#### Maximum Mean Discrepancy Gradient Flow

Michael Arbel<sup>1</sup> Anna Korba<sup>1</sup> Adil Salim<sup>2</sup> Arthur Gretton<sup>1</sup>

<sup>1</sup>Gatsby Computational Neuroscience Unit, UCL, London

<sup>2</sup>Visual Computing Center, KAUST, Saudi Arabia

Fields Institute September 2020

#### **Problem and Outline**

#### Problem:

- Transport mass from a starting probability distribution to a target distribution
- How? By finding a *continuous* path on the space of distributions, decreasing some loss (Wasserstein gradient flows)
- This work: Minimize the Maximum Mean Discrepancy (MMD) on the space of probability distributions.

*Application :* Insights on the theoretical properties of some large neural networks and alteration of the dynamics to improve convergence.

#### Outline

#### Background and motivation

Maximum Mean Discrepancy (Wasserstein) Gradient Flow

Convergence properties of the MMD gradient flow

A noise-injection algorithm for better convergence

#### Setting

Let  $\mathcal P$  the set of probability measures on  $\mathcal Z\subset \mathbb R^d$  with finite second moment :

$$\mathcal{P} = \{\mu \in \mathcal{P}(\mathcal{Z}), \ \int \|z\|^2 d\mu(z) < \infty\}$$

#### Setting

Let  $\mathcal P$  the set of probability measures on  $\mathcal Z \subset \mathbb R^d$  with finite second moment :

$$\mathcal{P} = \{\mu \in \mathcal{P}(\mathcal{Z}), \ \int \|z\|^2 d\mu(z) < \infty\}$$

The space  $\mathcal{P}$  is endowed with the Wassertein-2 distance from **Optimal transport** :

$$W_2^2(\nu,\mu) = \inf_{\pi \in \Pi(\nu,\mu)} \int_{\mathcal{Z} \times \mathcal{Z}} \left\| z - z' \right\|^2 d\pi(z,z') \qquad \forall 
u,\mu \in \mathcal{P}$$

where  $\Pi(\nu, \mu)$  is the set of possible couplings between  $\nu$  and  $\mu$ .

In other words  $\Pi(\nu, \mu)$  contains all possible distributions  $\pi$  on  $\mathcal{Z} \times \mathcal{Z}$  such that if  $(Z, Z') \sim \pi$  then  $Z \sim \nu$  and  $Z' \sim \mu$ .

#### Maximum Mean Discrepancy

• Let  $k : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  a positive, semi-definite kernel

$$k(z,z') = \langle \phi(z), \phi(z') 
angle_{\mathcal{H}}, \quad \phi : \mathcal{Z} \to \mathcal{H}$$

▶  $\mathcal{H}$  its RKHS (Reproducing Kernel Hilbert Space). *Recall:*  $\mathcal{H}$  is a Hilbert space with inner product  $\langle ., . \rangle_{\mathcal{H}}$  and norm  $\|.\|_{\mathcal{H}}$ . It satisfies the reproducing property:

$$\forall \quad f \in \mathcal{H}, \ z \in \mathcal{Z}, \quad f(z) = \langle f, k(z, .) \rangle_{\mathcal{H}}$$

Assume  $\mu \mapsto \int k(z, .) d\mu(z)$  injective (characteristic k).

#### Maximum Mean Discrepancy

• Let  $k : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$  a positive, semi-definite kernel

$$k(z,z') = \langle \phi(z), \phi(z') 
angle_{\mathcal{H}}, \quad \phi : \mathcal{Z} \to \mathcal{H}$$

▶  $\mathcal{H}$  its RKHS (Reproducing Kernel Hilbert Space). *Recall:*  $\mathcal{H}$  is a Hilbert space with inner product  $\langle ., . \rangle_{\mathcal{H}}$  and norm  $\|.\|_{\mathcal{H}}$ . It satisfies the reproducing property:

$$\forall \quad f \in \mathcal{H}, \ z \in \mathcal{Z}, \quad f(z) = \langle f, k(z, .) \rangle_{\mathcal{H}}$$

Assume  $\mu \mapsto \int k(z, .) d\mu(z)$  injective (characteristic *k*).

