



On the Convergence of Gradient Descent in GANs: MMD GAN As a Gradient Flow

Youssef Mroueh* and Truyen Nguyen*

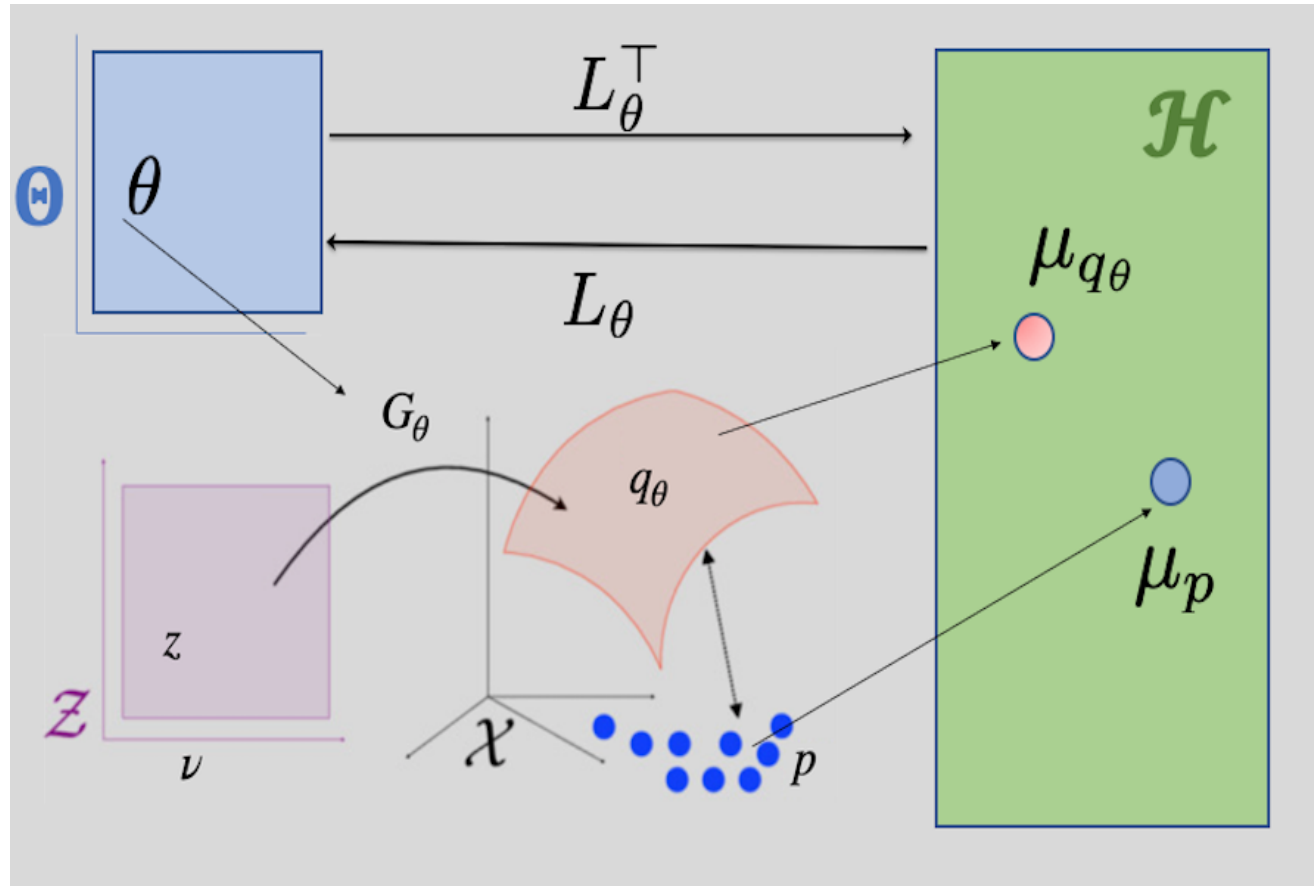
IBM Research & Akron University

<https://arxiv.org/abs/2011.02402>

Outline

- **Parametric Gradient** Regularization for MMD
- Assumptions for the convergence of MMD GAN both in the continuous and discrete case)
- **New Riemannian structure** and its gradient flow on the parametric statistical manifold
- Parametric Regularized **MMD GAN as a gradient flow** of the MMD functional w.r.t to the new Riemannian structure

