
On gradient regularizers for MMD GANs

Michael Arbel*
Gatsby Computational Neuroscience Unit
University College London
michael.n.arbel@gmail.com

Dougal J. Sutherland*
Gatsby Computational Neuroscience Unit
University College London
dougal@gmail.com

Mikołaj Bińkowski
Department of Mathematics
Imperial College London
mikbinkowski@gmail.com

Arthur Gretton
Gatsby Computational Neuroscience Unit
University College London
arthur.gretton@gmail.com

Abstract

We propose a principled method for gradient-based regularization of the critic of GAN-like models trained by adversarially optimizing the kernel of a Maximum Mean Discrepancy (MMD). We show that controlling the gradient of the critic is vital to having a sensible loss function, and devise a method to enforce exact, analytical gradient constraints at no additional cost compared to existing approximate techniques based on additive regularizers. The new loss function is provably continuous, and experiments show that it stabilizes and accelerates training, giving image generation models that outperform state-of-the-art methods on 160×160 CelebA and 64×64 unconditional ImageNet.

1 Introduction

There has been an explosion of interest in *implicit generative models* (IGMs) over the last few years, especially after the introduction of generative adversarial networks (GANs) [16]. These models allow approximate samples from a complex high-dimensional target distribution \mathbb{P} , using a model distribution \mathbb{Q}_θ , where estimation of likelihoods, exact inference, and so on are not tractable. GAN-type IGMs have yielded very impressive empirical results, particularly for image generation, far beyond the quality of samples seen from most earlier generative models [e.g. 18, 22, 23, 24, 38].

These excellent results, however, have depended on adding a variety of methods of regularization and other tricks to stabilize the notoriously difficult optimization problem of GANs [38, 42]. Some of this difficulty is perhaps because when a GAN is viewed as minimizing a discrepancy $\mathcal{D}_{\text{GAN}}(\mathbb{P}, \mathbb{Q}_\theta)$, its gradient $\nabla_\theta \mathcal{D}_{\text{GAN}}(\mathbb{P}, \mathbb{Q}_\theta)$ does not provide useful signal to the generator if the target and model distributions are not absolutely continuous, as is nearly always the case [2].

An alternative set of losses are the integral probability metrics (IPMs) [36], which can give credit to models \mathbb{Q}_θ “near” to the target distribution \mathbb{P} [3, 8, Section 4 of 15]. IPMs are defined in terms of a *critic function*: a “well behaved” function with large amplitude where \mathbb{P} and \mathbb{Q}_θ differ most. The IPM is the difference in the expected critic under \mathbb{P} and \mathbb{Q}_θ , and is zero when the distributions agree. The Wasserstein IPMs, whose critics are made smooth via a Lipschitz constraint, have been particularly successful in IGMs [3, 14, 18]. But the Lipschitz constraint must hold uniformly, which can be hard to enforce. A popular approximation has been to apply a gradient constraint only in expectation [18]: the critic’s gradient norm is constrained to be small on points chosen uniformly between \mathbb{P} and \mathbb{Q} .

Another class of IPMs used as IGM losses are the Maximum Mean Discrepancies (MMDs) [17], as in [13, 28]. Here the critic function is a member of a reproducing kernel Hilbert space (except in [50], who learn a deep approximation to an RKHS critic). Better performance can be obtained,

*These authors contributed equally.

however, when the MMD kernel is not based directly on image pixels, but on learned features of images. Wasserstein-inspired gradient regularization approaches can be used on the MMD critic when learning these features: [27] uses weight clipping [3], and [5, 7] use a gradient penalty [18].

The recent Sobolev GAN [33] uses a similar constraint on the expected gradient norm, but phrases it as estimating a Sobolev IPM rather than loosely approximating Wasserstein. This expectation can be taken over the same distribution as [18], but other measures are also proposed, such as $(\mathbb{P} + \mathbb{Q}_\theta) / 2$. A second recent approach, the spectrally normalized GAN [32], controls the Lipschitz constant of the critic by enforcing the spectral norms of the weight matrices to be 1. Gradient penalties also benefit GANs based on f -divergences [37]: for instance, the spectral normalization technique of [32] can be applied to the critic network of an f -GAN. Alternatively, a gradient penalty can be defined to approximate the effect of blurring \mathbb{P} and \mathbb{Q}_θ with noise [40], which addresses the problem of non-overlapping support [2]. This approach has recently been shown to yield locally convergent optimization in some cases with non-continuous distributions, where the original GAN does not [30].

In this paper, we introduce a novel regularization for the MMD GAN critic of [5, 7, 27], which *directly targets generator performance*, rather than adopting regularization methods intended to approximate Wasserstein distances [3, 18]. The new MMD regularizer derives from an approach widely used in semi-supervised learning [10, Section 2], where the aim is to define a classification function f which is positive on \mathbb{P} (the positive class) and negative on \mathbb{Q}_θ (negative class), in the absence of labels on many of the samples. The decision boundary between the classes is assumed to be in a region of low density for both \mathbb{P} and \mathbb{Q}_θ : f should therefore be flat where \mathbb{P} and \mathbb{Q}_θ have support (areas with constant label), and have a larger slope in regions of low density. Bousquet et al. [10] propose as their regularizer on f a sum of the variance and a density-weighted gradient norm.

We adopt a related penalty on the MMD critic, with the difference that we only apply the penalty on \mathbb{P} : thus, the critic is flatter where \mathbb{P} has high mass, but does not vanish on the generator samples from \mathbb{Q}_θ (which we optimize). In excluding \mathbb{Q}_θ from the critic function constraint, we also avoid the concern raised by [32] that a critic depending on \mathbb{Q}_θ will change with the current minibatch – potentially leading to less stable learning. The resulting discrepancy is no longer an integral probability metric: it is asymmetric, and the critic function class depends on the target \mathbb{P} being approximated.

We first discuss in Section 2 how MMD-based losses can be used to learn implicit generative models, and how a naive approach could fail. This motivates our new discrepancies, introduced in Section 3. Section 4 demonstrates that these losses outperform state-of-the-art models for image generation.

2 Learning implicit generative models with MMD-based losses

An IGM is a model \mathbb{Q}_θ which aims to approximate a target distribution \mathbb{P} over a space $\mathcal{X} \subseteq \mathbb{R}^d$. We will define \mathbb{Q}_θ by a *generator* function $G_\theta : \mathcal{Z} \rightarrow \mathcal{X}$, implemented as a deep network with parameters θ , where \mathcal{Z} is a space of latent codes, say \mathbb{R}^{128} . We assume a fixed distribution on \mathcal{Z} , say $Z \sim \text{Uniform}([-1, 1]^{128})$, and call \mathbb{Q}_θ the distribution of $G_\theta(Z)$. We will consider learning by minimizing a discrepancy \mathcal{D} between distributions, with $\mathcal{D}(\mathbb{P}, \mathbb{Q}_\theta) \geq 0$ and $\mathcal{D}(\mathbb{P}, \mathbb{P}) = 0$, which we call our *loss*. We aim to minimize $\mathcal{D}(\mathbb{P}, \mathbb{Q}_\theta)$ with stochastic gradient descent on an estimator of \mathcal{D} .

In the present work, we will build losses \mathcal{D} based on the Maximum Mean Discrepancy,

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) = \sup_{f: \|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}}[f(Y)], \quad (1)$$

an integral probability metric where the critic class is the unit ball within \mathcal{H}_k , the reproducing kernel Hilbert space with a kernel k . The optimization in (1) admits a simple closed-form optimal critic, $f^*(t) \propto \mathbb{E}_{X \sim \mathbb{P}}[k(X, t)] - \mathbb{E}_{Y \sim \mathbb{Q}}[k(Y, t)]$. There is also an unbiased, closed-form estimator of MMD_k^2 with appealing statistical properties [17] – in particular, its sample complexity is *independent* of the dimension of \mathcal{X} , compared to the exponential dependence [52] of the Wasserstein distance

$$\mathcal{W}(\mathbb{P}, \mathbb{Q}) = \sup_{f: \|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}}[f(Y)]. \quad (2)$$

The MMD is *continuous in the weak topology* for any bounded kernel with Lipschitz embeddings [46, Theorem 3.2(b)], meaning that if \mathbb{P}_n converges in distribution to \mathbb{P} , $\mathbb{P}_n \xrightarrow{D} \mathbb{P}$, then $\text{MMD}(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$. (\mathcal{W} is continuous in the slightly stronger Wasserstein topology [51, Definition 6.9]; $\mathbb{P}_n \xrightarrow{\mathcal{W}} \mathbb{P}$ implies

$\mathbb{P}_n \xrightarrow{D} \mathbb{P}$, and the two notions coincide if \mathcal{X} is bounded.) Continuity means the loss can provide better signal to the generator as \mathbb{Q}_θ approaches \mathbb{P} , as opposed to e.g. Jensen-Shannon where the loss could be constant until suddenly jumping to 0 [e.g. 3, Example 1]. The MMD is also *strict*, meaning it is zero iff $\mathbb{P} = \mathbb{Q}_\theta$, for *characteristic* kernels [45]. The Gaussian kernel yields an MMD both continuous in the weak topology and strict. Thus in principle, one need not conduct any alternating optimization in an IGM at all, but merely choose generator parameters θ to minimize MMD_k .

Despite these appealing properties, using simple pixel-level kernels leads to poor generator samples [8, 13, 28, 48]. More recent MMD GANs [5, 7, 27] achieve better results by using a parameterized *family* of kernels, $\{k_\psi\}_{\psi \in \Psi}$, in the Optimized MMD loss previously studied by [44, 46]:

$$\mathcal{D}_{\text{MMD}}^\Psi(\mathbb{P}, \mathbb{Q}) := \sup_{\psi \in \Psi} \text{MMD}_{k_\psi}(\mathbb{P}, \mathbb{Q}). \quad (3)$$

We primarily consider kernels defined by some fixed kernel K on top of a learned low-dimensional representation $\phi_\psi : \mathcal{X} \rightarrow \mathbb{R}^s$, i.e. $k_\psi(x, y) = K(\phi_\psi(x), \phi_\psi(y))$, denoted $k_\psi = K \circ \phi_\psi$. In practice, K is a simple characteristic kernel, e.g. Gaussian, and ϕ_ψ is usually a deep network with output dimension say $s = 16$ [7] or even $s = 1$ (in our experiments). If ϕ_ψ is powerful enough, this choice is sufficient; we need not try to ensure each k_ψ is characteristic, as did [27].

Proposition 1. *Suppose $k = K \circ \phi_\psi$, with K characteristic and $\{\phi_\psi\}$ rich enough that for any $\mathbb{P} \neq \mathbb{Q}$, there is a $\psi \in \Psi$ for which $\phi_{\psi\#\mathbb{P}} \neq \phi_{\psi\#\mathbb{Q}}$.² Then if $\mathbb{P} \neq \mathbb{Q}$, $\mathcal{D}_{\text{MMD}}^\Psi(\mathbb{P}, \mathbb{Q}) > 0$.*

Proof. Let $\hat{\psi} \in \Psi$ be such that $\phi_{\hat{\psi}}(\mathbb{P}) \neq \phi_{\hat{\psi}}(\mathbb{Q})$. Then, since K is characteristic,

$$\mathcal{D}_{\text{MMD}}^\Psi(\mathbb{P}, \mathbb{Q}) = \sup_{\psi \in \Psi} \text{MMD}_K(\phi_{\psi\#\mathbb{P}}, \phi_{\psi\#\mathbb{Q}}) \geq \text{MMD}_K(\phi_{\hat{\psi}\#\mathbb{P}}, \phi_{\hat{\psi}\#\mathbb{Q}}) > 0. \quad \square$$

To estimate $\mathcal{D}_{\text{MMD}}^\Psi$, one can conduct alternating optimization to estimate a $\hat{\psi}$ and then update the generator according to $\text{MMD}_{k_{\hat{\psi}}}$, similar to the scheme used in GANs and WGANs. (This form of estimator is justified by an envelope theorem [31], although it is invariably biased [7].) Unlike \mathcal{D}_{GAN} or \mathcal{W} , fixing a $\hat{\psi}$ and optimizing the generator still yields a sensible distance $\text{MMD}_{k_{\hat{\psi}}}$.

Early attempts at minimizing $\mathcal{D}_{\text{MMD}}^\Psi$ in an IGM, though, were unsuccessful [48, footnote 7]. This could be because for some kernel classes, $\mathcal{D}_{\text{MMD}}^\Psi$ is stronger than Wasserstein or MMD.

Example 1 (DiracGAN [30]). *We wish to model a point mass at the origin of \mathbb{R} , $\mathbb{P} = \delta_0$, with any possible point mass, $\mathbb{Q}_\theta = \delta_\theta$ for $\theta \in \mathbb{R}$. We use a Gaussian kernel of any bandwidth, which can be written as $k_\psi = K \circ \phi_\psi$ with $\phi_\psi(x) = \psi x$ for $\psi \in \Psi = \mathbb{R}$ and $K(a, b) = \exp(-\frac{1}{2}(a-b)^2)$. Then*

$$\text{MMD}_{k_\psi}^2(\delta_0, \delta_\theta) = 2 \left(1 - \exp\left(-\frac{1}{2}\psi^2\theta^2\right)\right), \quad \mathcal{D}_{\text{MMD}}^\Psi(\delta_0, \delta_\theta) = \begin{cases} \sqrt{2} & \theta \neq 0 \\ 0 & \theta = 0 \end{cases}.$$

Considering $\mathcal{D}_{\text{MMD}}^\Psi(\delta_0, \delta_{1/n}) = \sqrt{2} \not\rightarrow 0$, even though $\delta_{1/n} \xrightarrow{\mathcal{W}} \delta_0$, shows that the Optimized MMD distance is not continuous in the weak or Wasserstein topologies.

This also causes optimization issues. Figure 1 (a) shows gradient vector fields in parameter space, $v(\theta, \psi) \propto (-\nabla_\theta \text{MMD}_{k_\psi}^2(\delta_0, \delta_\theta), \nabla_\psi \text{MMD}_{k_\psi}^2(\delta_0, \delta_\theta))$. Some sequences following v (e.g. A) converge to an optimal solution $(0, \psi)$, but some (B) move in the wrong direction, and others (C) are stuck because there is essentially no gradient. Figure 1 (c, red) shows that the optimal $\mathcal{D}_{\text{MMD}}^\Psi$ critic is very sharp near \mathbb{P} and \mathbb{Q} ; this is less true for cases where the algorithm converged.

We can avoid these issues if we ensure a bounded Lipschitz critic:³

Proposition 2. *Assume the critics $f_\psi(x) = (\mathbb{E}_{X \sim \mathbb{P}} k_\psi(X, x) - \mathbb{E}_{Y \sim \mathbb{Q}} k_\psi(Y, x)) / \text{MMD}_{k_\psi}(\mathbb{P}, \mathbb{Q})$ are uniformly bounded and have a common Lipschitz constant: $\sup_{x \in \mathcal{X}, \psi \in \Psi} |f_\psi(x)| < \infty$ and $\sup_{\psi \in \Psi} \|f_\psi\|_{\text{Lip}} < \infty$. In particular, this holds when $k_\psi = K \circ \phi_\psi$ and*

$$\sup_{a \in \mathbb{R}^s} K(a, a) < \infty, \quad \|K(a, \cdot) - K(b, \cdot)\|_{\mathcal{H}_K} \leq L_K \|a - b\|_{\mathbb{R}^s}, \quad \sup_{\psi \in \Psi} \|\phi_\psi\|_{\text{Lip}} \leq L_\phi < \infty.$$

Then $\mathcal{D}_{\text{MMD}}^\Psi$ is continuous in the weak topology: if $\mathbb{P}_n \xrightarrow{D} \mathbb{P}$, then $\mathcal{D}_{\text{MMD}}^\Psi(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$.

² $f\#\mathbb{P}$ denotes the *pushforward* of a distribution: if $X \sim \mathbb{P}$, then $f(X) \sim f\#\mathbb{P}$.

³ [27, Theorem 4] makes a similar claim to Proposition 2, but its proof was incorrect: it tries to uniformly bound $\text{MMD}_{k_\psi} \leq \mathcal{W}^2$, but the bound used is for a Wasserstein in terms of $\|k_\psi(x, \cdot) - k_\psi(y, \cdot)\|_{\mathcal{H}_{k_\psi}}$.

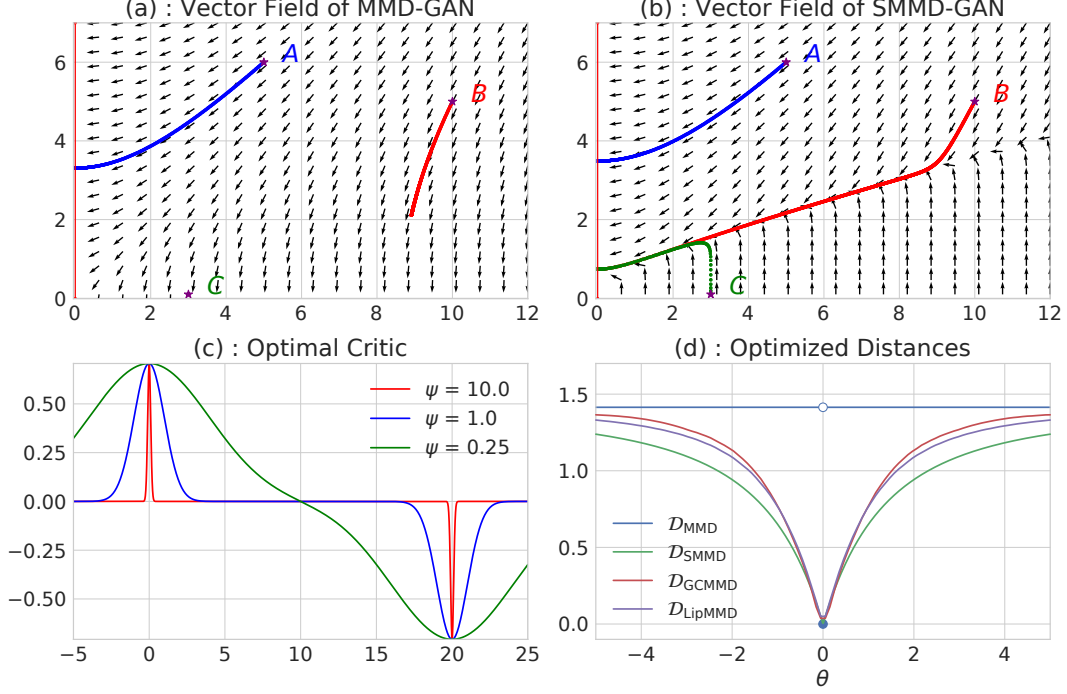


Figure 1: The setting of Example 1. (a, b): parameter-space gradient fields for the MMD and the SMMD (Section 3.3); the horizontal axis is θ , and the vertical $1/\psi$. (c): optimal MMD critics for $\theta = 20$ with different kernels. (d): the MMD and the distances of Section 3 optimized over ψ .