Maximum Mean Discrepancy ([Gretton et al., 2012]) defines a distance on  $\mathcal{P}$  (probability distributions on  $\mathcal{Z}$ ):

$$MMD(\mu,\nu) = \|f_{\mu,\nu}\|_{\mathcal{H}}, \text{ where}$$

$$\underbrace{f_{\mu,\nu}(.) = \int k(z,.)d\mu(z) - \int k(z,.)d\nu(z)}_{\text{"witness function"}}$$

#### MMD functional

For a target distribution  $\nu^*$  (fixed), for any  $\nu \in \mathcal{P}$ :

$$\begin{aligned} \mathcal{F}(\nu) &= \frac{1}{2} MMD^{2}(\nu, \nu^{*}) \\ &= \frac{1}{2} \|f_{\nu, \nu^{*}}\|_{\mathcal{H}}^{2} \\ &= \frac{1}{2} \int k(z, z') d\nu(z) d\nu(z') + \frac{1}{2} \int k(z, z') d\nu^{*}(z) d\nu^{*}(z') \\ &- \int k(z, z') d\nu(z) d\nu^{*}(z') \end{aligned}$$

*Proof* : use the reproducing property with  $f_{\nu,\nu^*}(.) = \int k(z,.)d\nu(z) - \int k(z,.)d\nu^*(z)$ 

Appear as a loss when optimizing some large neural networks.

Consider the following regression problem:

 $(x, y) \sim data$ 



Consider the following regression problem:



►  $\phi_{Z_i}(x) = w_i g(x, \theta_i),$   $Z_i = (w_i, \theta_i) \in \mathbb{R} \times \mathbb{R}^d$  $\phi_{Z_i}$ : non linearity

Example:

$$\phi_Z(x) = wg(ax + b)$$

where  $g : \mathbb{R} \to \mathbb{R}$ (sigmoid  $g(z) = 1/(1 + e^{-z})$ , RelU (g(z) = max(0, z)...)

Finite dimensional non-convex optimization (regression setting):

$$\min_{Z_1,...,Z_N\in\mathcal{Z}} \mathcal{F}\left(\frac{1}{N}\sum_{i=1}^N \delta_{Z_i}\right)$$



Finite dimensional non-convex optimization (regression setting):

$$\min_{Z_1,...,Z_N\in\mathcal{Z}} \mathcal{F}\left(\frac{1}{N}\sum_{i=1}^N \delta_{Z_i}\right)$$

 $(x, y) \sim data$  $\phi_{Z_N}$  $x_4$  $x_3$ ŵ  $\phi_{Z_1}$  $x_2$  $\left( \phi_{Z_{2}} \right)$  $x_1$  $\phi_{Z_1}$  $\min_{Z_1,\ldots,Z_N} \mathbb{E}_{data}[\|y - \frac{1}{N}\sum^N \phi_{Z_i}(x)\|^2]$ 

 Optimization using gradient descent (GD):

$$Z_{i}^{t+1} = Z_{i}^{t} - \gamma \nabla_{Z_{i}} \mathcal{F} \left( \frac{1}{N} \sum_{i=1}^{N} \delta_{Z_{i}^{t}} \right)$$

Finite dimensional non-convex optimization (regression setting):

$$\min_{Z_1,...,Z_N\in\mathcal{Z}} \mathcal{F}\left(\frac{1}{N}\sum_{i=1}^N \delta_{Z_i}\right)$$

 $(x, y) \sim data$  $\phi_{Z_N}$  $x_4$  $x_3$ ŵ  $\phi_{Z_1}$  $x_2$  $\left( \phi_{Z_{2}} \right)$  $x_1$  $\phi_{Z_1}$  $\min_{Z_1,\ldots,Z_N} \mathbb{E}_{data}[\|y - \frac{1}{N}\sum^N \phi_{Z_i}(x)\|^2]$ 

 Optimization using gradient descent (GD):

$$Z_{i}^{t+1} = Z_{i}^{t} - \gamma \nabla_{Z_{i}} \mathcal{F} \left( \frac{1}{N} \sum_{i=1}^{N} \delta_{Z_{i}^{t}} \right)$$

 Hard to describe the dynamics of GD!