Parametric Regularized MMD



Parametric Regularized MMD

$$\text{MMD}_{\alpha,\beta}(p, q_\theta)^2 = \sup_{f \in \mathcal{H}} \left\{ \mathbb{E}_p f(x) - \mathbb{E}_{q_\theta} f(x) - \frac{\alpha}{2} \left\| \int \nabla_\theta f(G_\theta(z)) d\nu(z) \right\|_{\mathbb{R}^p}^2 - \frac{\beta}{2} \|f\|_{\mathcal{H}}^2 \right\}$$

$$d\theta_t = L_{\theta_t}(f_t)dt$$

choice of f_t

$$(\alpha D(\theta_t) + \beta I) f_t = \mu_p - \mu_{q_{\theta_t}}$$

$q_{\theta_t} = (G_{\theta_t})_\#(\nu)$

Witness
function of
 $\text{MMD}_{\{\alpha, \beta\}}$

Parametric Regularized MMD GAN

$$d\theta_t = L_{\theta_t}(f_t)dt$$

choice of f_t $(\alpha D(\theta_t) + \beta I)f_t = \mu_p - \mu_{q_{\theta_t}}$

Witness
function of
 $\text{MMD}_{\{\alpha, \beta\}}$

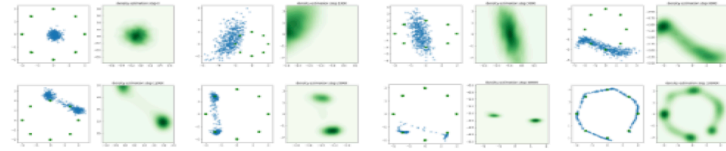
$$q_{\theta_t} = (G_{\theta_t})_{\#}(\nu)$$

Theorem 1 (Parametric Regularized Flows Decrease the *MMD* Distance).
Assume that $\alpha, \beta > 0$. Then the dynamic defined by the witness function of the parametric regularized MMD decreases the functional $\mathcal{F}(q_{\theta})$:

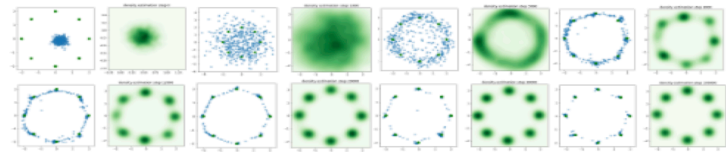
$$\frac{d\mathcal{F}(q_{\theta_t})}{dt} = -\frac{2}{\alpha} \left[\mathcal{F}(q_{\theta_t}) - \beta \text{MMD}_{\alpha, \beta}(p, q_{\theta_t})^2 \right] \leq 0. \quad (1)$$

Moreover, we have $\frac{d\mathcal{F}(q_{\theta_t})}{dt} < 0$ if and only if $D_{\theta_t} \mu_{p-q_{\theta_t}} \neq 0$.

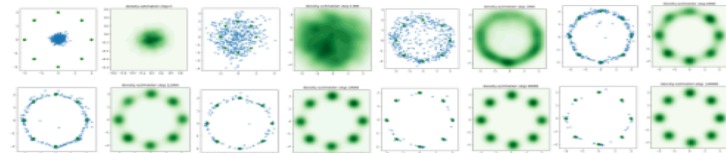
Experiments



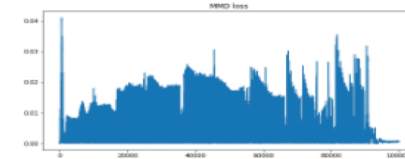
(a) Fixed Kernel: Trajectories of Flows for $\alpha = 0$: Euclidean Gradients Flows



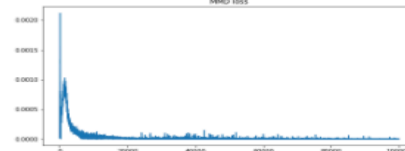
(c) Fixed Kernel: Trajectories of Flows for $\alpha = 100$: Kernelized Gradients Flows w.r.t. $d_{\alpha,\beta}$