Proof. The main result is [12, Corollary 11.3.4]. To show the claim for $k_\psi = K \circ \phi_\psi$, note that $|f_\psi(x) - f_\psi(y)| \leq \|f_\psi\|_{\mathcal{H}_{k_\psi}} \|k_\psi(x, \cdot) - k_\psi(y, \cdot)\|_{\mathcal{H}_{k_\psi}}$, which since $\|f_\psi\|_{\mathcal{H}_{k_\psi}} = 1$ is

$$\|K(\phi_\psi(x), \cdot) - K(\phi_\psi(y), \cdot)\|_{\mathcal{H}_K} \leq L_K \|\phi_\psi(x) - \phi_\psi(y)\|_{\mathbb{R}^s} \leq L_K L_\phi \|x - y\|_{\mathbb{R}^d}. \quad \square$$

Indeed, if we put a box constraint on ψ [27] or regularize the gradient of the critic function [7], the resulting MMD GAN generally matches or outperforms WGAN-based models. Unfortunately, though, an additive gradient penalty doesn't substantially change the vector field of Figure 1 (a), as shown in Figure 5 (Appendix B). We will propose distances with much better convergence behavior.

3 New discrepancies for learning implicit generative models

Our aim here is to introduce a discrepancy that can provide useful gradient information when used as an IGM loss. Proofs of results in this section are deferred to Appendix A.

3.1 Lipschitz Maximum Mean Discrepancy

Proposition 2 shows that an MMD-like discrepancy can be continuous under the weak topology even when optimizing over kernels, if we directly restrict the critic functions to be Lipschitz. We can easily define such a distance, which we call the Lipschitz MMD: for some $\lambda > 0$,

$$\text{LipMMD}_{k,\lambda}(\mathbb{P}, \mathbb{Q}) := \sup_{f \in \mathcal{H}_k : \|f\|_{\text{Lip}} + \lambda \|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} [f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}} [f(Y)]. \quad (4)$$

For a universal kernel k , we conjecture that $\lim_{\lambda \rightarrow 0} \text{LipMMD}_{k,\lambda}(\mathbb{P}, \mathbb{Q}) \rightarrow \mathcal{W}(\mathbb{P}, \mathbb{Q})$. But for any k and λ , LipMMD is upper-bounded by \mathcal{W} , as (4) optimizes over a smaller set of functions than (2). Thus $\mathcal{D}_{\text{LipMMD}}^{\Psi,\lambda}(\mathbb{P}, \mathbb{Q}) := \sup_{\psi \in \Psi} \text{LipMMD}_{k_\psi,\lambda}(\mathbb{P}, \mathbb{Q})$ is also upper-bounded by \mathcal{W} , and hence is continuous in the Wasserstein topology. It also shows excellent empirical behavior on Example 1 (Figure 1 (d), and Figure 5 in Appendix B). But estimating $\text{LipMMD}_{k,\lambda}$, let alone $\mathcal{D}_{\text{LipMMD}}^{\Psi,\lambda}$, is in general extremely difficult (Appendix D), as finding $\|f\|_{\text{Lip}}$ requires optimization in the input space. Constraining the *mean* gradient rather than the *maximum*, as we will do next, is far more tractable.

3.2 Gradient-Constrained Maximum Mean Discrepancy

We define the Gradient-Constrained MMD for $\lambda > 0$ and using some measure μ as

$$\text{GCMMMD}_{\mu,k,\lambda}(\mathbb{P}, \mathbb{Q}) := \sup_{f \in \mathcal{H}_k : \|f\|_{S(\mu),k,\lambda} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} [f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}} [f(Y)], \quad (5)$$

$$\text{where } \|f\|_{S(\mu),k,\lambda}^2 := \|f\|_{L^2(\mu)}^2 + \|\nabla f\|_{L^2(\mu)}^2 + \lambda \|f\|_{\mathcal{H}_k}^2. \quad (6)$$

$\|\cdot\|_{L^2(\mu)}^2 = \int \|\cdot\|^2 \mu(dx)$ denotes the squared L^2 norm. Rather than directly constraining the Lipschitz constant, the second term $\|\nabla f\|_{L^2(\mu)}^2$ encourages the function f to be flat where μ has mass. In experiments we use $\mu = \mathbb{P}$, flattening the critic near the target sample. We add the first term following [10]: in one dimension and with μ uniform, $\|\cdot\|_{S(\mu),\cdot,0}$ is then an RKHS norm with the kernel $\kappa(x, y) = \exp(-\|x - y\|)$, which is also a Sobolev space. The correspondence to a Sobolev norm is lost in higher dimensions [53, Ch. 10], but we also found the first term to be beneficial in practice.

We can exploit some properties of \mathcal{H}_k to compute (5) analytically. Call the difference in kernel mean embeddings $\eta := \mathbb{E}_{X \sim \mathbb{P}} [k(X, \cdot)] - \mathbb{E}_{Y \sim \mathbb{Q}} [k(Y, \cdot)] \in \mathcal{H}_k$; recall $\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\eta\|_{\mathcal{H}_k}$.

Proposition 3. *Let $\hat{\mu} = \sum_{m=1}^M \delta_{X_m}$. Define $\eta(X) \in \mathbb{R}^M$ with m th entry $\eta(X_m)$, and $\nabla \eta(X) \in \mathbb{R}^{Md}$ with (m, i) th entry⁴ $\partial_i \eta(X_m)$. Then under Assumptions (A) to (D) in Appendix A.1,*

$$\text{GCMMMD}_{\hat{\mu},k,\lambda}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{\lambda} (\text{MMD}^2(\mathbb{P}, \mathbb{Q}) - \bar{P}(\eta))$$

$$\bar{P}(\eta) = \begin{bmatrix} \eta(X) \\ \nabla \eta(X) \end{bmatrix}^\top \left(\begin{bmatrix} K & G^\top \\ G & H \end{bmatrix} + M\lambda I_{M+Md} \right)^{-1} \begin{bmatrix} \eta(X) \\ \nabla \eta(X) \end{bmatrix},$$

where K is the kernel matrix $K_{m,m'} = k(X_m, X_{m'})$, G is the matrix of left derivatives⁵ $G_{(m,i),m'} = \partial_i k(X_m, X_{m'})$, and H that of derivatives of both arguments $H_{(m,i),(m',j)} = \partial_i \partial_{j+d} k(X_m, X_{m'})$.

As long as \mathbb{P} and \mathbb{Q} have integrable first moments, and μ has second moments, Assumptions (A) to (D) are satisfied e.g. by a Gaussian or linear kernel on top of a differentiable ϕ_ψ . We can thus estimate the GCMMMD based on samples from \mathbb{P} , \mathbb{Q} , and μ by using the empirical mean $\hat{\eta}$ for η .

This discrepancy indeed works well in practice: Appendix F.2 shows that optimizing our estimate of $\mathcal{D}_{\text{GCMMMD}}^{\mu,\Psi,\lambda} = \sup_{\psi \in \Psi} \text{GCMMMD}_{\mu,k_\psi,\lambda}$ yields a good generative model on MNIST. But the linear system of size $M + Md$ is impractical: even on 28×28 images and using a low-rank approximation, the model took days to converge. We therefore design a less expensive discrepancy in the next section.

The GCMMMD is related to some discrepancies previously used in IGM training. The Fisher GAN [34] uses only the variance constraint $\|f\|_{L^2(\mu)}^2 \leq 1$. The Sobolev GAN [33] constrains $\|\nabla f\|_{L^2(\mu)}^2 \leq 1$, along with a vanishing boundary condition on f to ensure a well-defined solution (although this was not used in the implementation, and can cause very unintuitive critic behavior; see Appendix C). The authors considered several choices of μ , including the WGAN-GP measure [18] and mixtures $(\mathbb{P} + \mathbb{Q}_\theta)/2$. Rather than enforcing the constraints in closed form as we do, though, these models used additive regularization. We will compare to the Sobolev GAN in experiments.

3.3 Scaled Maximum Mean Discrepancy

We will now derive a lower bound on the Gradient-Constrained MMD which retains many of its attractive qualities but can be estimated in time linear in the dimension d .

Proposition 4. *Make Assumptions (A) to (D). For any $f \in \mathcal{H}_k$, $\|f\|_{S(\mu),k,\lambda} \leq \sigma_{\mu,k,\lambda}^{-1} \|f\|_{\mathcal{H}_k}$, where*

$$\sigma_{\mu,k,\lambda} := 1 / \sqrt{\lambda + \int k(x, x) \mu(dx) + \sum_{i=1}^d \int \frac{\partial^2 k(y, z)}{\partial y_i \partial z_i} \Big|_{(y,z)=(x,x)} \mu(dx)}.$$

We then define the Scaled Maximum Mean Discrepancy based on this bound of Proposition 4:

$$\text{SMMD}_{\mu,k,\lambda}(\mathbb{P}, \mathbb{Q}) := \sup_{f : \sigma_{\mu,k,\lambda}^{-1} \|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} [f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}} [f(Y)] = \sigma_{\mu,k,\lambda} \text{MMD}_k(\mathbb{P}, \mathbb{Q}). \quad (7)$$

⁴We use (m, i) to denote $(m - 1)d + i$; thus $\nabla \eta(X)$ stacks $\nabla \eta(X_1), \dots, \nabla \eta(X_M)$ into one vector.

⁵We use $\partial_i k(x, y)$ to denote the partial derivative with respect to x_i , and $\partial_{i+d} k(x, y)$ that for y_i .

Because the constraint in the optimization of (7) is more restrictive than in that of (5), we have that $\text{SMMD}_{\mu,k,\lambda}(\mathbb{P}, \mathbb{Q}) \leq \text{GCMMD}_{\mu,k,\lambda}(\mathbb{P}, \mathbb{Q})$. The Sobolev norm $\|f\|_{S(\mu),\lambda}$, and a fortiori the gradient norm under μ , is thus also controlled for the SMMD critic. We also show in Appendix F.1 that $\text{SMMD}_{\mu,k,\lambda}$ behaves similarly to $\text{GCMMD}_{\mu,k,\lambda}$ on Gaussians.

If $k_\psi = K \circ \phi_\psi$ and $K(a, b) = g(-\|a - b\|^2)$, then $\sigma_{k,\mu,\lambda}^{-2} = \lambda + g(0) + 2|g'(0)| \mathbb{E}_\mu [\|\nabla \phi_\psi(X)\|_F^2]$. Or if K is linear, $K(a, b) = a^\top b$, then $\sigma_{k,\mu,\lambda}^{-2} = \lambda + \mathbb{E}_\mu [\|\phi_\psi(X)\|^2 + \|\nabla \phi_\psi(X)\|_F^2]$. Estimating these terms based on samples from μ is straightforward, giving a natural estimator for the SMMD.

Of course, if μ and k are fixed, the SMMD is simply a constant times the MMD, and so behaves in essentially the same way as the MMD. But optimizing the SMMD over a kernel family Ψ , $\mathcal{D}_{\text{SMMD}}^{\mu,\Psi,\lambda}(\mathbb{P}, \mathbb{Q}) := \sup_{\psi \in \Psi} \text{SMMD}_{\mu,k_\psi,\lambda}(\mathbb{P}, \mathbb{Q})$, gives a distance very different from $\mathcal{D}_{\text{MMD}}^\Psi$ (3).

Figure 1 (b) shows the vector field for the Optimized SMMD loss in Example 1, using the WGAN-GP measure $\mu = \text{Uniform}(0, \theta)$. The optimization surface is far more amenable: in particular the location C , which formerly had an extremely small gradient that made learning effectively impossible, now converges very quickly by first reducing the critic gradient until some signal is available. Figure 1 (d) demonstrates that $\mathcal{D}_{\text{SMMD}}^{\mu,\Psi,\lambda}$, like $\mathcal{D}_{\text{GCMMD}}^{\mu,\Psi,\lambda}$ and $\mathcal{D}_{\text{LipMMD}}^{\Psi,\lambda}$ but in sharp contrast to $\mathcal{D}_{\text{MMD}}^\Psi$, is continuous with respect to the location θ and provides a strong gradient towards 0.

We can establish that $\mathcal{D}_{\text{SMMD}}^{\mu,\Psi,\lambda}$ is continuous in the Wasserstein topology under some conditions:

Theorem 1. *Let $k_\psi = K \circ \phi_\psi$, with $\phi_\psi : \mathcal{X} \rightarrow \mathbb{R}^s$ a fully-connected L -layer network with Leaky-ReLU $_\alpha$ activations whose layers do not increase in width, and K satisfying mild smoothness conditions $Q_K < \infty$ (Assumptions (II) to (V) in Appendix A.2). Let Ψ^κ be the set of parameters where each layer’s weight matrices have condition number $\text{cond}(W^l) = \|W^l\| / \sigma_{\min}(W^l) \leq \kappa < \infty$. If μ has a density (Assumption (I)), then*

$$\mathcal{D}_{\text{SMMD}}^{\mu,\Psi^\kappa,\lambda}(\mathbb{P}, \mathbb{Q}) \leq \frac{Q_K \kappa^{L/2}}{\sqrt{d_L} \alpha^{L/2}} \mathcal{W}(\mathbb{P}, \mathbb{Q}).$$

Thus if $\mathbb{P}_n \xrightarrow{\mathcal{W}} \mathbb{P}$, then $\mathcal{D}_{\text{SMMD}}^{\mu,\Psi^\kappa,\lambda}(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$, even if μ is chosen to depend on \mathbb{P} and \mathbb{Q} .

Uniform bounds vs bounds in expectation Controlling $\|\nabla f_\psi\|_{L^2(\mu)}^2 = \mathbb{E}_\mu \|\nabla f_\psi(X)\|^2$ does not necessarily imply a bound on $\|f\|_{\text{Lip}} \geq \sup_{x \in \mathcal{X}} \|\nabla f_\psi(X)\|$, and so does not in general give continuity via Proposition 2. Theorem 1 implies that when the network’s weights are well-conditioned, it is sufficient to only control $\|\nabla f_\psi\|_{L^2(\mu)}^2$, which is far easier in practice than controlling $\|f\|_{\text{Lip}}$.

If we instead tried to directly control $\|f\|_{\text{Lip}}$ with e.g. spectral normalization (SN) [32], we could significantly reduce the expressiveness of the parametric family. In Example 1, constraining $\|\phi_\psi\|_{\text{Lip}} = 1$ limits us to only $\Psi = \{1\}$. Thus $\mathcal{D}_{\text{MMD}}^{\{1\}}$ is simply the MMD with an RBF kernel of bandwidth 1, which has poor gradients when θ is far from 0 (Figure 1 (c), blue). The Cauchy-Schwartz bound of Proposition 4 allows jointly adjusting the smoothness of k_ψ , and the critic f , while SN must control the two independently. Relatedly, limiting $\|\phi\|_{\text{Lip}}$ by limiting the Lipschitz norm of each layer could substantially reduce capacity, while $\|\nabla f_\psi\|_{L^2(\mu)}$ need not be decomposed by layer. Another advantage is that μ provides a data-dependent measure of complexity as in [10]: we do not needlessly prevent ourselves from using critics that behave poorly only far from the data.

Spectral parametrization When the generator is near a local optimum, the critic might identify only one direction on which \mathbb{Q}_θ and \mathbb{P} differ. If the generator parameterization is such that there is no local way for the generator to correct it, the critic may begin to single-mindedly focus on this difference, choosing redundant convolutional filters and causing the condition number of the weights to diverge. If this occurs, the generator will be motivated to fix this single direction while ignoring all other aspects of the distributions, after which it may become stuck. We can help avoid this collapse by using a critic parameterization that encourages diverse filters with higher-rank weight matrices. Miyato et al. [32] propose to parameterize the weight matrices as $W = \gamma \bar{W} / \|\bar{W}\|_{\text{op}}$, where $\|\bar{W}\|_{\text{op}}$ is the spectral norm of \bar{W} . This parametrization works particularly well with $\mathcal{D}_{\text{SMMD}}^{\mu,\Psi,\lambda}$: Figure 2 (b) shows the singular values of the second layer of a critic’s network (and Figure 9, in Appendix F.3, shows more layers), while Figure 2 (d) shows the evolution of the condition number during training. The conditioning of the weight matrix remains stable throughout training with spectral parametrization, while it worsens through training in the default case.

