

# Sampling through Optimization of Divergences

Anna Korba

ENSAE, CREST, Institut Polytechnique de Paris

Journée SIGMA - MODE 2024

# Outline

- 1 Introduction
- 2 Sampling as Optimization
- 3 Choice of the Divergence
- 4 Optimization and Quantization errors
- 5 Quantization

# Why sampling?

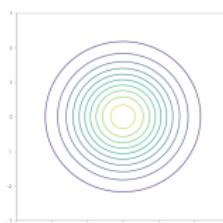
Suppose you are interested in some target probability distribution on  $\mathbb{R}^d$ , denoted  $\mu^*$ , and you have access only to partial information, e.g.:

- ① its unnormalized density (as in Bayesian inference)
- ② a discrete approximation  $\frac{1}{m} \sum_{k=1}^m \delta_{x_i} \approx \mu^*$  (e.g. i.i.d. samples, iterates of MCMC algorithms...)

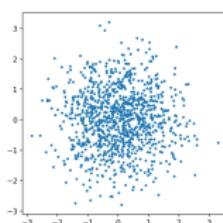
**Problem:** approximate  $\mu^* \in \mathcal{P}(\mathbb{R}^d)$  by a finite set of  $n$  points  $x_1, \dots, x_n$ , e.g. to compute functionals  $\int_{\mathbb{R}^d} f(x) d\mu^*(x)$ .

The quality of the set can be measured by the integral error:

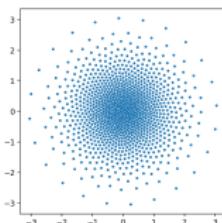
$$\left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \int_{\mathbb{R}^d} f(x) d\mu^*(x) \right|.$$



a Gaussian density



i.i.d. samples.



Some discrete approximation

# Example 1: Bayesian inference

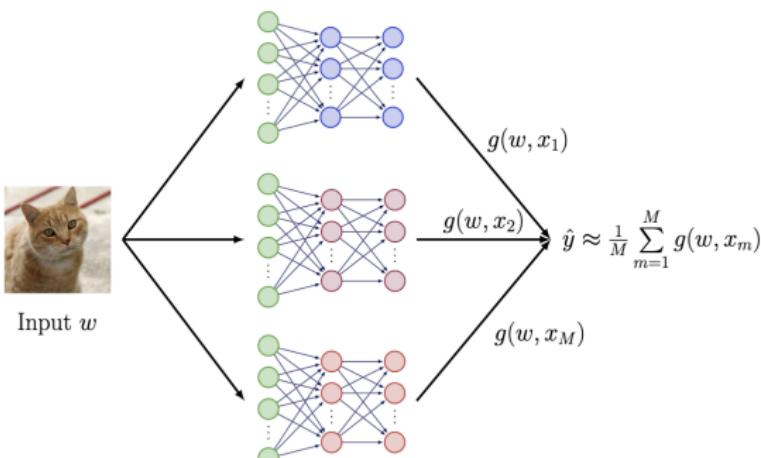
We want to sample from

$$\mu^*(x) \propto \exp(-V(x)), \quad V(\theta) = \underbrace{\sum_{i=1}^m \|y_i - g(w_i, x)\|^2}_{\text{loss on labeled data } (w_i, y_i)_{i=1}^m} + \frac{\|x\|^2}{2}.$$

Ensemble prediction for a new input  $w$ :

$$\hat{y} = \underbrace{\int_{\mathbb{R}^d} g(w, x) d\mu^*(x)}_{\text{"Bayesian model averaging"}}$$

Predictions of models parametrized by  $x \in \mathbb{R}^d$   
are reweighted by  $\mu^*(x)$ .

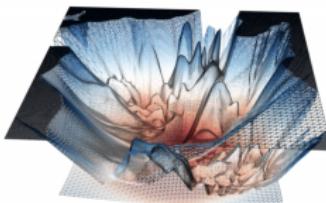


# Difficult cases (in practice and in theory)

Recall that

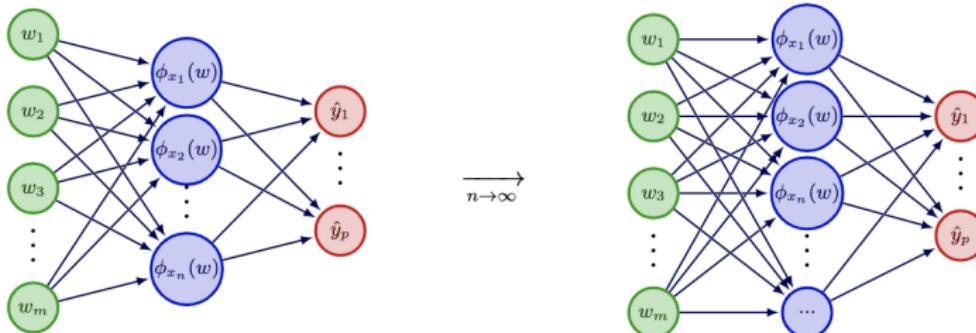
$$\mu^*(x) \propto \exp(-V(x)), \quad V(\theta) = \underbrace{\sum_{i=1}^m \|y_i - g(w_i, x)\|^2}_{\text{loss}} + \frac{\|x\|^2}{2}.$$

- if  $V$  is convex (e.g.  $g(w, x) = \langle w, x \rangle$ ) many sampling methods (e.g. Langevin Monte Carlo) are known to work quite well [Durmus and Moulines, 2016, Vempala and Wibisono, 2019]
- but if its not (e.g.  $g(w, x)$  is a neural network), the situation is much more delicate



A highly nonconvex loss surface, as is common in deep neural nets. From <https://www.telesens.co/2019/01/16/neural-network-loss-visualization>.