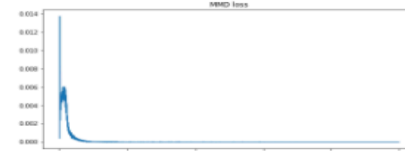
(e) Learned Kernel: Trajectories of Flows for $\alpha = 100$: Kernelized Gradients Flows w.r.t. $d_{\alpha,\beta}$



(b) MMD Loss



(d) MMD Loss



(f) MMD Loss

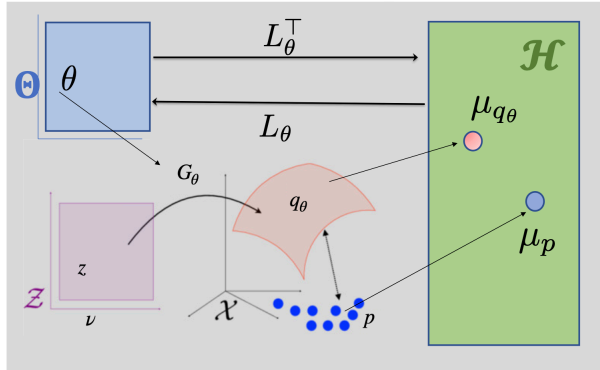
- Target Distribution is a mixture of Gaussian
- (a) Unregularized MMD GAN goes through cycles , **does not converge**
- (c) MMD GAN regularized with Parametric Gradients with fixed Kernel: **convergent**
- (e) MMD GAN regularized with Parametric Gradients with Learned Kernel: **convergent**

Youssef Mroueh* and Truyen Nguyen*
IBM Research & Akron University

Outline

- **Parametric Gradient** Regularization for MMD
- Assumptions for the convergence of MMD GAN both in the continuous and discrete case)
- **New Riemannian structure** and its gradient flow on the parametric statistical manifold
- Parametric Regularized **MMD GAN as a gradient flow** of the MMD functional w.r.t to the new Riemannian structure

Parametric Regularized MMD I



$$\mathcal{F}(q_\theta) = \|\mu_p - \mu_{q_\theta}\|_{\mathcal{H}}$$

$$L_\theta : \mathcal{H} \rightarrow \mathbb{R}^p \quad L_\theta(f) := \int J_\theta G_\theta(z) \nabla f(G_\theta(z)) \nu(dz)$$

$$f \rightarrow L_\theta(f)$$

$$L_\theta^\top : \mathbb{R}^p \rightarrow \mathcal{H}, \quad L_\theta^\top(v) := \int \langle \nabla_\theta [k(G_\theta(z)), \cdot], v \rangle \nu(dz)$$

$$v \rightarrow L_\theta^\top(v)$$

Parametric Grammian $D_\theta : \mathcal{H} \rightarrow \mathcal{H}$ is defined by $D_\theta := L_\theta^\top L_\theta$.

Parametric Regularized MMD II

$$\text{MMD}_{\alpha,\beta}(p, q_\theta)^2 = \sup_{f \in \mathcal{H}} \left\{ \mathbb{E}_p f(x) - \mathbb{E}_{q_\theta} f(x) - \frac{\alpha}{2} \left\| \int \nabla_\theta f(G_\theta(z)) \nu(dz) \right\|_{\mathbb{R}^p}^2 - \frac{\beta}{2} \|f\|_{\mathcal{H}}^2 \right\}$$

$$\text{MMD}_{\alpha,\beta}(p, q_\theta)^2 = \sup_{f \in \mathcal{H}} \left\{ \langle f, \mu_p - \mu_{q_\theta} \rangle_{\mathcal{H}} - \frac{1}{2} \langle f, (\alpha D_\theta + \beta I) f \rangle_{\mathcal{H}} \right\}$$

$$\text{MMD}_{\alpha,\beta}^2(p, q_\theta) = \frac{1}{2} \langle \mu_p - \mu_{q_\theta}, (\alpha D_\theta + \beta I)^{-1} \mu_p - \mu_{q_\theta} \rangle_{\mathcal{H}}$$

