

# WASSERSTEIN GRADIENT FLOWS OF MOREAU ENVELOPES OF $f$ -DIVERGENCES IN REPRODUCING KERNEL HILBERT SPACES

joint work with



Sebastian Neumayer, TU Chemnitz



Gabriele Steidl, TU Berlin



Nicolaj Rux, TU Berlin

**Goal.** Recover  $\nu \in \mathcal{P}(\mathbb{R}^d)$  from samples by minimizing  $f$ -divergence  $D_{f,\nu}$  to  $\nu$ , e.g.  $\text{KL}(\cdot | \nu)$ .

**Problem.** Only samples  $\rightsquigarrow$  empirical measures, but

$$\mu \not\ll \nu \implies D_{f,\nu}(\mu) = \infty.$$

weak convergence

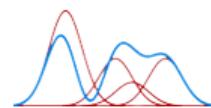
**Our Solution.** Regularize  $D_{f,\nu}: \mathcal{M}(\mathbb{R}^d) \rightarrow [0, \infty]$ .



$$m: \mathcal{M}(\mathbb{R}^d) \rightarrow \mathcal{H}_K, \quad \mu \mapsto \int_{\mathbb{R}^d} K(x, \cdot) d\mu(x)$$

pointwise convergence

$$“D_{f,\nu} \circ m^{-1}” = G_{f,\nu}: \mathcal{H}_K \rightarrow [0, \infty]$$



2. Moreau envelope regularization

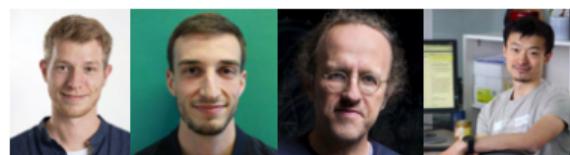
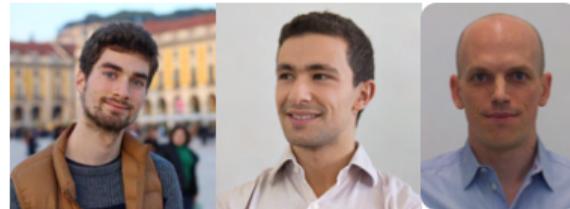
$${}^\lambda G_{f,\nu}(m(\mu)) = \min_{\sigma \in \mathcal{M}_+(\mathbb{R}^d)} D_{f,\nu}(\sigma) + \frac{1}{2\lambda} \|m(\sigma) - m(\mu)\|_{\mathcal{H}_K}^2, \quad \lambda > 0.$$

We prove existence & uniqueness of  $W_2$  gradient flows of  $({}^\lambda G_{f,\nu}) \circ m$ .

Simulate particle flows =  $W_2$  gradient flows starting at empirical measure

## LITERATURE REVIEW OF PRIOR WORK

- KALE functional = MMD-regularized KL divergence  
[Glaser, Arbel, Gretton. NeurIPS'21]  
No Moreau envelope interpretation.
- Kernel methods of moments =  $f$ -divergence-regularized MMD  
[Kremer, Nemmour, Schölkopf, Zhu. ICML'23]  
Doesn't cover all  $f$ -divergences.
- $(f, \Gamma)$ -divergence = Pasch-Hausdorff envelope of  $f$ -divergences.  
[Birrell, Dupuis, Katsoulakis, Pantazis, Rey-Bellet, JMLR'23]  
Yields only Lipschitz, not differentiable functional.
- $W_1$ -Moreau envelope of  $f$ -divergences [Terjék. ICML'21]  
No RKHS, which makes optimization finite-dimensional, hence  
tractable.



1. RKHS & MMD

2. Moreau envelopes

3.  $f$ -divergences

4. MMD-Moreau envelopes  
of  $f$ -divergences

5. Wasserstein gradient flow

6. WGF of MMD-Moreau envelopes of  $f$ -divergences

# REPRODUCING KERNEL HILBERT SPACES

“Kernel trick”: embed data into high-dimensional Hilbert space.