## Example 2 : Regression with infinite width NN



$$\min_{(x_i)_{i=1}^n \in \mathbb{R}^d} \mathbb{E}_{(w,y) \sim P_{\text{data}}} \left[ \underbrace{\left\| y - \frac{1}{n} \sum_{i=1}^n \phi_{x_i}(w) \right\|^2}_{\hat{y}} \right] \xrightarrow{n \rightarrow \infty} \min_{\mu \in \mathcal{P}(\mathbb{R}^d)} \underbrace{\mathbb{E}_{(w,y) \sim P_{\text{data}}} \left[ \left\| y - \int_{\mathbb{R}^d} \phi_x(w) d\mu(x) \right\|^2 \right]}_{\mathcal{F}(\mu)}$$

Define the target distribution  $\mu^* \in \arg \min \mathcal{F}(\mu)$ . Optimising the neural network  $\iff$  approximating  $\mu^*$  [Chizat and Bach, 2018, Mei et al., 2018].

If  $y(w) = \frac{1}{m} \sum_{i=1}^m \phi_{x_i}(w)$  is generated by a neural network (as in the student-teacher network setting), then  $\mathcal{F}$  can be identified to an MMD [Arbel et al., 2019].

# Outline

- 1 Introduction
- 2 Sampling as Optimization
- 3 Choice of the Divergence
- 4 Optimization and Quantization errors
- 5 Quantization

# Sampling as optimization over probability distributions

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

The sampling task can be recast as an optimization problem:

$$\mu^* = \arg \min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} D(\mu|\mu^*) := \mathcal{F}(\mu),$$

where  $D$  is a **dissimilarity functional**, for instance:

- a f-divergence:  $\int f\left(\frac{\mu}{\mu^*}\right) d\mu^*$ ,  $f$  convex,  $f(1) = 0$
- an integral probability metric:  $\sup_{f \in \mathcal{G}} \left| \int f d\mu - \int f d\mu^* \right|$
- an optimal transport distance...

Starting from an initial distribution  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^d)$ , one can then consider the **Wasserstein-2\* gradient flow** of  $\mathcal{F}$  over  $\mathcal{P}_2(\mathbb{R}^d)$  to transport  $\mu_0$  to  $\mu^*$ .

---

\*  $W_2^2(\nu, \mu) = \inf_{s \in \Gamma(\nu, \mu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 ds(x, y)$ , where  $\Gamma(\nu, \mu)$  = couplings between  $\nu, \mu$ .

# Wasserstein gradient flows (WGF) [Ambrosio et al., 2008]

The first variation of  $\mu \mapsto \mathcal{F}(\mu)$  evaluated at  $\mu \in \mathcal{P}(\mathbb{R}^d)$  is the unique function  $\frac{\partial \mathcal{F}(\mu)}{\partial \mu} : \mathbb{R}^d \rightarrow \mathbb{R}$  s. t. for any  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ ,  $\nu - \mu \in \mathcal{P}(\mathbb{R}^d)$ :

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} (\mathcal{F}(\mu + \epsilon(\nu - \mu)) - \mathcal{F}(\mu)) = \int_{\mathbb{R}^d} \frac{\partial \mathcal{F}(\mu)}{\partial \mu}(x) (d\nu - d\mu)(x).$$

The family  $\mu : [0, \infty] \rightarrow \mathcal{P}_2(\mathbb{R}^d)$ ,  $t \mapsto \mu_t$  is a **Wasserstein gradient flow** of  $\mathcal{F}$  if:

$$\frac{\partial \mu_t}{\partial t} = \nabla \cdot (\mu_t \nabla_{W_2} \mathcal{F}(\mu_t)),$$

where  $\nabla_{W_2} \mathcal{F}(\mu) := \nabla \frac{\partial \mathcal{F}(\mu)}{\partial \mu}$  denotes the **Wasserstein gradient** of  $\mathcal{F}$ .