4 Experiments

We evaluated unsupervised image generation on three datasets: CIFAR-10 [26] (60 000 images, 32×32), CelebA [29] (202 599 face images, resized and cropped to 160×160 as in [7]), and the more challenging ILSVRC2012 (ImageNet) dataset [41] (1 281 167 images, resized to 64×64). Code for all of these experiments is available at github.com/MichaelArbel/Scaled-MMD-GAN.

Losses All models are based on a scalar-output critic network $\phi_\psi : \mathcal{X} \rightarrow \mathbb{R}$, except MMDGAN-GP where $\phi_\psi : \mathcal{X} \rightarrow \mathbb{R}^{16}$ as in [7]. The WGAN and Sobolev GAN use a critic $f = \phi_\psi$, while the GAN uses a discriminator $D_\psi(x) = 1/(1 + \exp(-\phi_\psi(x)))$. The MMD-based methods use a kernel $k_\psi(x, y) = \exp(-(\phi_\psi(x) - \phi_\psi(y))^2/2)$, except for MMDGAN-GP which uses a mixture of RQ kernels as in [7]. Increasing the output dimension of the critic or using a different kernel didn't substantially change the performance of our proposed method. We also consider SMMD with a linear top-level kernel, $k(x, y) = \phi_\psi(x)\phi_\psi(y)$; because this becomes essentially identical to a WGAN (Appendix E), we refer to this method as SWGAN. SMMD and SWGAN use $\mu = \mathbb{P}$; Sobolev GAN uses $\mu = (\mathbb{P} + \mathbb{Q})/2$ as in [33]. We choose λ and an overall scaling to obtain the losses:

$$\text{SMMD: } \frac{\widehat{\text{MMD}}_{k_\psi}^2(\mathbb{P}, \mathbb{Q}_\theta)}{1 + 10 \mathbb{E}_{\mathbb{P}} [\|\nabla \phi_\psi(X)\|_F^2]}, \text{ SWGAN: } \frac{\mathbb{E}_{\mathbb{P}} [\phi_\psi(X)] - \mathbb{E}_{\mathbb{Q}_\theta} [\phi_\psi(X)]}{\sqrt{1 + 10 (\mathbb{E}_{\mathbb{P}} [|\phi_\psi(X)|^2] + \mathbb{E}_{\mathbb{P}} [\|\nabla \phi_\psi(X)\|_F^2])}}.$$

Architecture For CIFAR-10, we used the CNN architecture proposed by [32] with a 7-layer critic and a 4-layer generator. For CelebA, we used a 5-layer DCGAN discriminator and a 10-layer ResNet generator as in [7]. For ImageNet, we used a 10-layer ResNet for both the generator and discriminator. In all experiments we used 64 filters for the smallest convolutional layer, and double it at each layer (CelebA/ImageNet) or every other layer (CIFAR-10). The input codes for the generator are drawn from Uniform $([-1, 1]^{128})$. We consider two parameterizations for each critic: a standard one where the parameters can take any real value, and a spectral parametrization (denoted SN-) as above [32]. Models without explicit gradient control (SN-GAN, SN-MMDGAN, SN-MMGAN-L2, SN-WGAN) fix $\gamma = 1$, for spectral normalization; others learn γ , using a spectral parameterization.

Training All models were trained for 150 000 generator updates on a single GPU, except for ImageNet where the model was trained on 3 GPUs simultaneously. To limit communication overhead we averaged the MMD estimate on each GPU, giving the block MMD estimator [54]. We always used 64 samples per GPU from each of \mathbb{P} and \mathbb{Q} , and 5 critic updates per generator step. We used initial learning rates of 0.0001 for CIFAR-10 and CelebA, 0.0002 for ImageNet, and decayed these rates using the KID adaptive scheme of [7]: every 2 000 steps, generator samples are compared to those from 20 000 steps ago, and if the relative KID test [9] fails to show an improvement three consecutive times, the learning rate is decayed by 0.8. We used the Adam optimizer [25] with $\beta_1 = 0.5$, $\beta_2 = 0.9$.

Evaluation To compare the sample quality of different models, we considered three different scores based on the Inception network [49] trained for ImageNet classification, all using default parameters in the implementation of [7]. The *Inception Score (IS)* [42] is based on the entropy of predicted labels; higher values are better. Though standard, this metric has many issues, particularly on datasets other than ImageNet [4, 7, 20]. The *FID* [20] instead measures the similarity of samples from the generator and the target as the Wasserstein-2 distance between Gaussians fit to their intermediate representations. It is more sensible than the IS and becoming standard, but its estimator is strongly biased [7]. The *KID* [7] is similar to FID, but by using a polynomial-kernel MMD its estimates enjoy better statistical properties and are easier to compare. (A similar score was recommended by [21].)

Results Table 1a presents the scores for models trained on both CIFAR-10 and CelebA datasets. On CIFAR-10, SN-SWGAN and SN-SMMDGAN performed comparably to SN-GAN. But on CelebA, SN-SWGAN and SN-SMMDGAN dramatically outperformed the other methods with the same architecture in all three metrics. It also trained faster, and consistently outperformed other methods over multiple initializations (Figure 2 (a)). It is worth noting that SN-SWGAN far outperformed WGAN-GP on both datasets. Table 1b presents the scores for SMMDGAN and SN-SMMDGAN trained on ImageNet, and the scores of pre-trained models using BGAN [6] and SN-GAN [32].⁶ The

⁶These models are courtesy of the respective authors and also trained at 64×64 resolution. SN-GAN used the same architecture as our model, but trained for 250 000 generator iterations; BS-GAN used a similar 5-layer ResNet architecture and trained for 74 epochs, comparable to SN-GAN.

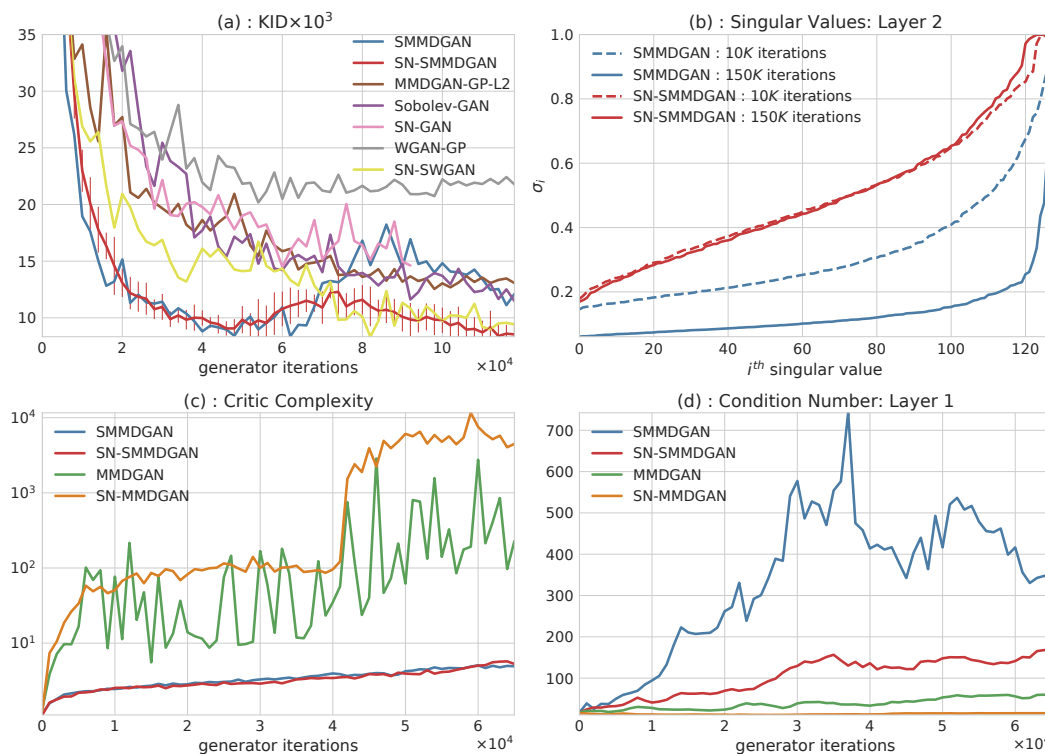


Figure 2: The training process on CelebA. (a) KID scores. We report a final score for SN-GAN slightly before its sudden failure mode; MMDGAN and SN-MMDGAN were unstable and had scores around 100. (b) Singular values of the second layer, both early (dashed) and late (solid) in training. (c) $\sigma_{\mu,k,\lambda}^{-2}$ for several MMD-based methods. (d) The condition number in the first layer through training. SN alone does not control $\sigma_{\mu,k,\lambda}$, and SMMD alone does not control the condition number.

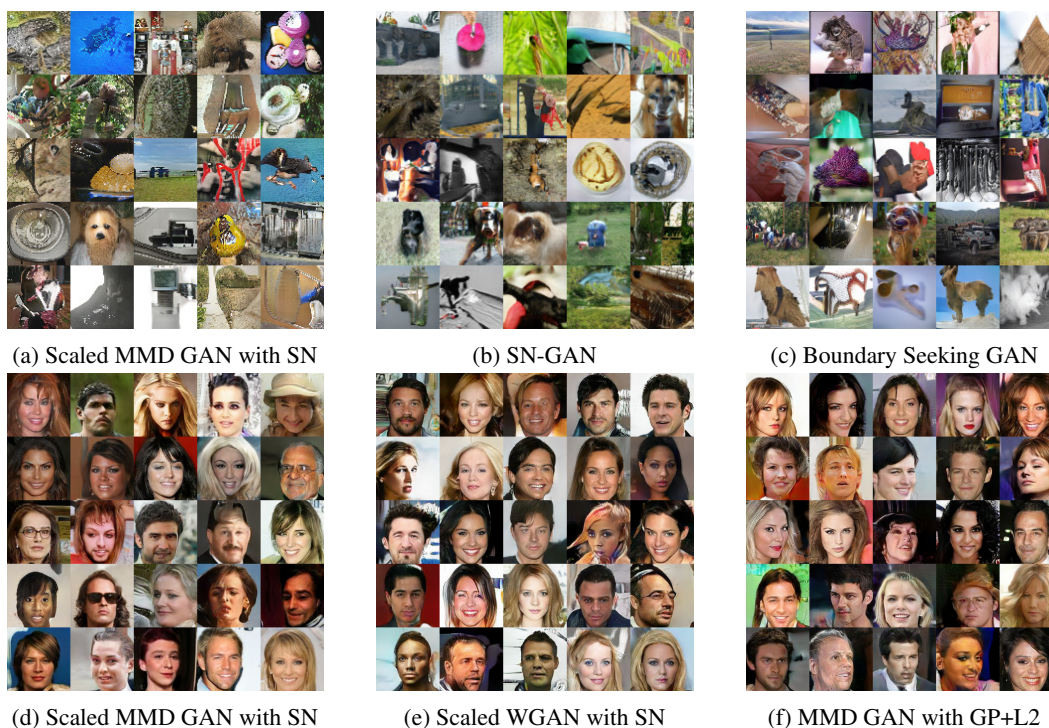


Figure 3: Samples from various models. Top: 64 \times 64 ImageNet; bottom: 160 \times 160 CelebA.

Table 1: Mean (standard deviation) of score estimates, based on 50 000 samples from each model.

(a) CIFAR-10 and CelebA.

Method	CIFAR-10			CelebA		
	IS	FID	KID $\times 10^3$	IS	FID	KID $\times 10^3$
WGAN-GP	6.9 \pm 0.2	31.1 \pm 0.2	22.2 \pm 1.1	2.7 \pm 0.0	29.2 \pm 0.2	22.0 \pm 1.0
MMDGAN-GP-L2	6.9 \pm 0.1	31.4 \pm 0.3	23.3 \pm 1.1	2.6 \pm 0.0	20.5 \pm 0.2	13.0 \pm 1.0
Sobolev-GAN	7.0 \pm 0.1	30.3 \pm 0.3	22.3 \pm 1.2	2.9\pm0.0	16.4 \pm 0.1	10.6 \pm 0.5
SMMDGAN	7.0 \pm 0.1	31.5 \pm 0.4	22.2 \pm 1.1	2.7 \pm 0.0	18.4 \pm 0.2	11.5 \pm 0.8
SN-GAN	7.2\pm0.1	26.7 \pm 0.2	16.1\pm0.9	2.7 \pm 0.0	22.6 \pm 0.1	14.6 \pm 1.1
SN-SWGAN	7.2\pm0.1	28.5 \pm 0.2	17.6\pm1.1	2.8 \pm 0.0	14.1 \pm 0.2	7.7 \pm 0.5
SN-SMMDGAN	7.3\pm0.1	25.0\pm0.3	16.6\pm2.0	2.8 \pm 0.0	12.4\pm0.2	6.1\pm0.4

(b) ImageNet.

Method	IS	FID	KID $\times 10^3$
BGAN	10.7 \pm 0.4	43.9 \pm 0.3	47.0 \pm 1.1
SN-GAN	11.2\pm0.1	47.5 \pm 0.1	44.4 \pm 2.2
SMMDGAN	10.7 \pm 0.2	38.4 \pm 0.3	39.3 \pm 2.5
SN-SMMDGAN	10.9 \pm 0.1	36.6\pm0.2	34.6\pm1.6

proposed methods substantially outperformed both methods in FID and KID scores. Figure 3 shows samples on ImageNet and CelebA; Appendix F.4 has more.

Spectrally normalized WGANs / MMDGANs To control for the contribution of the spectral parametrization to the performance, we evaluated variants of MMDGANs, WGANs and Sobolev-GAN using spectral normalization (in Table 2, Appendix F.3). WGAN and Sobolev-GAN led to unstable training and didn’t converge at all (Figure 11) despite many attempts. MMDGAN converged on CIFAR-10 (Figure 11) but was unstable on CelebA (Figure 10). The gradient control due to SN is thus probably too loose for these methods. This is reinforced by Figure 2 (c), which shows that the expected gradient of the critic network is much better-controlled by SMMD, even when SN is used. We also considered variants of these models with a learned γ while also adding a gradient penalty and an L_2 penalty on critic activations [7, footnote 19]. These generally behaved similarly to MMDGAN, and didn’t lead to substantial improvements. We ran the same experiments on CelebA, but aborted the runs early when it became clear that training was not successful.

Rank collapse We occasionally observed the failure mode for SMMD where the critic becomes low-rank, discussed in Section 3.3, especially on CelebA; this failure was obvious even in the training objective. Figure 2 (b) is one of these examples. Spectral parametrization seemed to prevent this behavior. We also found one could avoid collapse by reverting to an earlier checkpoint and increasing the RKHS regularization parameter λ , but did not do this for any of the experiments here.

5 Conclusion

We studied gradient regularization for MMD-based critics in implicit generative models, clarifying how previous techniques relate to the $\mathcal{D}_{\text{MMD}}^\Psi$ loss. Based on these insights, we proposed the Gradient-Constrained MMD and its approximation the Scaled MMD, a new loss function for IGMs that controls gradient behavior in a principled way and obtains excellent performance in practice.

One interesting area of future study for these distances is their behavior when used to diffuse particles distributed as \mathbb{Q} towards particles distributed as \mathbb{P} . Mroueh et al. [33, Appendix A.1] began such a study for the Sobolev GAN loss; [35] proved convergence and studied discrete-time approximations.

Another area to explore is the geometry of these losses, as studied by Bottou et al. [8], who showed potential advantages of the Wasserstein geometry over the MMD. Their results, though, do not address any distances based on optimized kernels; the new distances introduced here might have interesting geometry of their own.

References

- [1] B. Amos and J. Z. Kolter. “OptNet: Differentiable Optimization as a Layer in Neural Networks.” *ICML*. 2017. arXiv: [1703.00443](#).
- [2] M. Arjovsky and L. Bottou. “Towards Principled Methods for Training Generative Adversarial Networks.” *ICLR*. 2017. arXiv: [1701.04862](#).
- [3] M. Arjovsky, S. Chintala, and L. Bottou. “Wasserstein Generative Adversarial Networks.” *ICML*. 2017. arXiv: [1701.07875](#).
- [4] S. Barratt and R. Sharma. *A Note on the Inception Score*. 2018. arXiv: [1801.01973](#).
- [5] M. G. Bellemare, I. Danihelka, W. Dabney, S. Mohamed, B. Lakshminarayanan, S. Hoyer, and R. Munos. *The Cramer Distance as a Solution to Biased Wasserstein Gradients*. 2017. arXiv: [1705.10743](#).
- [6] D. Berthelot, T. Schumm, and L. Metz. *BEGAN: Boundary Equilibrium Generative Adversarial Networks*. 2017. arXiv: [1703.10717](#).
- [7] M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton. “Demystifying MMD GANs.” *ICLR*. 2018. arXiv: [1801.01401](#).
- [8] L. Bottou, M. Arjovsky, D. Lopez-Paz, and M. Oquab. “Geometrical Insights for Implicit Generative Modeling.” *Braverman Readings in Machine Learning: Key Ideas from Inception to Current State*. Ed. by L. Rozonoer, B. Mirkin, and I. Muchnik. LNAI Vol. 11100. Springer, 2018, pp. 229–268. arXiv: [1712.07822](#).
- [9] W. Bounliphone, E. Belilovsky, M. B. Blaschko, I. Antonoglou, and A. Gretton. “A Test of Relative Similarity For Model Selection in Generative Models.” *ICLR*. 2016. arXiv: [1511.04581](#).
- [10] O. Bousquet, O. Chapelle, and M. Hein. “Measure Based Regularization.” *NIPS*. 2004.
- [11] A. Brock, T. Lim, J. M. Ritchie, and N. Weston. “Neural Photo Editing with Introspective Adversarial Networks.” *ICLR*. 2017. arXiv: [1609.07093](#).
- [12] R. M. Dudley. *Real Analysis and Probability*. 2nd ed. Cambridge University Press, 2002.
- [13] G. K. Dziugaite, D. M. Roy, and Z. Ghahramani. “Training generative neural networks via Maximum Mean Discrepancy optimization.” *UAI*. 2015. arXiv: [1505.03906](#).
- [14] A. Genevay, G. Peyré, and M. Cuturi. “Learning Generative Models with Sinkhorn Divergences.” *AISTATS*. 2018. arXiv: [1706.00292](#).
- [15] T. Gneiting and A. E. Raftery. “Strictly proper scoring rules, prediction, and estimation.” *JASA* 102.477 (2007), pp. 359–378.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. “Generative Adversarial Nets.” *NIPS*. 2014. arXiv: [1406.2661](#).
- [17] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. J. Smola. “A Kernel Two-Sample Test.” *JMLR* 13 (2012).
- [18] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville. “Improved Training of Wasserstein GANs.” *NIPS*. 2017. arXiv: [1704.00028](#).
- [19] A. Güngör. “Some bounds for the product of singular values.” *International Journal of Contemporary Mathematical Sciences* (2007).
- [20] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, G. Klambauer, and S. Hochreiter. “GANs Trained by a Two Time-Scale Update Rule Converge to a Nash Equilibrium.” *NIPS*. 2017. arXiv: [1706.08500](#).
- [21] G. Huang, Y. Yuan, Q. Xu, C. Guo, Y. Sun, F. Wu, and K. Weinberger. *An empirical study on evaluation metrics of generative adversarial networks*. 2018. arXiv: [1806.07755](#).
- [22] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz. “Multimodal Unsupervised Image-to-Image Translation.” *ECCV*. 2018. arXiv: [1804.04732](#).
- [23] Y. Jin, K. Zhang, M. Li, Y. Tian, H. Zhu, and Z. Fang. *Towards the Automatic Anime Characters Creation with Generative Adversarial Networks*. 2017. arXiv: [1708.05509](#).
- [24] T. Karras, T. Aila, S. Laine, and J. Lehtinen. “Progressive Growing of GANs for Improved Quality, Stability, and Variation.” *ICLR*. 2018. arXiv: [1710.10196](#).
- [25] D. Kingma and J. Ba. “Adam: A Method for Stochastic Optimization.” *ICLR*. 2015. arXiv: [1412.6980](#).
- [26] A. Krizhevsky. *Learning Multiple Layers of Features from Tiny Images*. 2009.