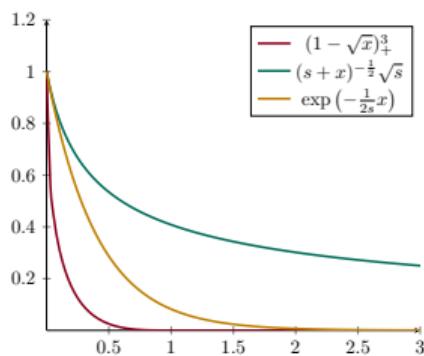
$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  **symmetric, positive definite.**

We consider **radial** kernels  $K(x, y) = \phi(\|x - y\|_2^2)$  with

$\phi \in \mathcal{C}^\infty((0, \infty)) \cap \mathcal{C}^2([0, \infty)), (-1)^k \phi^{(k)}(r) \geq 0, \forall k \in \mathbb{N}, r > 0.$

↪ **reproducing kernel Hilbert space (RKHS)**

$\mathcal{H}_K := \overline{\text{span}}(\{K(x, \cdot) : x \in \mathbb{R}^d\})$ . Key property:  $h \mapsto h(x)$  cts.



Examples (with parameter  $s > 0$ ).

- Gaussian  $\phi(r) = \exp(-\frac{1}{2s}r)$
- inverse multiquadric  $\phi(r) := (s + r)^{-\frac{1}{2}}$
- spline  $\phi(r) = \max(0, (1 - \sqrt{r})^{s+2})$ .

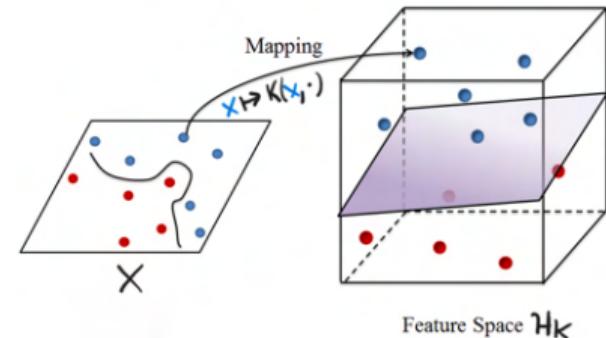


FIG. 1: “Kernel trick”.

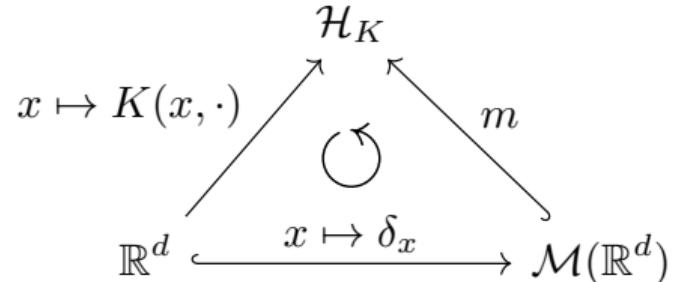
Source: [songy.net/posts/story-of-basis-and-kernels-part-2/](http://songy.net/posts/story-of-basis-and-kernels-part-2/)

Nonexamples.

- Laplace  $\phi(r) = \exp(-\frac{1}{2s}\sqrt{r})$   
(not smooth enough)
- $K(x, y) = \|x\| + \|y\| - \|x - y\|$   
(not radial)

“Kernel trick for signed measures”  $\mu \in \mathcal{M}(\mathbb{R}^d)$  (instead of points): **kernel mean embedding (KME)**

$$m: \mathcal{M}(\mathbb{R}^d) \rightarrow \mathcal{H}_K, \quad \mu \mapsto \int_{\mathbb{R}^d} K(x, \cdot) d\mu(x).$$



We require  $m$  to be injective ( $\mathcal{H}_K$  “characteristic”)  $\iff \mathcal{H}_K \subset \mathcal{C}_0(\mathbb{R}^d)$  dense.

↔ Instead of measures, compare their embeddings in  $\mathcal{H}_K$ : **maximum mean discrepancy (MMD)**

$$d_K: \mathcal{M}(\mathbb{R}^d) \times \mathcal{M}(\mathbb{R}^d) \rightarrow [0, \infty), \quad (\mu, \nu) \mapsto \|m(\mu - \nu)\|_{\mathcal{H}_K}.$$

$m$  injective  $\implies d_K$  is a metric, but  $(\mathcal{M}(\mathbb{R}^d), d_K)$  is not complete.

Easy to evaluate, e.g. for discrete measures since

$$d_K(\mu, \nu)^2 = \int_{\mathbb{R}^d \times \mathbb{R}^d} K(x, y) d(\mu - \nu)(x) d(\mu - \nu)(y) \quad \forall \mu, \nu \in \mathcal{M}(\mathbb{R}^d).$$

1. RKHS & MMD

2. Moreau envelopes

3.  $f$ -divergences

4. MMD-Moreau envelopes  
of  $f$ -divergences

5. Wasserstein gradient flow

6. WGF of MMD-Moreau envelopes of  $f$ -divergences

# REGULARIZATION IN CONVEX ANALYSIS - MOREAU ENVELOPES

Let  $(H, \langle \cdot, \cdot \rangle, \|\cdot\|)$  Hilbert space,  $f \in \Gamma_0(H)$ , i.e.  $f: H \rightarrow (-\infty, \infty]$  convex lower semicontinuous,  $\text{dom}(f) := \{x \in H : f(x) < \infty\} \neq \emptyset$ .

