

Gradient Flows on the Maximum Mean Discrepancy

Arthur Gretton

Gatsby Computational Neuroscience Unit,
Google Deepmind

First International Conference on Prob-
abilistic Numerics (Probnum 2025)

Outline

MMD and MMD flow

- Introduction to MMD as an integral probability metric
- Connection with neural net training
- Wasserstein-2 Gradient Flow on the MMD
- Convergence: adaptive kernel
 - Neural Net implementation
 - Interpolation to χ^2

Arbel, Korba, Salim, G., Maximum Mean Discrepancy Gradient Flow (NeurIPS 2019)
Galashov, De Bortoli, G., Deep MMD Gradient Flow without adversarial training
(ICLR 2025)

Chen, Mustafi, Glaser, Korba. G, Sriperumbudur (De)-regularized Maximum
Mean Discrepancy Gradient Flow (submitted JMLR)

Outline

MMD and MMD flow

- Introduction to MMD as an integral probability metric
- Connection with neural net training
- Wasserstein-2 Gradient Flow on the MMD
- Convergence: adaptive kernel
 - Neural Net implementation
 - Interpolation to χ^2

Main motivation: gradient flow when the target distribution represented by samples

- A different kind of particle flow to diffusion models
- Neural network training dynamics

Arbel, Korba, Salim, G., Maximum Mean Discrepancy Gradient Flow (NeurIPS 2019)
Galashov, De Bortoli, G., Deep MMD Gradient Flow without adversarial training
(ICLR 2025)

Chen, Mustafi, Glaser, Korba, G, Sriperumbudur (De)-regularized Maximum Mean Discrepancy Gradient Flow (submitted JMLR)

The MMD, and MMD flow

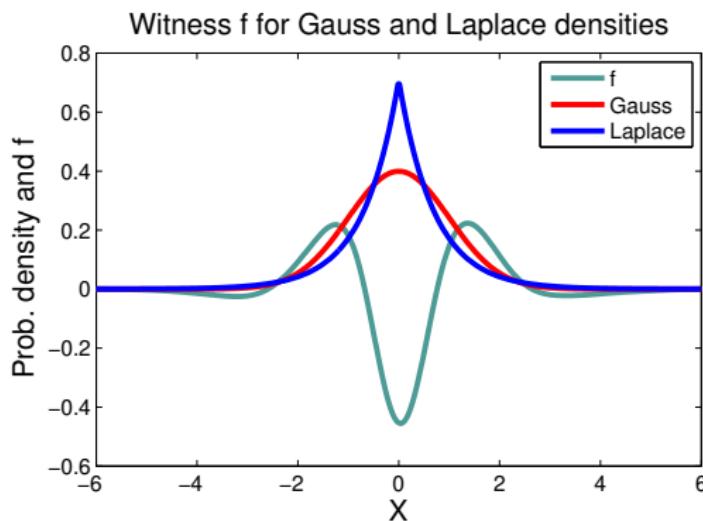
The MMD: an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$MMD(P, Q) := \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

$$f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}}$$

$$\langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}} = k(x, x')$$



The MMD: an integral probability metric

Maximum mean discrepancy: smooth function for P vs Q

$$MMD(P, Q) := \sup_{\|f\| \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

$$f(x) = \langle f, \varphi(x) \rangle_{\mathcal{H}}$$

$$\langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}} = k(x, x')$$

For characteristic RKHS \mathcal{H} , $MMD(P, Q) = 0$ iff $P = Q$

Other choices for witness function class:

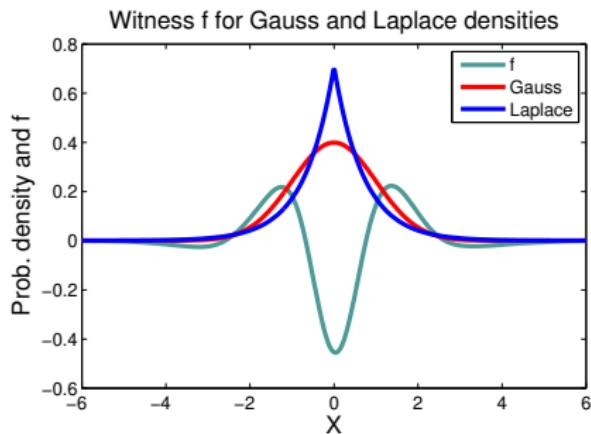
- Bounded continuous [Dudley, 2002]
- Bounded variation 1 (Kolmogorov metric) [Müller, 1997]
- Bounded Lipschitz (Wasserstein distances) [Dudley, 2002]

The MMD and witness in closed form

The MMD:

$$MMD(P, Q)$$

$$= \sup_{\|f\|_{\mathcal{H}} \leq 1} [E_P f(X) - E_Q f(Y)]$$



The MMD and witness in closed form

The MMD:

$$MMD(P, Q)$$

$$= \sup_{\|f\|_{\mathcal{H}} \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

$$= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{H}}$$

use

$$\begin{aligned}\mathbb{E}_P f(X) &= \mathbb{E}_P \langle \varphi(X), f \rangle_{\mathcal{H}} \\ &= \langle \mathbb{E}_P [\varphi(X)], f \rangle_{\mathcal{H}} \\ &= \langle \mu_P, f \rangle_{\mathcal{H}}\end{aligned}$$

The MMD and witness in closed form

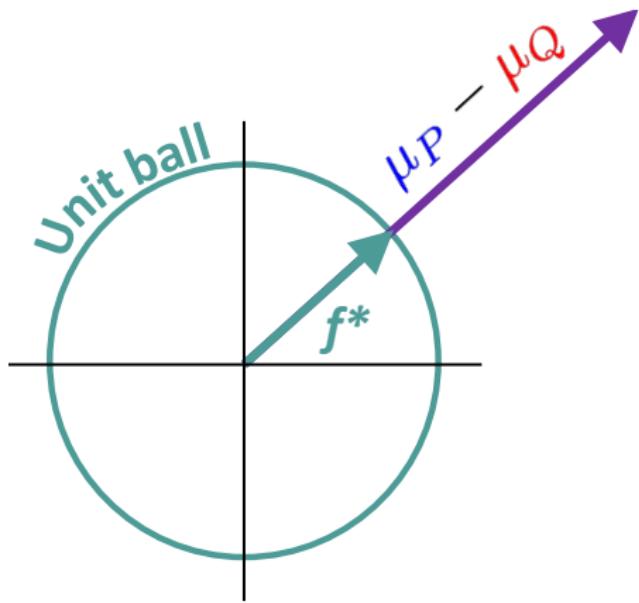
The MMD:

$$MMD(P, Q)$$

$$= \sup_{\|f\|_{\mathcal{H}} \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$

$$= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{H}}$$

$$= \|\mu_P - \mu_Q\|_{\mathcal{H}}$$



$$f^* = \frac{\mu_P - \mu_Q}{\|\mu_P - \mu_Q\|}$$

The MMD and witness in closed form

The MMD:

$$\begin{aligned} MMD(\mathcal{P}, \mathcal{Q}) &= \sup_{\|\mathbf{f}\|_{\mathcal{H}} \leq 1} [\mathbb{E}_{\mathcal{P}} f(\mathcal{X}) - \mathbb{E}_{\mathcal{Q}} f(\mathcal{Y})] \\ &= \sup_{\|\mathbf{f}\|_{\mathcal{H}} \leq 1} \langle \mathbf{f}, \boldsymbol{\mu}_{\mathcal{P}} - \boldsymbol{\mu}_{\mathcal{Q}} \rangle_{\mathcal{H}} \\ &= \|\boldsymbol{\mu}_{\mathcal{P}} - \boldsymbol{\mu}_{\mathcal{Q}}\|_{\mathcal{H}} \end{aligned}$$

$$\begin{aligned} \mathbf{f}^*(x) &\propto \langle \boldsymbol{\mu}_{\mathcal{P}} - \boldsymbol{\mu}_{\mathcal{Q}}, \varphi(x) \rangle_H \\ &= \mathbb{E}_{\mathcal{P}} k(\mathcal{X}, x) - \mathbb{E}_{\mathcal{Q}} k(\mathcal{Y}, x) \end{aligned}$$

The MMD and witness in closed form

The MMD:

$$\begin{aligned}MMD(P, Q) &= \sup_{\|\mathbf{f}\|_{\mathcal{H}} \leq 1} [\mathbb{E}_P f(\mathbf{X}) - \mathbb{E}_Q f(\mathbf{Y})] \\&= \sup_{\|\mathbf{f}\|_{\mathcal{H}} \leq 1} \langle \mathbf{f}, \mu_P - \mu_Q \rangle_{\mathcal{H}} \\&= \|\mu_P - \mu_Q\|_{\mathcal{H}}\end{aligned}$$

In terms of kernels:

$$\begin{aligned}MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{H}}^2 \\&= \underbrace{\mathbb{E}_P k(\mathbf{x}, \mathbf{x}')}_{(a)} + \underbrace{\mathbb{E}_Q k(\mathbf{y}, \mathbf{y}')}_{(a)} - 2 \underbrace{\mathbb{E}_{P, Q} k(\mathbf{x}, \mathbf{y})}_{(b)}\end{aligned}$$

(a)= within distrib. similarity, (b)= cross-distrib. similarity.

MMD Flow (NeurIPS 19)

arXiv > stat > arXiv:1906.04370

Statistics > Machine Learning

[Submitted on 11 Jun 2019 ([v1](#)), last revised 3 Dec 2019 (this version, v2)]

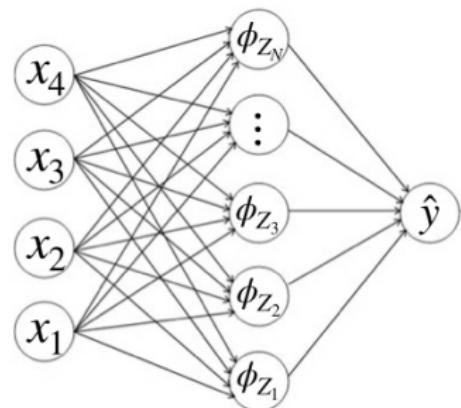
Maximum Mean Discrepancy Gradient Flow

[Michael Arbel](#), [Anna Korba](#), [Adil Salim](#), Arthur Gretton



Motivation: Neural Net training

$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} [\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2]$$

$$\min_{Z_1, \dots, Z_N \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{n} \sum_{i=1}^n \delta_{Z_i} \right)$$

Optimization using gradient descent:

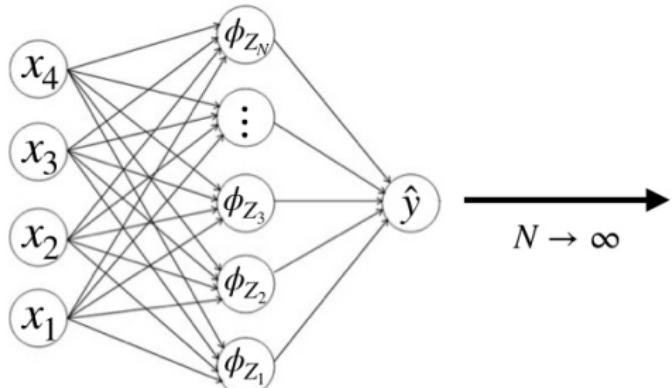
$$Z_i^{t+1} = Z_i^t - \gamma \nabla_{Z_i} \mathcal{L} \left(\frac{1}{n} \sum_{i=1}^n \delta_{Z_i^t} \right)$$

Chizat, Bach. "On the global convergence of gradient descent for over-parameterized models using optimal transport", NeurIPS (2018)

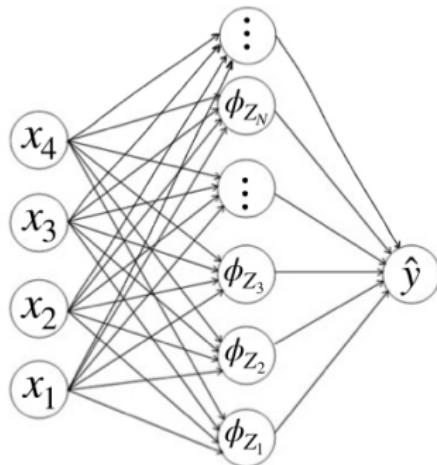
Motivation: Neural Net training

$$\min_{Z_1, \dots, Z_n \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{n} \sum_{i=1}^n \delta_{Z_i} \right) \xrightarrow{n \rightarrow \infty} \min_{\nu \in \mathcal{P}} \mathcal{L}(\nu)$$

$(x, y) \sim data$



$$N \rightarrow \infty$$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} [\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2] \xrightarrow{N \rightarrow \infty} \min_{\nu \in \mathcal{P}} \mathbb{E}_{data} [\|y - \mathbb{E}_{Z \sim \nu}[\phi_Z(x)]\|^2]$$

Chizat, Bach. "On the global convergence of gradient descent for over-parameterized models using optimal transport", NeurIPS (2018)

Motivation: Neural Net training

From previous slide:

$$\min_{\nu \in \mathcal{P}} \mathcal{L}(\nu) := \mathbb{E}_{(x,y)}[\|y - \mathbb{E}_{Z \sim \nu}[\phi_Z(x)]\|^2]$$

Chizat, Bach. "On the global convergence of gradient descent for over-parameterized models using optimal transport", NeurIPS (2018)

Motivation: Neural Net training

From previous slide:

$$\min_{\nu \in \mathcal{P}} \mathcal{L}(\nu) := \mathbb{E}_{(x,y)} [\|y - \mathbb{E}_{Z \sim \nu} [\phi_Z(x)]\|^2]$$

Connection to the MMD:

- Assume well-specified setting, $y(x) = \mathbb{E}_{U \sim \nu^*} [\phi_U(x)]$
- Random feature formulation,

$$\mathcal{L}(\nu) = \mathbb{E}_x \left[\|\mathbb{E}_{U \sim \nu^*} [\phi_U(x)] - \mathbb{E}_{Z \sim \nu} [\phi_Z(x)]\|^2 \right] = MMD^2(\nu, \nu^*)$$

- The kernel is: $k(U, Z) = \mathbb{E}_x [\phi_U(x)^\top \phi_Z(x)].$

Chizat, Bach. "On the global convergence of gradient descent for over-parameterized models using optimal transport", NeurIPS (2018)

Intuition: MMD as “force field” on ν

Assume henceforth

$$\nu, \nu^* \in \mathcal{P}_2(\mathbb{R}^d) := \left\{ \mu \in \mathcal{P}(\mathbb{R}^d) : \int \|x\|^2 d\mu(x) < \infty \right\}.$$

MMD as free energy: target ν^* , current distribution ν

$$\mathcal{F}(\nu) := \frac{1}{2} MMD^2(\nu^*, \nu) = \underbrace{\frac{1}{2} \mathbb{E}_{\nu} k(\mathbf{x}, \mathbf{x}')}_{\text{interaction}} + \underbrace{\frac{1}{2} \mathbb{E}_{\nu^*} k(\mathbf{y}, \mathbf{y}')}_{\text{constant}} - \underbrace{\mathbb{E}_{\nu, \nu^*} k(\mathbf{x}, \mathbf{y})}_{\text{confinement}}$$

[A] Ambrosio, Gigli, and Savaré. Gradient flows: in metric spaces and in the space of probability measures. (2008)

Intuition: MMD as “force field” on ν

Assume henceforth

$$\nu, \nu^* \in \mathcal{P}_2(\mathbb{R}^d) := \left\{ \mu \in \mathcal{P}(\mathbb{R}^d) : \int \|x\|^2 d\mu(x) < \infty \right\}.$$

MMD as free energy: target ν^* , current distribution ν

$$\mathcal{F}(\nu) := \frac{1}{2} MMD^2(\nu^*, \nu) = \frac{1}{2} \underbrace{\mathbb{E}_{\nu} k(\mathbf{x}, \mathbf{x}')}_{\text{interaction}} + \frac{1}{2} \underbrace{\mathbb{E}_{\nu^*} k(\mathbf{y}, \mathbf{y}')}_{\text{constant}} - \underbrace{\mathbb{E}_{\nu, \nu^*} k(\mathbf{x}, \mathbf{y})}_{\text{confinement}}$$

Consider $\{\mathbf{y}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \nu^*$ and $\{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \nu$.