- [27] C.-L. Li, W.-C. Chang, Y. Cheng, Y. Yang, and B. Póczos. “MMD GAN: Towards Deeper Understanding of Moment Matching Network.” *NIPS*. 2017. arXiv: [1705.08584](#).
- [28] Y. Li, K. Swersky, and R. Zemel. “Generative Moment Matching Networks.” *ICML*. 2015. arXiv: [1502.02761](#).
- [29] Z. Liu, P. Luo, X. Wang, and X. Tang. “Deep learning face attributes in the wild.” *ICCV*. 2015. arXiv: [1411.7766](#).
- [30] L. Mescheder, A. Geiger, and S. Nowozin. “Which Training Methods for GANs do actually Converge?” *ICML*. 2018. arXiv: [1801.04406](#).
- [31] P. Milgrom and I. Segal. “Envelope theorems for arbitrary choice sets.” *Econometrica* 70.2 (2002), pp. 583–601.
- [32] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. “Spectral Normalization for Generative Adversarial Networks.” *ICLR*. 2018. arXiv: [1802.05927](#).
- [33] Y. Mroueh, C.-L. Li, T. Sercu, A. Raj, and Y. Cheng. “Sobolev GAN.” *ICLR*. 2018. arXiv: [1711.04894](#).
- [34] Y. Mroueh and T. Sercu. “Fisher GAN.” *NIPS*. 2017. arXiv: [1705.09675](#).
- [35] Y. Mroueh, T. Sercu, and A. Raj. *Regularized Kernel and Neural Sobolev Descent: Dynamic MMD Transport*. 2018. arXiv: [1805.12062](#).
- [36] A. Müller. “Integral Probability Metrics and their Generating Classes of Functions.” *Advances in Applied Probability* 29.2 (1997), pp. 429–443.
- [37] S. Nowozin, B. Cseke, and R. Tomioka. “f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization.” *NIPS*. 2016. arXiv: [1606.00709](#).
- [38] A. Radford, L. Metz, and S. Chintala. “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks.” *ICLR*. 2016. arXiv: [1511.06434](#).
- [39] J. R. Retherford. “Review: J. Diestel and J. J. Uhl, Jr., Vector measures.” *Bull. Amer. Math. Soc.* 84.4 (July 1978), pp. 681–685.
- [40] K. Roth, A. Lucchi, S. Nowozin, and T. Hofmann. “Stabilizing Training of Generative Adversarial Networks through Regularization.” *NIPS*. 2017. arXiv: [1705.09367](#).
- [41] O. Russakovsky et al. *ImageNet Large Scale Visual Recognition Challenge*. 2014. arXiv: [1409.0575](#).
- [42] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. “Improved Techniques for Training GANs.” *NIPS*. 2016. arXiv: [1606.03498](#).
- [43] J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [44] B. K. Sriperumbudur, K. Fukumizu, A. Gretton, G. R. G. Lanckriet, and B. Schölkopf. “Kernel choice and classifiability for RKHS embeddings of probability distributions.” *NIPS*. 2009.
- [45] B. K. Sriperumbudur, K. Fukumizu, and G. R. G. Lanckriet. “Universality, Characteristic Kernels and RKHS Embedding of Measures.” *JMLR* 12 (2011), pp. 2389–2410. arXiv: [1003.0887](#).
- [46] B. Sriperumbudur. “On the optimal estimation of probability measures in weak and strong topologies.” *Bernoulli* 22.3 (2016), pp. 1839–1893. arXiv: [1310.8240](#).
- [47] I. Steinwart and A. Christmann. *Support Vector Machines*. Springer, 2008.
- [48] D. J. Sutherland, H.-Y. Tung, H. Strathmann, S. De, A. Ramdas, A. Smola, and A. Gretton. “Generative Models and Model Criticism via Optimized Maximum Mean Discrepancy.” *ICLR*. 2017. arXiv: [1611.04488](#).
- [49] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. “Rethinking the Inception Architecture for Computer Vision.” *CVPR*. 2016. arXiv: [1512.00567](#).
- [50] T. Unterthiner, B. Nessler, C. Seward, G. Klambauer, M. Heusel, H. Ramsauer, and S. Hochreiter. “Coulomb GANs: Provably Optimal Nash Equilibria via Potential Fields.” *ICLR*. 2018. arXiv: [1708.08819](#).
- [51] C. Villani. *Optimal Transport: Old and New*. Springer, 2009.
- [52] J. Weed and F. Bach. “Sharp asymptotic and finite-sample rates of convergence of empirical measures in Wasserstein distance.” *Bernoulli* (forthcoming). arXiv: [1707.00087](#).
- [53] H. Wendland. *Scattered Data Approximation*. Cambridge University Press, 2005.
- [54] W. Zaremba, A. Gretton, and M. B. Blaschko. “B-tests: Low Variance Kernel Two-Sample Tests.” *NIPS*. 2013. arXiv: [1307.1954](#).

A Proofs

We first review some basic properties of Reproducing Kernel Hilbert Spaces. We consider here a separable RKHS \mathcal{H} with basis $(e_i)_{i \in I}$, where I is either finite if \mathcal{H} is finite-dimensional, or $I = \mathbb{N}$ otherwise. We also assume that the reproducing kernel k is continuously twice differentiable.

We use a slightly nonstandard notation for derivatives: $\partial_i f(x)$ denotes the i th partial derivative of f evaluated at x , and $\partial_i \partial_{j+d} k(x, y)$ denotes $\frac{\partial^2 k(a,b)}{\partial a_i \partial b_j} \Big|_{(a,b)=(x,y)}$.

Then the following reproducing properties hold for any given function f in \mathcal{H} [47, Lemma 4.34]:

$$f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}} \quad (8)$$

$$\partial_i f(x) = \langle f, \partial_i k(x, \cdot) \rangle_{\mathcal{H}}. \quad (9)$$

We say that an operator $A : \mathcal{H} \mapsto \mathcal{H}$ is Hilbert-Schmidt if $\|A\|_{HS}^2 = \sum_{i \in I} \|Ae_i\|_{\mathcal{H}}^2$ is finite. $\|A\|_{HS}$ is called the Hilbert-Schmidt norm of A . The space of Hilbert-Schmidt operators itself a Hilbert space with the inner product $\langle A, B \rangle_{HS} = \sum_{i \in I} \langle Ae_i, Be_i \rangle_{\mathcal{H}}$. Moreover, we say that an operator A is trace-class if its trace norm is finite, i.e. $\|A\|_1 = \sum_{i \in I} \langle e_i, (A^* A)^{\frac{1}{2}} e_i \rangle_{\mathcal{H}} < \infty$. The outer product $f \otimes g$ for $f, g \in \mathcal{H}$ gives an $\mathcal{H} \rightarrow \mathcal{H}$ operator such that $(f \otimes g)v = \langle g, v \rangle_{\mathcal{H}} f$ for all v in \mathcal{H} .

Given two vectors f and g in \mathcal{H} and a Hilbert-Schmidt operator A we have the following properties:

- (i) The outer product $f \otimes g$ is a Hilbert-Schmidt operator with Hilbert-Schmidt norm given by: $\|f \otimes g\|_{HS} = \|f\|_{\mathcal{H}} \|g\|_{\mathcal{H}}$.
- (ii) The inner product between two rank-one operators $f \otimes g$ and $u \otimes v$ is $\langle f \otimes g, u \otimes v \rangle_{HS} = \langle f, u \rangle_{\mathcal{H}} \langle g, v \rangle_{\mathcal{H}}$.
- (iii) The following identity holds: $\langle f, Ag \rangle_{\mathcal{H}} = \langle f \otimes g, A \rangle_{HS}$.

Define the following covariance-type operators:

$$D_x = k(x, \cdot) \otimes k(x, \cdot) + \sum_{i=1}^d \partial_i k(x, \cdot) \otimes \partial_i k(x, \cdot) \quad D_\mu = \mathbb{E}_{X \sim \mu} D_X \quad D_{\mu, \lambda} = D_\mu + \lambda I; \quad (10)$$

these are useful in that, using (8) and (9), $\langle f, D_x g \rangle = f(x)g(x) + \sum_{i=1}^d \partial_i f(x) \partial_i g(x)$.

A.1 Definitions and estimators of the new distances

We will need the following assumptions about the distributions \mathbb{P} and \mathbb{Q} , the measure μ , and the kernel k :

- (A) \mathbb{P} and \mathbb{Q} have integrable first moments.
- (B) $\sqrt{k(x, x)}$ grows at most linearly in x : for all x in \mathcal{X} , $\sqrt{k(x, x)} \leq C(\|x\| + 1)$ for some constant C .
- (C) The kernel k is twice continuously differentiable.
- (D) The functions $x \mapsto k(x, x)$ and $x \mapsto \partial_i \partial_{i+d} k(x, x)$ for $1 \leq i \leq d$ are μ -integrable.

When $k = K \circ \phi_\psi$, Assumption (B) is automatically satisfied by a K such as the Gaussian; when K is linear, it is true for a quite general class of networks ϕ_ψ [7, Lemma 1].

We will first give a form for the Gradient-Constrained MMD (5) in terms of the operator (10):

Proposition 5. *Under Assumptions (A) to (D), the Gradient-Constrained MMD is given by*

$$\text{GCMMD}_{\mu, k, \lambda}(\mathbb{P}, \mathbb{Q}) = \sqrt{\langle \eta, D_{\mu, \lambda}^{-1} \eta \rangle_{\mathcal{H}}}. \quad (11)$$

Proof of Proposition 5. Let f be a function in \mathcal{H} . We will first express the squared λ -regularized Sobolev norm of f (6) as a quadratic form in \mathcal{H} . Recalling the reproducing properties of (8) and (9), we have:

$$\|f\|_{S(\mu), k, \lambda}^2 = \int \langle f, k(x, \cdot) \rangle_{\mathcal{H}}^2 \mu(dx) + \sum_{i=1}^d \int \langle f, \partial_i k(x, \cdot) \rangle_{\mathcal{H}}^2 \mu(dx) + \lambda \|f\|_{\mathcal{H}}^2.$$

Using Property (ii) and the operator (10), one further gets

$$\|f\|_{S(\mu),k,\lambda}^2 = \int \langle f \otimes f, D_x \rangle_{\text{HS}} \mu(dx) + \lambda \|f\|_{\mathcal{H}}^2.$$

Under Assumption (D), and using Lemma 6, one can take the integral inside the inner product, which leads to $\|f\|_{S(\mu),k,\lambda}^2 = \langle f \otimes f, D_\mu \rangle_{\text{HS}} + \lambda \|f\|_{\mathcal{H}}^2$. Finally, using Property (iii) it follows that

$$\|f\|_{S(\mu),k,\lambda}^2 = \langle f, D_{\mu,\lambda} f \rangle_{\mathcal{H}}.$$

Under Assumptions (A) and (B), Lemma 6 applies, and it follows that $k(x, \cdot)$ is also Bochner integrable under \mathbb{P} and \mathbb{Q} . Thus

$$\mathbb{E}_{\mathbb{P}}[\langle f, k(x, \cdot) \rangle_{\mathcal{H}}] - \mathbb{E}_{\mathbb{Q}}[\langle f, k(x, \cdot) \rangle_{\mathcal{H}}] = \langle f, \mathbb{E}_{\mathbb{P}}[k(x, \cdot)] - \mathbb{E}_{\mathbb{Q}}[k(x, \cdot)] \rangle_{\mathcal{H}} = \langle f, \eta \rangle_{\mathcal{H}},$$

where η is defined as this difference in mean embeddings.

Since $D_{\mu,\lambda}$ is symmetric positive definite, its square-root $D_{\mu,\lambda}^{\frac{1}{2}}$ is well-defined and is also invertible. For any $f \in \mathcal{H}$, let $g = D_{\mu,\lambda}^{\frac{1}{2}} f$, so that $\langle f, D_{\mu,\lambda} f \rangle_{\mathcal{H}} = \|g\|_{\mathcal{H}}^2$. Note that for any $g \in \mathcal{H}$, there is a corresponding $f = D_{\mu,\lambda}^{-\frac{1}{2}} g$. Thus we can re-express the maximization problem in (5) in terms of g :

$$\begin{aligned} \text{GCMMD}_{\mu,k,\lambda}(\mathbb{P}, \mathbb{Q}) &:= \sup_{\substack{f \in \mathcal{H} \\ \langle f, D_{\mu,\lambda} f \rangle_{\mathcal{H}} \leq 1}} \langle f, \eta \rangle_{\mathcal{H}} = \sup_{\substack{g \in \mathcal{H} \\ \|g\|_{\mathcal{H}} \leq 1}} \langle D_{\mu,\lambda}^{-\frac{1}{2}} g, \eta \rangle_{\mathcal{H}} \\ &= \sup_{\substack{g \in \mathcal{H} \\ \|g\|_{\mathcal{H}} \leq 1}} \langle g, D_{\mu,\lambda}^{-\frac{1}{2}} \eta \rangle_{\mathcal{H}} = \|D_{\mu,\lambda}^{-\frac{1}{2}} \eta\|_{\mathcal{H}} = \sqrt{\langle \eta, D_{\mu,\lambda}^{-1} \eta \rangle_{\mathcal{H}}}. \quad \square \end{aligned}$$

Proposition 5, though, involves inverting the infinite-dimensional operator $D_{\mu,\lambda}$ and thus doesn't directly give us a computable estimator. Proposition 3 solves this problem in the case where μ is a discrete measure:

Proposition 3. *Let $\hat{\mu} = \sum_{m=1}^M \delta_{X_m}$ be an empirical measure of M points. Let $\eta(X) \in \mathbb{R}^M$ have m th entry $\eta(X_m)$, and $\nabla \eta(X) \in \mathbb{R}^{M+d}$ have (m, i) th entry⁷ $\partial_i \eta(X_m)$. Then under Assumptions (A) to (D), the Gradient-Constrained MMD is*

$$\begin{aligned} \text{GCMMD}_{\hat{\mu},k,\lambda}^2(\mathbb{P}, \mathbb{Q}) &= \frac{1}{\lambda} (\text{MMD}^2(\mathbb{P}, \mathbb{Q}) - \bar{P}(\eta)) \\ \bar{P}(\eta) &= \begin{bmatrix} \eta(X) \\ \nabla \eta(X) \end{bmatrix}^{\text{T}} \left(\begin{bmatrix} K & G^{\text{T}} \\ G & H \end{bmatrix} + M\lambda I_{M+d} \right)^{-1} \begin{bmatrix} \eta(X) \\ \nabla \eta(X) \end{bmatrix}, \end{aligned}$$

where K is the kernel matrix $K_{m,m'} = k(X_m, X_{m'})$, G is the matrix of left derivatives $G_{(m,i),m'} = \partial_i k(X_m, X_{m'})$, and H that of derivatives of both arguments $H_{(m,i),(m',j)} = \partial_i \partial_j k(X_m, X_{m'})$.