witness function: $(\alpha D_\theta + \beta I) f = \mu_p - \mu_{q_\theta}$

Parametric Regularized MMD GAN

$$d\theta_t = L_{\theta_t}(f_t) dt$$

choice of f_t $(\alpha D(\theta_t) + \beta I) f_t = \mu_p - \mu_{q_{\theta_t}}$

$$q_{\theta_t} = (G_{\theta_t})_\#(\nu)$$

Theorem 1 (Parametric Regularized Flows Decrease the MMD Distance). Assume that $\alpha, \beta > 0$. Then the dynamic defined by the witness function of the parametric regularized MMD decreases the functional $\mathcal{F}(q_\theta)$:

$$\frac{d\mathcal{F}(q_{\theta_t})}{dt} = -\frac{2}{\alpha} [\mathcal{F}(q_{\theta_t}) - \beta \text{MMD}_{\alpha,\beta}(p, q_{\theta_t})^2] \leq 0. \quad (1)$$

Moreover, we have $\frac{d\mathcal{F}(q_{\theta_t})}{dt} < 0$ if and only if $D_{\theta_t} \mu_p - q_{\theta_t} \neq 0$.

Riemannian Structure and Gradient Flow

Definition 1 (Regularized MMD on a Statistical Manifold). Let $\alpha, \beta > 0$. Define

$$d_{\alpha,\beta}(q_{\theta_0}, q_{\theta_1})^2 = \min_{\theta_t, f_t} \int_0^1 \left(\alpha \|D_{\theta_t} f_t\|_{\mathcal{H}}^2 + \beta \langle f_t, D_{\theta_t} f_t \rangle_{\mathcal{H}} \right) dt,$$

$$\partial_t \theta_t = L_{\theta_t} f_t, \quad f_t \in \mathcal{H}, \quad \theta_{t=0} = \theta_0, \quad \theta_{t=1} = \theta_1.$$

$$\mathcal{F}(q_\theta) = H(\mu_{q_\theta})$$

$$\partial_{\theta_t} [H(\mu_{q_\theta})] = \langle h_{\theta_t}, \partial_{\theta_t} [\mu_{q_\theta}] \rangle_{\mathcal{H}}$$

Consider

$$\partial_t \theta_t = -\text{grad}_{d_{\alpha,\beta}} \mathcal{F}(q_{\theta_t}) = -L_{\theta_t} u_t, \quad (1)$$

where

$$(\alpha D_{\theta_t} + \beta I) u_t = h_{\theta_t}. \quad (2)$$

MMD GAN as A Gradient Flow w.r.t $d_{\{\alpha\beta\}}$

- Parametric Gradient Regularized MMD GAN update

$$d\theta_t = L_{\theta_t}(f_t) dt$$

$$\text{choice of } f_t \quad (\alpha D(\theta_t) + \beta I) f_t = \mu_p - \mu_{q_{\theta_t}}$$

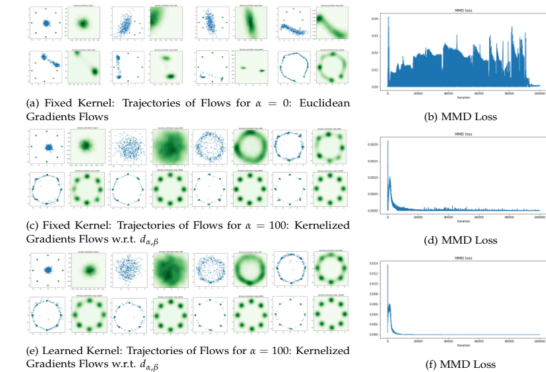
$$q_{\theta_t} = (G_{\theta_t})_\#(\nu)$$



$$\partial_t \theta_t = -\text{grad}_{d_{\alpha,\beta}} \frac{1}{2} \text{MMD}^2(p, q_{\theta_t})$$

- Parametric Gradient Flow of MMD w.r.t $d_{\{\alpha\beta\}}$

Experiments



- Target Distribution is a mixture of Gaussian
- (a) Unregularized MMD GAN goes through cycles, **does not converge**
- (c) MMD GAN regularized with Parametric Gradients with fixed Kernel: **convergent**
- (d) MMD GAN regularized with Parametric Gradients with Learned Kernel: **convergent**

Questions ,

mroueh@us.ibm.com
tn8@uakron.edu