Force on a particle \mathbf{z} :

$$-\sum_j \nabla_{\mathbf{z}} k(\mathbf{z}, \mathbf{x}_j) + \sum_j \nabla_{\mathbf{z}} k(\mathbf{z}, \mathbf{y}_j) = -\nabla_{\mathbf{z}} \hat{f}_{\nu^*, \nu_t}(z)$$

Can we formalize this?

[A] Ambrosio, Gigli, and Savaré. Gradient flows: in metric spaces and in the space of probability measures. (2008)

Wasserstein gradient flows

Tangent space of $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ at μ is $h \in L^2(\mu)$ where $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$.

Define $\nabla_{W_2}\mathcal{F}(\mu)$ of \mathcal{F} at μ using Taylor expansion

$$\mathcal{F}((\text{Id} + \epsilon h)_{\# \mu}) = \mathcal{F}(\mu) + \epsilon \langle \nabla_{W_2}\mathcal{F}(\mu), h \rangle_{L^2(\mu)} + o(\epsilon) \quad (1)$$

[A] Ambrosio, Gigli, and Savaré. Gradient flows: in metric spaces and in the space of probability measures. (2008)

Wasserstein gradient flows

Tangent space of $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ at μ is $h \in L^2(\mu)$ where $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$.

Define $\nabla_{W_2}\mathcal{F}(\mu)$ of \mathcal{F} at μ using Taylor expansion

$$\mathcal{F}((\text{Id} + \epsilon h)_{\#}\mu) = \mathcal{F}(\mu) + \epsilon \langle \nabla_{W_2}\mathcal{F}(\mu), h \rangle_{L^2(\mu)} + o(\epsilon) \quad (1)$$

The gradient flow is then:

$$\partial_t \nu_t = \text{div}(\nu_t \nabla_{W_2}\mathcal{F}(\nu_t))$$

[A] Ambrosio, Gigli, and Savaré. Gradient flows: in metric spaces and in the space of probability measures. (2008)

Wasserstein gradient flows

Tangent space of $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ at μ is $h \in L^2(\mu)$ where $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$.
Define $\nabla_{W_2}\mathcal{F}(\mu)$ of \mathcal{F} at μ using Taylor expansion

$$\mathcal{F}((\text{Id} + \epsilon h)_{\#}\mu) = \mathcal{F}(\mu) + \epsilon \langle \nabla_{W_2}\mathcal{F}(\mu), h \rangle_{L^2(\mu)} + o(\epsilon) \quad (1)$$

The gradient flow is then:

$$\partial_t \nu_t = \text{div}(\nu_t \nabla_{W_2}\mathcal{F}(\nu_t))$$

Under reasonable assumptions [A. Theorem 10.4.13]

$$\nabla_{W_2}\mathcal{F}(\mu) = \nabla \mathcal{F}'(\mu).$$

where **first variation** in direction ξ :

$$\mathcal{F}(\mu + \epsilon \xi) = \mathcal{F}(\mu) + \epsilon \int \mathcal{F}'(\mu)(x) d\xi(x) + o(\epsilon) \quad \mu + \epsilon \xi \in \mathcal{P}_2(\mathbb{R}^d) \quad (2)$$

[A] Ambrosio, Gigli, and Savaré. Gradient flows: in metric spaces and in the space of probability measures. (2008)

Wasserstein gradient flow on MMD

First variation of $\frac{1}{2} MMD^2(\nu^\star, \nu) =: \mathcal{F}(\nu)$

$$\mathcal{F}'(\nu)(z) := f_{\nu^\star, \nu}(z) = 2(\mathbb{E}_{U \sim \nu^\star}[k(U, z)] - \mathbb{E}_{U \sim \nu}[k(U, z)])$$

The W_2 gradient flow of the MMD:

$$\partial_t \nu_t = \operatorname{div}(\nu_t \nabla_{W_2} \mathcal{F}(\nu_t)) = \operatorname{div}(\nu_t \nabla f_{\nu^\star, \nu_t})$$

Ambrosio, Gigli, and Savaré. Gradient flows: in metric spaces and in the space of probability measures. (2008, Ch. 10)

Mroueh, Sercu, and Raj. Sobolev Descent. (AISTATS, 2019)

Arbel, Korba, Salim, G. (NeurIPS 2019)

Wasserstein gradient flow on MMD

First variation of $\frac{1}{2} MMD^2(\nu^\star, \nu) =: \mathcal{F}(\nu)$

$$\mathcal{F}'(\nu)(z) := f_{\nu^\star, \nu}(z) = 2(\mathbb{E}_{U \sim \nu^\star}[k(U, z)] - \mathbb{E}_{U \sim \nu}[k(U, z)])$$

The W_2 gradient flow of the MMD:

$$\partial_t \nu_t = \operatorname{div}(\nu_t \nabla_{W_2} \mathcal{F}(\nu_t)) = \operatorname{div}(\nu_t \nabla f_{\nu^\star, \nu_t})$$

McKean-Vlasov dynamics for particles (existence and uniqueness under **Assumption A**):

$$dZ_t = -\nabla_{Z_t} f_{\nu^\star, \nu_t}(Z_t) dt, \quad Z_0 \sim \nu_0$$

Assumption A: $k(x, x) \leq K$, for all $x \in \mathbb{R}^d$, $\sum_{i=1}^d \|\partial_i k(x, \cdot)\|^2 \leq K_{1d}$ and $\sum_{i,j=1}^d \|\partial_i \partial_j k(x, \cdot)\|^2 \leq K_{2d}$, d indicates scaling with dimension.

Ambrosio, Gigli, and Savaré. Gradient flows: in metric spaces and in the space of probability measures. (2008, Ch. 10)

Mroueh, Sercu, and Raj. Sobolev Descent. (AISTATS, 2019)

Arbel, Korba, Salim, G. (NeurIPS 2019)

Wasserstein gradient flow on the MMD

Forward Euler scheme [A, Section 2.2]:

$$\begin{aligned}\nu_{n+1} &= (I - \gamma \nabla f_{\nu^*, \nu_t})_{\#} \nu_n \\ Z_{n+1} &= Z_n - \gamma \nabla_{Z_n} f_{\nu^*, \nu_n}(Z_n), \quad Z_0 \sim \nu_0, \quad Z_n \sim \nu_n\end{aligned}$$

Under **Assumption A**, ν_n approaches ν_t as $\gamma \rightarrow 0$

Wasserstein gradient flow on the MMD

Forward Euler scheme [A, Section 2.2]:

$$\begin{aligned}\nu_{n+1} &= (I - \gamma \nabla f_{\nu^*, \nu_t})_{\#} \nu_n \\ Z_{n+1} &= Z_n - \gamma \nabla_{Z_n} f_{\nu^*, \nu_n}(Z_n), \quad Z_0 \sim \nu_0, \quad Z_n \sim \nu_n\end{aligned}$$

Under **Assumption A**, ν_n approaches ν_t as $\gamma \rightarrow 0$

Consistency? Does ν_t converge to ν^* as $t \rightarrow \infty$?

[A] Arbel, Korba, Salim, G. (NeurIPS 2019)

Consistency

Can we use geodesic (displacement) convexity?

- A geodesic ρ_t between ν_1 and ν_2 is given by the transport map $T_{\nu_1}^{\nu_2} : \mathbb{R}^d \rightarrow \mathbb{R}^d$:

$$\rho_t = ((1-t)\text{Id} + tT_{\nu_1}^{\nu_2})_{\# \nu_1}$$

Consistency

Can we use geodesic (displacement) convexity?

- A geodesic ρ_t between ν_1 and ν_2 is given by the transport map $T_{\nu_1}^{\nu_2} : \mathbb{R}^d \rightarrow \mathbb{R}^d$:

$$\rho_t = ((1-t)\text{Id} + tT_{\nu_1}^{\nu_2})_{\# \nu_1}$$

- A functional \mathcal{F} is displacement convex if:

$$\mathcal{F}(\rho_t) \leq (1-t)\mathcal{F}(\nu_1) + t\mathcal{F}(\nu_2)$$

Consistency

Can we use geodesic (displacement) convexity?

- A geodesic ρ_t between ν_1 and ν_2 is given by the transport map $T_{\nu_1}^{\nu_2} : \mathbb{R}^d \rightarrow \mathbb{R}^d$:

$$\rho_t = ((1-t)\text{Id} + tT_{\nu_1}^{\nu_2})_{\# \nu_1}$$

- A functional \mathcal{F} is displacement convex if:

$$\mathcal{F}(\rho_t) \leq (1-t)\mathcal{F}(\nu_1) + t\mathcal{F}(\nu_2)$$

MMD is not displacement convex in general
(it is always mixture convex¹).

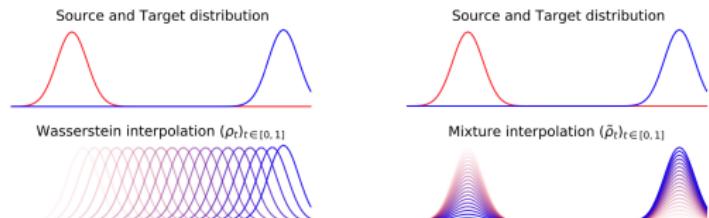
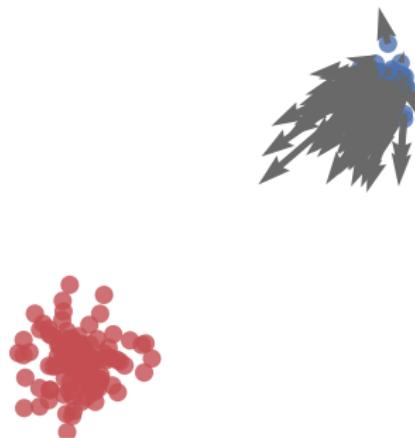


Figure from Korba, Salim, ICML 2022 Tutorial, "Sampling as First-Order Optimization over a space of probability measures"

1. $\mathcal{F}(t\nu_1 + (1-t)\nu_2) \leq t\mathcal{F}(\nu_1) + (1-t)\mathcal{F}(\nu_2) \quad \forall t \in [0, 1]).$