It can be implemented by the deterministic process:

$$\frac{dx_t}{dt} = -\nabla_{W_2} \mathcal{F}(\mu_t)(x_t), \quad \text{where } x_t \sim \mu_t$$

# Particle system/Gradient descent approximating the WGF

**Space/time discretization** : Introduce a particle system  $x_0^1, \dots, x_0^n \sim \mu_0$ , a step-size  $\gamma$ , and at each step:

$$x_{l+1}^i = x_l^i - \gamma \nabla_{W_2} \mathcal{F}(\hat{\mu}_l)(x_l^i) \quad \text{for } i = 1, \dots, n, \text{ where } \hat{\mu}_l = \frac{1}{n} \sum_{i=1}^n \delta_{x_l^i}$$

In particular, the algorithm above simply corresponds to gradient descent.  
We consider several questions:

- what can we say as time goes to infinity ? (**optimization error**)  
     $\Rightarrow$  heavily linked with the geometry (convexity, smoothness in the Wasserstein sense) of the loss
- (for minimizers) what can we say as the number of particles grow ? (**"quantization" error**)

# Outline

- 1 Introduction
- 2 Sampling as Optimization
- 3 Choice of the Divergence
- 4 Optimization and Quantization errors
- 5 Quantization

# Loss function for the unnormalized densities - the KL

Many possibilities for the choice of  $D(\cdot|\mu^*)$  among Wasserstein distances,  $f$ -divergences, Integral Probability Metrics...

For instance,  $D$  could be the Kullback-Leibler divergence:

$$\text{KL}(\mu|\mu^*) = \begin{cases} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\mu^*}(x)\right) d\mu(x) & \text{if } \mu \ll \mu^* \\ +\infty & \text{otherwise.} \end{cases}$$

The KL as an objective is convenient when the unnormalized density of  $\mu^*$  is known since it **does not depend on the normalization constant!**

Indeed writing  $\mu^*(x) = e^{-V(x)}/Z$  we have:

$$\text{KL}(\mu|\mu^*) = \int_{\mathbb{R}^d} \log\left(\frac{\mu}{e^{-V}}(x)\right) d\mu(x) + \log(Z).$$

**But, it is not convenient when we have a discrete approximation of  $\mu^*$ .  
Also, we cannot evaluate it for discrete  $\mu$ .**

# KL Gradient flow in practice

- The gradient flow of the KL can be implemented via the Probability Flow (ODE):

$$d\tilde{x}_t = -\nabla \log \left( \frac{\mu_t}{\mu^*} \right) (\tilde{x}_t) dt \quad (1)$$

or the Langevin diffusion (SDE):

$$dx_t = \nabla \log \mu^*(x_t) dt + \sqrt{2} dB_t \quad (2)$$

(they share the same marginal  $(\mu_t)_{t \geq 0}$ )

- (2) can be discretized in time as Langevin Monte Carlo (LMC)  
[Roberts and Tweedie, 1996]

$$x_{m+1} = x_m + \gamma \nabla \log \mu^*(x_m) + \sqrt{2\gamma} \epsilon_m, \quad \epsilon_m \sim \mathcal{N}(0, \text{Id}_{\mathbb{R}^d}).$$

- (1) can be approximated by a particle system (e.g. SVGD  
[Liu, 2017, He et al., 2022])
- however MCMC methods suffer an integral approximation error of order  $\mathcal{O}(n^{-1/2})$  if we use  $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  ( $x_i$  iterates of MCMC)  
[Łatuszyński et al., 2013]

# Another f-divergence?

- The chi-square divergence:

$$\chi^2(\mu \parallel \mu^*) := \begin{cases} \int \left( \frac{d\mu}{d\mu^*} - 1 \right)^2 d\mu^* & \mu \ll \mu^* \\ +\infty & \text{otherwise.} \end{cases}$$

not convenient neither when  $\mu^*$ 's unnormalized density is known, or if we have a discrete approximation.

- $\chi^2$ -gradient requires the normalizing constant of  $\mu^*$ :  $\nabla \frac{\mu}{\mu^*}$
- However, the GF of chi-square has interesting properties  
[Chewi et al., 2020, Craig et al., 2022] (e.g. linear convergence under Poincaré inequality on  $\mu^*$ , whereas the KL flow requires log-Sobolev)

# Losses for the discrete case

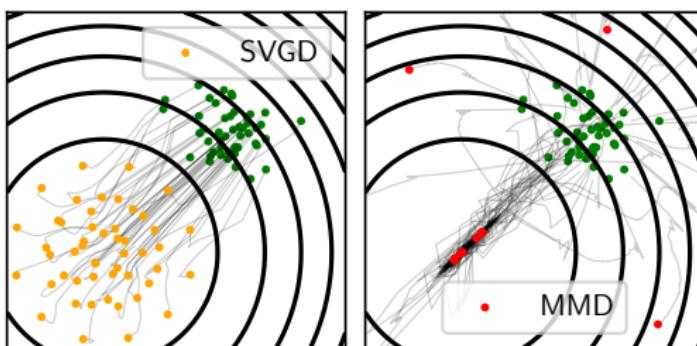
D could be the MMD (Maximum Mean Discrepancy):

$$\begin{aligned} \text{MMD}^2(\mu, \mu^*) &= \sup_{f \in \mathcal{H}_k, \|f\|_{\mathcal{H}_k} \leq 1} \left| \int f d\mu - \int f d\mu^* \right| \\ &\quad \iint_{\mathbb{R}^d} k(x, y) d\mu(x) d\mu(y) \\ &\quad + \iint_{\mathbb{R}^d} k(x, y) d\mu^*(x) d\mu^*(y) - 2 \iint_{\mathbb{R}^d} k(x, y) d\mu(x) d\mu^*(y). \end{aligned}$$

where  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a p.s.d. kernel (e.g.  $k(x, y) = e^{-\|x-y\|^2}$ ).