Finite dimensional non-convex optimization (regression setting):

$$\min_{Z_1,...,Z_N\in\mathcal{Z}} \mathcal{F}\left(\frac{1}{N}\sum_{i=1}^N \delta_{Z_i}\right)$$



 Optimization using gradient descent (GD):

$$Z_{i}^{t+1} = Z_{i}^{t} - \gamma \nabla_{Z_{i}} \mathcal{F}\left(\frac{1}{N} \sum_{i=1}^{N} \delta_{Z_{i}^{t}}\right)$$

- Hard to describe the dynamics of GD!
- Idea: look at the distribution of the Z<sub>i</sub>'s

### Infinite width regime

Infinite dimensional convex optimization [Chizat and Bach, 2018],

[Mei et al., 2018]



#### Minimization of the MMD : the well-specified case

Assume 
$$\exists \nu^* \in \mathcal{P}$$
,  $\mathbb{E}[y|X = x] = \mathbb{E}_{Z \sim \nu^*}[\phi_Z(x)]$ .

#### Minimization of the MMD : the well-specified case

Assume 
$$\exists \nu^* \in \mathcal{P}$$
,  $\mathbb{E}[y|X = x] = \mathbb{E}_{Z \sim \nu^*}[\phi_Z(x)]$ .

Then:  

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

$$\lim_{\nu \in \mathcal{P}} \mathbb{E}[\|\mathbb{E}_{Z \sim \nu^{*}}[\phi_{Z}(x)] - \mathbb{E}_{Z \sim \nu}[\phi_{Z}(x)]\|^{2}]$$

$$\lim_{\nu \in \mathcal{P}} \mathbb{E}_{Z \sim \nu^{*}}[k(Z, Z')] + \mathbb{E}_{Z \sim \nu}[k(Z, Z')] - 2\mathbb{E}_{Z \sim \nu^{*}}[k(Z, Z')]$$

$$\lim_{Z' \sim \nu} with \ k(Z, Z') = \mathbb{E}_{x \sim data}[\phi_{Z}(x)^{T}\phi_{Z'}(x)]$$

$$\lim_{\nu \in \mathcal{P}} MMD^{2}(\nu, \nu^{*})$$

#### Minimization of the MMD : the well-specified case

Assume 
$$\exists \nu^* \in \mathcal{P}$$
,  $\mathbb{E}[y|X = x] = \mathbb{E}_{Z \sim \nu^*}[\phi_Z(x)]$ .

Then:  

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

$$\lim_{\nu \in \mathcal{P}} \mathbb{E}[\|\mathbb{E}_{Z \sim \nu^{*}}[\phi_{Z}(x)] - \mathbb{E}_{Z \sim \nu}[\phi_{Z}(x)]\|^{2}]$$

$$\lim_{\nu \in \mathcal{P}} \mathbb{E}_{Z \sim \nu^{*}}[k(Z, Z')] + \mathbb{E}_{Z \sim \nu}[k(Z, Z')] - 2\mathbb{E}_{Z \sim \nu^{*}}[k(Z, Z')]$$

$$\lim_{Z' \sim \nu} with \ k(Z, Z') = \mathbb{E}_{x \sim data}[\phi_{Z}(x)^{T}\phi_{Z'}(x)]$$

$$\lim_{\nu \in \mathcal{P}} MMD^{2}(\nu, \nu^{*})$$

Optimizing the parameters of a NN  $\Leftrightarrow$  minimization of the MMD on  $\mathcal{P}$  in the population limit ( $N \to \infty$ ).



Background and motivation

Maximum Mean Discrepancy (Wasserstein) Gradient Flow

Convergence properties of the MMD gradient flow

A noise-injection algorithm for better convergence

We consider

$$\min_{\nu \in \mathcal{P}} \mathcal{F}(\nu) \text{ where } \mathcal{F}(\nu) = \frac{1}{2} \textit{MMD}^2(\nu, \nu^*)$$

 Gradient descent dynamics in this setting takes the form of a PDE (gradient flow on P)

$$rac{\partial 
u_t}{\partial t} = \textit{div}(
u_t 
abla rac{\partial \mathcal{F}(
u_t)}{\partial t}) = \textit{div}(
u_t 
abla f_{
u_t,
u^*})$$

Can be obtained as the limit when  $\tau \rightarrow 0$  of the **JKO** scheme [Jordan et al., 1998] :

$$\nu(n+1) = \operatorname*{argmin}_{\nu \in \mathcal{P}} \mathcal{F}(\nu) + \frac{1}{2\tau} W_2^2(\nu, \nu(n))$$

Density of particles following a Mc-Kean Vlasov dynamic :

$$\frac{dZ_t}{dt} = -\nabla_{Z_t} f_{\nu_t,\nu^*}(Z_t), \qquad Z_t \sim \nu_t$$
  
where  $\nabla_{Z_t} f_{\nu_t,\nu^*} = \int \nabla k(Z,Z_t) d\nu_t(Z) - \int \nabla k(Z,Z_t) d\nu^*(Z).$ 

