Bateman, H. and Erdélyi, A. (1954). *Tables of integral transforms*. McGraw-Hill.
- Bonneel, N., Rabin, J., Peyré, G., and Pfister, H. (2014). Sliced and radon wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51.
- Chang, W.-C., Li, C.-L., Mroueh, Y., and Yang, Y. (2020). Kernel stein generative modeling.
- Chizat, L. and Bach, F. (2020). Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on Learning Theory*, pages 1305–1338. PMLR.
- Cho, Y. and Saul, L. K. (2009). Kernel methods for deep learning. In *Advances in Neural Information Processing Systems 22*, pages 342–350. Curran Associates, Inc.
- Chwialkowski, K., Strathmann, H., and Gretton, A. (2016). A kernel test of goodness of fit. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2606–2615. PMLR.
- Domingo-Enrich, C., Bietti, A., Vanden-Eijnden, E., and Bruna, J. (2021). On energy-based models with overparametrized shallow neural networks. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 2771–2782. PMLR.
- Dziugaite, G. K., Roy, D. M., and Ghahramani, Z. (2015). Training generative neural networks via maximum mean discrepancy optimization. *UAI*.
- Efthimiou, C. and Frye, C. (2014). *Spherical Harmonics in p Dimensions*.
- Eldan, R. and Shamir, O. (2016). The power of depth for feedforward neural networks. In *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 907–940. PMLR.
- Genevay, A., Peyre, G., and Cuturi, M. (2018). Learning generative models with sinkhorn divergences. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1608–1617. PMLR.
- Ghorbani, B., Mei, S., Misiakiewicz, T., and Montanari, A. (2019). Limitations of lazy training of two-layers neural network. In *NeurIPS*.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *NIPS*.
- Gorham, J. and Mackey, L. (2015). Measuring sample quality with stein’s method. In *Advances in Neural Information Processing Systems*, pages 226–234.
- Gorham, J. and Mackey, L. (2017). Measuring sample quality with kernels. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 1292–1301. PMLR.

- Grathwohl, W., Wang, K.-C., Jacobsen, J.-H., Duvenaud, D., and Zemel, R. (2020). Learning the stein discrepancy for training and evaluating energy-based models without sampling. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 3732–3747.
- Gretton, A., Borgwardt, K. M., Rasch, M., Schölkopf, B., and Smola, A. J. (2007). A kernel method for the two-sample-problem. In *Advances in neural information processing systems*, pages 513–520.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012). A kernel two-sample test. *Journal of Machine Learning Research*, 13(25):723–773.
- Huggins, J. and Mackey, L. (2018). Random feature stein discrepancies. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Jacot, A., Gabriel, F., and Hongler, C. (2018). Neural tangent kernel: Convergence and generalization in neural networks. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 31, pages 8571–8580. Curran Associates, Inc.
- Kammler, D. (2000). *A first course in Fourier analysis*. Prentice Hall.
- Kolouri, S., Nadjahi, K., Simsekli, U., Badeau, R., and Rohde, G. (2019). Generalized sliced wasserstein distances. In *Advances in Neural Information Processing Systems*, volume 32.
- Kolouri, S., Zou, Y., and Rohde, G. K. (2016). Sliced wasserstein kernels for probability distributions. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5258–5267.
- Li, C.-L., Chang, W.-C., Cheng, Y., Yang, Y., and Poczos, B. (2017). Mmd gan: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems*, volume 30.
- Li, Y., Swersky, K., and Zemel, R. (2015). Generative moment matching networks. In *ICML*.
- Lin, T., Zheng, Z., Chen, E. Y., Cuturi, M., and Jordan, M. I. (2021). On projection robust optimal transport: Sample complexity and model misspecification.
- Liu, Q., Lee, J., and Jordan, M. (2016). A kernelized stein discrepancy for goodness-of-fit tests. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 276–284, New York, New York, USA. PMLR.
- Liu, Q. and Wang, D. (2016). Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in Neural Information Processing Systems*, volume 29, pages 2378–2386. Curran Associates, Inc.
- Mohamed, S. and Lakshminarayanan, B. (2017). Learning in implicit generative models. In *ICLR*.
- Nadjahi, K., Durmus, A., Chizat, L., Kolouri, S., Shahrampour, S., and Simsekli, U. (2020). Statistical and topological properties of sliced probability divergences. In *Advances in Neural Information Processing Systems*, volume 33, pages 20802–20812.
- Niles-Weed, J. and Rigollet, P. (2019). Estimation of wasserstein distances in the spiked transport model.
- Ongie, G., Willett, R., Soudry, D., and Srebro, N. (2019). A function space view of bounded norm infinite width relu nets: The multivariate case. In *International Conference on Learning Representations (ICLR 2020)*.
- Paty, F.-P. and Cuturi, M. (2019). Subspace robust wasserstein distances.
- Rezende, D. J. and Mohamed, S. (2015). Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*.

- Roux, N. L. and Bengio, Y. (2007). Continuous neural networks. In *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, volume 2 of *Proceedings of Machine Learning Research*, pages 404–411, San Juan, Puerto Rico.
- Santos, C. N. D., Mroueh, Y., Padhi, I., and Dognin, P. (2019). Learning implicit generative models by matching perceptual features. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4460–4469.
- Savarese, P., Evron, I., Soudry, D., and Srebro, N. (2019). How do infinite width bounded norm networks look in function space? In *Conference on Learning Theory*.
- Stein, C. (1972). A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 2: Probability Theory*, pages 583–602.
- Venturi, L., Jelassi, S., Ozuch, T., and Bruna, J. (2021). Depth separation beyond radial functions.
- Wei, C., Lee, J. D., Liu, Q., and Ma, T. (2020). Regularization matters: Generalization and optimization of neural nets v.s. their induced kernel.
- Williams, F., Trager, M., Silva, C., Panozzo, D., Zorin, D., and Bruna, J. (2019). Gradient dynamics of shallow univariate relu networks. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Woodworth, B., Gunasekar, S., Lee, J. D., Moroshko, E., Savarese, P., Golan, I., Soudry, D., and Srebro, N. (2020). Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*.

A \mathcal{F}_1 and \mathcal{F}_2 IPMs

Lemma 1. *The \mathcal{F}_1 IPM can be written as*

$$d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) = \sup_{\theta \in \mathbb{S}^d} \left| \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right|.$$

Proof.

$$\begin{aligned} d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) &= \sup_{f \in \mathcal{B}_{\mathcal{F}_1}} \mathbb{E}_{x \sim \mu} f(x) - \mathbb{E}_{x \sim \nu} f(x) \\ &= \sup_{\|\gamma\|_{\text{TV}} \leq 1} \mathbb{E}_{x \sim \mu} \int_{\mathbb{S}^d} \sigma(\langle x, \theta \rangle) d\gamma(\theta) - \mathbb{E}_{x \sim \nu} \int_{\mathbb{S}^d} \sigma(\langle x, \theta \rangle) d\gamma(\theta) \\ &= \sup_{\|\mu\|_{\text{TV}} \leq 1} \int_{\mathbb{S}^d} \left(\int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right) d\gamma(\theta) \\ &= \sup_{\theta \in \mathbb{S}^d} \left| \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right| \end{aligned}$$

In the last equality we have used that the set $\{\mu \in \mathcal{M}(\mathbb{S}^d) \mid \|\mu\|_{\text{TV}} \leq 1\}$ can be seen as the convex hull of $\{\delta_\theta \mid \theta \in \mathbb{S}^d\}$, which means that optimizing a convex function over one set and the other yields the same optimal value. \square

Lemma 2. *The \mathcal{F}_2 IPM can be written as*

$$d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) = \int_{\mathbb{S}^d} \left(\int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right)^2 d\tau(\theta).$$

Proof.

$$\begin{aligned} d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) &= \left(\sup_{f \in \mathcal{B}_{\mathcal{F}_2}} \mathbb{E}_{x \sim \mu} f(x) - \mathbb{E}_{x \sim \nu} f(x) \right)^2 = \left(\sup_{\|f\|_{\mathcal{F}_2} \leq 1} \int_K \langle f, k(x, \cdot) \rangle_{\mathcal{F}_2} d(\mu - \nu)(x) \right)^2 \\ &= \left(\sup_{\|f\|_{\mathcal{F}_2} \leq 1} \left\langle f, \int_K k(x, \cdot) d(\mu - \nu)(x) \right\rangle_{\mathcal{F}_2} \right)^2 = \left\| \int_K k(x, \cdot) d(\mu - \nu)(x) \right\|_{\mathcal{F}_2}^2 \\ &= \iint_{K \times K} k(x, y) d(\mu - \nu)(x) d(\mu - \nu)(y) = \iint_{K \times K} k(x, y) d(\mu - \nu)(x) d(\mu - \nu)(y) \\ &= \iint_{K \times K} \int_{\mathbb{S}^d} \sigma(\langle x, \theta \rangle) \sigma(\langle y, \theta \rangle) d\tau(\theta) d(\mu - \nu)(x) d(\mu - \nu)(y) \\ &= \int_{\mathbb{S}^d} \left(\int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right)^2 d\tau(\theta). \end{aligned} \tag{18}$$

\square

B Are the \mathcal{F}_1 and \mathcal{F}_2 IPMs distances?

In the following we consider unitary Fourier transforms with angular frequency: for $f \in L^1(\mathbb{R}^{d+1})$, we have $\hat{f}(\omega) = \frac{1}{(2\pi)^{(d+1)/2}} \int_{\mathbb{R}^{d+1}} f(x) e^{-i\langle \omega, x \rangle} dx$ and if $\hat{f} \in L^1(\mathbb{R}^{d+1})$, then $f(x) = \frac{1}{(2\pi)^{(d+1)/2}} \int_{\mathbb{R}^d} \hat{f}(\omega) e^{-i\langle \omega, x \rangle} d\omega$. We denote the space of tempered distributions on \mathbb{R}^{d+1} as $\mathcal{S}'(\mathbb{R}^{d+1})$, i.e., as the dual of the space $\mathcal{S}(\mathbb{R}^{d+1})$ of Schwartz functions, which are functions in $\mathcal{C}^\infty(\mathbb{R}^{d+1})$ whose derivatives of any order decay faster than polynomials of all orders. Functions that grow no faster than polynomials can be embedded in $\mathcal{S}'(\mathbb{R}^{d+1})$ by defining $g(\varphi) := \int_{\mathbb{R}^{d+1}} \varphi(x) g(x) dx$ for any

$\varphi \in \mathcal{S}(\mathbb{R}^{d+1})$. The Fourier transform of a tempered distribution T can be defined as the tempered distribution \hat{T} that acts on $\varphi \in \mathcal{S}(\mathbb{R}^{d+1})$ as $\langle \hat{T}, \varphi \rangle = \langle T, \hat{\varphi} \rangle$. Fourier transforms of two-layer neural networks have been used in prior works, e.g. [Venturi et al. \(2021\)](#).

Lemma 3. *Let $\hat{\sigma} \in \mathcal{S}'(\mathbb{R})$ be the Fourier transform of the activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ in the sense of tempered distributions. Let $g(\theta) = \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x)$. The Fourier transform of $g \in \mathcal{S}'(\mathbb{R}^{d+1})$ in the sense of tempered distributions is the tempered distribution \hat{g} defined as*

$$\langle \hat{g}, \varphi \rangle = (2\pi)^{d/2} \int_K \langle \hat{\sigma}, \varphi(\cdot x) \rangle d(\mu - \nu)(x)$$

for any $\varphi \in \mathcal{S}(\mathbb{R}^{d+1})$.

Proof. The Fourier transform of the tempered distribution T_x defined as $\langle T_x, \varphi \rangle = \int_{\mathbb{R}^{d+1}} \varphi(\theta) \sigma(\langle x, \theta \rangle) d\theta$ is \hat{T}_x defined as

$$\begin{aligned} \langle \hat{T}_x, \varphi \rangle &= \int_{\mathbb{R}^{d+1}} \left(\frac{1}{(2\pi)^{(d+1)/2}} \int_{\mathbb{R}^{d+1}} \varphi(\theta) e^{-i\langle \omega, \theta \rangle} d\theta \right) \sigma(\langle x, \omega \rangle) d\omega \\ &= \int_{\text{span}(x)} \int_{\text{span}(x)^\perp} \left(\frac{1}{(2\pi)^{(d+1)/2}} \int_{\text{span}(x)^\perp} \left(\int_{\text{span}(x)} \varphi(\theta) e^{-i\langle \omega_x, \theta_x \rangle} d\theta_x \right) e^{-i\langle \omega_{x^\perp}, \theta_{x^\perp} \rangle} d\theta_{x^\perp} \right) d\omega_{x^\perp} \sigma(\langle x, \omega_x \rangle) d\omega_x \\ &= \frac{1}{(2\pi)^{1/2}} \int_{\text{span}(x)} (2\pi)^{d/2} \int_{\text{span}(x)} \varphi(\theta_x) e^{-i\langle \omega_x, \theta_x \rangle} d\theta_x \sigma(\langle x, \omega_x \rangle) d\omega_x \\ &= \frac{1}{(2\pi)^{1/2}} \int_{\mathbb{R}} (2\pi)^{d/2} \int_{\mathbb{R}} \varphi(tx) e^{-i\omega t} dt \sigma(\omega) d\omega \\ &= (2\pi)^{d/2} \langle \hat{\sigma}, \varphi(\cdot x) \rangle \end{aligned}$$

Here, the first equality holds because by definition, $\langle \hat{T}_x, \varphi \rangle = \langle T_x, \hat{\varphi} \rangle$. In the second equality, we rewrite $\mathbb{R}^{d+1} = \text{span}(x) + \text{span}(x)^\perp$ and we use Fubini's theorem twice. In the third equality we make the following argument: denoting $h(\theta_{x^\perp}, \omega_x) = \int_{\text{span}(x)} \varphi(\theta_{x^\perp} + \theta_x) e^{-i\langle \omega_x, \theta_x \rangle} d\theta_x$, we have that $\int_{\text{span}(x)^\perp} \left(\int_{\text{span}(x)} h(\theta_{x^\perp}, \omega_x) e^{-i\langle \omega_{x^\perp}, \theta_{x^\perp} \rangle} d\theta_{x^\perp} \right) d\omega_{x^\perp} = (2\pi)^{d/2} \int_{\text{span}(x)^\perp} \hat{h}(\omega_{x^\perp}, \omega_x) d\omega_{x^\perp} = (2\pi)^d h(0, \omega_x) = (2\pi)^d \int_{\text{span}(x)} \varphi(\theta_x) e^{-i\langle \omega_x, \theta_x \rangle} d\theta_x$.