**It is convenient when we have a discrete approximation of  $\mu^*$  (to approximate integrals).**

# Why we care about the loss



**Figure:** Toy example with 2D standard Gaussian. The green points represent the initial positions of the particles. The light grey curves correspond to their trajectories under the different  $v_{\mu_t}$ .

**Gradient flow of the KL to a Gaussian**  $\mu^*(x) \propto e^{-\frac{\|x\|^2}{2}}$  is well-behaved, but not the MMD.

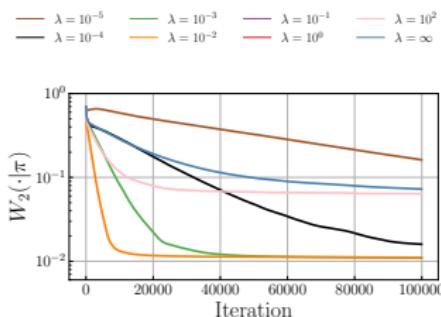
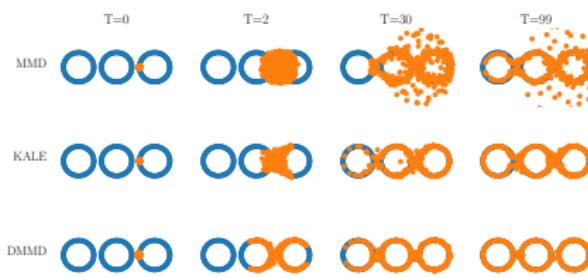
# A proposal<sup>†</sup>: Interpolate between MMD and $\chi^2$

"De-Regularized MMD" leverages the variational formulation of  $\chi^2$ :

$$\text{DMMD}(\mu||\mu^*) = (1 + \lambda) \left\{ \max_{h \in \mathcal{H}} \int h d\mu - \int h d\mu^* - \frac{1}{4} \int h^2 d\mu^* - \frac{1}{4} \lambda \|h\|_{\mathcal{H}}^2 \right\}$$

**Can be written in closed-form if  $\mu, \mu^*$  are discrete** (depends on  $\lambda$  and kernel matrices over samples of  $\mu, \mu^*$ ). It is a divergence for any  $\lambda$ .

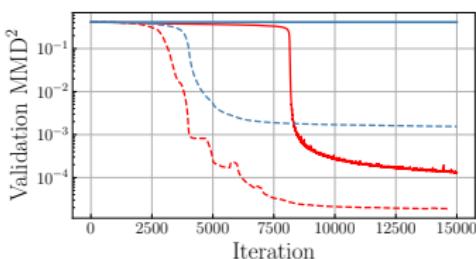
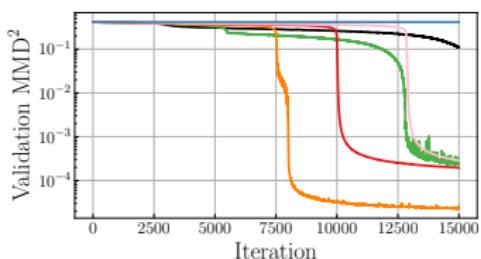
A similar idea was proposed for the KL, yielding Kale divergence [Glaser et al., 2021] but was not closed-form.



<sup>†</sup>with H. Chen, A. Gretton, P. Glaser (UCL), A. Mustafi, B. Sriperumbudur (CMU)

# Student-teacher networks experiment<sup>†</sup>

—  $\lambda = 10^{-5}$     —  $\lambda = 10^{-3}$     —  $\lambda = 10^{-1}$     —  $\lambda = 10^2$   
 —  $\lambda = 10^{-4}$     —  $\lambda = 10^{-2}$     —  $\lambda = 10^0$     —  $\lambda = \infty$   
 — DMMD    — DMMD (Noise)  
 — MMD    — MMD (Noise)



- the teacher network  $w \mapsto \Psi_{\mu^*}(w)$  is given by  $M$  particles  $(\xi_1, \dots, \xi_M)$  which are fixed during training  $\implies \mu = \frac{1}{M} \sum_{j=1}^M \delta_{\xi_j}$
- the student network  $w \mapsto \Psi_\mu(w)$  has  $n$  particles  $(x_1, \dots, x_n)$  that are initialized randomly  $\implies \mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

$$\min_{\mu} \mathbb{E}_{w \sim P_{data}} [(\Psi_{\mu^*}(w) - \Psi_{\mu}(w))^2]$$

$$\iff \min \text{MMD}(\mu, \mu^*) \text{ with } k(x, x') = \mathbb{E}_{w \sim P_{data}} [\phi_{x'}(w)\phi_x(w)].$$

<sup>†</sup>Same setting as [Arbel et al., 2019].