#### Example : Student-Teacher network

Satisfies the "well-specified" assumption !  $(\exists \nu^*, \mathbb{E}[y|X = x] = \mathbb{E}_{Z \sim \nu^*}[\phi_Z(x)])$ 

- ► the output of the Teacher network is deterministic and given by  $y = \int \phi_Z(x) d\nu^*(Z)$  where  $\nu^* = \frac{1}{M} \sum_{m=1}^M \delta_{U^m}$
- ► Student network parametrized by  $\nu_0 = \frac{1}{N} \sum_{n=1}^{N} \delta_{Z_0^n}$  tries to learn the mapping  $x \mapsto \int \phi_Z(x) d\nu^*(Z)$ .



Gradient descent on each parameter  $n \in \{1, ..., N\}$ :

$$z_{t+1}^n = z_t^n - \gamma \mathbb{E}_{x \sim data} \left[ \left( \frac{1}{N} \sum_{n'=1}^N \phi_{z_t^{n'}}(x) - \frac{1}{M} \sum_{m=1}^M \phi_{u^m}(x) \right) \nabla_{z_t^n} \phi_{z_t^n}(x) \right],$$

Re-arranging terms and recalling that  $k(Z, U) = \mathbb{E}_{x \sim data}[\phi_Z(x)^T \phi_U(x)]$ , the update becomes:

$$z_{t+1}^{n} = z_{t}^{n} - \gamma \underbrace{\left(\frac{1}{N} \sum_{n'=1}^{N} \nabla_{2} k(z_{t}^{n'}, z_{t}^{n}) - \frac{1}{M} \sum_{m=1}^{M} \nabla_{2} k(u^{m}, z_{t}^{n})\right)}_{\nabla f_{\nu^{\star}, \nu_{t}}(z_{t}^{n})}$$

The above equation is a time-discretized version of the gradient flow of the MMD.



Background and motivation

Maximum Mean Discrepancy (Wasserstein) Gradient Flow

Convergence properties of the MMD gradient flow

A noise-injection algorithm for better convergence

 $\mathcal{F}(\nu_t)$  along the MMD flow  $\frac{\partial \nu_t}{\partial t} = div(\nu_t \nabla \frac{\partial \mathcal{F}(\nu_t)}{\partial t})$ ?

 $\mathcal{F}(\nu_t)$  along the MMD flow  $\frac{\partial \nu_t}{\partial t} = div(\nu_t \nabla \frac{\partial \mathcal{F}(\nu_t)}{\partial t})$ ?

A functional *F* is (λ)-geodesically convex if it is convex along W<sub>2</sub> geodesics, i.e. if for any t ∈ [0, 1]:

 $\mathcal{F}(\rho(t)) \le (1-t)\mathcal{F}(\rho(0)) + t\mathcal{F}(\rho(1)) - t(1-t)\frac{\lambda}{2}W_2^2(\rho(0),\rho(1))^2$ 

where  $\rho(t) = ((1 - t)I + tT^{\rho(1)}_{\rho(0)})_{\#}\rho(0)$ 

 $\mathcal{F}(\nu_t)$  along the MMD flow  $\frac{\partial \nu_t}{\partial t} = div(\nu_t \nabla \frac{\partial \mathcal{F}(\nu_t)}{\partial t})$ ?

A functional *F* is (λ)-geodesically convex if it is convex along W<sub>2</sub> geodesics, i.e. if for any t ∈ [0, 1]:

 $\mathcal{F}(\rho(t)) \le (1-t)\mathcal{F}(\rho(0)) + t\mathcal{F}(\rho(1)) - t(1-t)\frac{\lambda}{2}W_2^2(\rho(0),\rho(1))^2$ 

where  $\rho(t) = ((1 - t)I + tT^{\rho(1)}_{\rho(0)})_{\#}\rho(0)$ 

 $\mathcal{F}(\nu_t)$  along the MMD flow  $\frac{\partial \nu_t}{\partial t} = div(\nu_t \nabla \frac{\partial \mathcal{F}(\nu_t)}{\partial t})$ ?