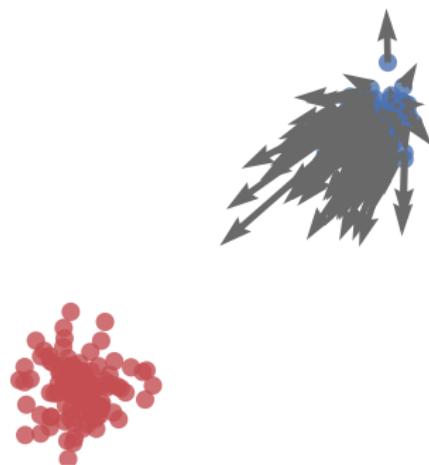
MMD flow in practice

- Data
- Particles



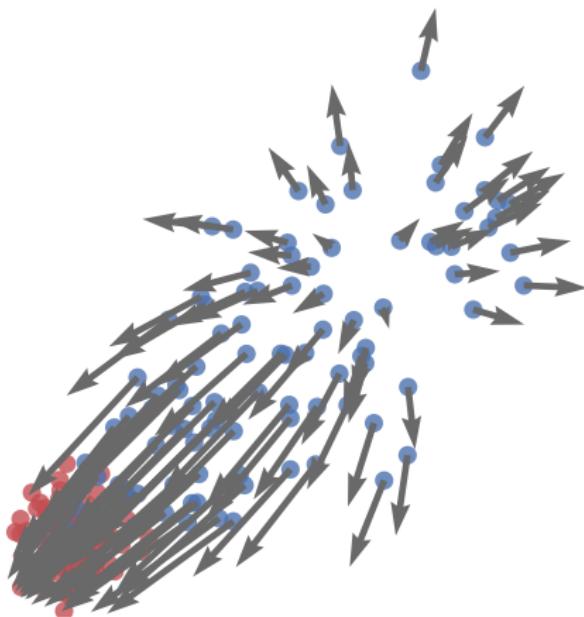
MMD flow in practice

- Data
- Particles

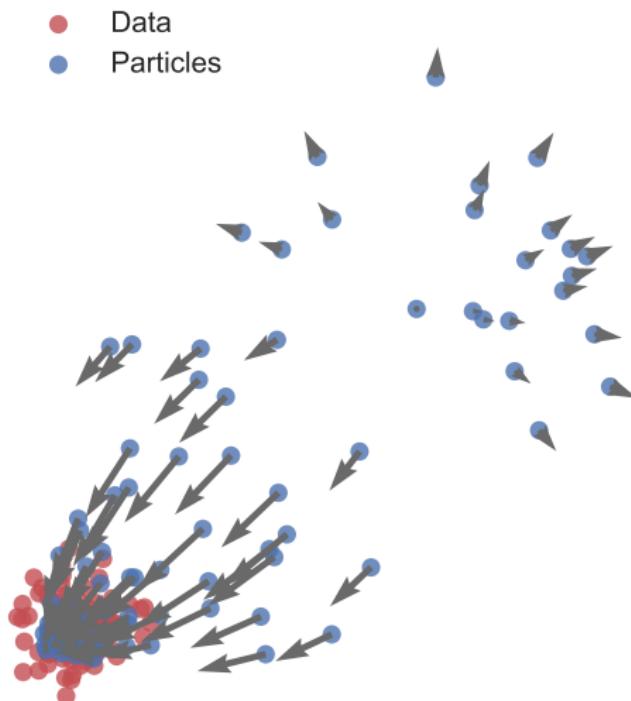


MMD flow in practice

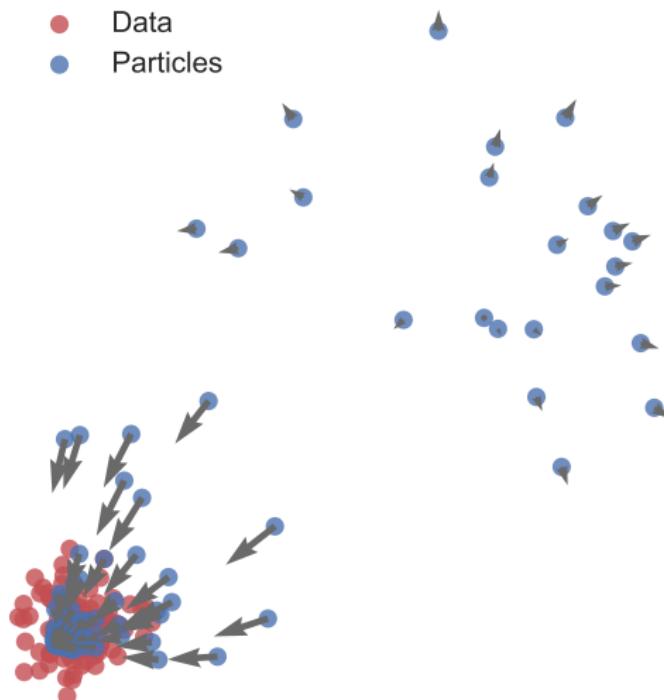
- Data
- Particles



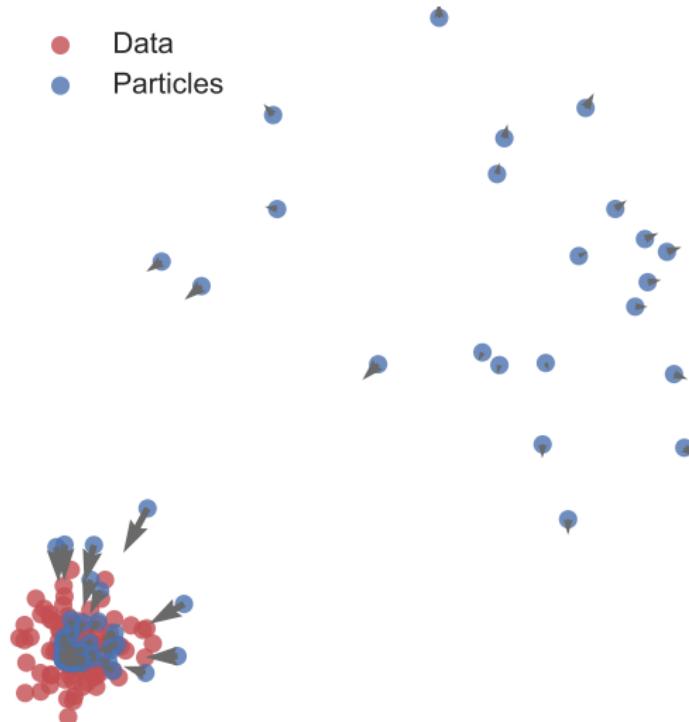
MMD flow in practice



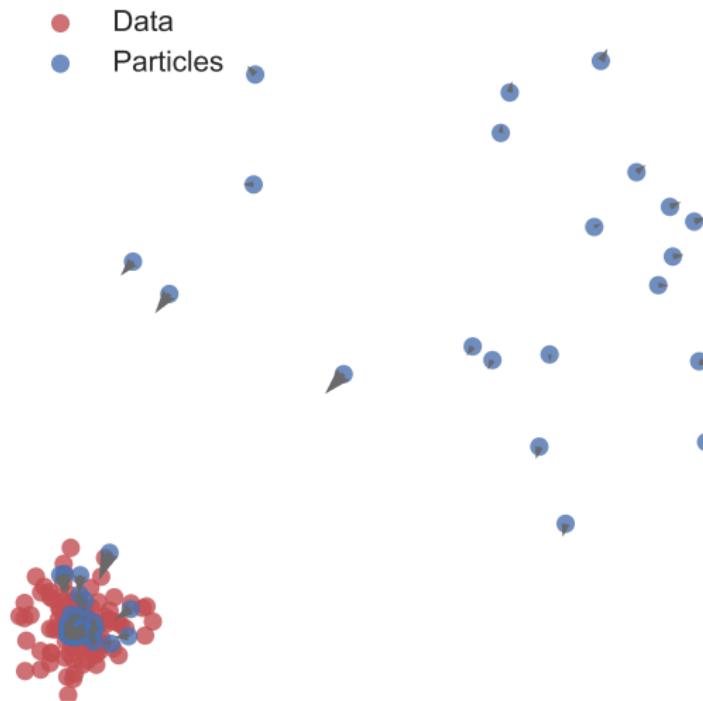
MMD flow in practice



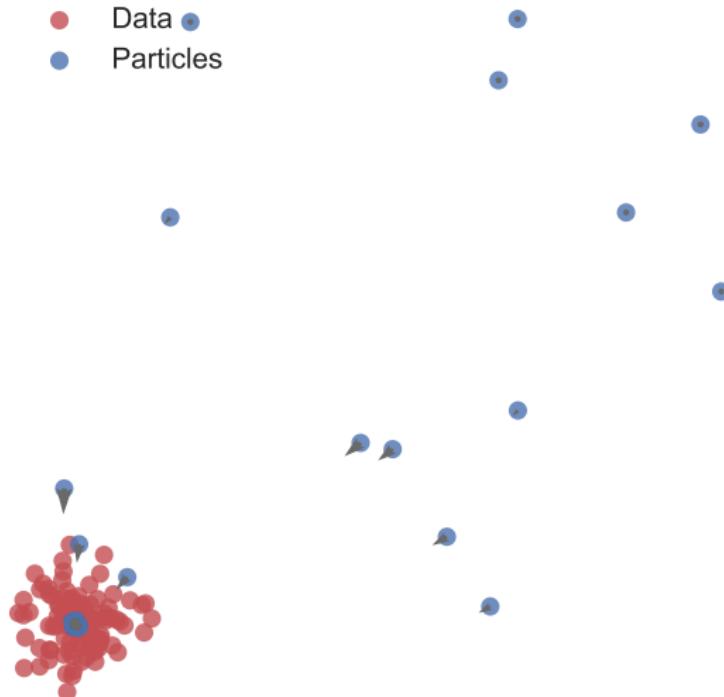
MMD flow in practice



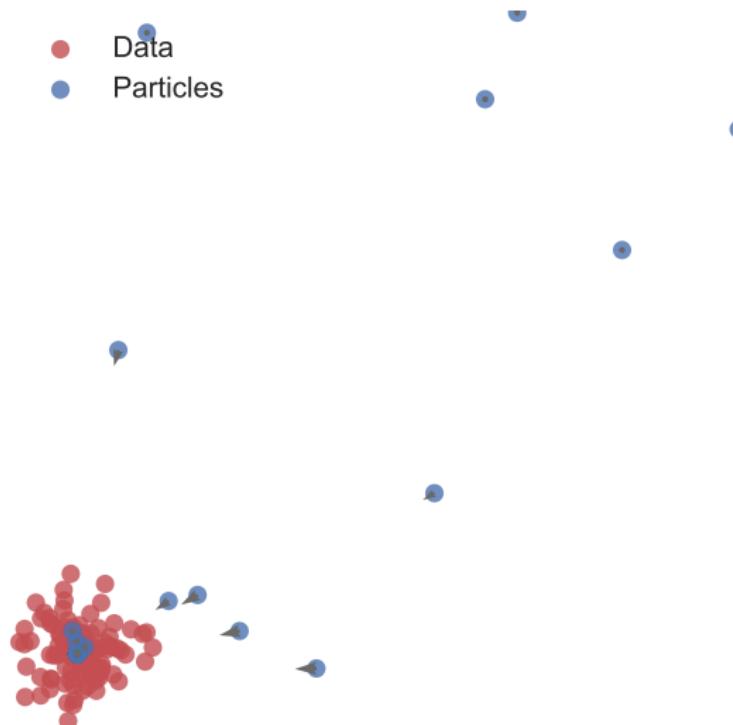
MMD flow in practice



MMD flow in practice

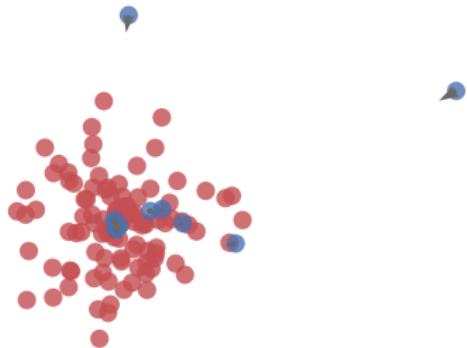


MMD flow in practice

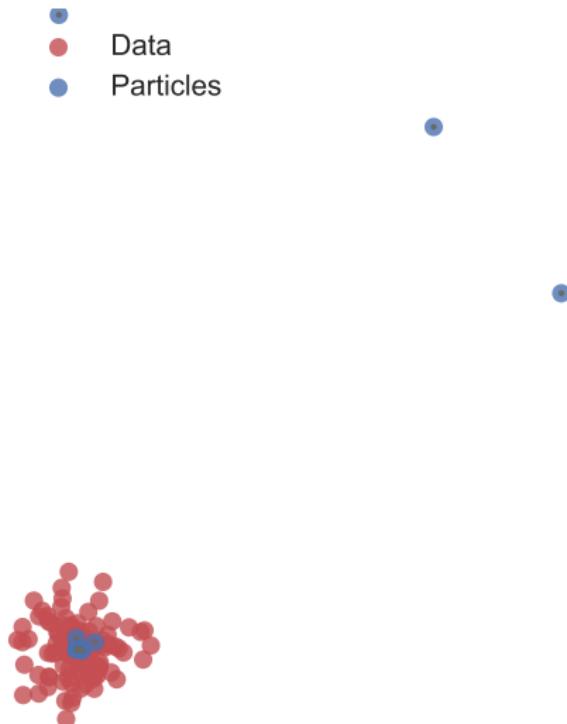


MMD flow in practice

- Data
- Particles



MMD flow in practice



Empirical observations

Some observations:

- Almost all particles tend to collapse at the center of mass m of the target ν^* , i.e.: ($\nu_t \simeq \delta_m$)
 - However, the loss stops decreasing: $\nabla f_{\nu^*, \nu_t}(z) \simeq 0$ for z on the support of ν_t (and is small when far from ν^*)...
 - ...and in general, $\nabla f_{\nu^*, \nu_t}(z) \neq 0$ outside the support of ν_t .

Empirical observations

Some observations:

- Almost all particles tend to collapse at the center of mass m of the target ν^* , i.e.: ($\nu_t \simeq \delta_m$)
 - However, the loss stops decreasing: $\nabla f_{\nu^*, \nu_t}(z) \simeq 0$ for z on the support of ν_t (and is small when far from ν^*)...
 - ...and in general, $\nabla f_{\nu^*, \nu_t}(z) \neq 0$ outside the support of ν_t .

Idea: Adapt the kernel according to distance of ν_t to ν^* .

- “Broad” kernel when distributions far apart,
- “narrow” kernel when they are close.

Noise injection in NeurIPS 2019 was a first attempt.

Noise injection for convergence

Noise injection: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!

- Sample $u_t \sim \mathcal{N}(0, 1)$ and β_t is the noise level:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$$

- Similar to continuation methods,¹ but extended to interacting particles.
- Different from entropic regularization:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t) + \beta_t u_t$$

- Blur RKHS kernel with t -dependent Gaussian noise

Noise injection for convergence

Noise injection: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!

- Sample $u_t \sim \mathcal{N}(0, 1)$ and β_t is the noise level:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$$

- Similar to continuation methods,¹ but extended to interacting particles.
- Different from entropic regularization:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t) + \beta_t u_t$$

- Blur RKHS kernel with t -dependent Gaussian noise

Noise injection for convergence

Noise injection: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!

- Sample $u_t \sim \mathcal{N}(0, 1)$ and β_t is the noise level:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$$

- Similar to continuation methods,¹ but extended to interacting particles.
- Different from entropic regularization:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t) + \beta_t u_t$$

- Blur RKHS kernel with t -dependent Gaussian noise

Noise injection for convergence

Noise injection: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!

- Sample $u_t \sim \mathcal{N}(0, 1)$ and β_t is the noise level:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$$

- Similar to continuation methods,¹ but extended to interacting particles.
- Different from entropic regularization:

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t) + \beta_t u_t$$

- Blur RKHS kernel with t -dependent Gaussian noise

Noise injection: consistency

Recall: $Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$

Tradeoff for β_t

- Large β_t : $\nu_{t+1} - \nu_t$ not a descent direction any more:
 $\mathcal{F}(\nu_{t+1}) > \mathcal{F}(\nu_t)$
- Small β_t : does not converge

Noise injection: consistency

Recall: $Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$

Tradeoff for β_t

- Large β_t : $\nu_{t+1} - \nu_t$ not a descent direction any more:
 $\mathcal{F}(\nu_{t+1}) > \mathcal{F}(\nu_t)$
- Small β_t : does not converge

Need β_t such that:

$$\mathcal{F}(\nu_{t+1}) - \mathcal{F}(\nu_t) \leq -C\gamma \mathbb{E}_{\substack{X_t \sim \nu_t \\ u_t \sim \mathcal{N}(0,1)}} [\|\nabla f_{\nu^*, \nu_t}(X_t + \beta_t u_t)\|^2]$$

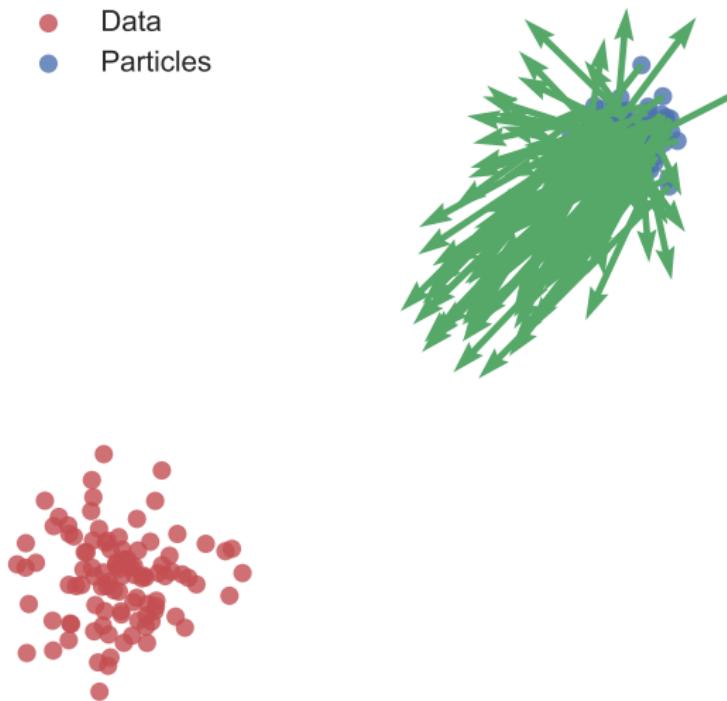
$$\sum_i^t \beta_i^2 \xrightarrow{t \rightarrow \infty} \infty$$