Notice that we can write g as a tempered distribution as $g = \int_K T_x d(\mu - \nu)(x)$. Thus, by linearity of the Fourier transform, we have that

$$\langle \hat{g}, \varphi \rangle = (2\pi)^{d/2} \int_K \langle \hat{\sigma}, \varphi(\cdot x) \rangle d(\mu - \nu)(x)$$

for any $\varphi \in \mathcal{S}(\mathbb{R}^{d+1})$. □

We compute $\hat{\sigma}$ for the specific case in which $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is an α -positive homogeneous activation function, i.e. $\sigma(x) = a(x)_+^\alpha + b(-x)_+^\alpha$ for some $a, b \in \mathbb{R}$ (equation (3)). It is known ([Bateman and Erdélyi, 1954](#); [Kammler, 2000](#)) that the Fourier transform of the Heaviside step function $u : \mathbb{R} \rightarrow \mathbb{R}$, defined as $u(x) = 1$ if $x \geq 0$ and $u(x) = 0$ if $x < 0$, is $\hat{u}(\omega) = \sqrt{\frac{\pi}{2}} \left(\text{p.v.} \left[\frac{1}{i\pi\omega} \right] + \delta(\omega) \right)$. Here $\text{p.v.} \left[\frac{1}{\omega} \right] \in \mathcal{S}'(\mathbb{R})$ is a Cauchy principal value, defined as $\text{p.v.} \left[\frac{1}{\omega} \right] (\varphi) = \lim_{\epsilon \rightarrow 0} \int_{\mathbb{R} \setminus [-\epsilon, \epsilon]} \frac{1}{\omega} \varphi(\omega) d\omega$ for any $\varphi \in \mathcal{S}(\mathbb{R})$.

Moreover, for any tempered distribution $f \in \mathcal{S}'(\mathbb{R})$, for $\alpha \geq 0$ integer, the Fourier transform of $x^\alpha f(x)$ is $i^\alpha \frac{d^\alpha \hat{f}(\omega)}{d\omega^\alpha}$, where the derivative of a tempered distribution is defined in the weak sense:

$\langle \frac{df}{d\omega}, \varphi \rangle = -\langle f, \frac{d\varphi}{d\omega} \rangle$. Since $\sigma(x) = (a - (-1)^\alpha b)(x)_+^\alpha + (-1)^\alpha b x^\alpha$, we have that

$$\begin{aligned}\hat{\sigma}(\omega) &= (a - (-1)^\alpha b) i^\alpha \sqrt{\frac{\pi}{2}} \frac{d^\alpha}{d\omega^\alpha} \left(\text{p.v.} \left[\frac{1}{i\pi\omega} \right] + \delta(\omega) \right) + (-1)^\alpha b i^\alpha \frac{d^\alpha}{d\omega^\alpha} \delta(\omega) \\ &= A \frac{d^\alpha}{d\omega^\alpha} \left(\text{p.v.} \left[\frac{1}{i\pi\omega} \right] \right) + B \frac{d^\alpha}{d\omega^\alpha} \delta(\omega),\end{aligned}\tag{19}$$

where $A = i^{\alpha-1} \frac{\alpha!}{\sqrt{2\pi}} (a - (-1)^\alpha b)$ and $B = i^\alpha \sqrt{\frac{\pi}{2}} (a - (-1)^\alpha b) + (-i)^\alpha b$.

Lemma 4 (Riesz-Markov theorem). *Let X be a locally compact Hausdorff space. For any continuous linear functional ψ on $C_0(X)$, there is a unique regular countably additive complex Borel measure μ on X such that $\forall f \in C_0(X)$, $\psi(f) = \int_X f(x) d\mu(x)$. The norm of ψ as a linear functional is the total variation of μ , that is $\|\psi\| = |\mu|(X)$.*

Theorem 6. *Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an α -positive homogeneous activation function of the form (3) and assume that $(-1)^\alpha a - b \neq 0$. Then, for $K = \mathbb{R}^d \times \{1\}$ both the \mathcal{F}_1 and the \mathcal{F}_2 IPMs are distances.*

Proof. By Lemma 1 and Lemma 2 we have $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) = \sup_{\theta \in \mathbb{S}^d} \left| \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right|$ and $d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) = \int_{\mathbb{S}^d} \left(\int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right)^2 d\tau(\theta)$, which means that both are distances if the function $g|_{\mathbb{S}^d} : \mathbb{S}^d \rightarrow \mathbb{R}$ defined as $g|_{\mathbb{S}^d}(\theta) = \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x)$ is different from zero in the L^2 sense when $\mu \neq \nu$. Since σ is α -positive homogeneous, the α -positive homogeneous extension $g : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$ of $g|_{\mathbb{S}^d}$ fulfills

$$g(\theta) = \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) = \|\theta\|_2^\alpha \int_K \sigma(\langle x, \theta/\|\theta\|_2 \rangle) d(\mu - \nu)(x) = \|\theta\|_2^\alpha g(\theta/\|\theta\|_2).$$

Thus, $g|_{\mathbb{S}^d}$ is different from zero in the L^2 sense if and only if g is. And g is different from zero if and only if its Fourier transform \hat{g} in the sense of tempered distributions is different from zero (this follows from $\langle \hat{g}, \varphi \rangle = \langle g, \hat{\varphi} \rangle$ for all $\varphi \in \mathcal{S}(\mathbb{R})$). By Lemma 3, we have that

$$\langle \hat{g}, \varphi \rangle = (2\pi)^{d/2} \int_K \langle \hat{\sigma}, \varphi(\cdot x) \rangle d(\mu - \nu)(x).\tag{20}$$

By (19), we have that

$$\langle \hat{\sigma}, \varphi(\cdot x) \rangle = (-1)^\alpha A \left(\text{p.v.} \left[\frac{1}{t} \right] \right) \left(\frac{d^\alpha}{dt^\alpha} \varphi(tx) \right) + (-1)^\alpha B \frac{d^\alpha}{dt^\alpha} \varphi(tx) \Big|_{t=0},\tag{21}$$

which means that

$$\langle \hat{g}, \varphi \rangle = (-1)^\alpha (2\pi)^{d/2} \int_K \left(A \left(\text{p.v.} \left[\frac{1}{t} \right] \right) \left(\frac{d^\alpha}{dt^\alpha} \varphi(tx) \right) + B \frac{d^\alpha}{dt^\alpha} \varphi(tx) \Big|_{t=0} \right) d(\mu - \nu)(x).$$

Suppose that $\mu \neq \nu$. Since μ and ν are Borel measures, they are regular, and thus $\mu - \nu$ is also regular. By Lemma 4, $\mu - \nu$ can be identified univocally with an element of the dual space $C_0(\mathbb{R}^{d+1})'$ of the space $C_0(\mathbb{R}^{d+1})$ of continuous functions on \mathbb{R}^{d+1} that vanish at infinity. Since $\mu - \nu \neq 0$, there must exist $\varphi \in C_0(\mathbb{R}^{d+1})$ such that $\int_K \varphi(x) d(\mu - \nu)(x) \neq 0$. Multiplying by the indicator function of a well chosen compact set and using a mollifier sequence, we can further assume that $\varphi \in C_c^\infty(\mathbb{R}^{d+1}) \subseteq \mathcal{S}(\mathbb{R}^{d+1})$.

Now, let η be a $C_c^\infty(\mathbb{R})$ function such that $\int_{\mathbb{R}} \eta(t) dt = 1$ and the support of η is compact and contained in $[1/2, +\infty)$. We define the function $\psi \in C_c^\infty(\mathbb{R}^{d+1}) \subseteq \mathcal{S}(\mathbb{R}^{d+1})$ as

$$\psi(tx) = \alpha! t^{\alpha+1} \varphi(x) \eta(t), \quad \forall x \in \mathbb{R}^d \times \{1\}, \forall t \in \mathbb{R}.$$

Remark that

$$\frac{d^\alpha}{dt^\alpha} \psi(tx) \Big|_{t=0} = 0,\tag{22}$$

because ψ is equal to zero in a neighborhood of the origin. Also, for all $x \in K$,

$$\begin{aligned} (-1)^\alpha \left(\text{p.v.} \left[\frac{1}{t} \right] \right) \left(\frac{d^\alpha}{dt^\alpha} \psi(tx) \right) &= (-1)^\alpha \int_{\mathbb{R}} \frac{1}{t} \frac{d^\alpha}{dt^\alpha} \left(\alpha! t^{\alpha+1} \varphi(x) \eta(t) \right) dt \\ &= \int_{\mathbb{R}} \frac{1}{t^{\alpha+1}} t^{\alpha+1} \varphi(x) \eta(t) dt = \varphi(x) \int_{\mathbb{R}} \eta(t) dt = \varphi(x). \end{aligned} \quad (23)$$

In the first equality we have used that $\lim_{\epsilon \rightarrow 0} \int_{\mathbb{R} \setminus [-\epsilon, \epsilon]} \frac{1}{t} \frac{d^\alpha}{dt^\alpha} \psi(tx) dt = \int_{\mathbb{R}} \frac{1}{t} \frac{d^\alpha}{dt^\alpha} \psi(tx) dt$, again because ψ is equal to zero in a neighborhood of the origin.

Notice that since we have assumed that $(-1)^\alpha a - b \neq 0$, we have $A \neq 0$. Hence,

$$\begin{aligned} 0 &\neq (2\pi)^{d/2} A \int_K \varphi(x) d(\mu - \nu)(x) \\ &= (-1)^\alpha (2\pi)^{d/2} \int_K \left(A \left(\text{p.v.} \left[\frac{1}{t} \right] \right) \left(\frac{d^\alpha}{dt^\alpha} \psi(tx) \right) + B \frac{d^\alpha}{dt^\alpha} \psi(tx) \Big|_{t=0} \right) d(\mu - \nu)(x) = \langle \hat{g}, \psi \rangle \end{aligned}$$

In the first equality, we have used (23) and (22). The last equality follows from (20) and (21). We have constructed a function $\psi \in \mathcal{S}(\mathbb{R}^{d+1})$ for which \hat{g} does not evaluate to zero, implying that $\hat{g} \neq 0$ and concluding the proof. \square

Theorem 7. Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an α -positive homogeneous activation function of the general form (3). Then, for $K = \mathbb{S}^d$ both the \mathcal{F}_1 and the \mathcal{F}_2 IPMs are not distances.

Proof. Since by Lemma 1 and Lemma 2 we have $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) = \sup_{\theta \in \mathbb{S}^d} \left| \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right|$ and $d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) = \int_{\mathbb{S}^d} \left(\int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x) \right)^2 d\tau(\theta)$, it suffices to see that the function $g|_{\mathbb{S}^d}(\theta) = \int_K \sigma(\langle x, \theta \rangle) d(\mu - \nu)(x)$ can be zero in the L^2 sense for some pairs $\mu \neq \nu$. We want to study the kernel of the map $\mu \mapsto \int_K \sigma(\langle x, \cdot \rangle) d(\mu - \nu)(x)$. Notice that

$$\begin{aligned} \sigma(\langle x, \theta \rangle) + (-1)^\alpha \sigma(\langle -x, \theta \rangle) &= a(\langle x, \theta \rangle)_+^\alpha + b(-\langle x, \theta \rangle)_+^\alpha + (-1)^\alpha (a(\langle -x, \theta \rangle)_+^\alpha + b(-\langle -x, \theta \rangle)_+^\alpha) \\ &= (a + (-1)^\alpha b)((\langle x, \theta \rangle)_+^\alpha + (-1)^\alpha (-\langle x, \theta \rangle)_+^\alpha) = (a + (-1)^\alpha b) \langle x, \theta \rangle^\alpha. \end{aligned}$$

If $\gamma \in \mathcal{M}(K)$ is an even (i.e., $(x \mapsto -x)_\# \gamma = \gamma$) or odd (i.e., $(x \mapsto -x)_\# \gamma = -\gamma$) signed measure of the same parity as α , this implies that

$$\int_K \sigma(\langle x, \theta \rangle) d\gamma(x) = \frac{a + (-1)^\alpha b}{2} \int_K \langle x, \theta \rangle^\alpha d|\gamma|(x),$$

which is a polynomial of degree α on θ . Consider the linear map $L : \mathcal{M}(K) \mapsto C(\mathbb{S}^d)$ defined as $\gamma \mapsto \int_K \sigma(\langle x, \cdot \rangle) d\gamma(x)$. Since L restricted to the measures of the parity of α has an infinite-dimensional domain and a finite-dimensional image, it must have an infinite-dimensional kernel.