A functional *F* is (λ)-geodesically convex if it is convex along W<sub>2</sub> geodesics, i.e. if for any t ∈ [0, 1]:

 $\mathcal{F}(\rho(t)) \le (1-t)\mathcal{F}(\rho(0)) + t\mathcal{F}(\rho(1)) - t(1-t)\frac{\lambda}{2}W_2^2(\rho(0),\rho(1))^2$ 

where 
$$\rho(t) = ((1 - t)I + tT^{\rho(1)}_{\rho(0)})_{\#}\rho(0)$$

If F is λ-convex with λ > 0, all gradient flows of F converge to the unique minimizer of F [Carrillo et al., 2006]

 $\mathcal{F}(\nu_t)$  along the MMD flow  $\frac{\partial \nu_t}{\partial t} = div(\nu_t \nabla \frac{\partial \mathcal{F}(\nu_t)}{\partial t})$ ?

A functional *F* is (λ)-geodesically convex if it is convex along W<sub>2</sub> geodesics, i.e. if for any t ∈ [0, 1]:

 $\mathcal{F}(\rho(t)) \le (1-t)\mathcal{F}(\rho(0)) + t\mathcal{F}(\rho(1)) - t(1-t)\frac{\lambda}{2}W_2^2(\rho(0),\rho(1))^2$ 

where 
$$\rho(t) = ((1 - t)I + tT^{\rho(1)}_{\rho(0)})_{\#}\rho(0)$$

If F is λ-convex with λ > 0, all gradient flows of F converge to the unique minimizer of F [Carrillo et al., 2006]

**Our finding**: The MMD is not  $\lambda$ -convex with  $\lambda > 0$  in general.

# Convergence of the MMD GF - a Lojasiewicz inequality

$$\frac{d\mathcal{F}(\nu_t)}{dt} \leq -\mathcal{CF}(\nu_t)^2$$

Applying Gronwall's lemma results in:  $\mathcal{F}(\nu_t) = \mathcal{O}(\frac{1}{t})$ .

# Convergence of the MMD GF - a Lojasiewicz inequality

$$rac{{\mathsf d} {\mathcal F}(
u_t)}{{\mathsf d} t} \leq - {old C} {\mathcal F}(
u_t)^2$$

Applying Gronwall's lemma results in:  $\mathcal{F}(\nu_t) = \mathcal{O}(\frac{1}{t})$ .

on the right, it's the RKHS norm:

$$\mathcal{F}(\nu_t) = rac{1}{2} \|f_{\nu_t, \nu^*}\|_{\mathcal{H}}^2$$

on the left we have the weighted Sobolev semi-norm:

$$\frac{d\mathcal{F}(\nu_t)}{dt} = -\int \|\nabla f_{\nu_t,\nu^*}(x)\|^2 d\nu_t(x) = -\|f_{\nu_t,\nu^*}\|^2_{\dot{H}(\nu_t)}$$

Since:

$$rac{\partial 
u_t}{\partial t} = \textit{div}(
u_t 
abla \textit{f}_{
u_t,
u^*})$$

(dissipation of the MMD along the MMD flow)

It can be shown that:

$$\|f_{\nu_t,\nu^*}\|_{\mathcal{H}}^2 \le \|f_{\nu_t,\nu^*}\|_{\dot{H}(\nu_t)}\|\nu^* - \nu_t\|_{\dot{H}^{-1}(\nu_t)}$$

where on the r.h.s. we have the dual norm of  $H(\nu_t)$ :

$$\|\nu_t - \nu^*\|_{\dot{H}^{-1}(\nu_t)} = \sup_{g, \ \mathbb{E}_{Z \sim \nu_t}[\|\nabla g(Z)\|^2] \le 1} |\mathbb{E}_{Z \sim \nu_t}[g(Z)] - \mathbb{E}_{U \sim \nu^*}[g(U)]|$$

It can be shown that:

$$\|f_{\nu_t,\nu^*}\|_{\mathcal{H}}^2 \le \|f_{\nu_t,\nu^*}\|_{\dot{H}(\nu_t)}\|\nu^* - \nu_t\|_{\dot{H}^{-1}(\nu_t)}$$

where on the r.h.s. we have the dual norm of  $H(\nu_t)$ :

$$\|\nu_t - \nu^*\|_{\dot{H}^{-1}(\nu_t)} = \sup_{g, \ \mathbb{E}_{Z \sim \nu_t}[\|\nabla g(Z)\|^2] \le 1} |\mathbb{E}_{Z \sim \nu_t}[g(Z)] - \mathbb{E}_{U \sim \nu^*}[g(U)]|$$