Then [A, Proposition 8]

$$\mathcal{F}(\nu_t) \leq \mathcal{F}(\nu_0) e^{-C\gamma \sum_i^t \beta_i^2}.$$

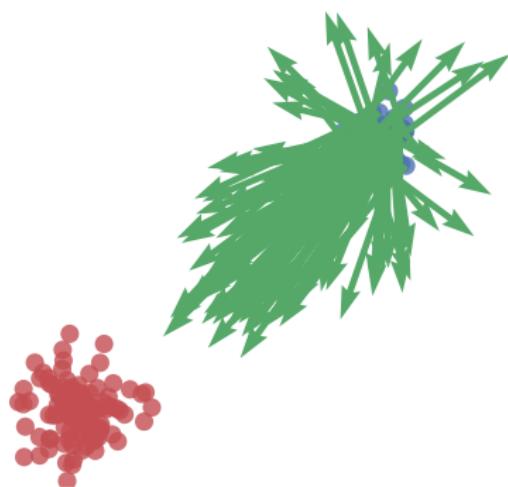
Noise injected MMD flow in practice

- Data
- Particles



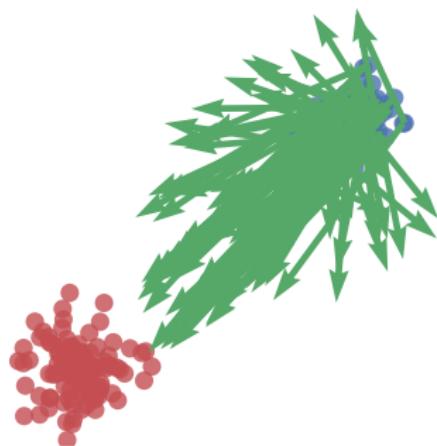
Noise injected MMD flow in practice

- Data
- Particles



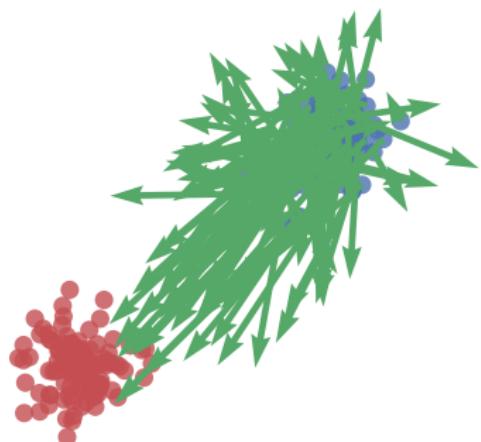
Noise injected MMD flow in practice

- Data
- Particles



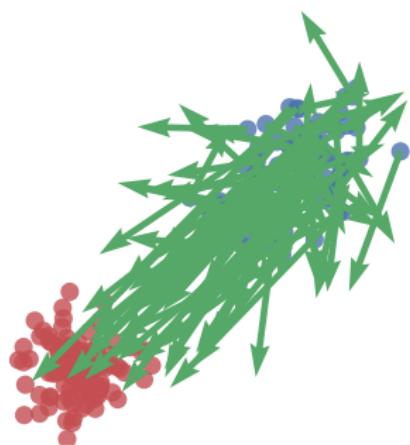
Noise injected MMD flow in practice

- Data
- Particles



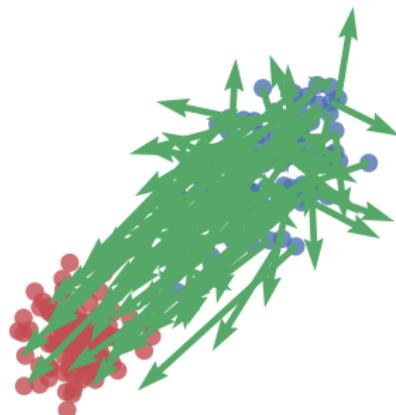
Noise injected MMD flow in practice

- Data
- Particles



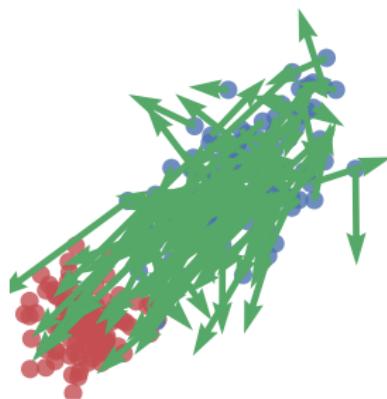
Noise injected MMD flow in practice

- Data
- Particles



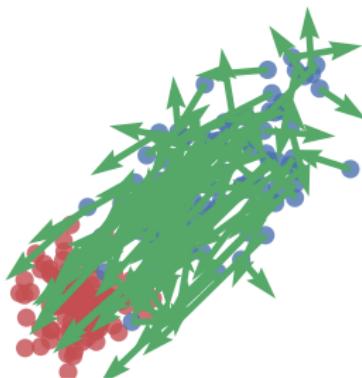
Noise injected MMD flow in practice

- Data
- Particles



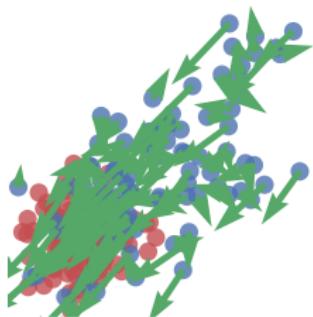
Noise injected MMD flow in practice

- Data
- Particles



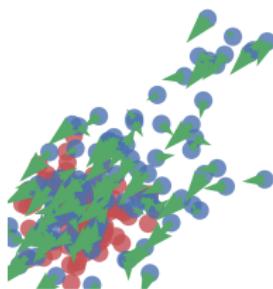
Noise injected MMD flow in practice

- Data
- Particles



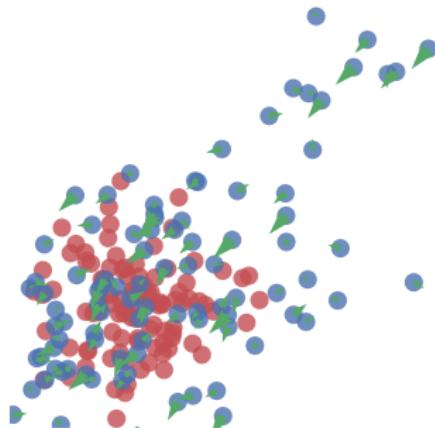
Noise injected MMD flow in practice

- Data
- Particles



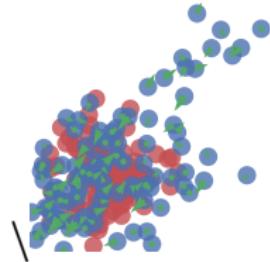
Noise injected MMD flow in practice

- Data
- Particles



Noise injected MMD flow in practice

- Data
- Particles



(De)-regularized MMD Gradient Flow (JMLR, submitted)

arXiv > stat > arXiv:2409.14980

Search..
Help | Ad

Statistics > Machine Learning

[Submitted on 23 Sep 2024]

(De)-regularized Maximum Mean Discrepancy Gradient Flow

Zonghao Chen, Aratrika Mustafi, Pierre Glaser, Anna Korba, Arthur Gretton, Bharath K. Sriperumbudur



χ^2 gradient flow has exponential convergence

Consider Wasserstein Gradient flow on χ^2 divergence,

$$\chi^2(\nu, \nu^*) = \left\| \frac{d\nu}{d\nu^*} - 1 \right\|_{L^2(\nu^*)}^2,$$

χ^2 gradient flow has exponential convergence

Consider Wasserstein Gradient flow on χ^2 divergence,

$$\chi^2(\nu, \nu^*) = \left\| \frac{d\nu}{d\nu^*} - 1 \right\|_{L^2(\nu^*)}^2,$$

Assume ν^* satisfies a Poincaré inequality,

$$\text{Var}_{\nu^*}[f] \leq C_P \mathbb{E}_{\nu^*} [\|\nabla f\|^2] \quad \forall f, \nabla f \in L^2(\nu^*)$$

E.g.: $C_P = \alpha$ if ν^* is α -log concave,

$$\nu^* \propto \exp(-V) \quad H V \succeq \alpha I$$

Detail: $\chi^2(\nu, \nu^*)$ satisfies a (modified) Polyak-Łojasiewicz inequality, a strict relaxation of strong convexity, when ν^* satisfies a Poincaré inequality.

χ^2 gradient flow has exponential convergence

Consider Wasserstein Gradient flow on χ^2 divergence,

$$\chi^2(\nu, \nu^*) = \left\| \frac{d\nu}{d\nu^*} - 1 \right\|_{L^2(\nu^*)}^2,$$

Assume ν^* satisfies a Poincaré inequality,

$$\text{Var}_{\nu^*}[f] \leq C_P \mathbb{E}_{\nu^*} [\|\nabla f\|^2] \quad \forall f, \nabla f \in L^2(\nu^*)$$

E.g.: $C_P = \alpha$ if ν^* is α -log concave,

$$\nu^* \propto \exp(-V) \quad H V \succeq \alpha I$$

Detail: $\chi^2(\nu, \nu^*)$ satisfies a (modified) Polyak-Łojasiewicz inequality, a strict relaxation of strong convexity, when ν^* satisfies a Poincaré inequality.

Convergence: Let $(\nu_t)_{t \geq 0}$ be the Wasserstein gradient flow of χ^2 . Then

$$\text{KL}(\nu_T, \nu^*) \leq \exp\left(-\frac{2T}{C_P}\right) \text{KL}(\nu_0, \nu^*) \quad \forall T \geq 0.$$

Can we interpolate between MMD and χ^2 ?

MMD vs the χ^2 divergence:

$$\chi^2(\nu, \nu^*) = \left\| \frac{d\nu}{d\nu^*} - 1 \right\|_{L^2(\nu^*)}^2,$$

$$\text{MMD}^2(\nu, \nu^*) = \left\| T_{\nu^*}^{\frac{1}{2}} \left(\frac{d\nu}{d\nu^*} - 1 \right) \right\|_{L^2(\nu^*)}^2$$

where

$$T_{\nu^*} f(\cdot) = \int k(x, \cdot) f(x) d\nu^*(x).$$

Can we interpolate between MMD and χ^2 ?

MMD vs the χ^2 divergence:

$$\chi^2(\nu, \nu^*) = \left\| \frac{d\nu}{d\nu^*} - 1 \right\|_{L^2(\nu^*)}^2,$$

$$\text{MMD}^2(\nu, \nu^*) = \left\| T_{\nu^*}^{\frac{1}{2}} \left(\frac{d\nu}{d\nu^*} - 1 \right) \right\|_{L^2(\nu^*)}^2$$

where

$$T_{\nu^*} f(\cdot) = \int k(x, \cdot) f(x) d\nu^*(x).$$

Deregularized MMD (DrMMD):

$$\text{DrMMD}^2(\nu, \nu^*) = (1 + \lambda) \left\| \left((T_{\nu^*} + \lambda \text{Id})^{-1} T_{\nu^*} \right)^{\frac{1}{2}} \left(\frac{d\nu}{d\nu^*} - 1 \right) \right\|_{L^2(\nu^*)}^2.$$

Can we interpolate between MMD and χ^2 ?

MMD vs the χ^2 divergence:

$$\chi^2(\nu, \nu^*) = \left\| \frac{d\nu}{d\nu^*} - 1 \right\|_{L^2(\nu^*)}^2,$$

$$\text{MMD}^2(\nu, \nu^*) = \left\| T_{\nu^*}^{\frac{1}{2}} \left(\frac{d\nu}{d\nu^*} - 1 \right) \right\|_{L^2(\nu^*)}^2$$

where

$$T_{\nu^*} f(\cdot) = \int k(x, \cdot) f(x) d\nu^*(x).$$

Deregularized MMD (DrMMD):

$$\text{DrMMD}^2(\nu, \nu^*) = (1 + \lambda) \left\| \left((T_{\nu^*} + \lambda \text{Id})^{-1} T_{\nu^*} \right)^{\frac{1}{2}} \left(\frac{d\nu}{d\nu^*} - 1 \right) \right\|_{L^2(\nu^*)}^2.$$

DrMMD interpolates between MMD and χ^2 :

$$\lim_{\lambda \rightarrow 0} \text{DrMMD}^2 = \chi^2, \quad \lim_{\lambda \rightarrow \infty} \text{DrMMD}^2 = \text{MMD}^2$$

for k bounded, continuous, c_0 -universal.

A kernel adaptation perspective

Eigendecomposition of kernel:

$$\varrho_i \psi_i(x) = [T_{\nu^*} \psi_i](x) = \int k(x, t) \psi_i(t) d\nu^*(t)$$

- MMD operator: $T_{\nu^*}^{\frac{1}{2}}$
- DrMMD operator: $(1 + \lambda)^{\frac{1}{2}} \left((T_{\nu^*} + \lambda \text{Id})^{-1} T_{\nu^*} \right)^{\frac{1}{2}}$

A kernel adaptation perspective

Eigendecomposition of kernel:

$$\varrho_i \psi_i(x) = [T_{\nu^*} \psi_i](x) = \int k(x, t) \psi_i(t) d\nu^*(t)$$

- MMD operator: $T_{\nu^*}^{\frac{1}{2}}$
- DrMMD operator: $(1 + \lambda)^{\frac{1}{2}} \left((T_{\nu^*} + \lambda \text{Id})^{-1} T_{\nu^*} \right)^{\frac{1}{2}}$

Interpret as **adjusting kernel eigenspectrum**:

$$(\text{MMD}) \quad \varrho_i \rightarrow (1 + \lambda) \frac{\varrho_i}{\varrho_i + \lambda} \quad (\text{DrMMD})$$

Wasserstein gradient of DrMMD

The Wasserstein gradient of DrMMD at ν is

$$(1 + \lambda) \nabla f_{\nu, \nu^*}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

where

$$f_{\nu, \nu^*} = (T_{\nu^*} + \lambda \text{Id})^{-1} T_{\nu^*} \left(\frac{d\nu}{d\nu^*} - 1 \right)$$

Wasserstein gradient of DrMMD

The Wasserstein gradient of DrMMD at ν is

$$(1 + \lambda) \nabla f_{\nu, \nu^*}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

where

$$\begin{aligned} f_{\nu, \nu^*} &= (T_{\nu^*} + \lambda \text{Id})^{-1} T_{\nu^*} \left(\frac{d\nu}{d\nu^*} - 1 \right) \\ &= (\Sigma_{\nu^*} + \lambda I)^{-1} \underbrace{(\int k(x, \cdot) d\nu - \int k(x, \cdot) d\nu^*)}_{\text{MMD witness}}. \end{aligned}$$

$\Sigma_{\nu^*} : \mathcal{H} \rightarrow \mathcal{H}$ is the covariance operator, defined

$$\langle f, \Sigma_{\nu^*} f \rangle_{\mathcal{H}} = \mathbb{E}_{\nu^*}[f(X)^2].$$

Wasserstein gradient of DrMMD

The Wasserstein gradient of DrMMD at ν is

$$(1 + \lambda) \nabla f_{\nu, \nu^*}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

where

$$\begin{aligned} f_{\nu, \nu^*} &= (T_{\nu^*} + \lambda \text{Id})^{-1} T_{\nu^*} \left(\frac{d\nu}{d\nu^*} - 1 \right) \\ &= (\Sigma_{\nu^*} + \lambda I)^{-1} \underbrace{(\int k(x, \cdot) d\nu - \int k(x, \cdot) d\nu^*)}_{\text{MMD witness}}. \end{aligned}$$

$\Sigma_{\nu^*} : \mathcal{H} \rightarrow \mathcal{H}$ is the covariance operator, defined

$$\langle f, \Sigma_{\nu^*} f \rangle_{\mathcal{H}} = \mathbb{E}_{\nu^*}[f(X)^2].$$

Relates to [kernel Fisher discriminant](#).

Mika, G. Rätsch, Weston, Schölkopf, and K.-R. Müller. Fisher discriminant analysis with kernels. (IEEE Neural Networks for Signal Processing 1999)

Harchaoui, Bach, and Moulines. Testing for homogeneity with kernel Fisher discriminant analysis (NeurIPS 2007)

Convergence of DrMMD flow

Suppose ν^* satisfies a Poincaré inequality with constant C_P ($= \alpha^{-1}$).

Suppose $\exists q_t \in L^2(\nu^*)$ such that

$$\frac{d\nu_t}{d\nu^*} - 1 = T_{\nu^*}^r q_t \quad \|q_t\|_{L^2(\nu^*)} < Q \quad T_{\nu^*} f(\cdot) = \int k(x, \cdot) f(x) d\nu^*(x).$$

Larger $r \rightarrow$ more regular trajectory.

Convergence of DrMMD flow

Suppose ν^* satisfies a Poincaré inequality with constant C_P ($= \alpha^{-1}$).