- For the case α odd, if $\gamma \in \mathcal{M}(K)$ is an odd measure belonging to the kernel of L with total variation norm $\|\gamma\|_{\text{TV}} = 2$ and such that $\gamma = \gamma_+ - \gamma_-$ with γ_+, γ_- non-negative, then choosing $\mu = \gamma_+$ and $\nu = \gamma_-$, we have that $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) = d_{\mathcal{B}_{\mathcal{F}_2}}(\mu, \nu) = 0$.
- For the case α even, let $\gamma \in \mathcal{M}(K)$ be an even measure belonging to the kernel of L with total variation norm $\|\gamma\|_{\text{TV}} = 2$. We must have that $\int_K d\gamma = 0$, because denoting by τ the uniform distribution over \mathbb{S}^d , we have

$$0 = \int_{\mathbb{S}^d} L\gamma(\theta) d\tau(\theta) = \int_{\mathbb{S}^d} \int_K \sigma(\langle x, \theta \rangle) d\gamma(x) d\tau(\theta) = \int_{\mathbb{S}^d} \sigma(\langle x', \theta \rangle) d\tau(\theta) \int_K d\gamma,$$

for all $x' \in K$, and $\int_{\mathbb{S}^d} \sigma(\langle x', \theta \rangle) d\tau(\theta)$ is a strictly positive quantity. In the last equality we used Fubini's theorem. Thus, the non-negative components γ_+, γ_- of the decomposition $\gamma = \gamma_+ - \gamma_-$ must fulfill $\int_K d\gamma_+ = \int_K d\gamma_- = 1$. Hence, choosing $\mu = \gamma_+$ and $\nu = \gamma_-$, we have that $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) = d_{\mathcal{B}_{\mathcal{F}_2}}(\mu, \nu) = 0$.

□

C Preliminaries on Legendre Polynomials and spherical harmonics

In the notation of [Sec. 5](#), $P_{k,d}(t)$ denotes the *Legendre polynomial* of degree k in dimension d .

It is known (equation (2.78), [Atkinson and Han \(2012\)](#)) that

$$\int_{-1}^1 P_{k,d}(t)^2 (1-t^2)^{\frac{d-3}{2}} dt = \frac{|\mathbb{S}^{d-1}|}{N_{k,d} |\mathbb{S}^{d-2}|} = \frac{\sqrt{\pi} \Gamma(\frac{d-1}{2})}{N_{k,d} \Gamma(\frac{d}{2})}, \quad (24)$$

where (equation (2.10), [Atkinson and Han \(2012\)](#))

$$N_{k,d} = \frac{(2k+d-2)(k+d-3)!}{k!(d-2)!} \quad (25)$$

is the dimension of the space of homogeneous harmonic polynomials of degree k in \mathbb{R}^d .

Also, the following uniform bound holds (equation (2.116), [Atkinson and Han \(2012\)](#)):

$$|P_{k,d}(t)| \leq 1, \quad k \geq 0, \quad d \geq 2, \quad t \in [-1, 1]. \quad (26)$$

The bound is achieved at $t = 1, -1$.

There is also a crucial link between Legendre polynomials and their derivatives (equation (2.90), [Atkinson and Han \(2012\)](#)):

$$P_{k,d}^{(j)}(t) = \frac{k!(k+j+d-3)! \Gamma(\frac{d-1}{2})}{2^j (k-j)!(k+d-3)! \Gamma(j + \frac{d-1}{2})} P_{k-j,d+2j}(t), \quad (27)$$

where $k \geq j$ and $d \geq 2$. Note that for $k < j$, $P_{k,d}^{(j)}(t) = 0$.

We recall two important equalities. Let $\{Y_{k,j} \mid 1 \leq j \leq N_{k,d}\}$ be an orthonormal basis of the space of homogeneous harmonic polynomials over \mathbb{R}^d of degree k , with real coefficients (some works like [Atkinson and Han \(2012\)](#) consider complex coefficients and all the results are unchanged up to complex conjugates). That is, $\int_{\mathbb{S}^d} Y_{k,j}(x) Y_{k,i}(x) d\tau(x) = \delta_{ij}$, where τ is the uniform probability measure over \mathbb{S}^{d-1} . Then, the addition theorem (Thm. 4.11, [Efthimiou and Frye \(2014\)](#)) states that

$$\sum_{j=1}^{N_{k,d}} Y_{k,j}(x) Y_{k,j}(y) = N_{k,d} P_{k,d}(\langle x, y \rangle)$$

The Funk-Hecke formula (Thm 2.22, [Atkinson and Han \(2012\)](#)) states that when $\int_{-1}^1 |f(t)|(1-t^2)^{\frac{d-3}{2}} dt < +\infty$, for any linear combination Y_k of $\{Y_{k,j} \mid 1 \leq j \leq N_{k,d}\}$ and for any $x \in \mathbb{S}^{d-1}$,

$$\int_{\mathbb{S}^{d-1}} f(\langle x, y \rangle) Y_k(y) d\tau(y) = \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} Y_k(x) \int_{-1}^1 P_{k,d}(t) f(t) (1-t^2)^{\frac{d-3}{2}} dt.$$

D Proofs of [Sec. 5](#)

Proposition 1. *If we choose $\gamma_{k,d} = 2 \left(\int_{\mathbb{S}^{d-1}} |L_{k,d}(x)| d\tau(x) \right)^{-1}$, then μ_d and ν_d are probability measures.*

Proof. Clearly, $\frac{d\mu_d}{d\lambda}(x) \geq 0$ and $\frac{d\nu_d}{d\lambda}(x) \geq 0$ for all $x \in \mathbb{S}^{d-1}$. Let $C_+ = \{x \in \mathbb{S}^{d-1} \mid L_{k,d}(x) > 0\}$ and $C_- = \{x \in \mathbb{S}^{d-1} \mid L_{k,d}(x) \leq 0\}$. For μ_d and ν_d to be probability measures, $\gamma_{k,d}$ must fulfill

$$1 = \gamma_{k,d} \int_{C_+} \frac{L_{k,d}(x)}{|\mathbb{S}^{d-1}|} d\lambda(x) \text{ and } 1 = -\gamma_{k,d} \int_{C_-} \frac{L_{k,d}(x)}{|\mathbb{S}^{d-1}|} d\lambda(x). \quad (28)$$

By equation (1.17) of [Atkinson and Han \(2012\)](#), if we parametrize \mathbb{S}^{d-1} as $x = te_d + (1-t^2)^{1/2}\xi_{(d-1)}$ with $t \in [-1, 1]$ and $\xi_{(d-1)} \in \mathbb{S}^{d-2}$, we have

$$d\lambda(x) = (1-t^2)^{\frac{d-3}{2}} dt d\lambda_{(d-2)}(\xi_{(d-1)}),$$

where $\lambda_{(d-2)}$ denotes the Hausdorff measure of \mathbb{S}^{d-2} . Hence,

$$\begin{aligned} \int_{\mathbb{S}^{d-1}} \frac{L_{k,d}(x)}{|\mathbb{S}^{d-1}|} d\lambda(x) &= \int_{\mathbb{S}^{d-2}} \int_{-1}^1 \frac{L_{k,d}(te_d + (1-t^2)^{1/2}\xi_{(d-1)})}{|\mathbb{S}^{d-1}|} (1-t^2)^{\frac{d-3}{2}} dt d\lambda_{(d-2)}(\xi_{(d-1)}) \\ &= \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 P_{k,d}(t) (1-t^2)^{\frac{d-3}{2}} dt = 0 \end{aligned} \quad (29)$$

The right-hand side is equal to zero because $P_{0,d} \equiv 1$, and the Legendre polynomials are orthogonal with respect to the scalar product with factor $(1-t^2)^{\frac{d-3}{2}}$ ([Sec. 5](#)). Equation (29) implies that $\int_{C_+} \frac{L_{k,d}(x)}{|\mathbb{S}^{d-1}|} d\lambda(x) = -\int_{C_-} \frac{L_{k,d}(x)}{|\mathbb{S}^{d-1}|} d\lambda(x)$, which means that conditions (28) are feasible. We also have that

$$\gamma_{k,d} = 2 \left(\int_{\mathbb{S}^{d-1}} \frac{|L_{k,d}(x)|}{|\mathbb{S}^{d-1}|} d\lambda(x) \right)^{-1} = 2 \left(\int_{\mathbb{S}^{d-1}} |L_{k,d}(x)| d\tau(x) \right)^{-1}.$$

□

Lemma 5. For μ_d, ν_d with densities given by (4), we have

$$d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d) = \gamma_{k,d} \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \frac{1}{\sqrt{N_{k,d}}} \left| \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right|$$

Proof. By [Lemma 2](#), we have

$$\begin{aligned} d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu_d, \nu_d) &= \int_{\mathbb{S}^{d-1}} \left(\int_{\mathbb{S}^{d-1}} \sigma(\langle x, \theta \rangle) d(\mu_d - \nu_d)(x) \right)^2 d\tau(\theta) \\ &= \int_{\mathbb{S}^{d-1}} \left(\int_{\mathbb{S}^{d-1}} \sigma(\langle x, \theta \rangle) \frac{\gamma_{k,d}}{|\mathbb{S}^{d-1}|} L_{k,d}(x) d\lambda(x) \right)^2 d\tau(\theta) \\ &= \gamma_{k,d}^2 \int_{\mathbb{S}^{d-1}} \left(\int_{\mathbb{S}^{d-1}} \sigma(\langle x, \theta \rangle) L_{k,d}(x) d\tau(x) \right)^2 d\tau(\theta) \end{aligned} \quad (30)$$

Now, we reproduce the argument of [Bach \(2017\)](#). If we define $g(\theta) = \int_{\mathbb{S}^{d-1}} \sigma(\langle x, \theta \rangle) L_{k,d}(x) d\tau(x)$, since $L_{k,d}(x)$ is a homogeneous harmonic polynomial of degree d , by the Funk-Hecke formula we can write

$$g(\theta) = \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} L_{k,d}(\theta) \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt = \lambda_{k,d} L_{k,d}(\theta), \quad (31)$$

where $\lambda_{k,d} = \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt$. Note as well that

$$\begin{aligned} &\int_{\mathbb{S}^{d-1}} L_{k,d}(\theta)^2 d\tau(\theta) \\ &= \frac{1}{|\mathbb{S}^{d-1}|} \int_{\mathbb{S}^{d-2}} \int_{-1}^1 L_{k,d}(te_d + (1-t^2)^{1/2}\xi_{(d-1)})^2 (1-t^2)^{\frac{d-3}{2}} dt d\lambda_{(d-2)}(\xi_{(d-1)}) \\ &= \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 P_{k,d}(t)^2 (1-t^2)^{\frac{d-3}{2}} dt = \frac{1}{N_{k,d}} \end{aligned} \quad (32)$$

In the first equality we used the same change of variables as in the proof of [Lemma 1](#), in the second equality we used that $P_{k,d}(t) = L_{k,d}(te_d + (1-t^2)^{1/2}\xi_{(d-1)})$ by definition, and the third equality

relies on equation (24). Using (31) and (32), the right-hand side of (30) becomes:

$$\begin{aligned} (\gamma_{k,d} \lambda_{k,d})^2 \int_{\mathbb{S}^{d-1}} L_{k,d}(\theta)^2 d\tau(\theta) &= (\gamma_{k,d} \lambda_{k,d})^2 \frac{1}{N_{k,d}} \\ &= \gamma_{k,d}^2 \left(\frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right)^2 \frac{1}{N_{k,d}} \end{aligned}$$

□

Lemma 6. For μ_d, ν_d with densities given by (4), we have

$$d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) = |\gamma_{k,d}| \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \left| \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right|$$

Proof. Using Lemma 1, we have

$$\begin{aligned} d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) &= \sup_{\theta \in \mathbb{S}^d} \left| \int \sigma(\langle x, \theta \rangle) d(\mu_d - \nu_d)(x) \right| = \sup_{\theta \in \mathbb{S}^d} \left| \int \sigma(\langle x, \theta \rangle) \frac{\gamma_{k,d}}{|\mathbb{S}^{d-1}|} L_{k,d}(x) d\lambda(x) \right| \\ &= \gamma_{k,d} \sup_{\theta \in \mathbb{S}^d} \left| \int \sigma(\langle x, \theta \rangle) L_{k,d}(x) d\tau(x) \right| \\ &= \gamma_{k,d} |\lambda_{k,d}| \sup_{\theta \in \mathbb{S}^d} |L_{k,d}(\theta)| = |\gamma_{k,d}| |\lambda_{k,d}| \\ &= |\gamma_{k,d}| \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \left| \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right|, \end{aligned}$$

In the fourth equality we used (31), in the fifth equality we used $\sup_{\theta \in \mathbb{S}^d} |L_{k,d}(\theta)| = 1$ by equation (26), and in the sixth equality we used the definition of $\lambda_{k,d}$. □

Theorem 2. Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an activation function that is bounded in $[-1, 1]$. For any $d \geq 2$ and $k \geq 1$, if we set $\gamma_{k,d}$ as in Proposition 1 we have that

$$d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) = \frac{2 \left| \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right|}{\int_{-1}^1 |P_{k,d}(t)| (1-t^2)^{\frac{d-3}{2}} dt}, \quad (5)$$

and

$$\frac{d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)}{d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)} = \sqrt{N_{k,d}} = \sqrt{\frac{(2k+d-2)(k+d-3)!}{k!(d-2)!}}, \quad (6)$$

where $N_{k,d}$ is the dimension of the space of spherical harmonics of order k over \mathbb{S}^{d-1} . That is,

$$\log \left(\frac{d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)}{d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)} \right) = \frac{1}{2} \left(k \log \left(\frac{k+d-3}{k} \right) + (d-2) \log \left(\frac{k+d-3}{d-2} \right) \right) + O(\log(k+d)). \quad (7)$$

Proof. Plugging the results of Lemma 5 and Lemma 6, we obtain

$$\frac{d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)}{d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)} = \frac{1}{\frac{1}{\sqrt{N_{k,d}}}} = \sqrt{N_{k,d}} = \sqrt{\frac{(2k+d-2)(k+d-3)!}{k!(d-2)!}}$$

The last equality follows from (25). Equation (7) follows from Stirling's approximation, which states that $\log n! = n \log n - n + O(\log n)$ and $\log(\Gamma(x)) = x \log x - x + O(\log x)$.