Before proving Proposition 3, we note the following interesting alternate form. Let \bar{e}_i be the i th standard basis vector for \mathbb{R}^{M+d} , and define $T : \mathcal{H} \rightarrow \mathbb{R}^{M+d}$ as the linear operator

$$T = \sum_{m=1}^M \bar{e}_m \otimes k(X_m, \cdot) + \sum_{m=1}^M \sum_{i=1}^d \bar{e}_{m+(m,i)} \otimes \partial_i k(X_m, \cdot).$$

Then $\begin{bmatrix} \eta(X) \\ \nabla \eta(X) \end{bmatrix} = T\eta$, and $\begin{bmatrix} K & G^{\text{T}} \\ G & H \end{bmatrix} = TT^*$. Thus we can write

$$\text{GCMMD}_{\hat{\mu},k,\lambda}^2 = \frac{1}{\lambda} \langle \eta, (I - T^*(TT^* + M\lambda I)^{-1}T) \eta \rangle_{\mathcal{H}}.$$

⁷We use (m, i) to denote $(m-1)d + i$; thus $\nabla \eta(X)$ stacks $\nabla \eta(X_1), \dots, \nabla \eta(X_M)$ into one vector.

Proof of Proposition 3. Let $g \in \mathcal{H}$ be the solution to the regression problem $D_{\mu,\lambda}g = \eta$:

$$\begin{aligned} & \frac{1}{M} \sum_{m=1}^M \left[g(X_m)k(X_m, \cdot) + \sum_{i=1}^d \partial_i g(X_m) \partial_i k(X_m, \cdot) \right] + \lambda g = \eta \\ g &= \frac{1}{\lambda} \eta - \frac{1}{\lambda M} \sum_{m=1}^M \left[g(X_m)k(X_m, \cdot) + \sum_{i=1}^d \partial_i g(X_m) \partial_i k(X_m, \cdot) \right]. \end{aligned} \quad (12)$$

Taking the inner product of both sides of (12) with $k(X_{m'}, \cdot)$ for each $1 \leq m' \leq M$ yields the following M equations:

$$g(X_{m'}) = \frac{1}{\lambda} \eta(X_{m'}) - \frac{1}{\lambda M} \sum_{m=1}^M \left[g(X_m) K_{m,m'} + \sum_{i=1}^d \partial_i g(X_m) G_{(m,i),m'} \right]. \quad (13)$$

Doing the same with $\partial_j k(X_{m'}, \cdot)$ gives Md equations:

$$\partial_j g(X_{m'}) = \frac{1}{\lambda} \partial_j \eta(X_{m'}) - \frac{1}{\lambda M} \sum_{m=1}^M \left[g(X_m) G_{(m',j),m} + \sum_{i=1}^d \partial_i g(X_m) H_{(m,i),(m',j)} \right]. \quad (14)$$

From (12), it is clear that g is a linear combination of the form:

$$g(x) = \frac{1}{\lambda} \eta(x) - \frac{1}{\lambda M} \sum_{m=1}^M \left[\alpha_m k(X_m, x) + \sum_{i=1}^d \beta_{m,i} \partial_i k(X_m, x) \right],$$

where the coefficients $\alpha := (\alpha_m = g(X_m))_{1 \leq m \leq M}$ and $\beta := (\beta_{m,i} = \partial_i g(X_m))_{\substack{1 \leq m \leq M \\ 1 \leq i \leq d}}$ satisfy the system of equations (13) and (14). We can rewrite this system as

$$\begin{bmatrix} K + M\lambda I_M & G^T \\ G & H + M\lambda I_{Md} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = M \begin{bmatrix} \eta(X) \\ \nabla \eta(X) \end{bmatrix},$$

where I_M, I_{Md} are the identity matrices of dimension M, Md . Since K and H must be positive semidefinite, an inverse exists. We conclude by noticing that

$$\text{GCMMD}_{\hat{\mu},k,\lambda}(\mathbb{P}, \mathbb{Q})^2 = \langle \eta, g \rangle_{\mathcal{H}} = \frac{1}{\lambda} \|\eta\|_{\mathcal{H}}^2 - \frac{1}{\lambda M} \sum_{m=1}^M \left[\alpha_m \eta(X_m) + \sum_{i=1}^d \beta_{m,i} \partial_i \eta(X_m) \right]. \quad \square$$

The following result was key to our definition of the SMMD in Section 3.3.

Proposition 4. *Under Assumptions (A) to (D), we have for all $f \in \mathcal{H}$ that*

$$\|f\|_{S(\mu),k,\lambda} \leq \sigma_{\mu,k,\lambda}^{-1} \|f\|_{\mathcal{H}_k},$$

where $\sigma_{k,\mu,\lambda} := 1/\sqrt{\lambda + \int k(x,x)\mu(dx) + \sum_{i=1}^d \int \partial_i \partial_{i+d} k(x,x)\mu(dx)}$.

Proof of Proposition 4. The key idea here is to use the Cauchy-Schwarz inequality for the Hilbert-Schmidt inner product. Letting $f \in \mathcal{H}$, $\|f\|_{S(\mu),k,\lambda}^2$ is

$$\begin{aligned} & \int f(x)^2 \mu(dx) + \int \|\nabla f(x)\|^2 \mu(dx) + \lambda \|f\|_{\mathcal{H}}^2 \\ & \stackrel{(a)}{=} \int \langle f, k(x, \cdot) \otimes k(x, \cdot) f \rangle_{\mathcal{H}} \mu(dx) + \sum_{i=1}^d \int \langle f, \partial_i k(x, \cdot) \otimes \partial_i k(x, \cdot) f \rangle_{\mathcal{H}} \mu(dx) + \lambda \|f\|_{\mathcal{H}}^2 \\ & \stackrel{(b)}{=} \int \langle f \otimes f, k(x, \cdot) \otimes k(x, \cdot) \rangle_{\text{HS}} \mu(dx) + \sum_{i=1}^d \int \langle f \otimes f, \partial_i k(x, \cdot) \otimes \partial_i k(x, \cdot) \rangle_{\text{HS}} \mu(dx) + \lambda \|f\|_{\mathcal{H}}^2 \\ & \stackrel{(c)}{\leq} \|f\|_{\mathcal{H}}^2 \left[\int k(x,x) \mu(dx) + \sum_{i=1}^d \int \partial_i \partial_{i+d} k(x,x) \mu(dx) + \lambda \right]. \end{aligned}$$

(a) follows from the reproducing properties (8) and (9) and Property (ii). (b) is obtained using Property (iii), while (c) follows from the Cauchy-Schwarz inequality and Property (i). \square

Lemma 6. Under Assumption **(D)**, D_x is Bochner integrable and its integral D_μ is a trace-class symmetric positive semi-definite operator with $D_{\mu,\lambda} = D + \lambda I$ invertible for any positive λ . Moreover, for any Hilbert-Schmidt operator A we have: $\langle A, D_\mu \rangle_{HS} = \int \langle A, D_x \rangle_{HS} \mu(dx)$.

Under Assumptions **(A)** and **(B)**, $k(x, \cdot)$ is Bochner integrable with respect to any probability distribution \mathbb{P} with finite first moment and the following relation holds: $\langle f, \mathbb{E}_{\mathbb{P}}[k(x, \cdot)] \rangle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}}[\langle f, k(x, \cdot) \rangle_{\mathcal{H}}]$ for all f in \mathcal{H} .

Proof. The operator D_x is positive self-adjoint. It is also trace-class, as by the triangle inequality

$$\begin{aligned} \|D_x\|_1 &\leq \|k(x, \cdot) \otimes k(x, \cdot)\|_1 + \sum_{i=1}^d \|\partial_i k(x, \cdot) \otimes \partial_i k(x, \cdot)\|_1 \\ &= \|k(x, \cdot)\|_{\mathcal{H}}^2 + \sum_{i=1}^d \|\partial_i k(x, \cdot)\|_{\mathcal{H}}^2 < \infty. \end{aligned}$$

By Assumption **(D)**, we have that $\int \|D_x\|_1 \mu(dx) < \infty$ which implies that D_x is μ -integrable in the Bochner sense [39, Definition 1 and Theorem 2]. Its integral D_μ is trace-class and satisfies $\|D_\mu\|_1 \leq \int \|D_x\|_1 \mu(dx)$. This allows to have $\langle A, D_\mu \rangle_{HS} = \int \langle A, D_x \rangle_{HS} \mu(dx)$ for all Hilbert-Schmidt operators A . Moreover, the integral preserves the symmetry and positivity. It follows that $D_{\mu,\lambda}$ is invertible.

The Bochner integrability of $k(x, \cdot)$ under a distribution \mathbb{P} with finite moment follows directly from Assumptions **(A)** and **(B)**, since $\int \|k(x, \cdot)\| \mathbb{P}(dx) \leq C \int (\|x\| + 1) \mathbb{P}(dx) < \infty$. This allows us to write $\langle f, \mathbb{E}_{\mathbb{P}}[k(x, \cdot)] \rangle_{\mathcal{H}} = \mathbb{E}_{\mathbb{P}}[\langle f, k(x, \cdot) \rangle_{\mathcal{H}}]$. \square

A.2 Continuity of the Optimized Scaled MMD in the Wasserstein topology

To prove Theorem 1, we will first need some new notation.

We assume the kernel is $k = K \circ \phi_\psi$, i.e. $k_\psi(x, y) = K(\phi_\psi(x), \phi_\psi(y))$, where the representation function ϕ_ψ is a network $\phi_\psi(X) : \mathbb{R}^d \rightarrow \mathbb{R}^{d_L}$ consisting of L fully-connected layers:

$$\begin{aligned} h_\psi^0(X) &= X \\ h_\psi^l(X) &= W^l \sigma_{l-1}(h_\psi^{l-1}(X)) + b^l \quad \text{for } 1 \leq l \leq L \\ \phi_\psi(X) &= h_\psi^L(X). \end{aligned} \tag{15}$$

The intermediate representations $h_\psi^l(X)$ are of dimension d_l , the weights W^l are matrices in $\mathbb{R}^{d_l \times d_{l-1}}$, and biases b^l are vectors in \mathbb{R}^{d_l} . The elementwise activation function σ is given by $\sigma_0(x) = x$, and for $l > 0$ the activation σ_l is a leaky ReLU with leak coefficient $0 < \alpha < 1$:

$$\sigma_l(x) = \sigma(x) = \begin{cases} x & x > 0 \\ \alpha x & x \leq 0 \end{cases} \quad \text{for } l > 0. \tag{16}$$

The parameter ψ is the concatenation of all the layer parameters:

$$\psi = ((W^L, b^L), (W^{L-1}, b^{L-1}), \dots, (W^1, b^1)).$$

We denote by Ψ the set of all such possible parameters, i.e. $\Psi = \mathbb{R}^{d_L \times d_{L-1}} \times \mathbb{R}^{d_L} \times \dots \times \mathbb{R}^{d_1 \times d} \times \mathbb{R}^{d_1}$. Define the following restrictions of Ψ :

$$\Psi^\kappa := \{\psi \in \Psi \mid \forall 1 \leq l \leq L, \text{cond}(W^l) \leq \kappa\} \tag{17}$$

$$\Psi_1^\kappa := \{\psi \in \Psi^\kappa \mid \forall 1 \leq l \leq L, \|W^l\| = 1\}. \tag{18}$$

Ψ^κ is the set of those parameters such that W^l have a small condition number, $\text{cond}(W) = \sigma_{\max}(W)/\sigma_{\min}(W)$. Ψ_1^κ is the set of per-layer normalized parameters with a condition number bounded by κ .

Recall the definition of Scaled MMD, (7), where $\lambda > 0$ and μ is a probability measure:

$$\text{SMMD}_{\mu,k,\lambda}(\mathbb{P}, \mathbb{Q}) := \sigma_{\mu,k,\lambda} \text{MMD}_k(\mathbb{P}, \mathbb{Q})$$

$$\sigma_{\mu,k,\lambda} := 1/\sqrt{\lambda + \int k(x, x) \mu(dx) + \sum_{i=1}^d \int \partial_i \partial_{i+d} k(x, x) \mu(dx)}.$$

The Optimized SMMD over the restricted set Ψ^κ is given by:

$$\mathcal{D}_{\text{SMMD}}^{\mu, \Psi^\kappa, \lambda}(\mathbb{P}, \mathbb{Q}) := \sup_{\psi \in \Psi^\kappa} \text{SMMD}_{\mu, k_\psi, \lambda}.$$

The constraint to $\psi \in \Psi^\kappa$ is critical to the proof. In practice, using a spectral parametrization helps enforce this assumption, as shown in Figures 2 and 9. Other regularization methods, like orthogonal normalization [11], are also possible.

We will use the following assumptions:

- (I) μ is a probability distribution absolutely continuous with respect to the Lebesgue measure.
- (II) The dimensions of the weights are decreasing per layer: $d_{l+1} \leq d_l$ for all $0 \leq l \leq L-1$.
- (III) The non-linearity used is Leaky-ReLU, (16), with leak coefficient $\alpha \in (0, 1)$.
- (IV) The top-level kernel K is globally Lipschitz in the RKHS norm: there exists a positive constant $L_K > 0$ such that $\|K(a, \cdot) - K(b, \cdot)\| \leq L_K \|a - b\|$ for all a and b in \mathbb{R}^{d_L} .
- (V) There is some $\gamma_K > 0$ for which K satisfies

$$\nabla_b \nabla_c K(b, c) \Big|_{(b,c)=(a,a)} \succeq \gamma^2 I \quad \text{for all } a \in \mathbb{R}^{d_L}. \quad (19)$$

Assumption (I) ensures that the points where $\phi_\psi(X)$ is not differentiable are reached with probability 0 under μ . This assumption can be easily satisfied e.g. if we define μ by adding Gaussian noise to \mathbb{P} .

Assumption (II) helps ensure that the span of W^l is never contained in the null space of W^{l+1} . Using Leaky-ReLU as a non-linearity, Assumption (III), further ensures that the network ϕ_ψ is locally full-rank almost everywhere; this might not be true with ReLU activations, where it could be always 0. Assumptions (II) and (III) can be easily satisfied by design of the network.

Assumptions (IV) and (V) only depend on the top-level kernel K and are easy to satisfy in practice. In particular, they always hold for a smooth translation-invariant kernel, such as the Gaussian, as well as the linear kernel.

We are now ready to prove Theorem 1.

Theorem 1. *Under Assumptions (I) to (V),*

$$\mathcal{D}_{\text{SMMD}}^{\mu, \Psi^\kappa, \lambda}(\mathbb{P}, \mathbb{Q}) \leq \frac{L_K \kappa^{L/2}}{\gamma \sqrt{d_L} \alpha^{L/2}} \mathcal{W}(\mathbb{P}, \mathbb{Q}),$$

which implies that if $\mathbb{P}_n \xrightarrow{\mathcal{W}} \mathbb{P}$, then $\mathcal{D}_{\text{SMMD}}^{\mu, \Psi^\kappa, \lambda}(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$.

Proof. Define the pseudo-distance corresponding to the kernel k_ψ

$$d_\psi(x, y) = \|k_\psi(x, \cdot) - k_\psi(y, \cdot)\|_{\mathcal{H}_\psi} = \sqrt{k_\psi(x, x) + k_\psi(y, y) - 2k_\psi(x, y)}.$$

Denote by $\mathcal{W}_{d_\psi}(\mathbb{P}, \mathbb{Q})$ the optimal transport metric between \mathbb{P} and \mathbb{Q} using the cost d_ψ , given by

$$\mathcal{W}_{d_\psi}(\mathbb{P}, \mathbb{Q}) = \inf_{\pi \in \Pi(\mathbb{P}, \mathbb{Q})} \mathbb{E}_{(X, Y) \sim \pi} [d_\psi(X, Y)].$$

where Π is the set of couplings with marginals \mathbb{P} and \mathbb{Q} . By Lemma 7,

$$\text{MMD}_\psi(\mathbb{P}, \mathbb{Q}) \leq \mathcal{W}_{d_\psi}(\mathbb{P}, \mathbb{Q}).$$

Recall that ϕ_ψ is Lipschitz, $\|\phi_\psi\|_{\text{Lip}} < \infty$, so along with Assumption (IV) we have that

$$d_\psi(x, y) \leq L_K \|\phi_\psi(x) - \phi_\psi(y)\| \leq L_K \|\phi_\psi\|_{\text{Lip}} \|x - y\|.$$

Thus

$$\mathcal{W}_{d_\psi}(\mathbb{P}, \mathbb{Q}) \leq \inf_{\pi \in \Pi(\mathbb{P}, \mathbb{Q})} \mathbb{E}_{(X, Y) \sim \pi} [L_K \|\phi_\psi\|_{\text{Lip}} \|X - Y\|] = L_K \|\phi_\psi\|_{\text{Lip}} \mathcal{W}(\mathbb{P}, \mathbb{Q}),$$

where \mathcal{W} is the standard Wasserstein distance (2), and so

$$\text{MMD}_\psi(\mathbb{P}, \mathbb{Q}) \leq L_K \|\phi_\psi\|_{\text{Lip}} \mathcal{W}(\mathbb{P}, \mathbb{Q}).$$

We have that $\partial_i \partial_{i+d} k(x, y) = [\partial_i \phi_\psi(x)]^\top \left[\nabla_a \nabla_b K(a, b) \Big|_{(a, b) = (\phi_\psi(x), \phi_\psi(y))} \right] [\partial_i \phi_\psi(y)]$, where the middle term is a $d_L \times d_L$ matrix and the outer terms are vectors of length d_L . Thus Assumption (V) implies that $\partial_i \partial_{i+d} k(x, x) \geq \gamma_K^2 \|\partial_i \phi_\psi(x)\|^2$, and hence

$$\sigma_{\mu, k, \lambda}^{-2} \geq \gamma_K^2 \mathbb{E}[\|\nabla \phi_\psi(X)\|_F^2]$$

so that

$$\text{SMMD}_\psi^2(\mathbb{P}, \mathbb{Q}) = \sigma_{\mu, k, \lambda}^2 \text{MMD}_\psi^2(\mathbb{P}, \mathbb{Q}) \leq \frac{L_K^2 \|\phi_\psi\|_{\text{Lip}}^2}{\gamma_K^2 \mathbb{E}[\|\nabla \phi_\psi(X)\|_F^2]} \mathcal{W}^2(\mathbb{P}, \mathbb{Q}).$$