For  $\varepsilon > 0$ , the  **$\varepsilon$ -Moreau envelope** of  $f$ ,

$${}^\varepsilon f: H \rightarrow \mathbb{R}, \quad x \mapsto \min \left\{ f(x') + \frac{1}{2\varepsilon} \|x - x'\|^2 : x' \in H \right\}$$

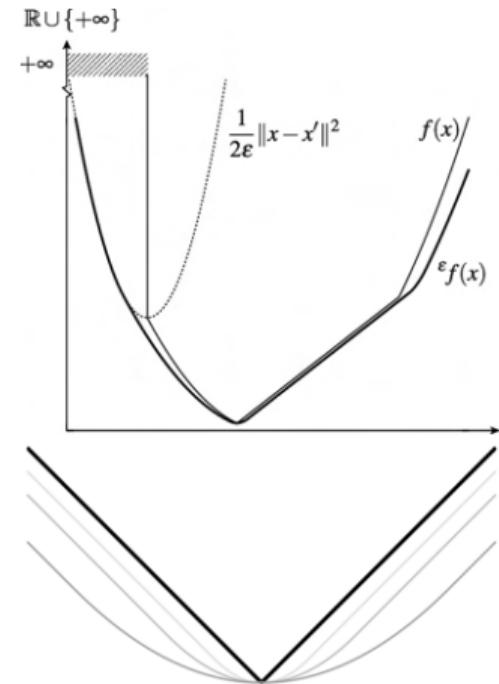
is convex, **differentiable** regularization of  $f$  preserving its **minimizers**.

Asymptotics:  ${}^\varepsilon f(x) \nearrow f(x)$  for  $\varepsilon \searrow 0$  and  ${}^\varepsilon f(x) \searrow \inf(f)$  for  $\varepsilon \rightarrow \infty$ .

$(\varepsilon, x) \mapsto {}^\varepsilon f(x)$  is viscosity solution of Hamilton-Jacobi equation:

$$\begin{cases} \partial_\varepsilon({}^\varepsilon f)(x) + \frac{1}{2} \|\nabla({}^\varepsilon f)(x)\|_2^2 = 0, \\ {}^0 f(x) \rightarrow f(x). \end{cases}$$

[Osher, Heaton, Fung, PNAS 120, **14**, 2023].



Moreau envelope of an extended-real-valued non-differentiable function (top) and of  $|\cdot|$  for different  $\varepsilon$  (bottom).