All that is left is checking that (5) holds. [Proposition 1](#) states that $\gamma_{k,d} = 2 \left(\int_{\mathbb{S}^{d-1}} |L_{k,d}(x)| d\tau(x) \right)^{-1} = 2 \left(\frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 |P_{k,d}(t)| (1-t^2)^{\frac{d-3}{2}} dt \right)^{-1}$ for μ_d and ν_d to be probability measures. Thus, [Lemma 6](#) implies that

$$d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) = \frac{2 \left| \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right|}{\int_{-1}^1 |P_{k,d}(t)| (1-t^2)^{\frac{d-3}{2}} dt}.$$

□

E Proofs of [Sec. 6](#)

Lemma 7. *Let ∇ denote the Riemannian gradient. We have that*

$$\int_{\mathbb{S}^d} \|\nabla L_{k,d}(x)\|^2 d\tau(x) = k(k+d-2) \frac{1}{N_{k,d}}$$

Proof. Through integration by parts, we have that

$$\int_{\mathbb{S}^d} \|\nabla L_{k,d}(x)\|^2 d\tau(x) = \int_{\mathbb{S}^d} L_{k,d}(x) (-\Delta) L_{k,d}(x) d\tau(x), \quad (33)$$

where Δ denotes the Laplace-Beltrami operator. Since the restriction of $L_{k,d}(x)$ to \mathbb{S}^{d-1} is a k -spherical harmonic and spherical harmonics are eigenfunctions of the Laplace-Beltrami operator (equation (3.19) of [Atkinson and Han \(2012\)](#)), we have

$$-\Delta L_{k,d}(x) = k(k+d-2) L_{k,d}(x).$$

Plugging this into (33) and using equality (24), we obtain

$$\int_{\mathbb{S}^d} \|\nabla L_{k,d}(x)\|^2 d\tau(x) = k(k+d-2) \int_{\mathbb{S}^d} L_{k,d}(x)^2 d\tau(x) = k(k+d-2) \frac{1}{N_{k,d}}.$$

□

Lemma 8. *Let $\hat{\nabla} L_{k,d}(x)$ be the Euclidean gradient of $L_{k,d} : \mathbb{R}^d \rightarrow \mathbb{R}$ and $\nabla L_{k,d}(x)$ be the Riemannian gradient of $L_{k,d} : \mathbb{S}^{d-1} \rightarrow \mathbb{R}$. Then,*

$$\hat{\nabla} L_{k,d}(x) = \nabla L_{k,d}(x) + k L_{k,d}(x) x.$$

Proof. By definition, for any $x \in \mathbb{S}^{d-1}$, $\nabla L_{k,d}(x)$ is the projection of $\hat{\nabla} L_{k,d}(x)$ to \mathbb{S}^{d-1} to $T_x \mathbb{S}^{d-1}$. That is,

$$\nabla L_{k,d}(x) = \hat{\nabla} L_{k,d}(x) - \langle \hat{\nabla} L_{k,d}(x), x \rangle x = \hat{\nabla} L_{k,d}(x) - \frac{\partial}{\partial r} \hat{\nabla} L_{k,d}(rx) \Big|_{r=1} x = \hat{\nabla} L_{k,d}(x) - k L_{k,d}(x) x.$$

In the last equality we used that $\frac{\partial}{\partial r} \hat{\nabla} L_{k,d}(rx) = \frac{k}{r} L_{k,d}(rx)$, which holds because $L_{k,d}$ is a homogeneous polynomial of degree k . □

Lemma 9. *Let $\hat{\nabla} L_{k,d}(x)$ be the Euclidean gradient of $L_{k,d} : \mathbb{R}^d \rightarrow \mathbb{R}$. Each component of $\hat{\nabla} L_{k,d}(x)$ is a $(k-1)$ -th spherical harmonic when restricted to \mathbb{S}^{d-1} .*

Proof. Spherical harmonics of degree k in dimension d can be characterized as the restrictions in \mathbb{S}^{d-1} of homogeneous harmonic polynomials of degree k in \mathbb{R}^d ([Atkinson and Han, 2012](#)), and harmonic functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ are those such that $\Delta f = \sum_{i=1}^d \partial_{ii} f = 0$. Notice that the i -th partial derivative of a homogeneous harmonic polynomial p of degree k is a homogeneous harmonic polynomial of degree $k-1$. That is because (i) the derivative of a homogeneous polynomial of

degree k is a homogeneous polynomial of degree $k - 1$ and (ii) by commutation of partial derivatives, we have

$$\Delta(\partial_i p) = \sum_{j=1}^{d+1} \partial_{jj} \partial_i p = \sum_{j=1}^{d+1} \partial_i \partial_{jj} p = \partial_i (\Delta p) = 0.$$

Thus, the restriction of $\partial_i p$ to \mathbb{S}^{d-1} is a $(k - 1)$ -th spherical harmonic. \square

Lemma 10. *let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an α -positive homogeneous activation function of the form (3).*

$$SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d) \geq |a + (-1)^{k+1} b| \gamma_{k,d} \lambda_{k,d}^{(\alpha+1)} \frac{k(d+k-3)}{\alpha+1}, \quad (34)$$

where

$$\lambda_{k,d}^{(\alpha+1)} = \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 P_{k,d}(t) (t)_+^{\alpha+1} (1-t^2)^{\frac{d-3}{2}} dt$$

Proof. For simplicity, we begin by considering the case $\sigma(x) = (x)_+^\alpha$. Remark that $\mu_d = \tau$, i.e., the uniform Borel probability measure over \mathbb{S}^{d-1} . We have

$$\begin{aligned} SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d) &= \sup_{h \in \mathcal{B}_{\mathcal{F}_1^d}} \mathbb{E}_{\mu_d} [\text{Tr}(\mathcal{A}_{\nu_d} h(x))] \\ &= \sup_{h \in \mathcal{B}_{\mathcal{F}_1^d}} \mathbb{E}_{\mu_d} [\text{Tr}(\mathcal{A}_{\nu_d} h(x)) - \text{Tr}(\mathcal{A}_{\mu_d} h(x))] \\ &= \sup_{h \in \mathcal{B}_{\mathcal{F}_1^d}} \mathbb{E}_{\mu_d} \left[\left(\nabla \log \left(\frac{d\nu_d}{d\tau}(x) \right) - \nabla \log \left(\frac{d\mu_d}{d\tau}(x) \right) \right)^\top h(x) \right] \\ &= \gamma_{k,d} \sup_{h \in \mathcal{B}_{\mathcal{F}_1^d}} \mathbb{E}_{\mu_d} [\nabla L_{k,d}(x)^\top h(x)] \\ &= \gamma_{k,d} \sup_{\substack{\|\mu_i\|_{\text{TV}} \leq 1 \\ \sum_i z_i^2 = 1}} \sum_{i=1}^d z_i \int_{\mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) \int_{\mathbb{S}^{d-1}} (\langle \theta, x \rangle)_+^\alpha d\mu_i(\theta) d\tau(x) \\ &= \gamma_{k,d} \sup_{\substack{\theta^{(i)} \in \mathbb{S}^{d-1} \\ \sum_i z_i^2 = 1}} \sum_{i=1}^d z_i \left| \int_{\mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) (\langle \theta^{(i)}, x \rangle)_+^\alpha d\tau(x) \right| \\ &= \gamma_{k,d} \sqrt{\sum_{i=1}^d \sup_{\theta^{(i)} \in \mathbb{S}^{d-1}} \left(\int_{\mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) (\langle \theta^{(i)}, x \rangle)_+^\alpha d\tau(x) \right)^2}. \end{aligned} \quad (35)$$

In the second equality, we have applied the Stein identity. The third equality relies on the definition of the Stein operator (equation (8)). In the fourth equality, we used that μ_d is uniform and ν_d has density given by (10). In the sixth equality we have used that for any function f and domain K , the supremum of $\int_K f d\mu$ over signed measures with total variation norm bounded by 1 is equal to $\sup_K f$. In the seventh equality we have used the Cauchy-Schwarz inequality. At this point, notice that by Lemma 8

$$\begin{aligned} \int_{\mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) (\langle \theta, x \rangle)_+^\alpha d\tau(x) &= \left(\int_{\mathbb{S}^{d-1}} \nabla L_{k,d}(x) (\langle \theta, x \rangle)_+^\alpha d\tau(x) \right)_i \\ &= \left(\int_{\mathbb{S}^{d-1}} (\hat{\nabla} L_{k,d}(x) - k L_{k,d}(x) x) (\langle \theta, x \rangle)_+^\alpha d\tau(x) \right)_i \end{aligned} \quad (36)$$

On the one hand, by the Funk-Hecke formula, since $\hat{\nabla} L_{k,d}(x)$ is a $(k - 1)$ -th spherical harmonic (Lemma 9), we have that for any $\theta \in \mathbb{S}^{d-1}$,

$$\begin{aligned} \int_{\mathbb{S}^{d-1}} \hat{\nabla} L_{k,d}(x) (\langle \theta, x \rangle)_+^\alpha d\tau(x) &= \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \hat{\nabla} L_{k,d}(\theta) \int_{-1}^1 P_{k-1,d}(t) (t)_+^\alpha (1-t^2)^{\frac{d-3}{2}} dt \\ &= \lambda_{k-1,d}^{(\alpha)} \hat{\nabla} L_{k,d}(\theta) \end{aligned} \quad (37)$$

where $\lambda_{k-1,d}^{(\alpha)}$ is defined accordingly.

On the other hand, since $\hat{\nabla}_\theta((\langle \theta, x \rangle)_+^{\alpha+1}) = (\alpha+1)(\langle \theta, x \rangle)_+^\alpha x$ for $\theta \in \mathbb{R}^d$, we have that for any $\theta \in \mathbb{R}^d$,

$$\begin{aligned}
& \int_{\mathbb{S}^{d-1}} L_{k,d}(x) x (\langle \theta, x \rangle)_+^\alpha d\tau(x) \\
&= \frac{1}{\alpha+1} \int_{\mathbb{S}^{d-1}} L_{k,d}(x) \hat{\nabla}_\theta((\langle \theta, x \rangle)_+^{\alpha+1}) d\tau(x) \\
&= \frac{1}{\alpha+1} \hat{\nabla}_\theta \int_{\mathbb{S}^{d-1}} L_{k,d}(x) (\langle \theta, x \rangle)_+^{\alpha+1} d\tau(x) \\
&= \frac{1}{\alpha+1} \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \hat{\nabla}_\theta (L_{k,d}(\theta) \|\theta\|^{\alpha+1-k}) \int_{-1}^1 P_{k,d}(t) (t)_+^{\alpha+1} (1-t^2)^{\frac{d-3}{2}} dt \\
&= \frac{1}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \hat{\nabla}_\theta (L_{k,d}(\theta) \|\theta\|^{\alpha+1-k}),
\end{aligned} \tag{38}$$

In the third equality, we used the Funk-Hecke formula, which says that for any $\theta \in \mathbb{S}^{d-1}$, $\int_{\mathbb{S}^{d-1}} L_{k,d}(x) (\langle \theta, x \rangle)_+^{\alpha+1} d\tau(x) = L_{k,d}(\theta) \int_{-1}^1 P_{k,d}(t) (t)_+^{\alpha+1} (1-t^2)^{\frac{d-3}{2}} dt$. To obtain the equality for a general $\theta \in \mathbb{R}^d$, we must use add the factor $\|\theta\|^{\alpha+1-k}$ so that the two sides have the same homogeneity parameter, yielding $\int_{\mathbb{S}^{d-1}} L_{k,d}(x) (\langle \theta, x \rangle)_+^{\alpha+1} d\tau(x) = L_{k,d}(\theta) \|\theta\|^{\alpha+1-k} \int_{-1}^1 P_{k,d}(t) (t)_+^{\alpha+1} (1-t^2)^{\frac{d-3}{2}} dt$.