Using Lemma 8, we can write $\phi_\psi(X) = \alpha(\psi) \phi_{\bar{\psi}}(X)$ with $\bar{\psi} \in \Psi_1^\kappa$. Then we have

$$\frac{\|\phi_\psi\|_{\text{Lip}}^2}{\mathbb{E}_\mu[\|\nabla \phi_\psi(X)\|_F^2]} = \frac{\alpha(\psi)^2 \|\phi_{\bar{\psi}}\|_{\text{Lip}}^2}{\alpha(\psi)^2 \mathbb{E}_\mu[\|\nabla \phi_{\bar{\psi}}(X)\|_F^2]} \leq \frac{1}{\mathbb{E}_\mu[\|\nabla \phi_{\bar{\psi}}(X)\|_F^2]},$$

where we used $\|\phi_{\bar{\psi}}\|_{\text{Lip}} \leq \prod_{l=1}^L \|\bar{W}^l\| = 1$. But by Lemma 9, for Lebesgue-almost all X , $\|\nabla \phi_{\bar{\psi}}(X)\|_F^2 \geq d_L (\alpha/\kappa)^L$. Using Assumption (I), this implies that

$$\frac{\|\phi_\psi\|_{\text{Lip}}^2}{\mathbb{E}_\mu[\|\nabla \phi_\psi(X)\|_F^2]} \leq \frac{1}{\mathbb{E}_\mu[\|\nabla \phi_{\bar{\psi}}(X)\|_F^2]} \leq \frac{\kappa^L}{d_L \alpha^L}.$$

Thus for any $\psi \in \Psi^\kappa$,

$$\text{SMMD}_\psi(\mathbb{P}, \mathbb{Q}) \leq \frac{L_K \kappa^{L/2}}{\gamma_K \sqrt{d_L} \alpha^{L/2}} \mathcal{W}(\mathbb{P}, \mathbb{Q}).$$

The desired bound on $\mathcal{D}_{\text{SMMD}}^{\mu, \Psi^\kappa, \lambda}$ follows immediately. \square

Lemma 7. Let $(x, y) \mapsto k(x, y)$ be the continuous kernel of an RKHS \mathcal{H} defined on a Polish space \mathcal{X} , and define the corresponding pseudo-distance $d_k(x, y) := \|k(x, \cdot) - k(y, \cdot)\|_{\mathcal{H}}$. Then the following inequality holds for any distributions \mathbb{P} and \mathbb{Q} on \mathcal{X} , including when the quantities are infinite:

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) \leq \mathcal{W}_{d_k}(\mathbb{P}, \mathbb{Q}).$$

Proof. Let \mathbb{P} and \mathbb{Q} be two probability distributions, and let $\Pi(\mathbb{P}, \mathbb{Q})$ be the set of couplings between them. Let $\pi^* \in \text{argmin}_{(X, Y) \sim \pi} [c_k(X, Y)]$ be an optimal coupling, which is guaranteed to exist [51, Theorem 4.1]; by definition $\mathcal{W}_{d_k}(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{(X, Y) \sim \pi^*} [d_k(X, Y)]$. When $\mathcal{W}_{d_k}(\mathbb{P}, \mathbb{Q}) = \infty$ the inequality trivially holds, so assume that $\mathcal{W}_{d_k}(\mathbb{P}, \mathbb{Q}) < \infty$.

Take a sample $(X, Y) \sim \pi^*$ and a function $f \in \mathcal{H}$ with $\|f\|_{\mathcal{H}} \leq 1$. By the Cauchy-Schwarz inequality,

$$\|f(X) - f(Y)\| \leq \|f\|_{\mathcal{H}} \|k(X, \cdot) - k(Y, \cdot)\|_{\mathcal{H}} \leq \|k(X, \cdot) - k(Y, \cdot)\|_{\mathcal{H}}.$$

Taking the expectation with respect to π^* , we obtain

$$\mathbb{E}_{\pi^*} [|f(X) - f(Y)|] \leq \mathbb{E}_{\pi^*} [\|k(X, \cdot) - k(Y, \cdot)\|_{\mathcal{H}}].$$

The right-hand side is just the definition of $\mathcal{W}_{d_k}(\mathbb{P}, \mathbb{Q})$. By Jensen's inequality, the left-hand side is lower-bounded by

$$|\mathbb{E}_{\pi^*} [f(X) - f(Y)]| = |\mathbb{E}_{X \sim \mathbb{P}} [f(X)] - \mathbb{E}_{Y \sim \mathbb{Q}} [f(Y)]|$$

since π^* has marginals \mathbb{P} and \mathbb{Q} . We have shown so far that for any $f \in \mathcal{H}$ with $\|f\|_{\mathcal{H}} \leq 1$,

$$|\mathbb{E}_{\mathbb{P}} [f(X)] - \mathbb{E}_{\mathbb{Q}} [f(Y)]| \leq \mathcal{W}_{c_k}(\mathbb{P}, \mathbb{Q});$$

the result follows by taking the supremum over f . \square

Lemma 8. Let $\psi = ((W^L, b^L), (W^{L-1}, b^{L-1}), \dots, (W^1, b^1)) \in \Psi^\kappa$. There exists a corresponding scalar $\alpha(\psi)$ and $\bar{\psi} = ((\bar{W}^L, \bar{b}^L), (\bar{W}^{L-1}, \bar{b}^{L-1}), \dots, (\bar{W}^1, \bar{b}^1)) \in \Psi_1^\kappa$, defined by (18), such that for all X ,

$$\phi_\psi(X) = \alpha(\psi) \phi_{\bar{\psi}}(X).$$

Proof. Set $\bar{W}^l = \frac{1}{\|W^l\|} W^l$, $\bar{b}^l = \frac{1}{\prod_{m=1}^l \|W^m\|} b^l$, and $\alpha(\psi) = \prod_{l=1}^L \|W^l\|$. Note that the condition number is unchanged, $\text{cond}(\bar{W}^l) = \text{cond}(W^l) \leq \kappa$, and $\|\bar{W}^l\| = 1$, so $\bar{\psi} \in \Phi_1^\kappa$. It is also easy to see from (16) that

$$h_{\bar{\psi}}^l(X) = \frac{1}{\prod_{m=1}^l \|W^m\|} h_\psi^l(X)$$

so that

$$\alpha(\psi) h_{\bar{\psi}}^L(X) = \frac{\prod_{l=1}^L \|W^l\|}{\prod_{l=1}^L \|W^l\|} h_\psi^L(X) = \phi_\psi(X). \quad \square$$

Lemma 9. Make Assumptions (II) and (III), and let $\psi \in \Psi_1^\kappa$. Then the set of inputs for which any intermediate activation is exactly zero,

$$\mathcal{N}_\psi = \bigcup_{l=1}^L \bigcup_{k=1}^{d_l} \left\{ X \in \mathbb{R}^d \mid (h_\psi^l(X))_k = 0 \right\},$$

has zero Lebesgue measure. Moreover, for any $X \notin \mathcal{N}_\psi$, $\nabla_X \phi_\psi(X)$ exists and

$$\|\nabla_X \phi_\psi(X)\|_F^2 \geq \frac{d_L \alpha^L}{\kappa^L}.$$

Proof. First, note that the network representation at layer l is piecewise affine. Specifically, define $M_X^l \in \mathbb{R}^{d_l}$ by, using Assumption (III),

$$(M_X^l)_k = \sigma'_l(h_k^l(X)) = \begin{cases} 1 & h_k^l(X) > 0 \\ \alpha & h_k^l(X) < 0 \end{cases};$$

it is undefined when any $h_k^l(X) = 0$, i.e. when $X \in \mathcal{N}_\psi$. Let $V_X^l := W^l \text{diag}(M_X^{l-1})$. Then

$$h_\psi^l(X) = W^l \sigma_{l-1}(h_\psi^{l-1}(X)) + b^l = V_X^l X + b^l,$$

and thus

$$h_\psi^l(X) = \underline{W}_X^l X + \underline{b}_X^l, \quad (20)$$

where $\underline{b}_X^0 = 0$, $\underline{b}_X^l = V_X^l \underline{b}^{l-1} + b^l$, and $\underline{W}_X^l = V_X^l V_X^{l-1} \dots V_X^1$, so long as $X \notin \mathcal{N}_\psi$.

Because $\psi \in \Psi_1^\kappa$, we have $\|W^l\| = 1$ and $\sigma_{\min}(W^l) \geq 1/\kappa$; also, $\|M_X^l\| \leq 1$, $\sigma_{\min}(M_X^l) \geq \alpha$. Thus $\|\underline{W}_X^l\| \leq 1$, and using Assumption (II) with Lemma 10 gives $\sigma_{\min}(\underline{W}_X^l) \geq (\alpha/\kappa)^l$. In particular, each \underline{W}_X^l is full-rank.

Next, note that \underline{b}_X^l and \underline{W}_X^l each only depend on X through the activation patterns M_X^l . Letting $H_X^l = (M_X^l, M_X^{l-1}, \dots, M_X^1)$ denote the full activation patterns up to level l , we can thus write

$$h_\psi^l(X) = \underline{W}^{H_X^l} X + \underline{b}^{H_X^l}.$$

There are only finitely many possible values for H_X^l ; we denote the set of such values as \mathcal{H}^l . Then we have that

$$\mathcal{N}_\psi \subseteq \bigcup_{l=0}^L \bigcup_{k=1}^{d_L} \bigcup_{H^l \in \mathcal{H}^l} \left\{ X \in \mathbb{R}^d \mid \underline{W}_k^{H^l} X + \underline{b}_k^{H^l} = 0 \right\}.$$

Because each $\underline{W}_k^{H^l}$ is of rank d_l , each set in the union is either empty or an affine subspace of dimension $d - d_l$. As each $d_l > 0$, each set in the finite union has zero Lebesgue measure, and \mathcal{N}_ψ also has zero Lebesgue measure.

We will now show that the activation patterns are piecewise constant, so that $\nabla_X h_\psi^l(X) = \underline{W}^{H_X^l}$ for all $X \notin \mathcal{N}_\psi$. Because $\psi \in \Psi_1^\kappa$, we have $\|h_\psi^l\|_{\text{Lip}} \leq 1$, and in particular

$$\left| (h_\psi^l(X))_k - (h_\psi^l(X'))_k \right| \leq \|X - X'\|.$$

Thus, take some $X \notin \mathcal{N}_\psi$, and find the smallest absolute value of its activations, $\epsilon = \min_{l=1,\dots,L} \min_{k=1,\dots,d_l} \left| (h_\psi^l(X))_k \right|$; clearly $\epsilon > 0$. For any X' with $\|X - X'\| < \epsilon$, we know that for all l and k ,

$$\text{sign} \left((h_\psi^l(X))_k \right) = \text{sign} \left((h_\psi^l(X'))_k \right),$$

implying that $H_X^l = H_{X'}^l$, as well as $X' \notin \mathcal{N}_\psi$. Thus for any point $X \notin \mathcal{N}_\psi$, $\nabla \phi_\psi(X) = \underline{W}^{H_X^L}$. Finally, we obtain

$$\|\nabla \phi_\psi(X)\|_F^2 = \|\underline{W}^{H_X^L}\|_F^2 \geq d_L \sigma_{\min} \left(\underline{W}^{H_X^L} \right)^2 \geq \frac{d_L \alpha^L}{\kappa^L}. \quad \square$$

Lemma 10. *Let $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$, with $m \geq n \geq p$. Then $\sigma_{\min}(AB) \geq \sigma_{\min}(A) \sigma_{\min}(B)$.*

Proof. A more general version of this result can be found in [19, Theorem 2]; we provide a proof here for completeness.

If B has a nontrivial null space, $\sigma_{\min}(B) = 0$ and the inequality holds. Otherwise, let \mathbb{R}_*^n denote $\mathbb{R}^n \setminus \{0\}$. Recall that for $C \in \mathbb{R}^{m \times n}$ with $m \geq n$,

$$\sigma_{\min}(C) = \sqrt{\lambda_{\min}(C^\top C)} = \sqrt{\inf_{x \in \mathbb{R}_*^n} \frac{x^\top C^\top C x}{x^\top x}} = \inf_{x \in \mathbb{R}_*^n} \frac{\|Cx\|}{\|x\|}.$$

Thus, as $Bx \neq 0$ for $x \neq 0$,

$$\begin{aligned} \sigma_{\min}(AB) &= \inf_{x \in \mathbb{R}_*^p} \frac{\|ABx\|}{\|x\|} = \inf_{x \in \mathbb{R}_*^p} \frac{\|ABx\| \|Bx\|}{\|Bx\| \|x\|} \\ &\geq \left(\inf_{x \in \mathbb{R}_*^p} \frac{\|ABx\|}{\|Bx\|} \right) \left(\inf_{x \in \mathbb{R}_*^p} \frac{\|Bx\|}{\|x\|} \right) \\ &\geq \left(\inf_{y \in \mathbb{R}_*^n} \frac{\|Ay\|}{\|y\|} \right) \left(\inf_{x \in \mathbb{R}_*^p} \frac{\|Bx\|}{\|x\|} \right) = \sigma_{\min}(A) \sigma_{\min}(B). \quad \square \end{aligned}$$

A.2.1 When some of the assumptions don't hold

Here we analyze through simple examples what happens when the condition number can be unbounded, and when Assumption **(II)**, about decreasing widths of the network, is violated.

Condition Number: We start by a first example where the condition number can be arbitrarily high. We consider a two-layer network on \mathbb{R}^2 , defined by

$$\phi_\alpha(X) = [1 \quad -1] \sigma(W_\alpha X) \quad W_\alpha = \begin{bmatrix} 1 & 1 \\ 1 & 1 + \alpha \end{bmatrix} \quad (21)$$

where $\alpha > 0$. As α approaches 0 the matrix W_α becomes singular which means that its condition number blows up. We are interested in analyzing the behavior of the Lipschitz constant of ϕ and the expected squared norm of its gradient under μ as α approaches 0.

One can easily compute the squared norm of the gradient of ϕ which is given by

$$\|\nabla \phi_\alpha(X)\|^2 = \begin{cases} \alpha^2 & X \in A_1 \\ \gamma^2 \alpha^2 & X \in A_2 \\ (1 - \gamma)^2 + (1 + \alpha - \gamma)^2 & X \in A_3 \\ (1 - \gamma)^2 + (\gamma \alpha + \gamma - 1)^2 & X \in A_4 \end{cases} \quad (22)$$

Here A_1, A_2, A_3 and A_4 are defined by (23) and are represented in Figure 4:

$$\begin{aligned}
A_1 &:= \{X \in \mathbb{R}^2 \mid X_1 + X_2 \geq 0 \quad X_1 + (1 + \alpha)X_2 \geq 0\} \\
A_2 &:= \{X \in \mathbb{R}^2 \mid X_1 + X_2 < 0 \quad X_1 + (1 + \alpha)X_2 < 0\} \\
A_3 &:= \{X \in \mathbb{R}^2 \mid X_1 + X_2 < 0 \quad X_1 + (1 + \alpha)X_2 \geq 0\} \\
A_4 &:= \{X \in \mathbb{R}^2 \mid X_1 + X_2 \geq 0 \quad X_1 + (1 + \alpha)X_2 < 0\}
\end{aligned} \tag{23}$$

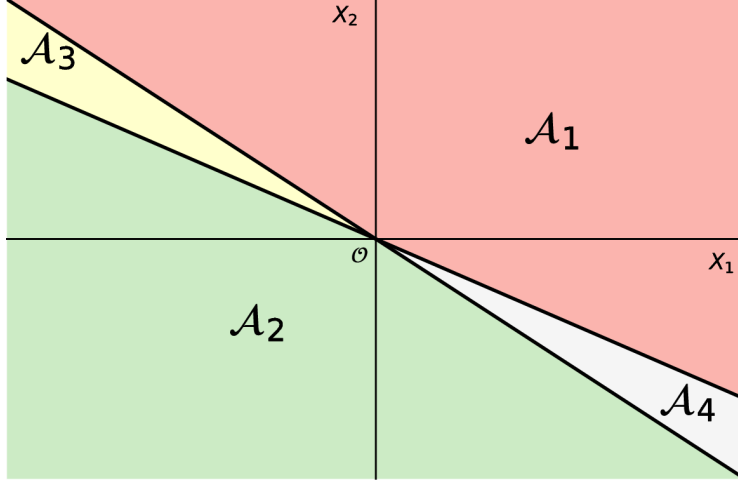


Figure 4: Decomposition of \mathbb{R}^2 into 4 regions A_1, A_2, A_3 and A_4 as defined in (23). As α approaches 0, the area of sets A_3 and A_4 becomes negligible.

It is easy to see that whenever μ has a density, the probability of the sets A_3 and A_4 goes to 0 as $\alpha \rightarrow 0$. Hence one can deduce that $\mathbb{E}_\mu[\|\nabla\phi_\alpha(X)\|^2] \rightarrow 0$ when $\alpha \rightarrow 0$. On the other hand, the squared Lipschitz constant of ϕ is given by $(1 - \gamma)^2 + (1 + \alpha - \gamma)^2$ which converges to $2(1 - \gamma)^2$. This shows that controlling the expectation of the gradient doesn't allow to effectively control the Lipschitz constant of ϕ .