# Another idea - "Mollified" discrepancies

Examples of mollifiers/kernels (Gaussian, Laplace, Riesz-s):

$$k_\epsilon^g(x) := \frac{\exp\left(-\frac{\|x\|_2^2}{2\epsilon^2}\right)}{Z^g(\epsilon)}, \quad k_\epsilon^l(x) := \frac{\exp\left(-\frac{\|x\|_2}{\epsilon}\right)}{Z^l(\epsilon)}, \quad k_\epsilon^s(x) := \frac{1}{(\|x\|_2^2 + \epsilon^2)^{s/2} Z^r(s, \epsilon)}$$



- Mollified chi-square [Li et al., 2022, Craig et al., 2022]:

$$\mathcal{E}_\epsilon(\mu) = \int \left( k_\epsilon * \frac{\mu}{\sqrt{\mu^*}} \right)(x) \frac{\mu}{\sqrt{\mu^*}}(x) dx \xrightarrow[\epsilon \rightarrow 0]{} \chi^2(\mu|\mu^*) + 1$$

- Mollified KL<sup>§</sup> [Craig and Bertozzi, 2016]:

$$\text{KL}(k_\epsilon * \mu | \mu^*) \xrightarrow[\epsilon \rightarrow 0]{} \text{KL}(\mu | \mu^*)$$

---

§ Also ongoing work with Tom Huix (CMAP).

# Outline

- 1 Introduction
- 2 Sampling as Optimization
- 3 Choice of the Divergence
- 4 Optimization and Quantization errors
- 5 Quantization

# (Geodesically)-convex and smooth losses

$\mathcal{F}$  is said to be  $\lambda$ -displacement convex if along  $W_2$  geodesics  $(\rho_t)_{t \in [0,1]}$ :

$$\mathcal{F}(\rho_t) \leq (1-t)\mathcal{F}(\rho_0) + t\mathcal{F}(\rho_1) - \frac{\lambda}{2}t(1-t)W_2^2(\rho_0, \rho_1) \quad \forall t \in [0, 1]. \quad (3)$$

The **Wasserstein Hessian** of a functional  $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$  at  $\mu$  is defined for any  $\psi \in \mathcal{C}_c^\infty(\mathbb{R}^d)$  as:

$$\text{Hess}_\mu \mathcal{F}(\psi, \psi) := \left. \frac{d^2}{dt^2} \right|_{t=0} \mathcal{F}(\mu_t)$$

where  $(\mu_t, v_t)_{t \in [0,1]}$  is a Wasserstein geodesic with  $\mu_0 = 0$ ,  $v_0 = \nabla \psi$ .

$$\mathcal{F} \text{ is } \lambda\text{-displacement convex} \iff \text{Hess}_\mu \mathcal{F}(\psi, \psi) \geq \lambda \|\nabla \psi\|_{L^2(\mu)}^2$$

(See [Villani, 2009, Proposition 16.2]). In an analog manner we can define **smooth functionals** as functionals with upper bounded Hessians.

# Guarantees for Wasserstein gradient descent

Consider Wasserstein gradient descent (Euler discretization of Wasserstein gradient flow)

$$\mu_{l+1} = (\text{Id} - \gamma \nabla \mathcal{F}'(\mu_l))_\# \mu_l$$

Assume  $\mathcal{F}$  is *M-smooth*. Then, we have a descent lemma:

$$\mathcal{F}(\mu_{l+1}) - \mathcal{F}(\mu_l) \leq -\gamma \left(1 - \frac{\gamma}{2} M\right) \|\nabla \mathcal{F}'(\mu_l)\|_{L^2(\mu_l)}^2.$$

Moreover, if  $\mathcal{F}$  is  *$\lambda$ -convex*, we have the global rate

$$\mathcal{F}(\mu_L) \leq \frac{W_2^2(\mu_0, \mu^*)}{2\gamma L} - \frac{\lambda}{L} \sum_{l=0}^L W_2^2(\mu_l, \mu^*).$$

(so the barrier term degrades with  $\lambda$ ).

# Some examples

- Let  $\mu^* \propto e^{-V}$ , we have [Wibisono, 2018]

$$\text{Hess}_\mu \text{KL}(\psi, \psi) = \int \left[ \langle H_V(x) \nabla \psi(x), \nabla \psi(x) \rangle + \|H\psi(x)\|_{HS}^2 \right] q(x) dx.$$

If  $V$  is  $m$ -strongly convex, then the KL is  $m$ -geo. convex; however it is not smooth (Hessian is unbounded). Similar story for  $\chi^2$ -square [Ohta and Takatsu, 2011].