And for $\theta \in \mathbb{S}^{d-1}$, we have

$$\begin{aligned}
\hat{\nabla}_\theta (L_{k,d}(\theta) \|\theta\|^{\alpha+1-k}) &= \hat{\nabla}_\theta L_{k,d}(\theta) \|\theta\|^{\alpha+1-k} + (\alpha+1-k) L_{k,d}(\theta) \|\theta\|^{\alpha-1-k} \theta \\
&= \hat{\nabla}_\theta L_{k,d}(\theta) + (\alpha+1-k) L_{k,d}(\theta) \theta
\end{aligned} \tag{39}$$

Thus, the right-hand side of (35) can be developed as

$$\gamma_{k,d} \sqrt{\sum_{i=1}^d \sup_{\theta^{(i)}} \left(\left(\lambda_{k-1,d}^{(\alpha)} - \frac{k}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \right) \hat{\nabla}_i L_{k,d}(\theta^{(i)}) - \frac{k(\alpha+1-k)}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} L_{k,d}(\theta^{(i)}) \theta_i^{(i)} \right)^2} \tag{40}$$

Bach (2017) (App. D.2) shows the following equality

$$\begin{aligned}
& \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 P_{k,d}(t) (t)_+^\alpha (1-t^2)^{\frac{d-3}{2}} dt \\
&= \begin{cases} 0 & \text{if } k \equiv \alpha \pmod{2}, k > \alpha \\ \frac{\Gamma(d/2)}{\sqrt{\pi}\Gamma((d-1)/2)} \frac{\alpha!(-1)^{(k-1-\alpha)/2}}{2^k} \frac{\Gamma(\frac{d-1}{2})\Gamma(k-\alpha)}{\Gamma(\frac{k}{2}-\frac{\alpha}{2}+\frac{1}{2})\Gamma(\frac{k}{2}+\frac{\alpha}{2}+\frac{\alpha}{2})} & \text{if } k \not\equiv \alpha \pmod{2}, k \geq \alpha+1 \end{cases} \tag{41}
\end{aligned}$$

Notice that in Bach (2017) the factor $\frac{d-1}{2\pi}$ is a typo, and should instead be $\frac{\Gamma(d/2)}{\sqrt{\pi}\Gamma((d-1)/2)}$. Using equality (41), we get that when $k \not\equiv \alpha \pmod{2}$ and $k \geq \alpha+1$,

$$\begin{aligned}
\lambda_{k,d}^{(\alpha)} &= \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} \int_{-1}^1 P_{k,d}(t) (t)_+^\alpha (1-t^2)^{\frac{d-3}{2}} dt \\
&= \frac{\Gamma(d/2)}{\sqrt{\pi}\Gamma((d-1)/2)} \frac{\alpha!(-1)^{(k-\alpha-1)/2}}{2^k} \frac{\Gamma((d-1)/2)\Gamma(k-\alpha)}{\Gamma(\frac{k-\alpha+1}{2})\Gamma(\frac{k+d+\alpha-1}{2})},
\end{aligned}$$

and $\lambda_{k,d}^{(\alpha)} = 0$ otherwise. Thus,

$$\frac{\lambda_{k-1,d}^{(\alpha)}}{\frac{k}{\alpha+1} \lambda_{k,d}^{(\alpha+1)}} = \frac{\alpha+1}{k} \frac{\frac{\alpha!}{2^{k-1}}}{\frac{(\alpha+1)!}{2^k}} \frac{\frac{1}{\Gamma(\frac{k+d+\alpha-2}{2})}}{\frac{1}{\Gamma(\frac{k+d+\alpha}{2})}} = \frac{k+d+\alpha-2}{k}$$

Hence, the arguments of the suprema in (40) can be rewritten as

$$\begin{aligned} & \frac{k}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \left(\left(\frac{k+d+\alpha-2}{k} - 1 \right) \hat{\nabla}_i L_{k,d}(\theta^{(i)}) - (\alpha+1-k) L_{k,d}(\theta^{(i)}) \theta_i^{(i)} \right) \\ &= \frac{k}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \left(\frac{d+\alpha-2}{k} (\nabla_i L_{k,d}(\theta^{(i)}) + k L_{k,d}(\theta^{(i)}) \theta_i^{(i)}) - (\alpha+1-k) L_{k,d}(\theta^{(i)}) \theta_i^{(i)} \right) \end{aligned} \quad (42)$$

If we substitute $i = d$ and $\theta^{(i)} = e_d$ in this expression, and use that $\nabla_d L_{k,d}(e_d) = 0$ (by the fact that $\nabla L_{k,d}(e_d) \in T_{e_d} \mathbb{S}^{d-1}$) and $L_{k,d}(e_d) = 1$ we obtain

$$\lambda_{k,d}^{(\alpha+1)} \left(\frac{(d+\alpha-2)k}{\alpha+1} + \frac{k(k-\alpha-1)}{\alpha+1} \right) = \lambda_{k,d}^{(\alpha+1)} \frac{k(d+k-3)}{\alpha+1}, \quad (43)$$

which means that (40) is lower-bounded by $\gamma_{k,d} \lambda_{k,d}^{(\alpha+1)} \frac{k(d+k-3)}{\alpha+1}$.

When $\sigma(x) = (-x)_+^\alpha$, we reproduce the same argument. Equation (37) becomes $\int_{\mathbb{S}^{d-1}} \hat{\nabla} L_{k,d}(x) (-\langle \theta, x \rangle)_+^\alpha d\tau(x) = \lambda_{k-1,d}^{(\alpha)} \hat{\nabla} L_{k,d}(-\theta) = (-1)^{k+1} \lambda_{k-1,d}^{(\alpha)} \hat{\nabla} L_{k,d}(\theta)$, where we have used that $L_{k,d}(-\theta) = (-1)^k L_{k,d}(\theta)$. Since $\hat{\nabla}_\theta((-\langle \theta, x \rangle)_+^{\alpha+1}) = -(\alpha+1)(-\langle \theta, x \rangle)_+^\alpha x$, equation (38) becomes $\int_{\mathbb{S}^{d-1}} L_{k,d}(x) x (-\langle \theta, x \rangle)_+^\alpha d\tau(x) = -\frac{1}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \hat{\nabla}_\theta(L_{k,d}(-\theta)) - \theta \|\theta\|^{\alpha+1-k}$. Since $L_{k,d}(-\theta) = (-1)^k L_{k,d}(\theta)$, we have that $(\hat{\nabla} L_{k,d})(-\theta) = (-1)^{k+1} \hat{\nabla} L_{k,d}(\theta)$. Thus, equation (39) becomes

$$\begin{aligned} -\hat{\nabla}_\theta(L_{k,d}(-\theta)) - \theta \|\theta\|^{\alpha+1-k} &= -\hat{\nabla}_\theta(L_{k,d}(-\theta)) - \theta \|\theta\|^{\alpha+1-k} - \hat{\nabla}_\theta(\|\theta\|^{\alpha+1-k}) L_{k,d}(-\theta) \\ &= (\hat{\nabla}_\theta L_{k,d})(-\theta) - \theta \|\theta\|^{\alpha+1-k} - (\alpha+1-k) \|\theta\|^{\alpha-1-k} \theta L_{k,d}(-\theta) \\ &= (-1)^{k+1} \left(\hat{\nabla} L_{k,d}(\theta) + (\alpha+1-k) \theta L_{k,d}(\theta) \right) \end{aligned}$$

Hence, for $\sigma(x) = (-x)_+^\alpha$ the expression (40) is unchanged, and the rest of the argument holds in the same way. When $\sigma(x) = a(x)_+^\alpha + b(-x)_+^\alpha$, the argument of the square root in expression (40) gets multiplied by $|a + (-1)^{k+1}b|$, and this factor is carried over for the rest of the argument. This concludes the proof. \square

Lemma 11. *Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an α -positive homogeneous activation function of the form (3). Then, $SD_{\mathcal{B}_{\mathbb{F}_2^d}}(\mu_d, \nu_d)$ is upper-bounded by*

$$|a + (-1)^{k+1}b| \gamma_{k,d} \lambda_{k,d}^{(\alpha+1)} \sqrt{\frac{2}{N_{k,d}} \left(k(k+d-2) \left(\frac{d+\alpha-2}{\alpha+1} \right)^2 + \left(\frac{k(d+k-3)}{\alpha+1} \right)^2 \right)}.$$

Proof. For simplicity, we begin by considering the case $\sigma(x) = (x)_+^\alpha$.

$$\begin{aligned}
\text{SD}_{\mathcal{B}_{\mathcal{F}_2^d}}(\mu_d, \nu_d) &= \sup_{h \in \mathcal{B}_{\mathcal{F}_2^d}} \mathbb{E}_{\mu_d} [\text{Tr}(\mathcal{A}_{\nu_d} h(x))] \\
&= \sup_{h \in \mathcal{B}_{\mathcal{F}_2^d}} \mathbb{E}_{\mu_d} [\text{Tr}(\mathcal{A}_{\nu_d} h(x)) - \text{Tr}(\mathcal{A}_{\mu_d} h(x))] \\
&= \sup_{h \in \mathcal{B}_{\mathcal{F}_2^d}} \mathbb{E}_{\mu_d} \left[\left(\nabla \log \left(\frac{d\nu_d}{d\tau}(x) \right) - \nabla \log \left(\frac{d\mu_d}{d\tau}(x) \right) \right)^\top h(x) \right] \\
&= \gamma_{k,d} \sup_{h \in \mathcal{B}_{\mathcal{F}_2^d}} \mathbb{E}_{\mu_d} [\nabla L_{k,d}(x)^\top h(x)] = \gamma_{k,d} \sup_{\substack{\|h_i\|_{\mathcal{F}_2} \leq 1 \\ \sum_i z_i^2 = 1}} z_i \mathbb{E}_{\mu_d} [\nabla_i L_{k,d}(x) \langle k(x, \cdot), h_i \rangle_{\mathcal{F}_2}] \\
&= \gamma_{k,d} \sup_{\substack{\|h_i\|_{\mathcal{F}_2} \leq 1 \\ \sum_i z_i^2 = 1}} z_i \left\langle \mathbb{E}_{\mu_d} [\nabla_i L_{k,d}(x) k(x, \cdot)], h_i \right\rangle_{\mathcal{F}_2} \\
&= \gamma_{k,d} \sqrt{\sum_{i=1}^d \left\| \mathbb{E}_{\mu_d} [\nabla_i L_{k,d}(x) k(x, \cdot)] \right\|_{\mathcal{F}_2}^2}.
\end{aligned}$$

And we can rewrite the right-hand side as

$$\begin{aligned}
&\gamma_{k,d} \sqrt{\sum_{i=1}^d \iint_{\mathbb{S}^{d-1} \times \mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) k(x, y) \nabla_i L_{k,d}(y) d\tau(x) d\tau(y)} \\
&= \gamma_{k,d} \sqrt{\sum_{i=1}^d \int_{\mathbb{S}^{d-1}} \left(\int_{\mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) (\langle x, \theta \rangle)_+^\alpha d\tau(x) \right)^2 d\tau(\theta)}. \tag{44}
\end{aligned}$$

At this point, we express $\int_{\mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) (\langle x, \theta \rangle)_+^\alpha d\tau(x)$ using the development in equations (36), (37), (38), (39):

$$\begin{aligned}
&\left(\int_{\mathbb{S}^{d-1}} \nabla_i L_{k,d}(x) (\langle \theta_i, x \rangle)_+^\alpha d\tau(x) \right)^2 \\
&= \left(\frac{k}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \right)^2 \left(\left(\frac{k+d+\alpha-2}{k} - 1 \right) \hat{\nabla}_i L_{k,d}(\theta) - (\alpha+1-k) L_{k,d}(\theta) \theta_i \right)^2 \\
&= \left(\frac{k}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \right)^2 \left(\frac{d+\alpha-2}{k} (\nabla_i L_{k,d}(\theta) + k L_{k,d}(\theta) \theta_i) - (\alpha+1-k) L_{k,d}(\theta) \theta_i \right)^2 \\
&\leq 2A^2 (\nabla_i L_{k,d}(\theta))^2 + 2B^2 (L_{k,d}(\theta) \theta_i)^2,
\end{aligned}$$

where, using the computations in the proof of Lemma 10,

$$\begin{aligned}
A &= \lambda_{k,d}^{(\alpha+1)} \frac{d+\alpha-2}{\alpha+1}, \\
B &= \lambda_{k,d}^{(\alpha+1)} \left(\frac{(d+\alpha-2)k}{\alpha+1} + \frac{k(k-\alpha-1)}{\alpha+1} \right) = \lambda_{k,d}^{(\alpha+1)} \frac{k(d+k-3)}{\alpha+1}.
\end{aligned}$$

Thus, the right-hand side of (44) is upper-bounded by:

$$\gamma_{k,d} \sqrt{2A^2 \int_{\mathbb{S}^{d-1}} \|\nabla L_{k,d}(\theta)\|^2 d\tau(\theta) + 2B^2 \int_{\mathbb{S}^{d-1}} L_{k,d}(\theta)^2 d\tau(\theta)} \tag{45}$$

We can use Lemma 7 to compute the first integral: $\int_{\mathbb{S}^d} \|\nabla L_{k,d}(x)\|^2 d\tau(x) = k(k+d-2) \frac{1}{N_{k,d}}$. And for the second integral we have $\int_{\mathbb{S}^{d-1}} L_{k,d}(\theta)^2 d\tau(\theta) = \frac{1}{N_{k,d}}$ by equation (32). Substituting

everything into (45) yields

$$\gamma_{k,d} \lambda_{k,d}^{(\alpha+1)} \sqrt{\frac{2}{N_{k,d}} \left(k(k+d-2) \left(\frac{d+\alpha-2}{\alpha+1} \right)^2 + \left(\frac{k(d+k-3)}{\alpha+1} \right)^2 \right)}. \quad (46)$$