Monotonicity of the dimensions: We would like to consider a second example where Assumption (II) doesn't hold. Consider the following two layer network defined by:

$$\phi(X) = [-1 \quad 0 \quad 1] \sigma(W_\beta X) \quad W_\beta := \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & \beta \end{bmatrix} \tag{24}$$

for $\beta > 0$. Note that W_β is a full rank matrix, but Assumption (II) doesn't hold. Depending on the sign of the components of $W_\beta X$ one has the following expression for $\|\nabla\phi_\alpha(X)\|^2$:

$$\|\nabla\phi_\alpha(X)\|^2 = \begin{cases} \beta^2 & X \in B_1 \\ \gamma^2 \beta^2 & X \in B_2 \\ \beta^2 & X \in B_3 \\ (1 - \gamma)^2 + \gamma^2 \beta^2 & X \in B_4 \\ (1 - \gamma)^2 + \beta^2 & X \in B_5 \\ \gamma^2 \beta^2 & X \in B_6 \end{cases} \tag{25}$$

where $(B_i)_{1 \leq i \leq 6}$ are defined by (26)

$$\begin{aligned}
B_1 &:= \{X \in \mathbb{R}^2 \mid X_1 \geq 0 \quad X_2 \geq 0\} \\
B_2 &:= \{X \in \mathbb{R}^2 \mid X_1 < 0 \quad X_2 < 0\} \\
B_3 &:= \{X \in \mathbb{R}^2 \mid X_1 \geq 0 \quad X_2 < 0 \quad X_1 + \beta X_2 \geq 0\} \\
B_4 &:= \{X \in \mathbb{R}^2 \mid X_1 \geq 0 \quad X_2 < 0 \quad X_1 + \beta X_2 < 0\} \\
B_5 &:= \{X \in \mathbb{R}^2 \mid X_1 > 0 \quad X_2 \geq 0 \quad X_1 + \beta X_2 \geq 0\} \\
B_6 &:= \{X \in \mathbb{R}^2 \mid X_1 > 0 \quad X_2 \geq 0 \quad X_1 + \beta X_2 < 0\}
\end{aligned} \tag{26}$$

The squared Lipschitz constant is given by $\|\phi\|_L^2(1-\gamma)^2 + \beta^2$ while the expected squared norm of the gradient of ϕ is given by:

$$\mathbb{E}_\mu[\|\phi(X)\|^2] = 3\beta^2(p(B_1 \cup B_3 \cup B_5) + \gamma^2 p(B_2 \cup B_4 \cup B_6)) + (1-\gamma)^2 p(B_4 \cup B_5). \quad (27)$$

Again the set $B_4 \cup B_5$ becomes negligible as β approaches 0 which implies that $\mathbb{E}_\mu[\|\phi(X)\|^2] \rightarrow 0$. On the other hand $\|\phi\|_L^2$ converges to $(1-\gamma)^2$. Note that unlike in the first example in (21), the matrix W_β has a bounded condition number. In this example, the columns of W_0 are all in the null space of $[-1 \ 0 \ 1]$, which implies $\nabla\phi_0(X) = 0$ for all $X \in \mathbb{R}^2$, even though all matrices have full rank.

B DiracGAN vector fields for more losses

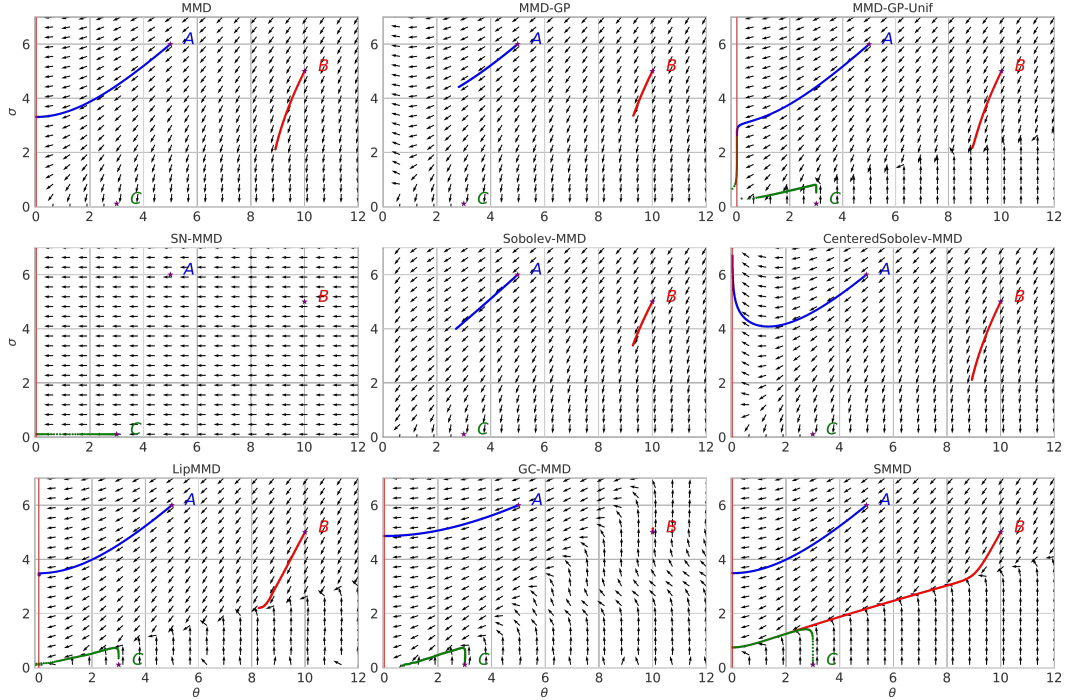


Figure 5: Vector fields for different losses with respect to the generator parameter θ and the feature representation parameter ψ ; the losses use a Gaussian kernel, and are shown in (28). Following [30], $\mathbb{P} = \delta_0$, $\mathbb{Q} = \delta_\theta$ and $\phi_\psi(x) = \psi x$. The curves show the result of taking simultaneous gradient steps in (θ, ψ) beginning from three initial parameter values.

Figure 5 shows parameter vector fields, like those in Figure 6, for Example 1 for a variety of different losses:

$$\begin{aligned}
\text{MMD:} & - \text{MMD}_\psi^2 \\
\text{MMD-GP:} & - \text{MMD}_\psi^2 + \lambda \mathbb{E}_{\mathbb{P}}[(\|\nabla f(X)\| - 1)^2] \\
\text{MMD-GP-Unif:} & - \text{MMD}_\psi^2 + \lambda \mathbb{E}_{\tilde{X} \simeq \mu^*}[(\|\nabla f(\tilde{X})\| - 1)^2] \\
\text{SN-MMD:} & - 2 \text{MMD}_1(\mathbb{P}, \mathbb{Q})^2 \\
\text{Sobolev-MMD:} & - \text{MMD}_\psi^2 + \lambda (\mathbb{E}_{(\mathbb{P}+\mathbb{Q})/2}[\|\nabla f(X)\|^2] - 1)^2 \\
\text{CenteredSobolev-MMD:} & - \text{MMD}_\psi^2 + \lambda (\mathbb{E}_{(\mathbb{P}+\mathbb{Q})/2}[\|\nabla f(X)\|^2])^2 \\
\text{LipMMD:} & - \text{LipMMD}_{k_\psi, \lambda}^2 \\
\text{GC-MMD:} & - \text{GCMMD}_{\mathcal{N}(0, 10^2), k_\psi, \lambda}^2 \\
\text{SMMD:} & - \text{SMMD}_{k_\psi, \mathbb{P}, \lambda}^2
\end{aligned} \quad (28)$$

The squared MMD between δ_0 and δ_θ under a Gaussian kernel of bandwidth $1/\psi$ and is given by $2(1 - e^{-\frac{\psi^2 \theta^2}{2}})$. MMD-GP-unif uses a gradient penalty as in [7] where each samples from μ^* is obtained by first sampling X and Y from \mathbb{P} and \mathbb{Q} and then sampling uniformly between X and Y . MMD-GP uses the same gradient penalty, but the expectation is taken under \mathbb{P} rather than μ^* . SN-MMD refers to MMD with spectral normalization; here this means that $\psi = 1$. Sobolev-MMD refers to the loss used in [33] with the quadratic penalty only. $\text{GCMMMD}_{\mu,k,\lambda}$ is defined by (5), with $\mu = \mathcal{N}(0, 10^2)$.

C Vector fields of Gradient-Constrained MMD and Sobolev GAN critics

Mroueh et al. [33] argue that *the gradient of the critic (...) defines a transportation plan for moving the distribution mass* (from generated to reference distribution) and present the solution of Sobolev PDE for 2-dimensional Gaussians. We observed that in this simple example the gradient of the Sobolev critic can be very high outside of the areas of high density, which is not the case with the Gradient-Constrained MMD. Figure 6 presents critic gradients in both cases, using $\mu = (\mathbb{P} + \mathbb{Q})/2$ for both.

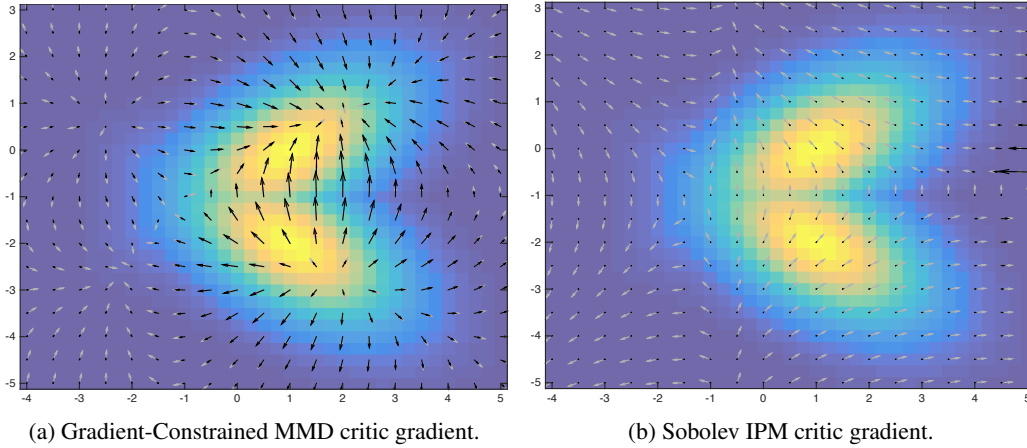


Figure 6: Vector fields of critic gradients between two Gaussians. The grey arrows show normalized gradients, i.e. gradient directions, while the black ones are the actual gradients. Note that for the Sobolev critic, gradients norms are orders of magnitudes higher on the right hand side of the plot than in the areas of high density of the given distributions.

This unintuitive behavior is most likely related to the vanishing boundary condition, assumed by Sobolev GAN. Solving the actual Sobolev PDE, we found that the Sobolev critic has very high gradients close to the boundary in order to match the condition; moreover, these gradients point in opposite directions to the target distribution.

D An estimator for Lipschitz MMD

We now describe briefly how to estimate the Lipschitz MMD in low dimensions. Recall that

$$\text{LipMMD}_{k,\lambda}(\mathbb{P}, \mathbb{Q}) = \sup_{f \in \mathcal{H}_k : \|f\|_{\text{Lip}}^2 + \lambda \|f\|_{\mathcal{H}_k}^2 \leq 1} \mathbb{E}_{X \sim \mathbb{P}}[f(X)] - \mathbb{E}_{X \sim \mathbb{Q}}[f(Y)].$$

For $f \in \mathcal{H}_k$, it is the case that

$$\|f\|_{\text{Lip}}^2 = \sup_{x \in \mathbb{R}^d} \|\nabla f(x)\|^2 = \sup_{x \in \mathbb{R}^d} \sum_{i=1}^d \langle \partial_i k(x, \cdot), f \rangle_{\mathcal{H}_k}^2 = \sup_{x \in \mathbb{R}^d} \left\langle f, \sum_{i=1}^d [\partial_i k(x, \cdot) \otimes \partial_i k(x, \cdot)] f \right\rangle_{\mathcal{H}_k}.$$

Thus we can approximate the constraint $\|f\|_{\text{Lip}}^2 + \lambda \|f\|_{\mathcal{H}_k}^2 \leq 1$ by enforcing the constraint on a set of m points $\{Z_i\}$ reasonably densely covering the region around the supports of \mathbb{P} and \mathbb{Q} , rather

than enforcing it at every point in \mathcal{X} . An estimator of the Lipschitz MMD based on $X \sim \mathbb{P}^{n_X}$ and $Y \sim \mathbb{Q}^{n_Y}$ is

$$\widehat{\text{LipMMD}}_{k,\lambda}(X, Y, Z) \approx \sup_{f \in \mathcal{H}_k} \frac{1}{n_X} \sum_{j=1}^{n_X} f(X_j) - \frac{1}{n_Y} \sum_{j=1}^{n_Y} f(Y_j) \quad (29)$$

$$\text{s.t. } \forall j, \|\nabla f(Z_j)\|^2 + \lambda \|f\|_{\mathcal{H}_k}^2 \leq 1.$$

By the generalized representer theorem, the optimal f for (29) will be of the form

$$f(\cdot) = \sum_{j=1}^{n_X} \alpha_j k(X_j, \cdot) + \sum_{j=1}^{n_Y} \beta_j k(Y_j, \cdot) + \sum_{i=1}^d \sum_{j=1}^m \gamma_{(i,j)} \partial_i k(Z_j, \cdot).$$

Writing $\delta = (\alpha, \beta, \gamma)$, the objective function is linear in δ ,

$$\left[\frac{1}{n_X} \quad \dots \quad \frac{1}{n_X} \quad -\frac{1}{n_Y} \quad \dots \quad -\frac{1}{n_Y} \quad 0 \quad \dots \quad 0 \right] \delta.$$

The constraints are quadratic, built from the following matrices, where the X and Y samples are concatenated together, as are the derivatives with each dimension of the Z samples:

$$K := \begin{bmatrix} k(X_1, X_1) & \dots & k(X_1, Y_{n_Y}) \\ \vdots & \ddots & \vdots \\ k(Y_{n_X}, X_1) & \dots & k(Y_{n_Y}, Y_{n_Y}) \end{bmatrix}$$

$$B := \begin{bmatrix} \partial_1 k(Z_1, X_1) & \dots & \partial_1 k(Z_1, Y_{n_Y}) \\ \vdots & \ddots & \vdots \\ \partial_d k(Z_m, X_1) & \dots & \partial_d k(Z_m, Y_{n_Y}) \end{bmatrix}$$

$$H := \begin{bmatrix} \partial_1 \partial_{1+d} k(Z_1, Z_1) & \dots & \partial_1 \partial_{d+d} k(Z_1, Z_m) \\ \vdots & \ddots & \vdots \\ \partial_d \partial_{1+d} k(Z_m, Z_1) & \dots & \partial_d \partial_{d+d} k(Z_m, Z_m) \end{bmatrix}.$$

Given these matrices, and letting $O_j = \sum_{i=1}^d e_{(i,j)} e_{(i,j)}^\top$ where $e_{(i,j)}$ is the (i, j) th standard basis vector in \mathbb{R}^{md} , we have that

$$\|f\|_{\mathcal{H}_k}^2 = \delta^\top \begin{bmatrix} K & B^\top \\ B & H \end{bmatrix} \delta \quad \|\nabla f(Z_j)\|^2 = \sum_{i=1}^d (\partial_i f(Z_j))^2 = \delta^\top \begin{bmatrix} B^\top O_j B & B^\top O_j H \\ H O_j B & H O_j H \end{bmatrix} \delta.$$

Thus the optimization problem (29) is a linear problem with convex quadratic constraints, which can be solved by standard convex optimization software. The approximation is reasonable only if we can effectively cover the region of interest with densely spaced $\{Z_i\}$; it requires a nontrivial amount of computation even for the very simple 1-dimensional toy problem of Example 1.

One advantage of this estimator, though, is that finding its derivative with respect to the input points or the kernel parameterization is almost free once we have computed the estimate, as long as our solver has computed the dual variables μ corresponding to the constraints in (29). We just need to exploit the envelope theorem and then differentiate the KKT conditions, as done for instance in [1]. The differential of (29) ends up being, assuming the optimum of (29) is at $\hat{\delta} \in \mathbb{R}^{n_X+n_Y+md}$ and $\hat{\mu} \in \mathbb{R}^m$,

$$d\widehat{\text{LipMMD}}_{k,\lambda}(X, Y, Z) = \hat{\delta}^\top \begin{bmatrix} dK \\ dB \end{bmatrix} \left[\frac{1}{n_X} \quad \dots \quad \frac{1}{n_X} \quad -\frac{1}{n_Y} \quad \dots \quad -\frac{1}{n_Y} \right]^\top - \sum_{j=1}^m \hat{\mu}_j \hat{\delta}^\top (dP_j) \hat{\delta}$$

$$dP_j := \begin{bmatrix} (dB)^\top O_j B + B^\top O_j (dH) & (dB)^\top O_j H + B^\top O_j (dH) \\ (dH) O_j B + H O_j (dB) & (dH) O_j H + H O_j (dH) \end{bmatrix} + \lambda \begin{bmatrix} dK & dB^\top \\ dB & dH \end{bmatrix}.$$

E Near-equivalence of WGAN and linear-kernel MMD GANs

For an MMD GAN-GP with kernel $k(x, y) = \phi(x)\phi(y)$, we have that

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) = |\mathbb{E}_{\mathbb{P}} \phi(x) - \mathbb{E}_{\mathbb{Q}} \phi(y)|$$

and the corresponding critic function is

$$\frac{\eta(t)}{\|\eta\|_{\mathcal{H}}} = \frac{\mathbb{E}_{X \sim \mathbb{P}} \phi(X)\phi(t) - \mathbb{E}_{Y \sim \mathbb{Q}} \phi(Y)\phi(t)}{|\mathbb{E}_{\mathbb{P}} \phi(X) - \mathbb{E}_{\mathbb{Q}} \phi(Y)|} = \text{sign}(\mathbb{E}_{X \sim \mathbb{P}} \phi(X) - \mathbb{E}_{Y \sim \mathbb{Q}} \phi(Y)) \phi(t).$$

Thus if we assume $\mathbb{E}_{X \sim \mathbb{P}} \phi(X) > \mathbb{E}_{Y \sim \mathbb{Q}} \phi(Y)$, as that is the goal of our critic training, we see that the MMD becomes identical to the WGAN loss, and the gradient penalty is applied to the same function.