- For a  $M$ -smooth kernel  $k$  [Arbel et al., 2019]

$$\begin{aligned} \text{Hess}_\mu \text{MMD}^2(\psi, \psi) &= \left\| \int \nabla \phi(x)^\top \nabla_1 k(x, \cdot) d\mu(x) \right\|_{\mathcal{H}}^2 + \\ &2 \int \nabla \phi(x)^\top \left( \int H_1 k(x, z) d\rho(z) - \int H_1 k(x, z) d\mu^*(z) \right) \nabla \phi(x) d\mu(x) \end{aligned}$$

It is  $M$ -smooth but not geodesically convex (Hessian lower bounded by a big negative constant)

- Partial results for other discrepancies (De-Regularized MMD, mollified chi-square [Li et al., 2022] and KL)

# Outline

- 1 Introduction
- 2 Sampling as Optimization
- 3 Choice of the Divergence
- 4 Optimization and Quantization errors
- 5 Quantization

# What is known

What can we say on  $\inf_{x_1, \dots, x_n} D(\mu_n | \mu^*)$  where  $\mu_n = \sum_{i=1}^n \delta_{x_i}$ ?

- Quantization rates for the Wasserstein distance  
[Kloeckner, 2012, Mérigot et al., 2021]

$$W_2(\mu_n, \mu^*) \sim O(n^{-\frac{1}{d}})$$

- Forward KL [Li and Barron, 1999]: for every  $g_P = \int k_\epsilon(\cdot - w) dP(w)$ ,

$$\arg \min_{\mu_n} \text{KL}(\mu^* | k_\epsilon \star \mu_n) \leq \text{KL}(\mu^* | g_P) + \frac{C_{\mu^*, P}^2 \gamma}{n}$$

where  $C_{\mu^*, P}^2 = \int \frac{\int k_\epsilon(x-m)^2 dP(m)}{(\int k_\epsilon(x-w) dP(w))^2} d\mu^*(x)$ , and  $\gamma = 4 \log(3\sqrt{e} + a)$  is a constant depending on  $\epsilon$  with  $a = \sup_{z, z' \in \mathbb{R}^d} \log(k_\epsilon(x-z)/k_\epsilon(x-z'))$ .

# Recent results

- For smooth and bounded kernels in [Xu et al., 2022] and  $\mu^*$  with exponential tails, we get using Koksma-Hlawka inequality

$$\min_{\mu_n} \text{MMD}(\pi, \mu_n) \leq C_d \frac{(\log n)^{\frac{5d+1}{2}}}{n}.$$

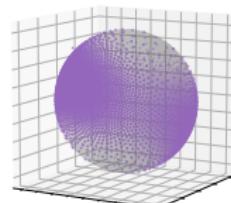
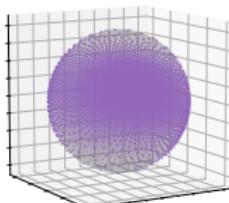
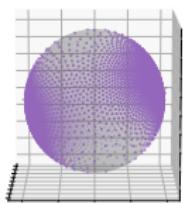
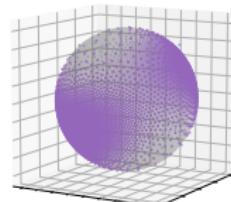
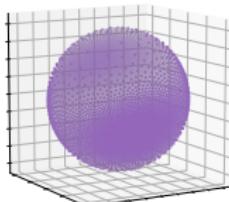
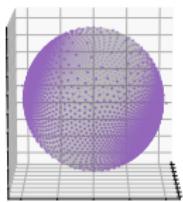
This bounds the integral error for  $f \in \mathcal{H}_k$  by Cauchy-Schwartz:

$$\left| \int_{\mathbb{R}^d} f(x) d\pi(x) - \int_{\mathbb{R}^d} f(x) d\mu(x) \right| \leq \|f\|_{\mathcal{H}_k} \text{MMD}(\mu, \pi).$$

- For the reverse KL (joint work with Tom Huix) we get (in the well-specified case) adapting the proof of [Li and Barron, 1999]:

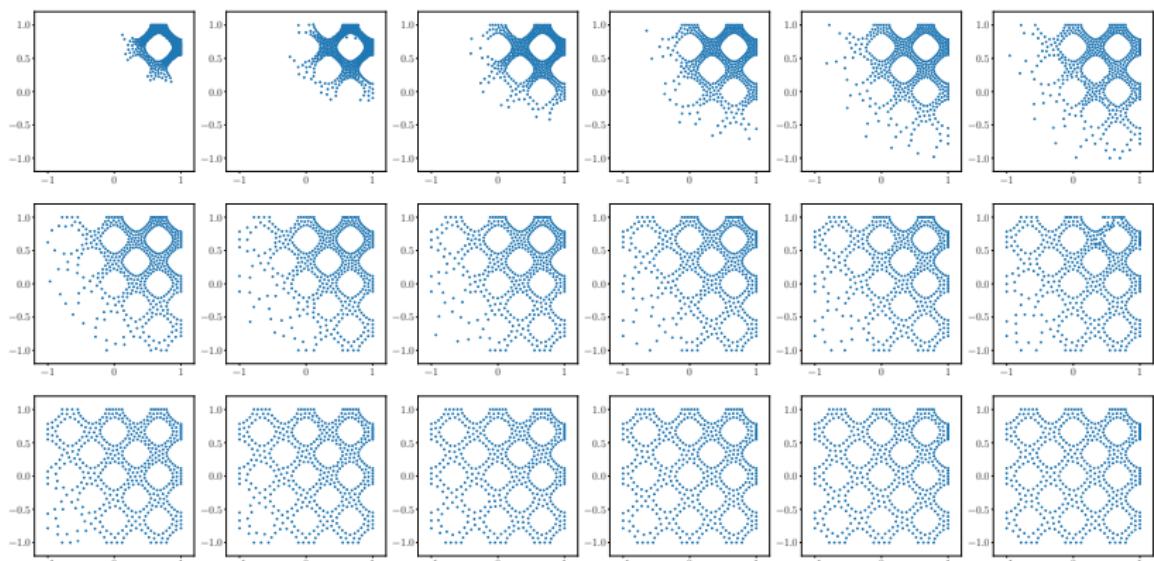
$$\min_{\mu_n} \text{KL}(k_\epsilon * \mu | \mu^*) \leq C_{\mu^*}^2 \frac{\log(n) + 1}{n}.$$