For the general case $\sigma(x) = a(x)_+^\alpha + b(-x)_+^\alpha$, we use arguments analogous to those of Lemma 10, and we obtain that the upper-bound (46) gets multiplied by a factor $|a + (-1)^{k+1}b|$. \square

Theorem 3. *Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an α -positive homogeneous activation function of the form (3) such that $a + (-1)^{k+1}b \neq 0$. For all $k \geq 1$, $d \geq 2$ we can choose $\gamma_{k,d} \in [-1, 1]$ such that $SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d) = 1$ and*

$$\frac{SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d)}{SD_{\mathcal{B}_{\mathcal{F}_2^d}}(\mu_d, \nu_d)} \geq \frac{\frac{k(d+k-3)}{\alpha+1}}{\sqrt{\frac{2}{N_{k,d}} \left(k(k+d-2) \left(\frac{d+\alpha-2}{\alpha+1} \right)^2 + \left(\frac{k(d+k-3)}{\alpha+1} \right)^2 \right)}} \quad (11)$$

That is,

$$\log \left(\frac{SD_{\mathcal{B}_{\mathcal{F}_1^d}}(\mu_d, \nu_d)}{SD_{\mathcal{B}_{\mathcal{F}_2^d}}(\mu_d, \nu_d)} \right) \geq \frac{1}{2} \left(k \log \left(\frac{k+d-3}{k} \right) + (d-2) \log \left(\frac{k+d-3}{d-2} \right) \right) + O(\log(k+d)) \quad (12)$$

Proof. We obtain (11) from Lemma 10 and Lemma 11. Taking the logarithm and using Stirling's approximation yields (12). The only relevant factor is $\log(\sqrt{N_{k,d}})$, as the other ones are $O(\log(k+d))$. \square

F Proofs of Sec. 7

Lemma 12 (Approximation of Lipschitz-continuous functions on the unit ball by \mathcal{F}_2 functions, Bach (2017)). *Let $\sigma(x) = (x)_+^\alpha$ be the α -th power of the ReLU activation function, where α is a non-negative integer. For δ greater than a constant depending only on d , for any function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ such that for all x, y such that for any $\|x\|_q \leq R$, $\|y\|_q \leq R$ we have $|f(x)| \leq \eta$ and $|f(x) - f(y)| \leq \eta R^{-1} \|x - y\|_q$, there exists $h \in \mathcal{F}_2(\mathbb{R}^d \times \{R\})$, such that $\|h\|_{\mathcal{F}_2} \leq \delta$ and*

$$\sup_{\|x\|_q \leq R} |h(x) - f(x)| \leq C(d, \alpha) \eta \left(\frac{R\delta}{\eta} \right)^{-\frac{1}{\alpha+(d-1)/2}} \log \left(\frac{R\delta}{\eta} \right)$$

Proof. See Proposition 6 of Bach (2017). Notice that in Bach (2017) the factor in the bound is $\left(\frac{\delta}{\eta} \right)^{-2/(d+1)} \log \left(\frac{\delta}{\eta} \right)$, while we have $\left(\frac{R\delta}{\eta} \right)^{-2/(d+1)} \log \left(\frac{R\delta}{\eta} \right)$. The R factor stems from the fact that we consider the neural network features to lie in \mathbb{S}^d , while Bach (2017) considers them in the hypersphere of radius R^{-1} . \square

Theorem 4. *Let $\delta > 0$ be larger than a certain constant depending on k and α . Let $\sigma(x) = (x)_+^\alpha$ be the α -th power of the ReLU activation function, where α is a non-negative integer. Let μ, ν be Borel probability measures with support included in $\{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\} \times \{1\}$. Let $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu)$ be as defined in (2) and $\overline{\mathcal{W}}_{1,k}$ as defined in (13). Then,*

$$\delta \overline{\mathcal{W}}_{1,k}(\mu, \nu) \geq \delta d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) \geq \overline{\mathcal{W}}_{1,k}(\mu, \nu) - 2C(k, \alpha) \delta^{-\frac{1}{\alpha+(k-1)/2}} \log(\delta), \quad (15)$$

where $C(k, \alpha)$ is a constant that depends only on k and α . If we optimize the lower bound in (15) with respect to δ , we obtain $\overline{\mathcal{W}}_{1,k}(\mu, \nu) \geq d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) \geq \tilde{\Omega}(\overline{\mathcal{W}}_{1,k}(\mu, \nu)^{\alpha + \frac{k+1}{2}})$ where $\tilde{\Omega}$ hides log factors.

Proof. We begin with the lower bound. Let U^* and f^* be the matrix in \mathcal{V}_k and function in $\text{Lip}_1(\mathbb{R}^k)$ where $\overline{\mathcal{W}}_{1,k}(\mu, \nu)$ is achieved (it is known in optimal transport that the supremum is achieved (Niles-Weed and Rigollet, 2019)), i.e.

$$\overline{\mathcal{W}}_{1,k}(\mu, \nu) = \mathbb{E}_{x \sim \mu}[f^*(U^*x)] - \mathbb{E}_{x \sim \nu}[f^*(U^*x)].$$

If μ (resp. ν) is supported in the closed unit ball in \mathbb{R}^d , $\forall x \in \text{supp}(\mu)$, $\|U^*x\|_2 \leq \|x\|_2 \leq 1$. Thus, all that matters is the restriction of f^* to the closed unit ball of \mathbb{R}^k . We can apply Lemma 12 with $R = 1, \eta = 1$, which yields the existence of $h \in \mathcal{F}_2(\mathbb{R}^k)$, such that $\|h\|_{\mathcal{F}_2} \leq \delta$ and

$$\sup_{x \in \mathbb{R}^k, \|x\|_2 \leq 1} |h(x) - f^*(x)| \leq C(k, \alpha) \delta^{-\frac{1}{\alpha + (k-1)/2}} \log(\delta),$$

when δ is larger than a constant depending on k and α . Thus,

$$\sup_{x \in \mathbb{R}^d, \|x\|_2 \leq 1} |h(U^*x) - f^*(U^*x)| \leq C(k, \alpha) \delta^{-\frac{1}{\alpha + (k-1)/2}} \log(\delta),$$

which implies that

$$|\overline{\mathcal{W}}_{1,k}(\mu, \nu) - (\mathbb{E}_{x \sim \mu}[h(U^*x)] - \mathbb{E}_{x \sim \nu}[h(U^*x)])| \leq 2C(k, \alpha) \delta^{-\frac{1}{\alpha + (k-1)/2}} \log(\delta).$$

Now, $h \circ U^*$ belongs to $\mathcal{F}_1(\mathbb{R}^d)$ by the argument of Section 4.6 of Bach (2017). Namely, if $h(x) = \int_{\mathbb{S}^k} (\langle \theta, (x, 1) \rangle)_+^\alpha d\mu_h(\theta)$, we can write

$$\begin{aligned} h(U^*x) &= \int_{\mathbb{S}^k} (\langle \theta, (U^*x, 1) \rangle)_+^\alpha d\mu_h(\theta) = \int_{\mathbb{S}^k} (\langle \theta_{1:k}, U^*x \rangle + \theta_{k+1})_+^\alpha d\mu_h(\theta) \\ &= \int_{\mathbb{S}^k} (\langle (U^*)^\top \theta_{1:k}, x \rangle + \theta_{k+1})_+^\alpha d\mu_h(\theta) = \int_{\mathbb{S}^k} (\langle ((U^*)^\top \theta_{1:k}, \theta_{k+1}), (x, 1) \rangle)_+^\alpha d\mu_h(\theta) \\ &= \int_{\mathbb{S}^k} (\langle \theta, (x, 1) \rangle)_+^\alpha d\tilde{\mu}_h(\theta), \end{aligned}$$

where $\tilde{\mu}_h$ is the pushforward of μ_h by the map $\theta \mapsto ((U^*)^\top \theta_{1:k}, \theta_{k+1})$. The last equality follows from the fact that $\|((U^*)^\top \theta_{1:k}, \theta_{k+1})\|_2^2 = \theta_{1:k}^\top U^* (U^*)^\top \theta_{1:k} + \theta_{k+1}^2 = \|\theta\|_2^2 = 1$. Moreover, this argument also shows that $h \circ U^*$ has \mathcal{F}_1 norm $\gamma_1(h \circ U^*) \leq \gamma_2(h) \leq \delta$. Hence, $\forall \mu, \nu \in \mathcal{P}(B_1(\mathbb{R}^d))$, for δ larger than a constant depending on k ,

$$\delta d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) \geq \overline{\mathcal{W}}_{1,k}(\mu, \nu) - 2C(k, \alpha) \delta^{-\frac{1}{\alpha + (k-1)/2}} \log(\delta).$$

The upper bound $\overline{\mathcal{W}}_{1,k}(\mu, \nu) \geq d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu)$ follows from

$$\begin{aligned} \overline{\mathcal{W}}_{1,k}(\mu, \nu) &= \max_{U \in \mathcal{V}_k} \sup_{f \in \text{Lip}_1(\mathbb{R}^k)} \mathbb{E}_{x \sim \mu}[f(Ux)] - \mathbb{E}_{x \sim \nu}[f(Ux)] \\ &\geq \max_{U \in \mathcal{V}_k} \sup_{f \in \mathcal{B}_{\mathcal{F}_1}(\mathbb{R}^k)} \mathbb{E}_{x \sim \mu}[f(Ux)] - \mathbb{E}_{x \sim \nu}[f(Ux)] \\ &= \sup_{f \in \mathcal{B}_{\mathcal{F}_1}(\mathbb{R}^k)} \mathbb{E}_{x \sim \mu}[f(x)] - \mathbb{E}_{x \sim \nu}[f(x)] = d_{\mathcal{B}_{\mathcal{F}_1}}(\mu, \nu) \end{aligned}$$

In the second to last inequality we used once again that for all $f \in \mathcal{F}_1(\mathbb{R}^k)$ such that $\|f\|_{\mathcal{F}_1(\mathbb{R}^k)} = 1$, we have $f \circ U \in \mathcal{F}_1(\mathbb{R}^d)$ and $\|f \circ U\|_{\mathcal{F}_1(\mathbb{R}^d)} = 1$. \square

Proposition 2. ($\tilde{\tau}$ as a rescaling of uniform measure) *The measure $d\tilde{\tau}(\sqrt{1-t^2}\xi, t) = \frac{1}{\pi}(1-t^2)^{-1/2} dt d\tau_{(d-1)}(\xi)$ is a probability measure. For comparison, the uniform measure over \mathbb{S}^d can be written as $d\tau(\sqrt{1-t^2}\xi, t) = \frac{\Gamma((d+1)/2)}{\sqrt{\pi}\Gamma(d/2)}(1-t^2)^{\frac{d-1}{2}} dt d\tau_{(d-1)}(\xi)$.*

Proof. The measure $d\tilde{\tau}(\sqrt{1-t^2}\xi, t) = \frac{1}{\pi}(1-t^2)^{-1/2} dt d\tau_{(d-1)}(\xi)$ is normalized because $\int_{\mathbb{S}^{d-1}} \int_{-1}^1 (1-t^2)^{-1/2} dt d\tau_{(d-1)}(\xi) = \int_{-1}^1 (1-t^2)^{-1/2} dt = \arcsin(1) - \arcsin(-1) = \pi$, where we used that $\int_{\mathbb{S}^{d-1}} d\tau_{(d-1)}(\xi) = 1$ by definition of $\tau_{(d-1)}$. The characterization of

the uniform measure τ follows from equation (1.17) of [Atkinson and Han \(2012\)](#): $d\tau(\theta) = \frac{|\mathbb{S}^{d-1}|}{|\mathbb{S}^d|} (1-t^2)^{\frac{d-1}{2}} dt d\tau_{(d-1)}(\xi) = \frac{\Gamma((d+1)/2)}{\sqrt{\pi}\Gamma(d/2)} (1-t^2)^{\frac{d-1}{2}} dt d\tau_{(d-1)}(\xi)$. For clarity, if we plug this change of variables into equation (1), we obtain that the F_2 kernel reads:

$$k(x, y) = \int_{\mathbb{S}^d} \sigma(\langle(x, 1), \theta\rangle) \sigma(\langle(y, 1), \theta\rangle) d\tau(\theta) = \frac{\Gamma((d+1)/2)}{\sqrt{\pi}\Gamma(d/2)} \int_{\mathbb{S}^{d-1}} \int_{-1}^1 \sigma(\langle(x, 1), (\sqrt{1-t^2}\xi_{(d)}, t)\rangle) \sigma(\langle(y, 1), (\sqrt{1-t^2}\xi, t)\rangle) (1-t^2)^{\frac{d-1}{2}} dt d\tau_{(d-1)}(\xi).$$

Notice that beyond the normalization factors, the main difference between \tilde{k} and k is the factor $(1-t^2)^{-1/2}$ instead of $(1-t^2)^{\frac{d-1}{2}}$. \square