Suppose $\exists q_t \in L^2(\nu^*)$ such that

$$\frac{d\nu_t}{d\nu^*} - 1 = T_{\nu^*}^r q_t \quad \|q_t\|_{L^2(\nu^*)} < Q \quad T_{\nu^*} f(\cdot) = \int k(x, \cdot) f(x) d\nu^*(x).$$

Larger $r \rightarrow$ more regular trajectory.

Near-global convergence of DrMMD flow ν_t :

$$\text{KL}(\nu_T, \nu^*) \leq \exp\left(-\frac{2}{C_P} T\right) \text{KL}(\nu_0, \nu^*) + \lambda^r C_P Q (\mathcal{J} + \mathcal{I}) \quad \forall T \geq 0$$

Fine print: assume

$$\left\| \nabla(\log \nu^*) \nabla \left(\frac{d\nu_t}{d\nu^*} \right) \right\|_{L^2(\nu^*)} \leq \mathcal{J} \quad \left\| \Delta \left(\frac{d\nu_t}{d\nu^*} \right) \right\|_{L^2(\nu^*)} \leq \mathcal{I}$$

Convergence of DrMMD flow

Suppose ν^* satisfies a Poincaré inequality with constant C_P ($= \alpha^{-1}$).

Suppose $\exists q_t \in L^2(\nu^*)$ such that

$$\frac{d\nu_t}{d\nu^*} - 1 = T_{\nu^*}^r q_t \quad \|q_t\|_{L^2(\nu^*)} < Q \quad T_{\nu^*} f(\cdot) = \int k(x, \cdot) f(x) d\nu^*(x).$$

Larger $r \rightarrow$ more regular trajectory.

Near-global convergence of DrMMD flow ν_t :

$$\text{KL}(\nu_T, \nu^*) \leq \exp\left(-\frac{2}{C_P}T\right) \text{KL}(\nu_0, \nu^*) + \lambda^r C_P Q (\mathcal{J} + \mathcal{I}) \quad \forall T \geq 0$$

Fine print: assume

$$\left\| \nabla(\log \nu^*) \nabla \left(\frac{d\nu_t}{d\nu^*} \right) \right\|_{L^2(\nu^*)} \leq \mathcal{J} \quad \left\| \Delta \left(\frac{d\nu_t}{d\nu^*} \right) \right\|_{L^2(\nu^*)} \leq \mathcal{I}$$

...so just make λ as small as possible?

DrMMD flow in practice

DrMMD gradient descent with step size γ :

$$\nu_{n+1} = (\text{Id} + \gamma(1 + \lambda) \nabla f_{\nu_n, \nu^*})_{\#} \nu_n.$$

DrMMD flow in practice

DrMMD gradient descent with step size γ :

$$\nu_{n+1} = (\text{Id} + \gamma(1 + \lambda)\nabla f_{\nu_n, \nu^*})_{\#}\nu_n.$$

Additional **discretization error**:

$$\text{total error} = \underbrace{\gamma \lambda^r C_P Q (\mathcal{J} + \mathcal{I})}_{\text{approximation}} + \underbrace{C_2 \gamma^2 \lambda^{-1} \chi^2(\nu_n, \nu^*)}_{\text{discretization}}$$

DrMMD flow in practice

DrMMD gradient descent with step size γ :

$$\nu_{n+1} = (\text{Id} + \gamma(1 + \lambda)\nabla f_{\nu_n, \nu^*})_{\#}\nu_n.$$

Additional *discretization error*:

$$\text{total error} = \underbrace{\gamma \lambda^r C_P Q(\mathcal{J} + \mathcal{I})}_{\text{approximation}} + \underbrace{C_2 \gamma^2 \lambda^{-1} \chi^2(\nu_n, \nu^*)}_{\text{discretization}}$$

Adaptive λ_n :

$$\lambda_n \propto \chi^2(\nu_n, \nu^*)^{\frac{1}{r+1}}$$

Larger λ_n at the start, smaller λ_n near convergence.

DrMMD flow in practice

DrMMD gradient descent with step size γ :

$$\nu_{n+1} = (\text{Id} + \gamma(1 + \lambda)\nabla f_{\nu_n, \nu^*})_{\#}\nu_n.$$

Additional **discretization error**:

$$\text{total error} = \underbrace{\gamma \lambda^r C_P Q(\mathcal{J} + \mathcal{I})}_{\text{approximation}} + \underbrace{C_2 \gamma^2 \lambda^{-1} \chi^2(\nu_n, \nu^*)}_{\text{discretization}}$$

Adaptive λ_n :

$$\lambda_n \propto \chi^2(\nu_n, \nu^*)^{\frac{1}{r+1}}$$

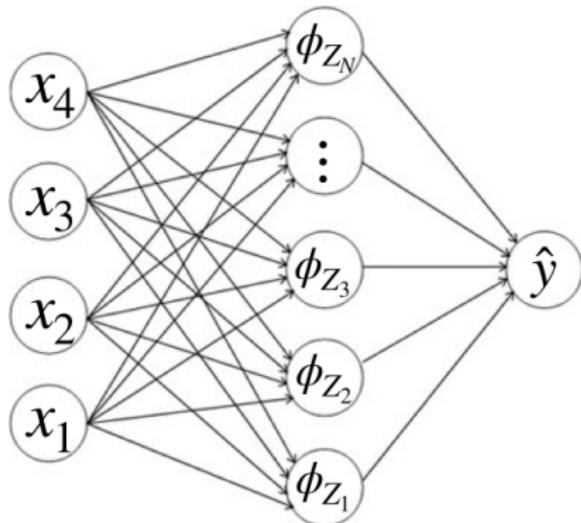
Larger λ_n at the start, smaller λ_n near convergence.

Iteration complexity:

- $\text{KL}(\nu_n, \nu^*) \leq \delta$ in $\mathcal{O}\left(\left(\frac{1}{\delta}\right)^{\frac{r+1}{r}} \log \frac{1}{\delta}\right)$ iterations.
- Langevin Monte Carlo (known density $\propto e^{-V(x)}$): $\mathcal{O}\left(\frac{1}{\delta} \log \frac{1}{\delta}\right)$

Neural net student-teacher setting

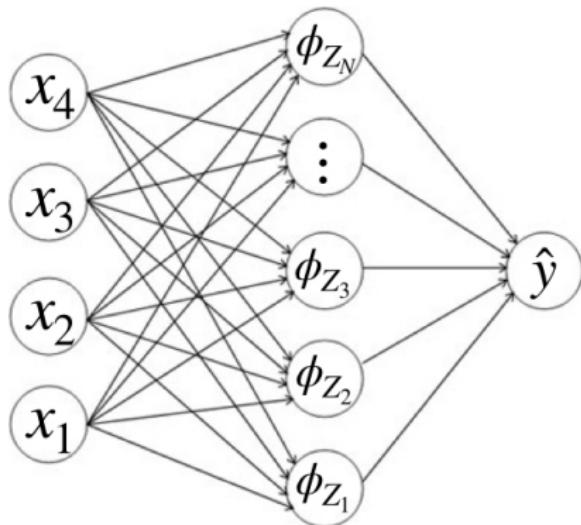
$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\left\| \frac{1}{M} \sum_m^M \phi_{U^m}(x) - \frac{1}{N} \sum_{n=1}^N \phi_{Z^n}(x) \right\|^2 \right]$$

Neural net student-teacher setting

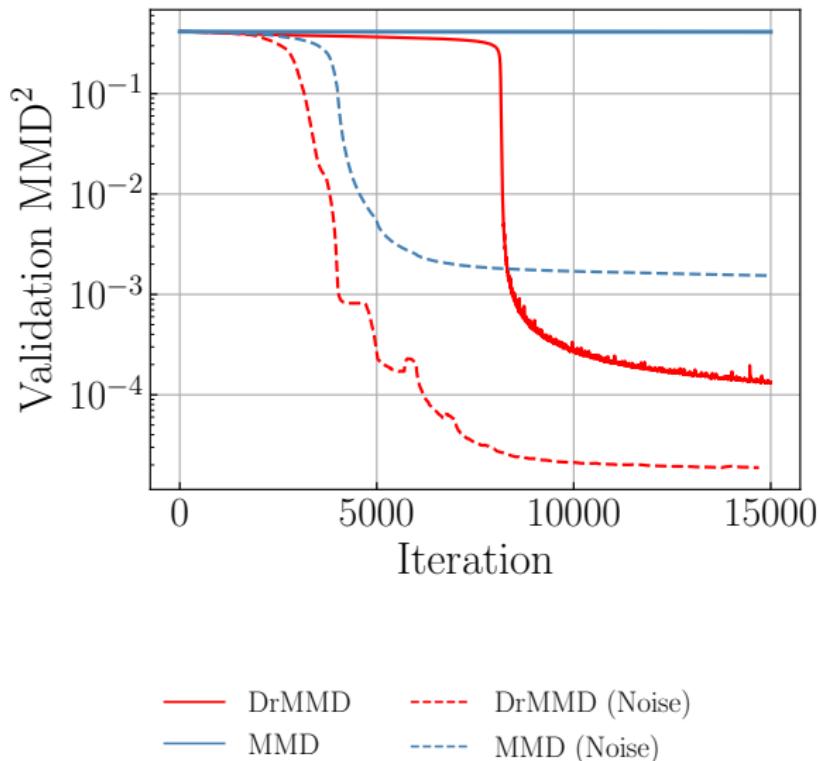
$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} MMD^2(\nu^*, \frac{1}{N} \sum_{n=1}^N \delta_{Z^n})$$

$$k(Z, Z') = \mathbb{E}_{data}[\phi_Z(x)\phi_{Z'}(x)]$$

Neural net student-teacher setting



Adaptive MMD Flow (ICLR 25)

arXiv > cs > arXiv:2405.06780

Computer Science > Machine Learning

[Submitted on 10 May 2024]

Deep MMD Gradient Flow without adversarial training

Alexandre Galashov, Valentin de Bortoli, Arthur Gretton



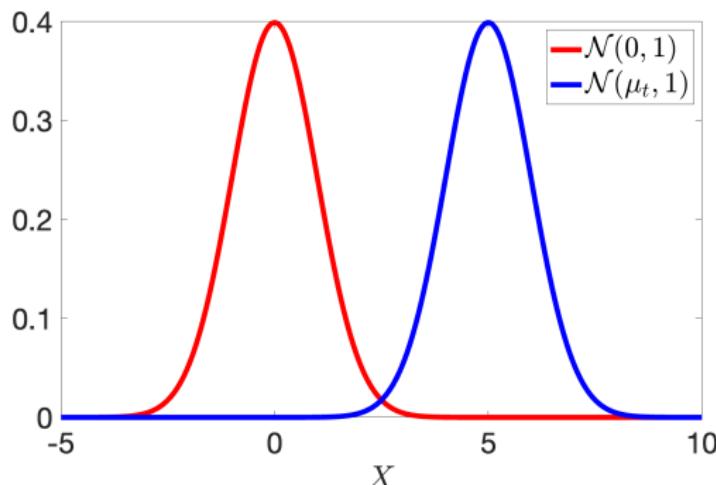
Will an adaptive kernel help?

Define the two measures:

$$\nu^* := \mathcal{N}(0, \sigma^2 \text{Id}) \quad \nu_t := \mathcal{N}(\mu_t, \sigma^2 \text{Id}).$$

Consider the family of MMDs:

$$\text{MMD}_\alpha^2(\nu^*, \nu_t) \quad \text{with} \quad k_\alpha(x, y) = \alpha^{-d} \exp[-\|x - y\|^2/(2\alpha^2)]$$



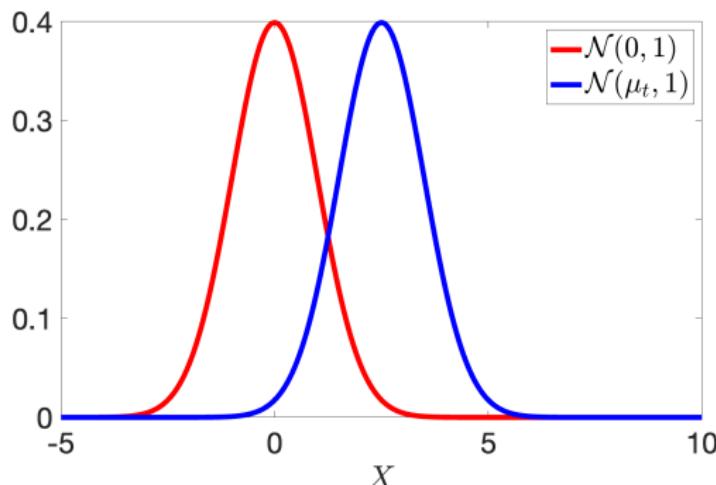
Will an adaptive kernel help?

Define the two measures:

$$\nu^* := \mathcal{N}(0, \sigma^2 \text{Id}) \quad \nu_t := \mathcal{N}(\mu_t, \sigma^2 \text{Id}).$$

Consider the family of MMDs:

$$\text{MMD}_\alpha^2(\nu^*, \nu_t) \quad \text{with} \quad k_\alpha(x, y) = \alpha^{-d} \exp[-\|x - y\|^2/(2\alpha^2)]$$



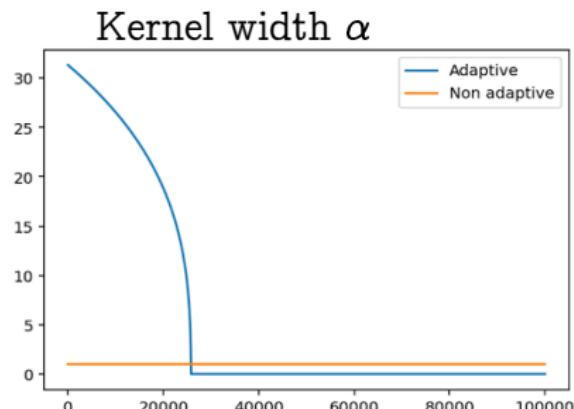
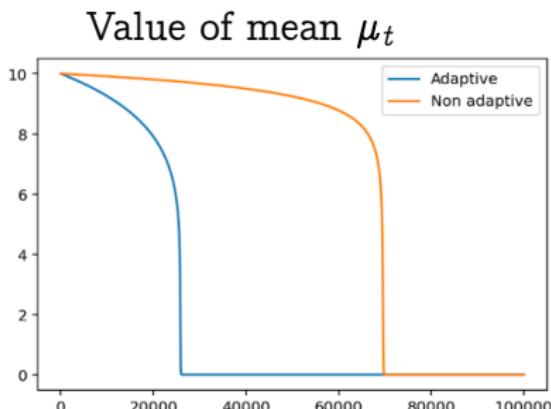
Will an adaptive kernel help?