# A numerical example from [Li et al., 2022]



**Figure:** Sampling from the von Mises-Fisher distribution obtained by constraining a 3-dimensional Gaussian to the unit sphere. The unit-sphere constraint is enforced using the dynamic barrier method and the shown results are obtained using MIED with Riesz kernel and  $s = 3$ . The six plots are views from six evenly spaced angles.

# Another numerical example from [Li et al., 2022]



**Figure:** Uniform sampling of the region  $\{(x, y) \in [-1, 1]^2 : (\cos(3\mu^*x) + \cos(3\mu^*y))^2 < 0.3\}$  using MIED with a Riesz mollifier ( $s = 3$ ) where the constraint is enforced using the dynamic barrier method. The plot in row  $i$  column  $j$  shows the samples at iteration  $100 + 200(6i + j)$ . The initial samples are drawn uniformly from the top-right square  $[0.5, 1.0]^2$ .

# Open questions

- Finite-particle/quantization guarantees are still missing for many losses or in the non-well specified case

$$D(\mu_n || \mu^*) \leq f(n, \mu^*)?$$

- How to improve the performance of the algorithms for highly non-log concave targets? e.g. through sequence of targets  $(\mu^*)_{t \in [0,1]}$  interpolating between  $\mu_0$  and  $\mu^*$ ?

References (with code):

- Maximum Mean Discrepancy Gradient Flow. Arbel, M., Korba, A., Salim, A., and Gretton, A. (Neurips 2019).
- Accurate quantization of measures via interacting particle-based optimization. Xu, L., Korba, A., and Slep̄cev, D. (ICML 2022).
- Sampling with mollified interaction energy descent. Li, L., Liu, Q., Korba, A., Yurochkin, M., and Solomon, J. (ICLR 2023).

# References I

-  Ambrosio, L., Gigli, N., and Savaré, G. (2008).  
*Gradient flows: in metric spaces and in the space of probability measures.*  
Springer Science & Business Media.
-  Arbel, M., Korba, A., Salim, A., and Gretton, A. (2019).  
Maximum mean discrepancy gradient flow.  
In *Advances in Neural Information Processing Systems*, pages 6481–6491.
-  Chewi, S., Le Gouic, T., Lu, C., Maunu, T., and Rigollet, P. (2020).  
Svsgd as a kernelized wasserstein gradient flow of the chi-squared divergence.  
*Advances in Neural Information Processing Systems*, 33:2098–2109.
-  Chizat, L. and Bach, F. (2018).  
On the global convergence of gradient descent for over-parameterized models using optimal transport.  
*Advances in neural information processing systems*, 31.
-  Chopin, N., Crucinio, F. R., and Korba, A. (2023).  
A connection between tempering and entropic mirror descent.  
*arXiv preprint arXiv:2310.11914*.
-  Craig, K. and Bertozzi, A. (2016).  
A blob method for the aggregation equation.  
*Mathematics of computation*, 85(300):1681–1717.
-  Craig, K., Elamvazhuthi, K., Haberland, M., and Turanova, O. (2022).  
A blob method method for inhomogeneous diffusion with applications to multi-agent control and sampling.  
*arXiv preprint arXiv:2202.12927*.
-  Durmus, A. and Moulines, E. (2016).  
Sampling from strongly log-concave distributions with the unadjusted langevin algorithm.  
*arXiv preprint arXiv:1605.01559*, 5.