Lemma 13. *Let $d_{\mathcal{B}_{\tilde{\mathcal{F}}_2}}$ be as defined in (16) and let $K = \{x \in \mathbb{R}^d | \|x\|_2 \leq 1\} \times \{1\}$. Then, for any $\mu, \nu \in \mathcal{P}(K)$, $d_{\mathcal{B}_{\tilde{\mathcal{F}}_2}}^2(\mu, \nu)$ is lower-bounded by*

$$\frac{1}{2\pi} \frac{5}{6\alpha 2^{\alpha/2}} \int_{\mathbb{S}^{d-1}} \sup_{\gamma \in [0, 2\pi]} \left| \int_K \left(\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle \right)_+^\alpha d(\mu - \nu)(x) \right|^3 d\tau_{(d-1)}(\xi_{(d)}).$$

Proof. Using the change of variables $t = \sin(\gamma)$, we have

$$\begin{aligned} d_{\mathcal{B}_{\tilde{\mathcal{F}}_2}}^2(\mu, \nu) &= \frac{1}{\pi} \int_{\mathbb{S}^{d-1}} \int_{-1}^1 \left(\int_K \left(\langle(x, 1), (\sqrt{1-t^2}\xi_{(d)}, t)\rangle \right)_+^\alpha d(\mu - \nu)(x, 1) \right)^2 (1-t^2)^{-1/2} dt d\tau_{(d-1)}(\xi_{(d)}) \\ &= \frac{1}{\pi} \int_{\mathbb{S}^{d-1}} \int_{-\pi/2}^{\pi/2} \left(\int_K \left(\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle \right)_+^\alpha d(\mu - \nu)(x, 1) \right)^2 d\gamma d\tau_{(d-1)}(\xi_{(d)}) \\ &= \frac{1}{2\pi} \int_{\mathbb{S}^{d-1}} \int_0^{2\pi} \left(\int_K \left(\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle \right)_+^\alpha d(\mu - \nu)(x, 1) \right)^2 d\gamma d\tau_{(d-1)}(\xi_{(d)}). \end{aligned} \tag{47}$$

We want to compute the Lipschitz constant of $\gamma \mapsto \int \left(\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle \right)_+^\alpha d(\mu - \nu)(x)$. For $\alpha \geq 1$, the derivative of this mapping is:

$$\int \alpha \left(\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle \right)_+^{\alpha-1} \langle(x, 1), (-\sin(\gamma)\xi_{(d)}, \cos(\gamma))\rangle d(\mu - \nu)(x, 1),$$

and its absolute value is upper-bounded by

$$\begin{aligned} 2\alpha |\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle|^{\alpha-1} |\langle(x, 1), (-\sin(\gamma)\xi_{(d)}, \cos(\gamma))\rangle| &\leq 2\alpha \|x, 1\|_2^{\alpha-1} \|x, 1\|_2 \\ &= 2\alpha \|x, 1\|_2^\alpha \leq \alpha 2^{\alpha/2}, \end{aligned}$$

where we used that $\|x\|_2 \leq 1$ for x in the support of μ or ν . Thus, if we denote $s = \sup_{\gamma \in [0, 2\pi]} \left| \int \left(\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle \right)_+^\alpha d(\mu - \nu)(x) \right|$, we have

$$\begin{aligned} \int_0^{2\pi} \left(\int \left(\langle(x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma))\rangle \right)_+^\alpha d(\mu - \nu)(x) \right)^2 d\gamma &\geq \int_0^{\frac{s}{\alpha 2^{\alpha/2}}} \left(s - \gamma \alpha 2^{\alpha/2} \right)^2 d\gamma \\ &= \int_0^{\frac{s}{\alpha 2^{\alpha/2}}} \left(s^2 - \gamma \alpha 2^{\alpha/2} s + \left(\gamma \alpha 2^{\alpha/2} \right)^2 \right) d\gamma = \frac{s^3}{\alpha 2^{\alpha/2}} - \frac{\alpha 2^{\alpha/2} s}{2} \left(\frac{s}{\alpha 2^{\alpha/2}} \right)^2 + \frac{\left(\alpha 2^{\alpha/2} \right)^2}{3} \left(\frac{s}{\alpha 2^{\alpha/2}} \right)^3 \\ &= \frac{5s^3}{6\alpha 2^{\alpha/2}}. \end{aligned}$$

Hence,

$$d_{\mathcal{B}_{\tilde{\mathcal{F}}_2}}^2(\mu, \nu) \geq \frac{1}{2\pi} \frac{5}{6\alpha 2^{\alpha/2}} \int_{\mathbb{S}^{d-1}} \sup_{\gamma \in [0, 2\pi]} \left| \int \left(\langle (x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma)) \rangle \right)_+^\alpha d(\mu - \nu)(x) \right|^3 d\tau_{(d-1)}(\xi_{(d)}).$$

□

Theorem 5. Let $\delta > 0$ be larger than a certain constant depending on k and α . Let $\sigma(x) = (x)_+^\alpha$ be the α -th power of the ReLu activation function, where α is a non-negative integer. Let μ, ν be Borel probability measures with support included in $\{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\} \times \{1\}$. Let $d_{\mathcal{B}_{\tilde{\mathcal{F}}_2}}$ be as defined in (16) and $\underline{\mathcal{W}}_{1,1}$ as defined in (14). Then,

$$\delta d_{\tilde{\mathcal{F}}_2}^{2/3}(\mu, \nu) \geq \left(\frac{5}{12\pi\alpha 2^{\alpha/2}} \right)^{1/3} \left(\underline{\mathcal{W}}_{1,1}(\mu, \nu) - 2C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta) \right). \quad (17)$$

and $\pi d_{\mathcal{B}_{\tilde{\mathcal{F}}_2}}^2(\mu, \nu) \leq \underline{\mathcal{W}}_{1,1}(\mu, \nu)$. If we optimize the lower bound in (17) with respect to δ , we obtain $d_{\tilde{\mathcal{F}}_2}^{2/3}(\mu, \nu) \geq \tilde{\Omega}(\underline{\mathcal{W}}_{1,1}(\mu, \nu)^{1+\alpha})$.

Proof. We begin with the lower bound (17). By the definition of the integral 1-dimensional projection robust Wasserstein distance and the fact that the Stiefel manifold for $k = 1$ is \mathcal{S}^{d-1} , we have

$$\underline{\mathcal{W}}_{1,1}(\mu, \nu) = \int_{\mathbb{S}^{d-1}} \mathcal{W}_1(u_\# \mu, u_\# \nu) d\tau(u),$$

where $u_\# \mu$ denotes the pushforward of μ by the map $\theta \mapsto \langle u, \theta \rangle$ and thus, $\mathcal{W}_1(u_\# \mu, u_\# \nu) = \min_{\pi \in \Gamma(\mu, \nu)} \int \|Ux - Uy\| d\pi(x, y)$. By the dual characterization of the 1-Wasserstein distance, for any $u \in \mathbb{S}^{d-1}$ we can write

$$\mathcal{W}_1(u_\# \mu, u_\# \nu) = \mathbb{E}_{x \sim \mu}[f_u^*(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[f_u^*(\langle u, x \rangle)]$$

for some function in $\text{Lip}_1(\mathbb{R})$. Using the same argument as in Theorem 4, Lemma 12 with $R = 1, \eta = 1$ yields the existence of $h_u \in \mathcal{F}_2(\mathbb{R})$ such that $\|h_u\|_{\mathcal{F}_2} \leq \delta$ and

$$\sup_{x \in \mathbb{R}, |x| \leq 1} |h_u(x) - f_u^*(x)| \leq C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta),$$

when δ is larger than a constant depending on k and α . Thus,

$$\sup_{x \in \mathbb{R}^d, \|x\|_2 \leq 1} |h_u(\langle u, x \rangle) - f_u^*(\langle u, x \rangle)| \leq C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta),$$

which implies that

$$\begin{aligned} & |\mathcal{W}_1(u_\# \mu, u_\# \nu) - (\mathbb{E}_{x \sim \mu}[h_u(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[h_u(\langle u, x \rangle)])| \leq 2C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta) \\ \implies & \mathbb{E}_{x \sim \mu}[h_u(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[h_u(\langle u, x \rangle)] \geq \mathcal{W}_1(u_\# \mu, u_\# \nu) - 2C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta). \end{aligned} \quad (48)$$

And since $h_u(y) = \int_{\mathbb{S}^1} \sigma(\langle \theta, (y, 1) \rangle) d\mu_{h_u}(\theta)$ for some $\mu_{h_u} \in \mathcal{M}(\mathbb{S}^1)$ such that $\|\mu_{h_u}\|_{\text{TV}} \leq \delta$, we have

$$\begin{aligned} \mathbb{E}_{x \sim \mu}[h_u(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[h_u(\langle u, x \rangle)] &= \int \int_{\mathbb{S}^1} (\langle \theta, (\langle u, x \rangle, 1) \rangle)_+^\alpha d\mu_{h_u}(\theta) d(\mu - \nu)(x) \\ &= \int_{\mathbb{S}^1} \int (\langle (\theta_1 u, \theta_2), (x, 1) \rangle)_+^\alpha d(\mu - \nu)(x) d\mu_{h_u}(\theta) \\ &\leq \delta \sup_{\theta \in \mathbb{S}^1} \left| \int (\langle (\theta_1 u, \theta_2), (x, 1) \rangle)_+^\alpha d(\mu - \nu)(x) \right| \\ &= \delta \sup_{\gamma \in [0, 2\pi]} \left| \int (\langle (\cos(\gamma)u, \sin(\gamma)), (x, 1) \rangle)_+^\alpha d(\mu - \nu)(x) \right|. \end{aligned} \quad (49)$$

Hence, (48) and (49) yield

$$\begin{aligned} & \delta \int_{\mathbb{S}^{d-1}} \sup_{\gamma \in [0, 2\pi]} \left| \int (\langle (\cos(\gamma)u, \sin(\gamma)), (x, 1) \rangle)_+^\alpha d(\mu - \nu)(x) \right| d\tau(\nu) \\ & \geq \mathcal{W}_{1,1}(\mu, \nu) - 2C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta). \end{aligned} \quad (50)$$

If we use the Hölder inequality in the left-hand side of (50), we obtain

$$\begin{aligned} & \delta \left(\int_{\mathbb{S}^{d-1}} \sup_{\gamma \in [0, 2\pi]} \left| \int (\langle (\cos(\gamma)u, \sin(\gamma)), (x, 1) \rangle)_+^\alpha d(\mu - \nu)(x) \right|^3 d\tau(\nu) \right)^{1/3} \\ & \geq \mathcal{W}_{1,1}(\mu, \nu) - 2C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta). \end{aligned} \quad (51)$$

By Lemma 13, we have

$$d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) \geq \frac{1}{2\pi} \frac{5}{6\alpha 2^{\alpha/2}} \int_{\mathbb{S}^{d-1}} \sup_{\gamma \in [0, 2\pi]} \left| \int (\langle (x, 1), (\cos(\gamma)\xi_{(d)}, \sin(\gamma)) \rangle)_+^\alpha d(\mu - \nu)(x) \right|^3 d\tau_{(d-1)}(\xi_{(d)}).$$

Hence, combining this bound with (51) we conclude that

$$\delta d_{\mathcal{F}_2}^{2/3}(\mu, \nu) \geq \left(\frac{5}{12\pi\alpha 2^{\alpha/2}} \right)^{1/3} \left(\mathcal{W}_{1,1}(\mu, \nu) - 2C(1, \alpha)\delta^{-\frac{1}{\alpha}} \log(\delta) \right).$$

The upper bound $\pi d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) \leq \mathcal{W}_{1,1}(\mu, \nu)$ follows from

$$\begin{aligned} \mathcal{W}_{1,1}(\mu, \nu) &= \int_{\mathbb{S}^{d-1}} \left(\sup_{f \in \text{Lip}_1(\mathbb{R})} \mathbb{E}_{x \sim \mu}[f(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[f(\langle u, x \rangle)] \right) d\tau(u) \\ &\geq \frac{1}{2} \int_{\mathbb{S}^{d-1}} \left(\sup_{f \in \text{Lip}_1(\mathbb{R})} \mathbb{E}_{x \sim \mu}[f(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[f(\langle u, x \rangle)] \right)^2 d\tau(u) \\ &\geq \frac{1}{2} \int_{\mathbb{S}^{d-1}} \left(\sup_{f \in \mathcal{B}_{\mathcal{F}_2}(\mathbb{R})} \mathbb{E}_{x \sim \mu}[f(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[f(\langle u, x \rangle)] \right)^2 d\tau(u) \\ &= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_{\mathbb{S}^1} \left(\int (\langle (\langle u, x \rangle, 1), \theta \rangle)_+^\alpha d(\mu - \nu)(x) \right)^2 d\tau_{(1)}(\theta) d\tau(u) \\ &= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \int_0^{2\pi} \left(\int (\langle (x, 1), (\cos(\gamma)u, \sin(\gamma)) \rangle)_+^\alpha d(\mu - \nu)(x) \right)^2 d\gamma d\tau(u) \\ &= \pi d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu, \nu) \end{aligned}$$

In the first inequality, we used that $|\sup_{f \in \text{Lip}_1(\mathbb{R})} \mathbb{E}_{x \sim \mu}[f(\langle u, x \rangle)] - \mathbb{E}_{x \sim \nu}[f(\langle u, x \rangle)]| \leq 2$ since $\|x\|_2 \leq 1$ for all x in the support of μ or ν . In the second inequality, we used that $\mathcal{B}_{\mathcal{F}_2}(\mathbb{R}) \subseteq \mathcal{B}_{\mathcal{F}_1}(\mathbb{R}) \subseteq \text{Lip}_1(\mathbb{R})$. The next equality follows from (18). The last equality is from (47). \square

G Experimental details

For the figures, the experiments were run with CPUs from a cluster, using a different 15GB RAM node for each dimension and repetition. The experiments for Figure 2 were the most computationally expensive taking about 40 hours to complete. We need to use a high amount of Monte Carlo samples from the measures to reduce the variance of the estimator, and samples are computationally expensive to obtain because the rejection rate for rejection sampling, which was the method we chose for simplicity, was high. The code would be faster if we had used MCMC methods to obtain the samples, but we are not too concerned about the speed because the only purpose is to plot figures, not to design an algorithm that can be implemented.