Choose kernel such that:

$$\alpha^* = \operatorname{argmax}_{\alpha \geq 0} \|\nabla_{\mu_t} \text{MMD}_\alpha^2(\nu^*, \nu_t)\|.$$

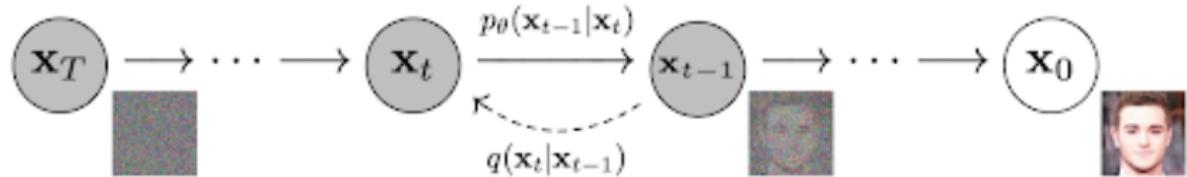
Then

$$\alpha^* = \text{ReLU}(\|\mu_t\|^2/(d+2) - 2\sigma^2)^{1/2}.$$



How to train an adaptive MMD (1)

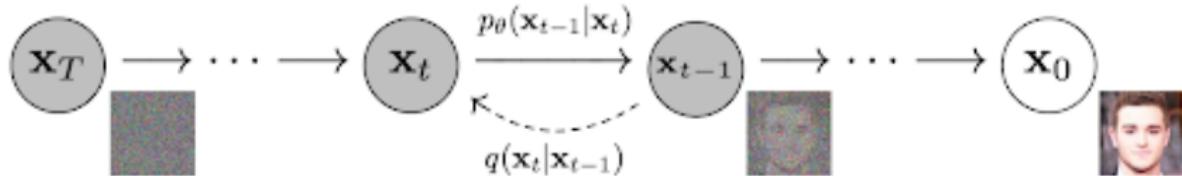
Diffusion:



Generate forward path $\tilde{\nu}_t$, $t \in [0, 1]$, such that $\tilde{\nu}_0 = \nu^*$, and $\tilde{\nu}_1 = N(0, \text{Id})$ is a Gaussian noise.

How to train an adaptive MMD (1)

Diffusion:



Generate forward path $\tilde{\nu}_t$, $t \in [0, 1]$, such that $\tilde{\nu}_0 = \nu^*$, and $\tilde{\nu}_1 = N(0, \text{Id})$ is a Gaussian noise.

Given samples $\tilde{x}_0 \sim \tilde{\nu}_0$, the samples $\tilde{x}_t | \tilde{x}_0$ are given by

$$\tilde{x}_t = \alpha_t \tilde{x}_0 + \beta_t \epsilon, \quad \epsilon \sim N(0, \text{Id}),$$

with $\alpha_0 = \beta_1 = 1$ and $\alpha_1 = \beta_0 = 0$.

- low t : \tilde{x}_t close to the original data \tilde{x}_0 ,
- high t : \tilde{x}_t close to a unit Gaussian

Schedule (α_t, β_t) is the variance-preserving one of Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole. Score-based generative modeling through stochastic differential equations (ICLR 2021)

How to train an adaptive MMD (2)

Time-dependent MMD **training loss**:

$$\mathcal{F}(\theta, t) := \frac{1}{2} \mathbb{E}_{\tilde{\nu}_t} k_{\theta, t}(\tilde{x}_t, \tilde{x}_t^l) + \mathbb{E}_{\tilde{\nu}_t, \nu^*} k_{\theta, t}(\tilde{x}_t, y)$$

with kernel

$$k_{\theta, t}(x, y) = \phi(x; t, \theta)^\top \phi(y; t, \theta)$$

and witness $f_{\nu^*, \tilde{\nu}_t}^{(\theta, t)}$.

How to train an adaptive MMD (2)

Time-dependent MMD **training loss**:

$$\mathcal{F}(\theta, t) := \frac{1}{2} \mathbb{E}_{\tilde{\nu}_t} k_{\theta, t}(\tilde{x}_t, \tilde{x}_t^l) + \mathbb{E}_{\tilde{\nu}_t, \nu^*} k_{\theta, t}(\tilde{x}_t, y)$$

with kernel

$$k_{\theta, t}(x, y) = \phi(x; t, \theta)^\top \phi(y; t, \theta)$$

and witness $f_{\nu^*, \tilde{\nu}_t}^{(\theta, t)}$.

Train θ by minimizing noise-conditional loss on **forward path**:

$$\mathcal{F}_{\text{tot}}(\theta, t) = \mathcal{F}(\theta, t) + \lambda_{\ell_2} \mathcal{F}_{\ell_2}(\theta, t) + \lambda_{\nabla} \mathcal{F}_{\nabla}(\theta, t),$$

$$\mathcal{F}_{\text{tot}}(\theta) = \mathbb{E}_{t \sim U[0,1]} [\mathcal{F}_{\text{tot}}(\theta, t)]$$

where

- $\mathcal{F}_{\ell_2}(\theta, t)$ is a “variance”-style penalty
- $\mathcal{F}_{\nabla}(\theta, t) = \frac{1}{N} \sum_{i=1}^N (\|\nabla f_{\nu^*, \tilde{\nu}_t}^{(\theta, t)}(\tilde{x}_{t,i})\|_2 - 1)^2$, is a gradient penalty

Gulrajani, Ahmed, Arjovsky, Dumoulin, Courville, Improved Training of Wasserstein GANs (NeurIPS 2017)

Binkowski, Sutherland, Arbel, G. (NeurIPS 2018)

Sample generation

Algorithm Noise-adaptive MMD gradient flow

Sample initial particles $Z \sim N(0, \text{Id})$

Set $\Delta t = (t_{\max} - t_{\min}) / T$

for $i = T$ to 0 do

 Set the noise level $t = i\Delta t$

 Set $Z_t^0 = Z$

 for $n = 0$ to $N_s - 1$ do

$Z_t^{n+1} = Z_t^n - \eta \nabla \textcolor{teal}{f}_{\nu^\star, \nu_t}^{(\theta^\star, t)}(Z_t^n)$

 end for

 Set $Z = Z_t^N$

end for

Output Z

Results

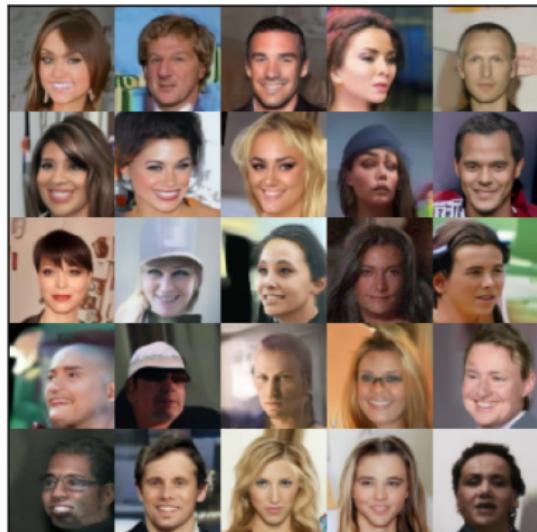
Table: Unconditional generation, CIFAR-10. MMD GAN (orig.), used mixed-RQ kernel. "Orig." – original paper, "impl." – our implementation.

Method	FID	IS	NFE
MMD GAN (orig.)	39.90	6.51	-
MMD GAN (impl.)	13.62	8.93	-
DDPM (orig.)	3.17	9.46	1000
DDPM (impl.)	5.19	8.90	100
Discriminator flows			
DGGF-KL	28.80	-	110
JKO-Flow	23.10	7.48	~ 150
GS-MMD-RK	55.00	-	86
DMMD (ours)	8.31	9.09	100
DMMD (ours)	7.74	9.12	250

DDPM from (Ho et al., 2020). Discriminator flows include two KL gradient flows trained adversarially: JKO-Flow (Fan et al., 2022) and Deep Generative Wasserstein Gradient Flows (DGGF-KL) (Heng et al., 2023). GS-MMD-RK is Generative Sliced MMD Flows with Riesz Kernels (Hertrich et al., 2024)

Images

CELEB-A (64x64)



LSUN Church (64x64)



Summary

- Gradient flows based on kernel dependence measures
- NeurIPS 2019, NeurIPS 2021, ICLR 2025, JMLR (submitted)

NeurIPS 2019:

 > stat > arXiv:1906.04370

Statistics > Machine Learning

[Submitted on 11 Jun 2019 (v1), last revised 3 Dec 2019 (this version, v2)]

Maximum Mean Discrepancy Gradient Flow

Michael Arbel, Anna Korba, Adil Salim, Arthur Gretton

NeurIPS 2021:

 > stat > arXiv:2106.08929

Statistics > Machine Learning

[Submitted on 16 Jun 2021 (v1), last revised 29 Oct 2021 (this version, v2)]

KALE Flow: A Relaxed KL Gradient Flow for Probabilities with Disjoint Support

Pierre Glaser, Michael Arbel, Arthur Gretton

Adaptive MMD (ICLR 25):

 > cs > arXiv:2405.06780

Computer Science > Machine Learning

[Submitted on 10 May 2024]

Deep MMD Gradient Flow without adversarial training

Alexandre Galashov, Valentin de Bortoli, Arthur Gretton

(De)regularized MMD
(JMLR, submitted):

 > stat > arXiv:2409.14980

Search...

Help | About

Statistics > Machine Learning

[Submitted on 23 Sep 2024]

(De)-regularized Maximum Mean Discrepancy Gradient Flow

Zonghao Chen, Aratrika Mustafi, Pierre Glaser, Anna Korba, Arthur Gretton, Bharath K. Sriperumbudur

Research support

Work supported by:

The Gatsby Charitable Foundation



Google Deepmind

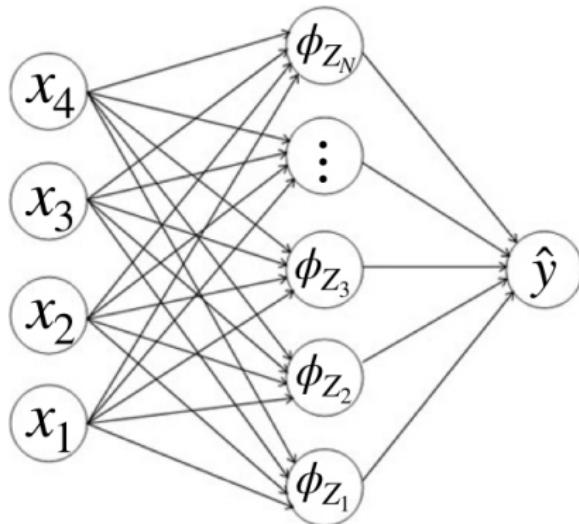


Questions?



Noise injection: neural net setting

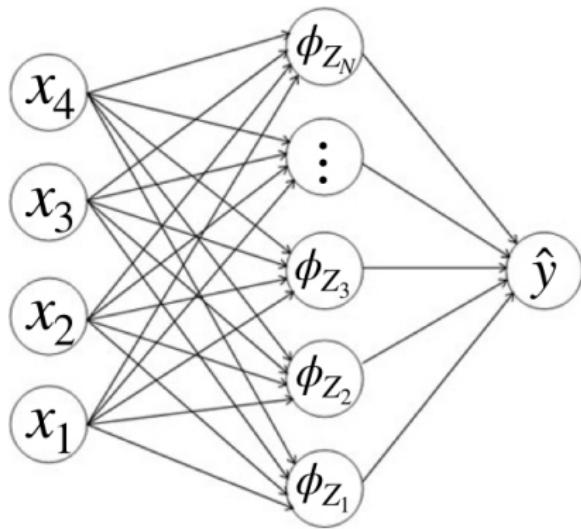
$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\left\| \frac{1}{M} \sum_m^M \phi_{U^m}(x) - \frac{1}{N} \sum_{n=1}^N \phi_{Z^n}(x) \right\|^2 \right]$$

Noise injection: neural net setting

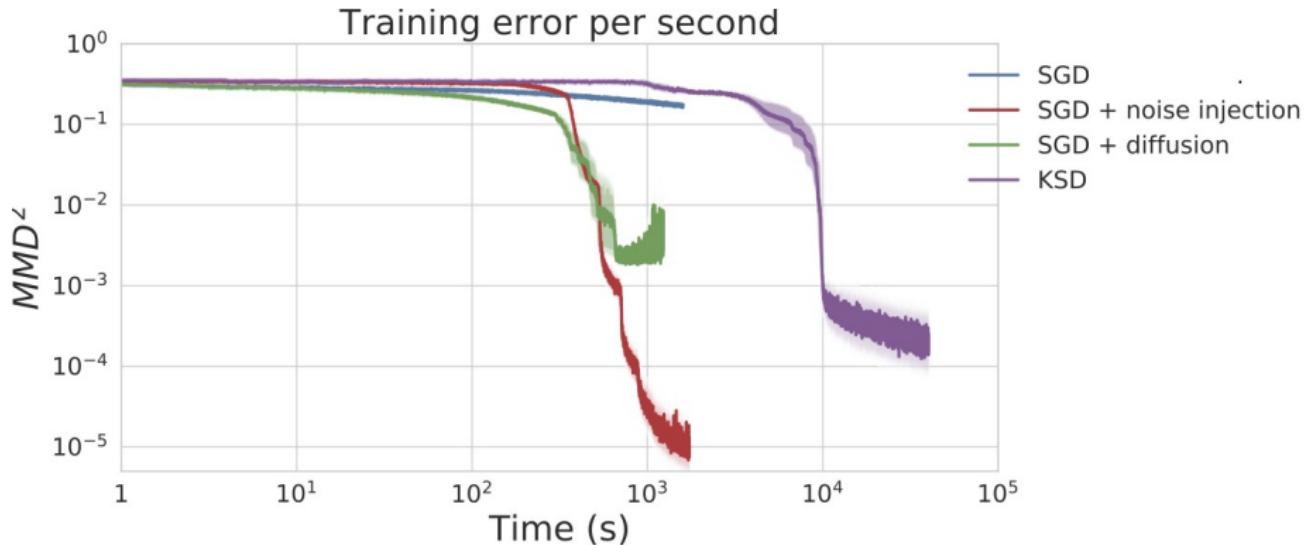
$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} MMD^2(\nu^*, \frac{1}{N} \sum_{n=1}^N \delta_{Z^n})$$

$$k(Z, Z') = \mathbb{E}_{data}[\phi_Z(x)\phi_{Z'}(x)]$$

Noise injection: neural net setting



KSD is Kernel Sobolev Discrepancy. Y. Mroueh, T. Sercu, and A. Raj. "Sobolev Descent." In: AISTATS. 2019.

The KALE, and KALE flow



Wasserstein gradient flow: KL lower bound

$$D_{KL}(\textcolor{blue}{P}, \textcolor{red}{Q}) = \int \log \left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)} \right) \textcolor{blue}{p}(z) dz$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Wasserstein gradient flow: KL lower bound

$$D_{KL}(\textcolor{blue}{P}, \textcolor{red}{Q}) = \int \log \left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)} \right) \textcolor{blue}{p}(z) dz$$
$$\geq \sup_{f \in \mathcal{H}} \mathbb{E}_{\textcolor{blue}{P}} f(\textcolor{blue}{X}) + 1 - \mathbb{E}_{\textcolor{red}{Q}} \underbrace{\exp(f(\textcolor{red}{Y}))}_{\phi^*(f(\textcolor{red}{Y})+1)}$$

This is a KL Approximate Lower-bound Estimator.

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

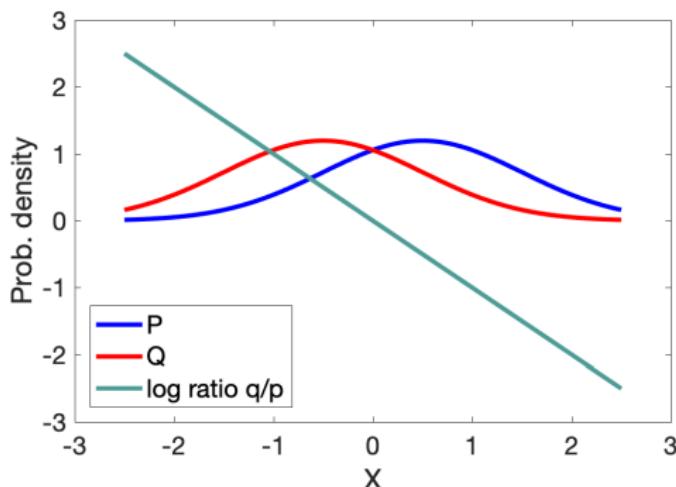
Wasserstein gradient flow: KL lower bound

$$D_{KL}(P, Q) = \int \log\left(\frac{p(z)}{q(z)}\right) p(z) dz$$
$$\geq \sup_{f \in \mathcal{H}} E_{Pf}(X) + 1 - E_Q \exp(f(Y))$$

Bound tight when:

$$f^\diamond(z) = \log \frac{q(z)}{p(z)}$$

if ratio defined.