# References II

-  Glaser, P., Arbel, M., and Gretton, A. (2021).  
Kale flow: A relaxed kl gradient flow for probabilities with disjoint support.  
*Advances in Neural Information Processing Systems*, 34:8018–8031.
-  He, Y., Balasubramanian, K., Sriperumbudur, B. K., and Lu, J. (2022).  
Regularized stein variational gradient flow.  
*arXiv preprint arXiv:2211.07861*.
-  Kloeckner, B. (2012).  
Approximation by finitely supported measures.  
*ESAIM: Control, Optimisation and Calculus of Variations*, 18(2):343–359.
-  Łatuszyński, K., Miasojedow, B., and Niemiro, W. (2013).  
Nonasymptotic bounds on the estimation error of mcmc algorithms.  
*Bernoulli*, 19(5A):2033–2066.
-  Li, J. and Barron, A. (1999).  
Mixture density estimation.  
*Advances in neural information processing systems*, 12.
-  Li, L., Liu, Q., Korba, A., Yurochkin, M., and Solomon, J. (2022).  
Sampling with mollified interaction energy descent.  
*arXiv preprint arXiv:2210.13400*.
-  Liu, Q. (2017).  
Stein variational gradient descent as gradient flow.  
In *Advances in neural information processing systems*, pages 3115–3123.

# References III

-  Mei, S., Montanari, A., and Nguyen, P.-M. (2018).  
A mean field view of the landscape of two-layer neural networks.  
*Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671.
-  Mérigot, Q., Santambrogio, F., and Sarazin, C. (2021).  
Non-asymptotic convergence bounds for wasserstein approximation using point clouds.  
*Advances in Neural Information Processing Systems*, 34:12810–12821.
-  Ohta, S.-i. and Takatsu, A. (2011).  
Displacement convexity of generalized relative entropies.  
*Advances in Mathematics*, 228(3):1742–1787.
-  Roberts, G. O. and Tweedie, R. L. (1996).  
Exponential convergence of langevin distributions and their discrete approximations.  
*Bernoulli*, pages 341–363.
-  Vempala, S. and Wibisono, A. (2019).  
Rapid convergence of the unadjusted langevin algorithm: Isoperimetry suffices.  
*Advances in neural information processing systems*, 32.
-  Villani, C. (2009).  
*Optimal transport: old and new*, volume 338.  
Springer.
-  Wibisono, A. (2018).  
Sampling as optimization in the space of measures: The langevin dynamics as a composite optimization problem.  
In *Conference on Learning Theory*, pages 2093–3027. PMLR.

# References IV

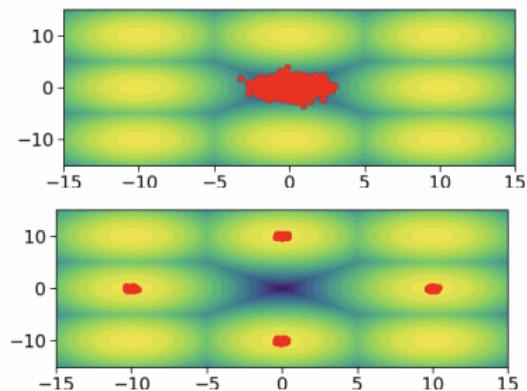


Xu, L., Korba, A., and Slepčev, D. (2022).

Accurate quantization of measures via interacting particle-based optimization.  
*International Conference on Machine Learning*.

# Mixture of Gaussians

Langevin Monte Carlo on a mixture of Gaussians does not manage to target all modes in reasonable time, even in low dimensions.

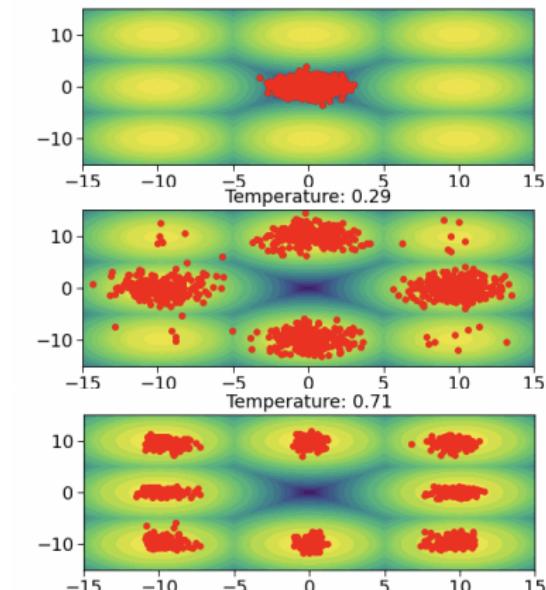


# Annealing

One possible fix : sequence of tempered targets as:

$$\mu_\beta^* \propto \mu_0^\beta (\mu^*)^{1-\beta}, \quad \beta \in [0, 1]$$

Can be seen as **discretized Fisher-Rao gradient flow**  
[Chopin et al., 2023].



## Other tempered path

"Convolutional path" ( $\beta \in [0, +\infty[$ ) frequently used in Diffusion Models

$$\mu_\beta^* = \frac{1}{\sqrt{1-\beta}} \mu_0 \left( \frac{\cdot}{\sqrt{1-\beta}} \right) * \frac{1}{\sqrt{\beta}} \mu^* \left( \frac{\cdot}{\sqrt{\beta}} \right)$$

(vs "geometric path"  $\mu_\beta^* \propto \mu_0^\beta (\mu^*)^{1-\beta}$ )