Details on Figure 2. To get the theoretical \mathcal{F}_1 IPM estimate (which is *still* an estimate, i.e. not a closed form expression), we use (6), which states that

$$d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) = \frac{2 \left| \int_{-1}^1 P_{k,d}(t) \sigma(t) (1-t^2)^{\frac{d-3}{2}} dt \right|}{\int_{-1}^1 |P_{k,d}(t)| (1-t^2)^{\frac{d-3}{2}} dt}.$$

To approximate this quantity, we observe that it can be expressed as $2\mathbb{E}_t[\sigma(t)\text{sign}(P_{k,d}(t))]$ when the distribution of t has a density proportional to $|P_{k,d}(t)|(1-t^2)^{\frac{d-3}{2}}$ restricted to $[-1, 1]$. We sample from this density using rejection sampling and obtain the desired estimate as the Monte Carlo estimate of $2\mathbb{E}_t[\sigma(t)\text{sign}(P_{k,d}(t))]$.

The empirical \mathcal{F}_1 IPM estimate in the left plot is computed by writing $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) = \sup_{\theta \in \mathbb{S}^{d-1}} \left| \int \sigma(\langle x, \theta \rangle) d(\mu_d - \nu_d)(x) \right|$ per Lemma 1. Since this supremum is attained at $\theta = e_d$ (see the proof of Lemma 6 in App. D), we rely on the Monte Carlo estimate $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d) \approx \frac{1}{M} \left| \sum_{i=1}^M \sigma(\langle x_i, e_d \rangle) - \sigma(\langle y_i, e_d \rangle) \right|$, where $(x_i)_{i=1}^M$ and $(y_i)_{i=1}^M$ are i.i.d. samples from μ_d and ν_d respectively. Analogously, the \mathcal{F}_2 IPM estimate in the left plot is computed by writing $d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu_d, \nu_d) = \int_{\mathbb{S}^{d-1}} \left(\int \sigma(\langle x, \theta \rangle) d(\mu_d - \nu_d)(x) \right)^2 d\tau(\theta)$ per Lemma 2. We use Monte Carlo estimates to approximate the integrals over $\mu_d - \nu_d$ and over τ , i.e.

$$d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu_d, \nu_d) \approx \frac{1}{N} \sum_{j=1}^N \left(\frac{1}{M} \sum_{i=1}^M \sigma(\langle x_i, \theta_j \rangle) - \sigma(\langle y_i, \theta_j \rangle) \right)^2.$$

In Figure 2 we used $M = 6000000$, $N = 10000$ and we obtained the samples $(x_i)_{i=1}^M$ and $(y_i)_{i=1}^M$ using rejection sampling. The curves for the empirical estimates in the left plot are obtained by running the Monte Carlo estimate 10 times; thick lines show the average, and error bars indicate the minimum and maximum values over the 10 repetitions. The empirical ratio in the right plot is obtained by dividing the \mathcal{F}_1 IPM estimate over the \mathcal{F}_2 IPM estimate, and its error bars are obtained by dividing the minimum value (resp. maximum) for the \mathcal{F}_1 IPM over the 10 repetitions by the maximum value (resp. minimum) for the \mathcal{F}_2 IPM.

Details on Figure 3. The \mathcal{F}_1 SD estimate in the left plot is computed using that $\text{SD}_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)$ is equal to

$$\gamma_{k,d} \frac{k}{\alpha+1} \lambda_{k,d}^{(\alpha+1)} \sqrt{\sum_{i=1}^d \sup_{\theta^{(i)} \in \mathbb{S}^{d-1}} \left(\frac{d+\alpha-2}{k} \hat{\nabla}_i L_{k,d}(\theta^{(i)}) - (\alpha+1-k) L_{k,d}(\theta^{(i)}) \theta_i^{(i)} \right)^2}.$$

by equations (35), (40) and (42). In (43) we lower-bound the supremum for $i = d$, which suffices for the lower bound in Lemma 10. However, we need a procedure to approximate the suprema for $i = 1, \dots, d$. By the fact that $L_{k,d}(x) = \|x\|^k P_{k,d}(\langle e_d, x \rangle / \|x\|)$ for any $x \in \mathbb{R}^d$ (see the second paragraph of Sec. 5), we have that for all $\theta^{(i)} \in \mathbb{S}^{d-1}$,

$$\begin{aligned} \hat{\nabla}_i L_{k,d}(\theta^{(i)}) &= \hat{\nabla}_i \left(\|\theta^{(i)}\|^k P_{k,d}(\langle e_d, \theta^{(i)} \rangle / \|\theta^{(i)}\|) \right) \\ &= \mathbb{1}_{i=d} P'_{k,d}(\langle e_d, \theta^{(i)} \rangle) - \langle e_d, \theta^{(i)} \rangle P'_{k,d}(\langle e_d, \theta^{(i)} \rangle) \theta_i^{(i)} + k P_{k,d}(\langle e_d, \theta^{(i)} \rangle) \theta_i^{(i)} \\ &= \frac{k(k+d-2)}{d-1} P_{k-1,d+2}(\langle e_d, \theta^{(i)} \rangle) (\mathbb{1}_{i=d} - \langle e_d, \theta^{(i)} \rangle \theta_i^{(i)}) + k P_{k,d}(\langle e_d, \theta^{(i)} \rangle) \theta_i^{(i)} \end{aligned}$$

In the last equality we have used (27). Thus, for $i \neq d$,

$$\begin{aligned} &\frac{d+\alpha-2}{k} \hat{\nabla}_i L_{k,d}(\theta^{(i)}) - (\alpha+1-k) L_{k,d}(\theta^{(i)}) \theta_i^{(i)} \\ &= -\frac{(d+\alpha-2)(k+d-2)}{d-1} \langle e_d, \theta^{(i)} \rangle P_{k-1,d+2}(\langle e_d, \theta^{(i)} \rangle) \theta_i^{(i)} + (d+k-3) P_{k,d}(\langle e_d, \theta^{(i)} \rangle) \theta_i^{(i)} \end{aligned} \tag{52}$$

We want to find $\theta^{(i)}$ that maximizes the absolute value of (52) within \mathbb{S}^{d-1} , which via the change of variables $t = \langle e_d, \theta^{(i)} \rangle$ is equivalent to minimizing the one-dimensional function

$$\left| -\frac{(d+\alpha-2)(k+d-2)}{d-1}t(1-t^2)^{1/2}P_{k-1,d+2}(t) + (d+k-3)(1-t^2)^{1/2}P_{k,d}(t) \right| \quad (53)$$

over $[-1, 1]$. Here, we have used that the absolute value of (52) is maximized when $\theta^{(i)} = \theta_i^{(i)} + \langle e_d, \theta^{(i)} \rangle e_d$, which implies that $\theta_i^{(i)} = \pm(1-t^2)^{1/2}$. We can optimize (53) over $[-1, 1]$ via brute force, since it is a one dimensional problem. On the other hand, when $i = d$ we have

$$\begin{aligned} & \frac{d+\alpha-2}{k} \nabla_i L_{k,d}(\theta^{(i)}) - (\alpha+1-k) L_{k,d}(\theta^{(i)}) \theta_i^{(i)} \\ &= \frac{(d+\alpha-2)(k+d-2)}{d-1} P_{k-1,d+2}(\langle e_d, \theta^{(i)} \rangle) (1 - \langle e_d, \theta^{(i)} \rangle \theta_i^{(i)}) + (d+k-3) P_{k,d}(\langle e_d, \theta^{(i)} \rangle) \theta_i^{(i)} \end{aligned}$$

We again the change of variables $t = \langle e_d, \theta^{(i)} \rangle$, which in this case implies that $t = \theta_i^{(i)}$. Thus, the problem to be solved for $i = d$ is

$$\left| \frac{(d+\alpha-2)(k+d-2)}{d-1} P_{k-1,d+2}(t)(1-t^2) + (d+k-3) P_{k,d}(t)t \right|.$$

The theoretical lower bound on the \mathcal{F}_1 SD is obtained directly by evaluating the right-hand side of (34).

The \mathcal{F}_2 SD estimate is obtained as a Monte-Carlo estimate of the right-hand side of (44). Namely, if $(\theta_j)_{j=1}^N$ and $(x_l)_{l=1}^M$ are uniform i.i.d. samples over \mathbb{S}^{d-1} ,

$$d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu_d, \nu_d) \approx \gamma_{k,d}^2 \sum_{i=1}^d \frac{1}{N} \sum_{j=1}^N \left(\frac{1}{M} \sum_{l=1}^M \nabla_i L_{k,d}(x_l) (\langle x_l, \theta_j \rangle)_+^\alpha \right)^2.$$

Details on Figure 4. Denoting by μ_d the standard d -variate Gaussian and by ν_d the d -variate Gaussian with unit variance in all directions except for e_d with variance 0.1, we have taken M samples $(x_i)_{i=1}^M$ of μ_d and M samples $(y_i)_{i=1}^M$ of ν_d . We used the same estimate for $d_{\mathcal{B}_{\mathcal{F}_1}}(\mu_d, \nu_d)$ and $d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)$ as in Figure 2, although in this case with bias term and with $M = 100000$ and $N = 10000$. To obtain an estimate $d_{\mathcal{B}_{\mathcal{F}_2}}(\mu_d, \nu_d)$, we sample N uniform samples $(\theta_i)_{i=1}^N$ from \mathbb{S}^{d-1} and N uniform samples (t_i) from $[-1, 1]$ and we compute

$$d_{\mathcal{B}_{\mathcal{F}_2}}^2(\mu_d, \nu_d) \approx \frac{1}{N} \sum_{j=1}^N \left(\frac{1}{M} \sum_{i=1}^M \sigma(\langle (x_i, 1), (\sqrt{1-t_j^2}\theta_j, t_j) \rangle) - \sigma(\langle (y_i, 1), (\sqrt{1-t_j^2}\theta_j, t_j) \rangle) \right)^2.$$

Let $\mathcal{W}_1((\theta)_{\#}\mu, (\theta)_{\#}\nu)$ be the one-dimensional Wasserstein distance between the projections of μ and ν_d to the one dimensional subspace spanned by $\theta \in \mathbb{S}^{d-1}$. Let $\hat{\mu}_d = \frac{1}{M} \sum_{i=1}^M \delta_{x_i}$ and $\hat{\nu}_d = \frac{1}{M} \sum_{i=1}^M \delta_{y_i}$. To estimate the max-sliced Wasserstein $\overline{\mathcal{W}}_{1,1}(\hat{\mu}_d, \hat{\nu}_d)$ we compute $\overline{\mathcal{W}}_{1,1}(\mu_d, \nu_d) \approx \mathcal{W}_1((e_d)_{\#}\hat{\mu}_d, (e_d)_{\#}\hat{\nu}_d)$ to e_d , because we know that in theory e_d is the direction of maximal discrepancy. The one-dimensional Wasserstein distance can be computed quickly.

To estimate the sliced Wasserstein distance $\underline{\mathcal{W}}_{1,1}(\mu_d, \nu_d)$ we sample N uniform samples $(\theta_i)_{i=1}^N$ from \mathbb{S}^{d-1} and we compute

$$\underline{\mathcal{W}}_{1,1}(\mu_d, \nu_d) \approx \frac{1}{N} \sum_{j=1}^N \mathcal{W}_1((\theta_j)_{\#}\hat{\mu}_d, (\theta_j)_{\#}\hat{\nu}_d).$$

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#)
 - (c) Did you discuss any potential negative societal impacts of your work? [\[No\]](#) The work is of theoretical nature and has no direct negative societal impacts. Our results give insight into generative modeling algorithms, which are useful for many relevant tasks such as image generation, image-to-image or text-to-image translation and video prediction. As always, we note that machine learning improvements like ours come in the form of “building machines to do X better”. For a sufficiently malicious or ill-informed choice of X, such as surveillance or recidivism prediction, almost any progress in machine learning might indirectly lead to a negative outcome, and our work is not excluded from that.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#)
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) We include the code in the supplementary material, and will make the repository available if/once the paper is accepted.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) We provide the most relevant details on the experiments in [Sec. 8](#), and a complete specification in [App. G](#).
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) And in the first paragraph of [Sec. 8](#) we explain: “The empirical estimates in the plots are detailed in [App. G](#). They are averaged over 10 repetitions; the error bars show the maximum and minimum.”
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) Discussed in the first paragraph of [App. G](#).
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[No\]](#) Our work does not use existing assets beyond the ideas of the works we cite.
 - (b) Did you mention the license of the assets? [\[N/A\]](#)
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#) We include our code in the supplementary material, if it may be considered as such.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)