It can be shown that:

$$\|f_{\nu_t,\nu^*}\|_{\mathcal{H}}^2 \le \|f_{\nu_t,\nu^*}\|_{\dot{H}(\nu_t)}\|\nu^* - \nu_t\|_{\dot{H}^{-1}(\nu_t)}$$

where on the r.h.s. we have the dual norm of  $H(\nu_t)$ :

$$\|\nu_t - \nu^*\|_{\dot{H}^{-1}(\nu_t)} = \sup_{g, \ \mathbb{E}_{Z \sim \nu_t}[\|\nabla g(Z)\|^2] \le 1} |\mathbb{E}_{Z \sim \nu_t}[g(Z)] - \mathbb{E}_{U \sim \nu^*}[g(U)]|$$

Assume that  $\|\nu_t - \nu^*\|_{\dot{H}^{-1}(\nu_t)} \leq C$  for all t, then

$$MMD^{2}(\nu_{t}, \nu^{*}) \leq rac{1}{MMD^{2}(\nu_{0}, \nu^{*}) + 4C^{-1}t}$$

It can be shown that:

$$\|f_{\nu_t,\nu^*}\|_{\mathcal{H}}^2 \le \|f_{\nu_t,\nu^*}\|_{\dot{H}(\nu_t)}\|\nu^* - \nu_t\|_{\dot{H}^{-1}(\nu_t)}$$

where on the r.h.s. we have the dual norm of  $H(\nu_t)$ :

$$\|\nu_t - \nu^*\|_{\dot{H}^{-1}(\nu_t)} = \sup_{g, \ \mathbb{E}_{Z \sim \nu_t}[\|\nabla g(Z)\|^2] \le 1} |\mathbb{E}_{Z \sim \nu_t}[g(Z)] - \mathbb{E}_{U \sim \nu^*}[g(U)]|$$

Assume that  $\|\nu_t - \nu^*\|_{\dot{H}^{-1}(\nu_t)} \leq C$  for all t, then

$$MMD^{2}(\nu_{t},\nu^{*}) \leq \frac{1}{MMD^{2}(\nu_{0},\nu^{*}) + 4C^{-1}t}$$

**Problem:** Depends on the whole sequence  $\nu_t$ ; Hard to verify in general [Peyre, 2018]; and we've seen failure cases in practice.







#### Convergence issues

► The condition we exhibited for global convergence may not hold and (*F*(*v*<sub>t</sub>))<sub>t</sub> might be stuck at a local minima.

$$\begin{aligned} \frac{d\mathcal{F}(\nu_t)}{dt} &= -\int \|\nabla f_{\nu_t,\nu^*}(x)\|^2 d\nu_t(x) \text{ at equilibrium} \\ &\implies \int \|\nabla f_{\nu^{\infty},\nu^*}(x)\|^2 d\nu^{\infty}(x) = 0 \end{aligned}$$

If  $\nu^{\infty}$  positive everywhere this implies  $f_{\nu^{\infty},\nu^*} = cte = 0^1$ But  $\nu^{\infty}$  might be singular...

Idea : Evaluate ∇f<sub>νt,ν\*</sub> outside of the support of νt to get a better signal!

 $<sup>^{1}\</sup>text{as}$  soon as the RKHS  $\mathcal H$  does not contain non-zero constant functions, e.g. for a gaussian kernel.



Background and motivation

Maximum Mean Discrepancy (Wasserstein) Gradient Flow

Convergence properties of the MMD gradient flow

A noise-injection algorithm for better convergence

#### Noise injection

At each iteration *n*, sample ξ<sub>n</sub> ~ N(0, 1) and β<sub>n</sub> is the noise level:

$$Z_{n+1} = Z_n - \gamma \nabla f_{\nu_n,\nu^*} (Z_n + \beta_n \xi_n)$$

<sup>2</sup>[Duchi et al., 2012] <sup>3</sup>[Mei et al., 2018]

#### Noise injection

At each iteration n, sample ξ<sub>n</sub> ~ N(0, 1) and β<sub>n</sub> is the noise level:

$$Z_{n+1} = Z_n - \gamma \nabla f_{\nu_n,\nu^*} (Z_n + \beta_n \xi_n)$$

 Similar to randomized smoothing<sup>2</sup>, but extended to interacting particles.