(MMD GANs, however, would typically train on the unbiased estimator of MMD^2 , giving a very slightly different loss function. [7] also applied the gradient penalty to η rather than the true critic $\eta/\|\eta\|_{\mathcal{H}}$.)

The SMMD with a linear kernel is thus analogous to applying the scaling operator to a WGAN; hence the name SWGAN.

F Additional experiments

F.1 Comparison of Gradient-Constrained MMD to Scaled MMD

Figure 7 shows the behavior of the MMD, the Gradient-Constrained SMMD, and the Scaled MMD when comparing Gaussian distributions. We can see that $\text{MMD} \propto \text{SMMD}$ and the Gradient-Constrained MMD behave similarly in this case, and that optimizing the SMMD and the Gradient-Constrained MMD is also similar. Optimizing the MMD would yield an essentially constant distance.

F.2 IGMs with Optimized Gradient-Constrained MMD loss

We implemented the estimator of Proposition 3 using the empirical mean estimator of η , and sharing samples for $\mu = \mathbb{P}$. To handle the large but approximately low-rank matrix system, we used an incomplete Cholesky decomposition [43, Algorithm 5.12] to obtain $R \in \mathbb{R}^{\ell \times M(1+d)}$ such that

$$\begin{bmatrix} K & G^T \\ G & H \end{bmatrix} \approx R^T R. \text{ Then the Woodbury matrix identity allows an efficient evaluation:}$$

$$(R^T R + M\lambda I)^{-1} = \frac{1}{M\lambda} (I - R(RR^T + M\lambda I)^{-1} R).$$

Even though only a small ℓ is required for a good approximation, and the full matrices K , G , and H need never be constructed, backpropagation through this procedure is slow and not especially GPU-friendly; training on CPU was faster. Thus we were only able to run the estimator on MNIST, and even that took days to conduct the optimization on powerful workstations.

The learned models, however, were reasonable. Using a DCGAN architecture, batches of size 64, and a procedure that otherwise agreed with the setup of Section 4, samples with and without spectral normalization are shown in Figures 8a and 8b. After the points in training shown, however, the same rank collapse as discussed in Section 4 occurred. Here it seems that spectral normalization may have delayed the collapse, but not prevented it. Figure 8c shows generator loss estimates through training, including the obvious peak at collapse; Figure 8d shows KID scores based on the MNIST-trained convnet representation [7], including comparable SMMD models for context. The fact that SMMD models converged somewhat faster than Gradient-Constrained MMD models here may be more related to properties of the estimator of Proposition 3 rather than the distances; more work would be needed to fully compare the behavior of the two distances.

F.3 Spectral normalization and Scaled MMD

Figure 9 shows the distribution of critic weight singular values, like Figure 2, at more layers. Figure 11 and Table 2 show results for the spectral normalization variants considered in the experiments. MMDGAN, with neither spectral normalization nor a gradient penalty, did surprisingly well in this case, though it fails badly in other situations.

Figure 9 compares the decay of singular values for layer of the critic’s network at both early and later stages of training in two cases: with or without the spectral parametrization. The model was trained on CelebA using SMMD. Figure 11 shows the evolution per iteration of Inception score,

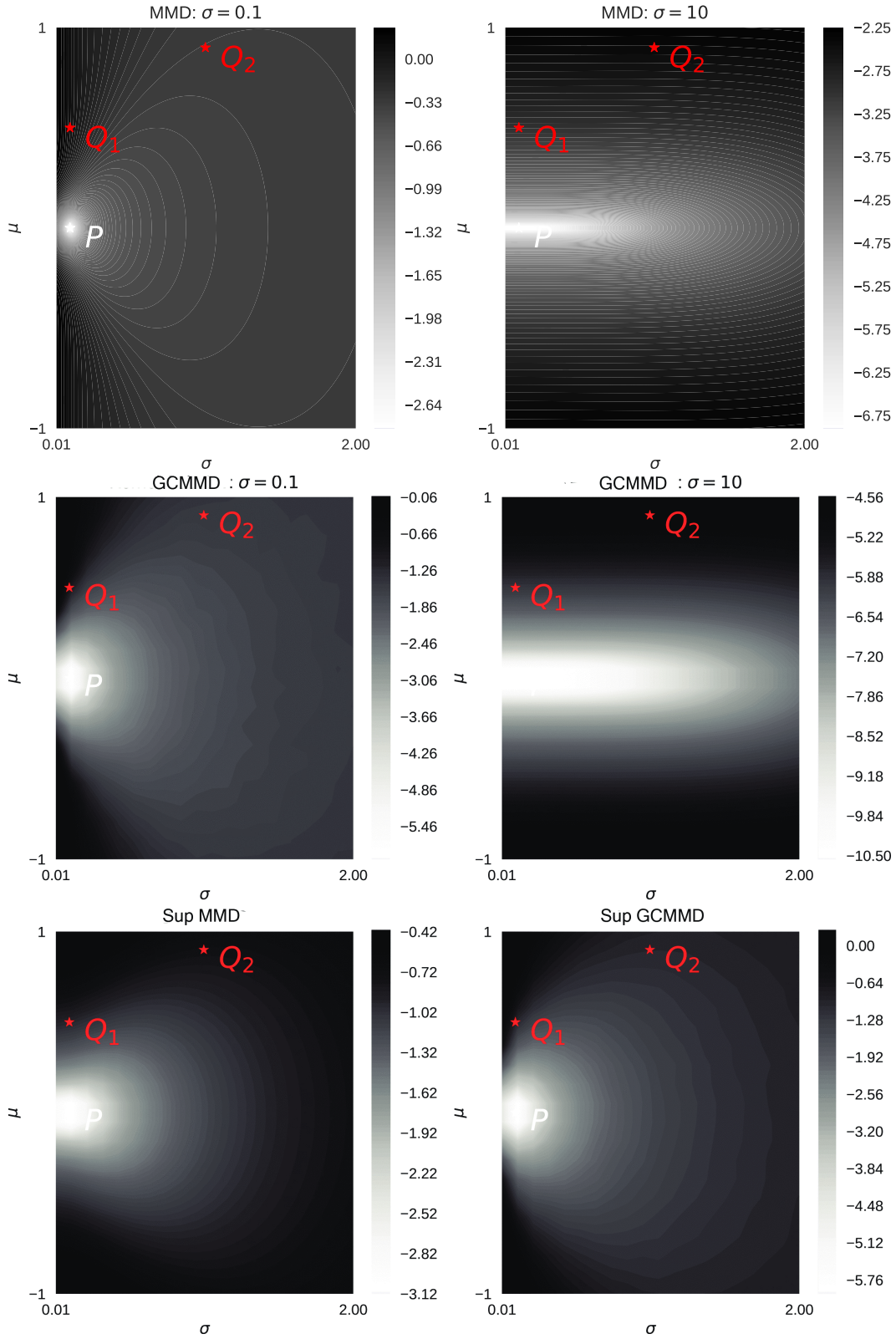


Figure 7: Plots of various distances between one dimensional Gaussians, where $P = \mathcal{N}(0, 0.1^2)$, and the colors show $\log \mathcal{D}(P, \mathcal{N}(\mu, \sigma^2))$. All distances use $\lambda = 1$. Top left: MMD with a Gaussian kernel of bandwidth $\psi = 0.1$. Top right: MMD with bandwidth $\psi = 10$. Middle left: Gradient-Constrained MMD with bandwidth $\psi = 0.1$. Middle right: Gradient-Constrained MMD with bandwidth $\psi = 10$. Bottom left: Optimized SSMD, allowing any $\psi \in \mathbb{R}$. Bottom right: Optimized Gradient-Constrained MMD.

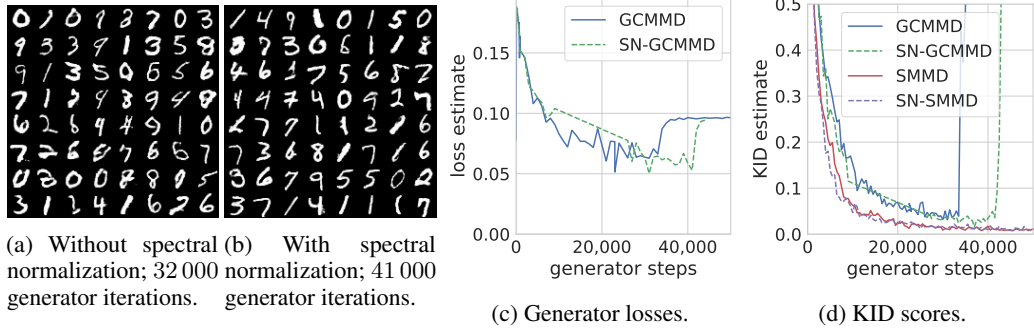


Figure 8: The MNIST models with Optimized Gradient-Constrained MMD loss.

FID and KID for Sobolev-GAN, MMDGAN and variants of MMDGAN and WGAN using spectral normalization. It is often the case that this parametrization alone is not enough to achieve good results.

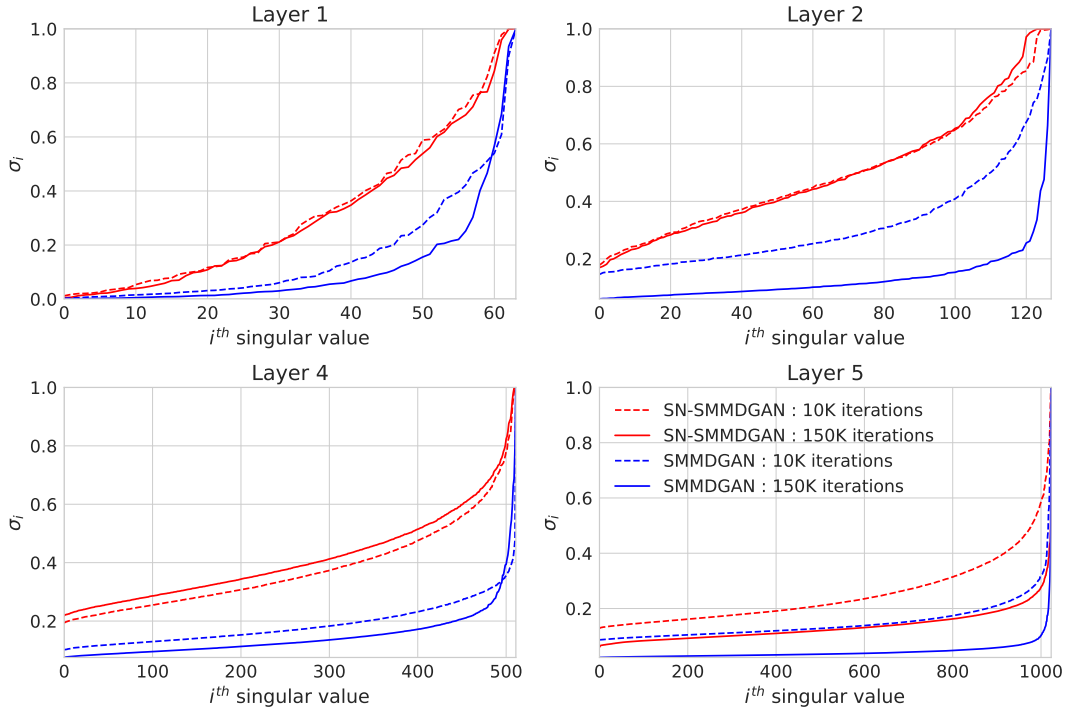


Figure 9: Singular values at different layers, for the same setup as Figure 2.

F.4 Additional samples

Figures 12 and 13 give extra samples from the models.

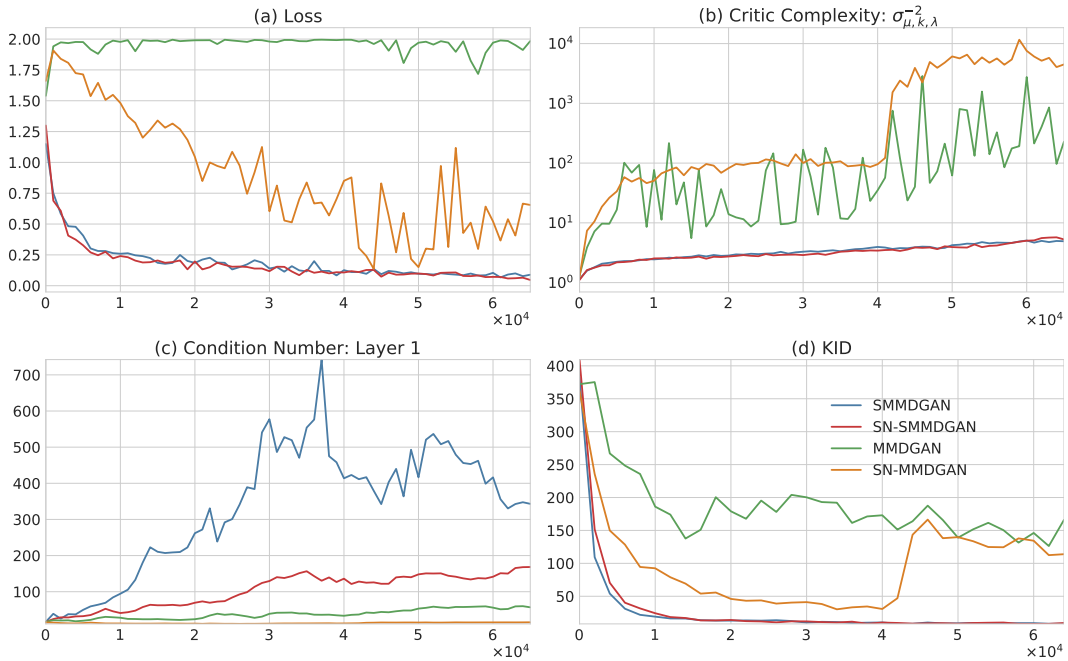


Figure 10: Evolution of various quantities per generator iteration on CelebA during training. 4 models are considered: (SMMDGAN, SN-SMMDGAN, MMDGAN, SN-MMDGAN). (a) Loss: $\text{SMMD}^2 = \sigma_{\mu,k,\lambda}^2 \text{MMD}_k^2$ for SMMDGAN and SN-SMMDGAN, and MMD_k^2 for MMDGAN and SN-MMDGAN. The loss saturates for MMDGAN (green); spectral normalization allows some improvement in loss, but training is still unstable (orange). SMMDGAN and SN-SMMDGAN both lead to stable, fast training (blue and red). (b) SMMD controls the critic complexity well, as expected (blue and red); SN has little effect on the complexity (orange). (c) Ratio of the highest singular value to the smallest for the first layer of the critic network: $\sigma_{\max}/\sigma_{\min}$. SMMD tends to increase the condition number of the weights during training (blue), while SN helps controlling it (red). (d) KID score during training: Only variants using SMMD lead to stable training in this case.

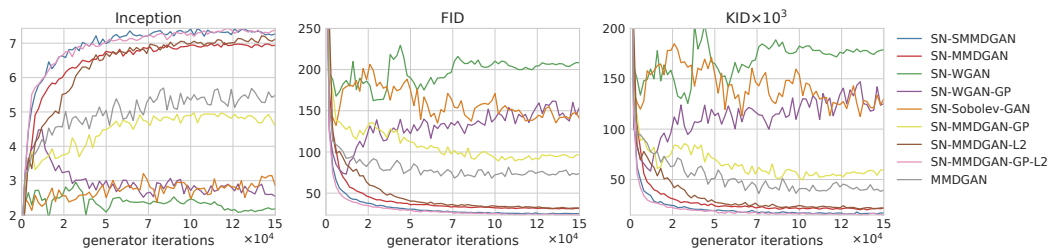


Figure 11: Evolution per iteration of different scores for variants of methods, mostly using spectral normalization, on CIFAR-10.

Table 2: Mean (standard deviation) of score evaluations on CIFAR-10 for different methods using Spectral Normalization.

Method	IS	FID	KID $\times 10^3$
MMDGAN	5.5 \pm 0.0	73.9 \pm 0.1	39.4 \pm 1.5
SN-WGAN	2.2 \pm 0.0	208.5 \pm 0.2	178.9 \pm 1.5
SN-WGAN-GP	2.5 \pm 0.0	154.3 \pm 0.2	125.3 \pm 0.9
SN-Sobolev-GAN	2.9 \pm 0.0	140.2 \pm 0.2	130.0 \pm 1.9
SN-MMDGAN-GP	4.6 \pm 0.1	96.8 \pm 0.4	59.5 \pm 1.4
SN-MMDGAN-L2	7.1 \pm 0.1	31.9 \pm 0.2	21.7 \pm 0.9
SN-MMDGAN	6.9 \pm 0.1	31.5 \pm 0.2	21.7 \pm 1.0
SN-MMDGAN-GP-L2	6.9 \pm 0.2	32.3 \pm 0.3	20.9 \pm 1.1
SN-SMMDGAN	7.3\pm0.1	25.0\pm0.3	16.6\pm2.0



Figure 12: Samples from a generator trained on ImageNet dataset using Scaled MMD with Spectral Normalization: SN-SMMDGAN.



(a) SNGAN

(b) SobolevGAN



(c) MMDGAN-GP-L2

(d) SN-SMMD GAN



(e) SN SWGAN

(f) SMMD GAN

Figure 13: Comparison of samples from different models trained on CelebA with 160×160 resolution.