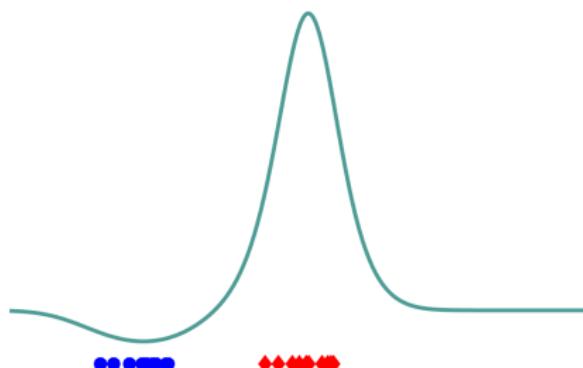
Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Wasserstein gradient flow: KL lower bound

$$D_{KL}(P, Q) = \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz$$

$$\geq \sup_{f \in \mathcal{H}} \mathbb{E}_{Pf}(X) + 1 - \mathbb{E}_Q \exp(f(Y))$$

$$KALE(Q, P; \mathcal{H}) = 0.18$$



Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

A different flow: KALE flow

arXiv > stat > arXiv:2106.08929

Statistics > Machine Learning

[Submitted on 16 Jun 2021 (v1), last revised 29 Oct 2021 (this version, v2)]

KALE Flow: A Relaxed KL Gradient Flow for Probabilities with Disjoint Support

Pierre Glaser, Michael Arbel, Arthur Gretton

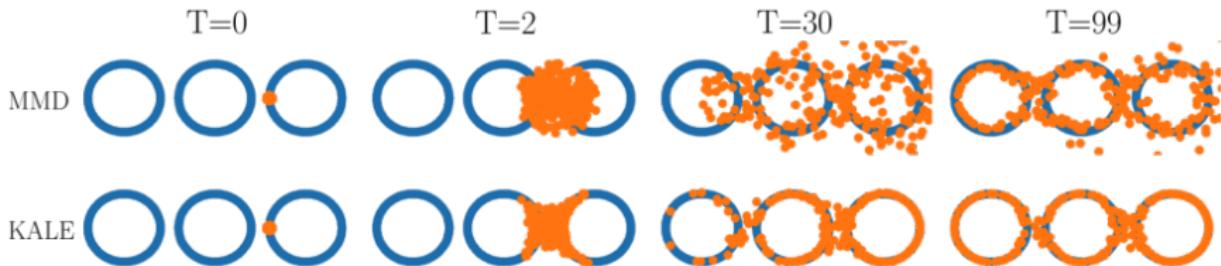


Figure 1: MMD and KALE flow trajectories for “three rings” target

Glaser, Arbel, G. (NeurIPS 2021)

The KALE, and KALE flow



The ϕ -divergences

Define the ϕ -divergence(f -divergence):

$$D_\phi(\textcolor{blue}{P}, \textcolor{red}{Q}) = \int \phi\left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \textcolor{red}{q}(z) dz$$

where ϕ is convex, lower-semicontinuous, $\phi(1) = 0$.

■ Example: $\phi(u) = u \log(u)$ gives KL divergence,

$$\begin{aligned} D_{KL}(\textcolor{blue}{P}, \textcolor{red}{Q}) &= \int \log\left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \textcolor{blue}{p}(z) dz \\ &= \int \left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \log\left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \textcolor{red}{q}(z) dz \end{aligned}$$

The ϕ -divergences

Define the ϕ -divergence(f -divergence):

$$D_\phi(\textcolor{blue}{P}, \textcolor{red}{Q}) = \int \phi\left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \textcolor{red}{q}(z) dz$$

where ϕ is convex, lower-semicontinuous, $\phi(1) = 0$.

■ Example: $\phi(u) = u \log(u)$ gives KL divergence,

$$\begin{aligned} D_{KL}(\textcolor{blue}{P}, \textcolor{red}{Q}) &= \int \log\left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \textcolor{blue}{p}(z) dz \\ &= \int \left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \log\left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \textcolor{red}{q}(z) dz \end{aligned}$$

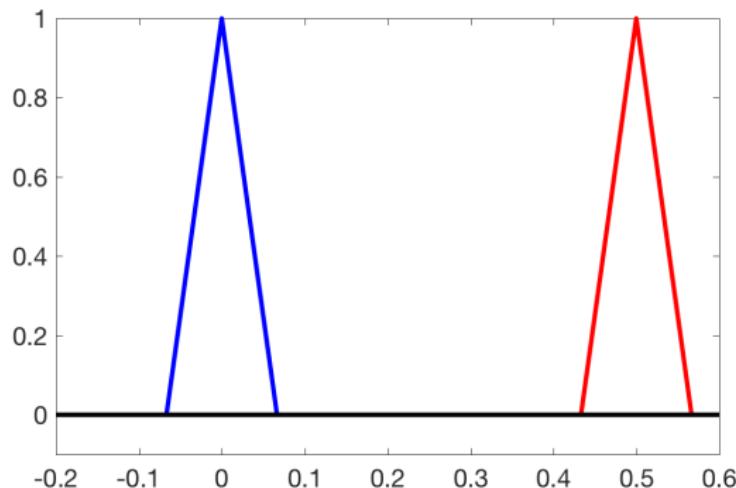
The challenge of disjoint support



Simple example: disjoint support.

Goodfellow et al. (NeurIPS 2014), Arjovsky and Bottou [ICLR 2017]

$$D_{KL}(P, Q) = \infty \quad D_{JS}(P, Q) = \log 2$$



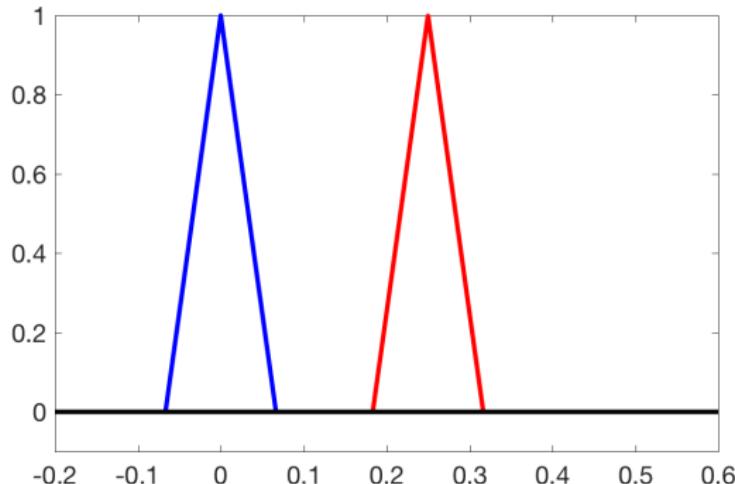
The challenge of disjoint support



Simple example: disjoint support.

Goodfellow et al. (NeurIPS 2014), Arjovsky and Bottou [ICLR 2017]

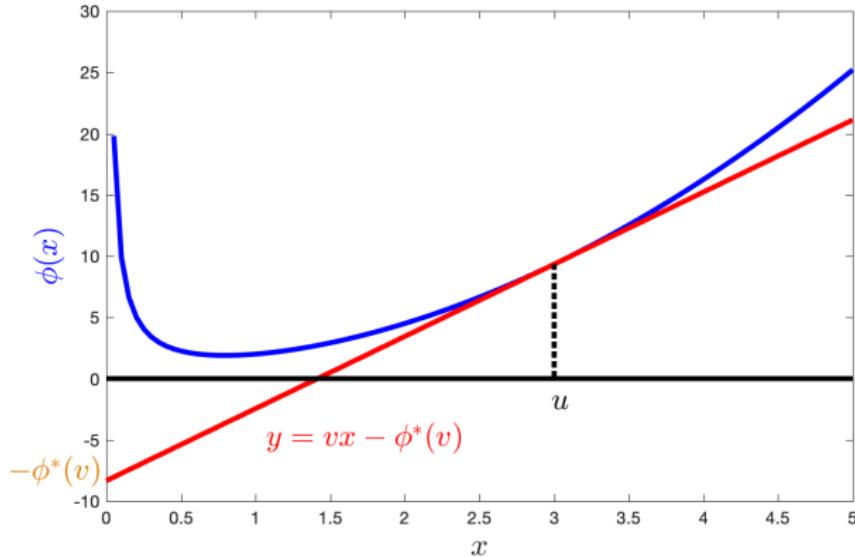
$$D_{KL}(P, Q) = \infty \quad D_{JS}(P, Q) = \log 2$$



ϕ -divergences in practice

Notation: the conjugate (Fenchel) dual

$$\phi^*(v) = \sup_{u \in \mathbb{R}} \{uv - \phi(u)\}.$$



- $\phi^*(v)$ is negative intercept of tangent to ϕ with slope v

ϕ -divergences in practice

Notation: the conjugate (Fenchel) dual

$$\phi^*(v) = \sup_{u \in \mathbb{R}} \{uv - \phi(u)\}.$$

- For a convex l.s.c. ϕ we have

$$\phi^{**}(x) = \phi(x) = \sup_{v \in \mathbb{R}} \{xv - \phi^*(v)\}$$

ϕ -divergences in practice

Notation: the conjugate (Fenchel) dual

$$\phi^*(v) = \sup_{u \in \mathbb{R}} \{uv - \phi(u)\}.$$

- For a convex l.s.c. ϕ we have

$$\phi^{**}(x) = \phi(x) = \sup_{v \in \mathbb{R}} \{xv - \phi^*(v)\}$$

- KL divergence:

$$\phi(x) = x \log(x) \quad \phi^*(v) = \exp(v - 1)$$

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$D_\phi(P, Q) = \int q(z)\phi\left(\frac{p(z)}{q(z)}\right) dz$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$\begin{aligned} D_\phi(P, Q) &= \int q(z)\phi\left(\frac{p(z)}{q(z)}\right) dz \\ &= \int q(z) \underbrace{\sup_{f_z} \left(\frac{p(z)}{q(z)} f_z - \phi^*(f_z) \right)}_{\phi\left(\frac{p(z)}{q(z)}\right)} \end{aligned}$$

$\phi^*(v)$ is dual of $\phi(x)$.

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$\begin{aligned} D_\phi(P, Q) &= \int q(z)\phi\left(\frac{p(z)}{q(z)}\right) dz \\ &= \int q(z) \sup_{f_z} \left(\frac{p(z)}{q(z)} f_z - \phi^*(f_z) \right) \\ &\geq \sup_{f \in \mathcal{H}} E_P f(X) - E_Q \phi^*(f(Y)) \end{aligned}$$

(restrict the function class)

A variational lower bound

A lower-bound ϕ -divergence approximation:

$$\begin{aligned} D_\phi(P, Q) &= \int q(z)\phi\left(\frac{p(z)}{q(z)}\right) dz \\ &= \int q(z) \sup_{f_z} \left(\frac{p(z)}{q(z)} f_z - \phi^*(f_z) \right) \\ &\geq \sup_{f \in \mathcal{H}} E_P f(X) - E_Q \phi^*(f(Y)) \end{aligned}$$

(restrict the function class)

Bound tight when:

$$f^\diamond(z) = \partial \phi \left(\frac{p(z)}{q(z)} \right)$$

if ratio defined.

Case of the KL

$$D_{KL}(P, Q) = \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Case of the KL

$$D_{KL}(P, Q) = \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz$$

$$\geq \sup_{f \in \mathcal{H}} \mathbb{E}_{Pf}(X) + 1 - \mathbb{E}_Q \underbrace{\exp(f(Y))}_{\phi^*(f(Y)+1)}$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Case of the KL

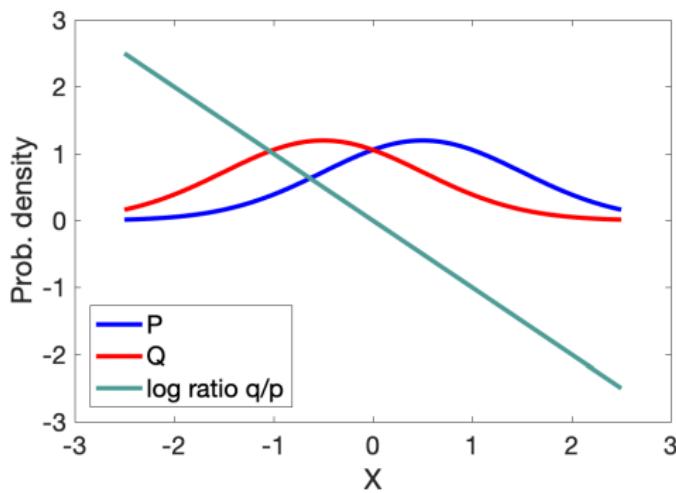
$$D_{KL}(P, Q) = \int \log \left(\frac{p(z)}{q(z)} \right) p(z) dz$$

$$\geq \sup_{f \in \mathcal{H}} E_{Pf}(X) + 1 - E_Q \exp(f(Y))$$

Bound tight when:

$$f^\diamond(z) = \log \frac{q(z)}{p(z)}$$

if ratio defined.



Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log\left(\frac{\textcolor{blue}{p}(z)}{\textcolor{red}{q}(z)}\right) \textcolor{blue}{p}(z) dz \\ &\geq \sup_{f \in \mathcal{H}} \mathbb{E}_{Pf}(X) + 1 - \mathbb{E}_Q \exp(f(Y)) \quad \begin{array}{l} x_i \stackrel{\text{i.i.d.}}{\sim} P \\ y_i \stackrel{\text{i.i.d.}}{\sim} Q \end{array} \\ &\approx \sup_{f \in \mathcal{H}} \left[\frac{1}{n} \sum_{j=1}^n f(\textcolor{blue}{x}_i) - \frac{1}{n} \sum_{i=1}^n \exp(\textcolor{teal}{f}(\textcolor{red}{y}_i)) \right] + 1 \end{aligned}$$

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log\left(\frac{p(z)}{q(z)}\right) p(z) dz \\ &\geq \sup_{f \in \mathcal{H}} E_P f(X) + 1 - E_Q \exp(f(Y)) \\ &\approx \sup_{f \in \mathcal{H}} \left[\frac{1}{n} \sum_{j=1}^n f(x_i) - \frac{1}{n} \sum_{i=1}^n \exp(f(y_i)) \right] + 1 \end{aligned}$$

This is a

KL

Approximate

Lower-bound

Estimator.

Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log\left(\frac{p(z)}{q(z)}\right) p(z) dz \\ &\geq \sup_{f \in \mathcal{H}} E_P f(X) + 1 - E_Q \exp(f(Y)) \\ &\approx \sup_{f \in \mathcal{H}} \left[\frac{1}{n} \sum_{j=1}^n f(x_i) - \frac{1}{n} \sum_{i=1}^n \exp(f(y_i)) \right] + 1 \end{aligned}$$

This is a

K

A

L

E

Case of the KL

$$\begin{aligned} D_{KL}(P, Q) &= \int \log\left(\frac{p(z)}{q(z)}\right) p(z) dz \\ &\geq \sup_{f \in \mathcal{H}} E_P f(X) + 1 - E_Q \exp(f(Y)) \\ &\approx \sup_{f \in \mathcal{H}} \left[\frac{1}{n} \sum_{j=1}^n f(x_i) - \frac{1}{n} \sum_{i=1}^n \exp(f(y_i)) \right] + 1 \end{aligned}$$

The KALE divergence

Nguyen, Wainwright, Jordan, IEEE Transactions on Information Theory (2010);
Nowozin, Cseke, Tomioka, NeurIPS (2016)

Empirical properties of KALE

$$KALE(P, Q; \mathcal{H}) = \sup_{f \in \mathcal{H}} E_P f(X) - E_Q \exp(f(Y)) + 1$$



$$f = \langle w, \phi(x) \rangle_{\mathcal{H}} \quad \mathcal{H} \text{ an RKHS}$$
$$\|w\|_{\mathcal{H}}^2 \quad \text{penalized}$$

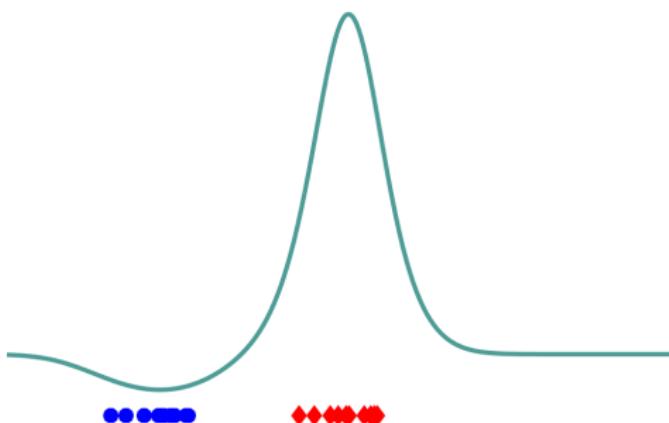
Empirical properties of KALE

$$KALE(P, Q; \mathcal{H}) = \sup_{f \in \mathcal{H}} E_P f(X) - E_Q \exp(f(Y)) + 1$$



$$f = \langle w, \phi(x) \rangle_{\mathcal{H}} \quad \mathcal{H} \text{ an RKHS}$$
$$\|w\|_{\mathcal{H}}^2 \text{ penalized}$$

$$KALE(Q, P; \mathcal{H}) = 0.18$$



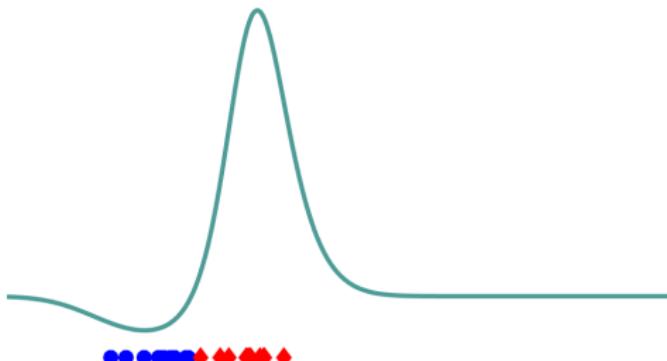
Empirical properties of KALE

$$KALE(P, Q; \mathcal{H}) = \sup_{f \in \mathcal{H}} E_P f(X) - E_Q \exp(f(Y)) + 1$$



$$f = \langle w, \phi(x) \rangle_{\mathcal{H}} \quad \mathcal{H} \text{ an RKHS}$$
$$\|w\|_{\mathcal{H}}^2 \quad \text{penalized}$$

$$KALE(Q, P; \mathcal{H}) = 0.12$$



Glaser, Arbel, G. "KALE Flow: A Relaxed KL Gradient Flow for Probabilities with Disjoint Support," (NeurIPS 2021, Section 2)

Topological properties of KALE (1)

Key requirements on \mathcal{H} and \mathcal{X} :

- Compact domain \mathcal{X} ,
- \mathcal{H} dense in the space $C(\mathcal{X})$ of continuous functions on \mathcal{X} wrt $\|\cdot\|_\infty$.
- If $f \in \mathcal{H}$ then $-f \in \mathcal{H}$ and $cf \in \mathcal{H}$ for $0 \leq c \leq C_{\max}$.

Theorem: $KALE(P, Q; \mathcal{H}) \geq 0$ and $KALE(P, Q; \mathcal{H}) = 0$ iff $P = Q$.

Zhang, Liu, Zhou, Xu, and He. "On the Discrimination-Generalization Tradeoff in GANs"
(ICLR 2018, Corollary 2.4; Theorem B.1)
Arbel, Liang, G. (ICLR 2021, Proposition 1)

Topological properties of KALE (1)

Key requirements on \mathcal{H} and \mathcal{X} :

- Compact domain \mathcal{X} ,
- \mathcal{H} dense in the space $C(\mathcal{X})$ of continuous functions on \mathcal{X} wrt $\|\cdot\|_\infty$.
- If $f \in \mathcal{H}$ then $-f \in \mathcal{H}$ and $cf \in \mathcal{H}$ for $0 \leq c \leq C_{\max}$.

Theorem: $KALE(P, Q; \mathcal{H}) \geq 0$ and $KALE(P, Q; \mathcal{H}) = 0$ iff $P = Q$.

\mathcal{H} dense in $C(\mathcal{X})$ for $\mathcal{X} \subset \mathbb{R}^d$ when:

$$\mathcal{H} = \text{span}\{\sigma(w^\top x + b) : [w, b] \in \Theta\}$$

$$\sigma(u) = \max\{u, 0\}^\alpha, \alpha \in \mathbb{N}, \text{ and } \{\lambda\theta : \lambda \geq 0, \theta \in \Theta\} = \mathbb{R}^{d+1}.$$

Zhang, Liu, Zhou, Xu, and He. "On the Discrimination-Generalization Tradeoff in GANs"
(ICLR 2018, Corollary 2.4; Theorem B.1)
Arbel, Liang, G. (ICLR 2021, Proposition 1)

Topological properties of KALE (2)

Additional requirement: all functions in \mathcal{H} Lipschitz in their inputs with constant L

Theorem: $\text{KALE}(\mathcal{P}, \mathcal{Q}^n; \mathcal{H}) \rightarrow 0$ iff $\mathcal{Q}^n \rightarrow \mathcal{P}$ under the weak topology.

Liu, Bousquet, Chaudhuri. “Approximation and Convergence Properties of Generative Adversarial Learning” (NeurIPS 2017); Arbel, Liang, G. (ICLR 2021, Proposition 1)

Topological properties of KALE (2)

Additional requirement: all functions in \mathcal{H} Lipschitz in their inputs with constant L

Theorem: $\text{KALE}(\mathbf{P}, \mathbf{Q}^n; \mathcal{H}) \rightarrow 0$ iff $\mathbf{Q}^n \rightarrow \mathbf{P}$ under the weak topology.

Partial proof idea:

$$\begin{aligned}\text{KALE}(\mathbf{P}, \mathbf{Q}; \mathcal{H}) &= \int \mathbf{f} d\mathbf{P} - \int \exp(\mathbf{f}) d\mathbf{Q} + 1 \\ &= - \int \mathbf{f}(x) d\mathbf{Q}(x) + \mathbf{f}(x') d\mathbf{P}(x') \\ &\quad - \int \underbrace{(\exp(\mathbf{f}) - \mathbf{f} - 1)}_{\geq 0} d\mathbf{Q} \\ &\leq \int \mathbf{f}(x') d\mathbf{P}(x') - \int \mathbf{f}(x) d\mathbf{Q}(x) \leq LW_1(\mathbf{P}, \mathbf{Q})\end{aligned}$$

Liu, Bousquet, Chaudhuri. "Approximation and Convergence Properties of Generative Adversarial Learning" (NeurIPS 2017); Arbel, Liang, G. (ICLR 2021, Proposition 1)

KALE vs KL vs MMD

A scaled KALE (non-degenerate for $\lambda = 0$ or $\lambda \rightarrow \infty$):

$$\begin{aligned} \text{KALE}_\lambda(\mathcal{P}, \mathcal{Q}; \mathcal{H}) &= (1 + \lambda) \sup_{f \in \mathcal{H}} \left[E_{\mathcal{P}} f(\mathcal{X}) - E_{\mathcal{Q}} \exp(f(\mathcal{Y})) \right. \\ &\quad \left. + 1 - \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 \right] \end{aligned}$$

MMD limit:

$$\lim_{\lambda \rightarrow +\infty} \text{KALE}_\lambda(\mathcal{P}, \mathcal{Q}; \mathcal{H}) = \frac{1}{2} \text{MMD}^2(\mathcal{P}, \mathcal{Q}).$$

KL limit (assuming $\log \frac{d\mathcal{P}}{d\mathcal{Q}} \in \mathcal{H}$):

$$\lim_{\lambda \rightarrow 0} \text{KALE}_\lambda(\mathcal{P}, \mathcal{Q}; \mathcal{H}) = \text{KL}(\mathcal{P}, \mathcal{Q}).$$

Glaser, Arbel, G. (NeurIPS 2021, Proposition 1)

Wasserstein gradient flow on KALE

First variation of the $KALE_\lambda(\nu, \nu^*)$

$$\frac{\partial KALE_\lambda}{\partial \nu}(\nu)(z) := (1 + \lambda) f_{\nu, \nu^*}(z)$$

where f_{ν, ν^*} is the solution of

$$f_{\nu, \nu^*} = \arg \max_{f \in \mathcal{H}} \{ \mathcal{K}(f, \nu) \},$$

where

$$\mathcal{K}(f, \nu) := E_\nu f(X) - E_{\nu^*} \exp(f(Y)) + 1 - \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

Wasserstein gradient flow on KALE

First variation of the $KALE_\lambda(\nu, \nu^*)$

$$\frac{\partial KALE_\lambda}{\partial \nu}(\nu)(z) := (1 + \lambda) f_{\nu, \nu^*}(z)$$

where f_{ν, ν^*} is the solution of

$$f_{\nu, \nu^*} = \arg \max_{f \in \mathcal{H}} \{ \mathcal{K}(f, \nu) \},$$

where

$$\mathcal{K}(f, \nu) := E_\nu f(X) - E_{\nu^*} \exp(f(Y)) + 1 - \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

Proof (idea):

$$\frac{\partial KALE_\lambda}{\partial \nu} = \frac{\partial \mathcal{K}(f_{\nu, \nu^*}, \nu)}{\partial \nu} + \underbrace{\left. \frac{\partial \mathcal{K}(f, \nu)}{\partial f} \right|_{f=f_{\nu, \nu^*}}}_{=0} \frac{\partial f_{\nu, \nu^*}}{\partial \nu}$$

as long as $\frac{\partial f_{\nu, \nu^*}}{\partial \nu}$ exists (via implicit function theorem)

Wasserstein gradient flow on KALE

The W_2 gradient flow of the KALE:

$$\partial_t \nu_t = -(1 + \lambda) \operatorname{div}(\nu_t \nabla f_{\nu_t, \nu^*}), \quad \nu_0 = P_0$$

where

$$f_{\nu, \nu^*} = \arg \max_f \mathcal{K}(f, \nu)$$

Glaser, Arbel, G. (NeurIPS 2021, Lemma 3)

Consistency (2)

Again, under the (strong!) assumption

$$\begin{aligned} S(\nu^* | \nu_t) &:= \sup_{g, \mathbb{E}_{Z \sim \nu_t} [\|\nabla g(Z)\|^2] \leq 1} |\mathbb{E}_{Z \sim \nu_t} [g(Z)] - \mathbb{E}_{U \sim \nu^*} [g(U)]| \\ &\leq C \end{aligned}$$

we have

$$\text{KALE}(\nu_t) \leq \frac{1}{\text{KALE}(\nu_0)^{-1} + C^{-1}t}$$

Once again, noise injection can be used (similar result to MMD flow).

Consistency (2)

Again, under the (strong!) assumption

$$\begin{aligned} S(\nu^* | \nu_t) &:= \sup_{g, \mathbb{E}_{Z \sim \nu_t} [\|\nabla g(Z)\|^2] \leq 1} |\mathbb{E}_{Z \sim \nu_t} [g(Z)] - \mathbb{E}_{U \sim \nu^*} [g(U)]| \\ &\leq C \end{aligned}$$

we have

$$\text{KALE}(\nu_t) \leq \frac{1}{\text{KALE}(\nu_0)^{-1} + C^{-1}t}$$

Once again, noise injection can be used (similar result to MMD flow). Compare with linear rate for Wasserstein-2 flow on KL when ν^* satisfies log-Sobolev inequality with constant ρ :

$$\frac{d}{dt} KL(\nu_t, \nu^*) \leq -2\rho KL(\nu_t, \nu^*)$$

Glaser, Arbel, G. (NeurIPS 2021, Proposition 3)

Consistency (2)

Dissipation inequalities:

- Rate by which \mathcal{F} decreases along the gradient flow [A, Proposition 2]

$$\frac{d\mathcal{F}(\nu_t)}{dt} = -\mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

- Assume the dissipation rate is controlled (path-dependent Lojasiewicz inequality)

$$\mathcal{F}(\nu_t) \leq C \mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

- From above, [A, Proposition 7]:

$$\mathcal{F}(\nu_t) \leq \frac{1}{\mathcal{F}(\nu_0)^{-1} + 2C^{-1}t}$$

[A] Arbel, Korba, Salim, G. (NeurIPS 2019)

Consistency (2)

Dissipation inequalities:

- Rate by which \mathcal{F} decreases along the gradient flow [A, Proposition 2]

$$\frac{d\mathcal{F}(\nu_t)}{dt} = -\mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

- Assume the dissipation rate is controlled (path-dependent Lojasiewicz inequality)

$$\mathcal{F}(\nu_t) \leq C \mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

- From above, [A, Proposition 7]:

$$\mathcal{F}(\nu_t) \leq \frac{1}{\mathcal{F}(\nu_0)^{-1} + 2C^{-1}t}$$

[A] Arbel, Korba, Salim, G. (NeurIPS 2019)

Consistency (2)

Dissipation inequalities:

- Rate by which \mathcal{F} decreases along the gradient flow [A, Proposition 2]

$$\frac{d\mathcal{F}(\nu_t)}{dt} = -\mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

- Assume the dissipation rate is controlled (path-dependent Lojasiewicz inequality)

$$\mathcal{F}(\nu_t) \leq C \mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

- From above, [A, Proposition 7]:

$$\mathcal{F}(\nu_t) \leq \frac{1}{\mathcal{F}(\nu_0)^{-1} + 2C^{-1}t}$$

[A] Arbel, Korba, Salim, G. (NeurIPS 2019)

Consistency (2)

Check: Lojasiewicz inequality for MMD?

- Does there exist $C > 0$ such that

$$\mathcal{F}(\nu_t) \leq C \mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

- By Cauchy-Schwarz in the RKHS, [A, eq. 16]

$$\mathcal{F}(\nu_t) =: \frac{1}{2} MMD^2(\nu_t, \nu^*) \leq S(\nu^* | \nu_t) \mathbb{E}_{\nu_t} [\|\nabla f_{\nu^*, \nu_t}\|^2]$$

where $S(\nu^* | \nu_t)$ is the Negative Sobolev Distance²

- Require $S(\nu^* | \nu_t) < C$ for entire sequence ν_t : hard to check in theory, fails in practice.

[A] Arbel, Korba, Salim, G. (NeurIPS 2019)