<sup>&</sup>lt;sup>2</sup>[Duchi et al., 2012] <sup>3</sup>[Mei et al., 2018]

#### Noise injection

At each iteration n, sample ξ<sub>n</sub> ~ N(0, 1) and β<sub>n</sub> is the noise level:

$$Z_{n+1} = Z_n - \gamma \nabla f_{\nu_n,\nu^*} (Z_n + \beta_n \xi_n)$$

- Similar to randomized smoothing<sup>2</sup>, but extended to interacting particles.
- Different from adding noise outside ("diffusion")

$$Z_{n+1} = Z_n - \gamma \nabla f_{\nu_n,\nu^*}(Z_n) + \beta_n \xi_n$$

which corresponds to an entropic regularization of the original loss <sup>3</sup>.

<sup>&</sup>lt;sup>2</sup>[Duchi et al., 2012]

<sup>&</sup>lt;sup>3</sup>[Mei et al., 2018]

#### Noise Injection: Theory (discrete time)

Tradeoff for the level of noise  $\beta_n$ 

• Too large  $\beta_n$ :  $\nu_{n+1}$  not a descent direction anymore:  $\mathcal{F}(\nu_{n+1}, \nu^*) > \mathcal{F}^2(\nu_n, \nu^*)$ 

$$\Longrightarrow \beta_n^2 \mathcal{F}^2(\nu_n) \le C_k \mathbb{E}_{\substack{Z_n \sim \nu_n \\ u_n \sim \mathcal{N}(0,1)}} [\|\nabla f_{\nu_n,\nu^*}(Z_n + \beta_n \xi_n)\|^2]$$
(1)

#### Noise Injection: Theory (discrete time)

Tradeoff for the level of noise  $\beta_n$ 

• Too large  $\beta_n$ :  $\nu_{n+1}$  not a descent direction anymore:  $\mathcal{F}(\nu_{n+1}, \nu^*) > \mathcal{F}^2(\nu_n, \nu^*)$ 

$$\Longrightarrow \beta_n^2 \mathcal{F}^2(\nu_n) \le C_k \mathbb{E}_{\substack{Z_n \sim \nu_n \\ u_n \sim \mathcal{N}(0,1)}} [\|\nabla f_{\nu_n,\nu^*}(Z_n + \beta_n \xi_n)\|^2]$$
(1)

• Too small  $\beta_n$ : Back to the failure mode:  $\nabla f_{\nu_n,\nu^*}(Z_n + \beta_n u_n) \simeq 0.$ 

$$\implies \sum_{n=1}^{N} \beta_n^2 \to \infty \tag{2}$$

Under (1) and (2) :

$$\mathcal{F}^2(\nu_N,\nu^*) \leq \mathcal{F}^2(\nu_0,\nu^*) \boldsymbol{e}^{-C_k\gamma(1-\gamma C'_k)\sum_{n=1}^N \beta_n^2}$$

#### Noise Injection: Experiments (Gaussians)

# Noise Injection: Experiments (Student-Teacher)

Methods:

- SGD
- SGD + Noise injection
- SGD + diffusion
- KSD <sup>4</sup>: SGD using a (regularized) Negative Sobolev distance as a loss function; also decreases the MMD.

<sup>&</sup>lt;sup>4</sup>"Kernel Sobolev Descent" [Mroueh et al., 2019]

### Noise Injection: Experiments



dimension = 50

# Noise Injection: Experiments



### Contributions and openings

- Provided a convergence criterion for the Wasserstein gradient flow of the MMD.
- Proposed a pertubation of the dynamics with a noise injection and showed it effectiveness on simple examples.
- new insights for training large neural networks.

Openings:

- Control of the weighted negative Sobolev norm?
- Stronger guarantees for the convergence for the noise injection algorithm.

Carrillo, J. A., McCann, R. J., and Villani, C. (2006). Contractions in the 2-wasserstein length space and thermalization of granular media. *Archive for Rational Mechanics and Analysis*, 179(2):217–263.

- Chizat, L. and Bach, F. (2018). On the global convergence of gradient descent for over-parameterized models using optimal transport. NIPS.
- Duchi, J. C., Bartlett, P. L., and Wainwright, M. J. (2012).
   Randomized smoothing for stochastic optimization.
   SIAM Journal on Optimization, 22(2):674–701.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. J. (2012).
   A kernel two-sample test. *JMLR*, 13.
- Jordan, R., Kinderlehrer, D., and Otto, F. (1998). The variational formulation of the fokker–planck equation. <sup>30/30</sup>