©Trygve U. Helgaker, Pontus Giselsson

1. RKHS & MMD

2. Moreau envelopes

3.  $f$ -divergences

4. MMD-Moreau envelopes  
of  $f$ -divergences

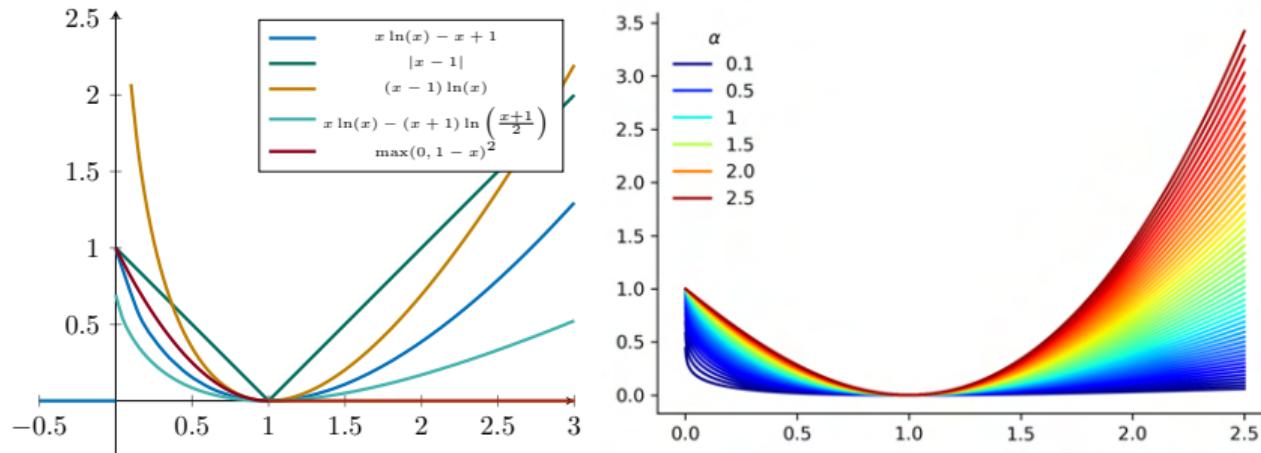
5. Wasserstein gradient flow

6. WGF of MMD-Moreau envelopes of  $f$ -divergences

## ENTROPY FUNCTIONS

We consider  $f \in \Gamma_0(\mathbb{R})$  with  $f|_{(-\infty, 0)} \equiv \infty$  and with **unique minimizer at 1**:  $f(1) = 0$  and positive recession constant  $f'_\infty := \lim_{t \rightarrow \infty} \frac{1}{t} f(t) > 0$ .

*Examples.*  $f_{\text{KL}}(x) := x \ln(x) - x + 1$  for  $x \geq 0$  yields the **Kullback-Leibler divergence** and  $f_\alpha(x) := \frac{1}{\alpha-1} (x^\alpha - \alpha x + \alpha - 1)$  the **Tsallis- $\alpha$  divergence**  $T_\alpha$  for  $\alpha > 0$ . In the limit:  $T_1 = \text{KL}$ .



Left: Examples of entropy functions, except the red. Right: The functions  $f_\alpha$  for  $\alpha \in [0.1, 2.5]$ .

**$f$ -divergence** of  $\mu = \rho\nu + \mu_s \in \mathcal{M}_+(\mathbb{R}^d)$  (unique Lebesgue decomposition) to  $\nu \in \mathcal{M}_+(\mathbb{R}^d)$

$$\begin{aligned} D_{f,\nu}(\rho\nu + \mu_s) &:= \int_{\mathbb{R}^d} f \circ \rho \, d\nu + f'_\infty \cdot \mu_s(\mathbb{R}^d) \quad (\infty \cdot 0 := 0) \\ &= \sup_{h \in \mathcal{C}_b(\mathbb{R}^d; \text{dom}(f^*))} \mathbb{E}_\mu[h] - \mathbb{E}_\nu[f^* \circ h], \quad \mathbb{E}_\sigma[h] := \int_{\mathbb{R}^d} h(x) \, d\sigma(x) \end{aligned}$$

The **convex conjugate** of  $f$  is  $f^*: \mathbb{R} \rightarrow (-\infty, \infty]$ ,  $s \mapsto \sup \{st - f(t) : t \geq 0\}$ .

## THEOREM (PROPERTIES OF $D_{f,\nu}$ )

$D_{f,\nu}: \mathcal{M}_+(\mathbb{R}^d) \rightarrow [0, \infty]$  is convex, weak\* lower semicontinuous. We have:  $D_{f,\nu}(\mu) = 0 \iff \mu = \nu$ .

We define the **MMD-regularized  $f$ -divergence** functional

$$D_{f,\nu}^\lambda(\mu) := \min \left\{ D_{f,\nu}(\sigma) + \frac{1}{2\lambda} d_K(\mu, \sigma)^2 : \sigma \in \mathcal{M}(\mathbb{R}^d) \right\}, \quad \lambda > 0, \mu \in \mathcal{M}(\mathbb{R}^d). \quad (1)$$

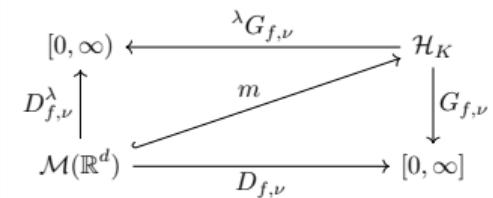
**THEOREM** (MOREAU ENVELOPE INTERPRETATION OF  $D_{f,\nu}^\lambda$  [NSSR24])

The  $\mathcal{H}_K$ -extension of  $D_{f,\nu}$ ,

$$G_{f,\nu}: \mathcal{H}_K \rightarrow [0, \infty], \quad h \mapsto \begin{cases} D_{f,\nu}(\mu), & \text{if } \exists \mu \in \mathcal{M}_+(\mathbb{R}^d) \text{ s.t. } h = m(\mu), \\ \infty, & \text{else.} \end{cases}$$

is convex, **lower semicontinuous** and its Moreau envelope concatenated with  $m$  is the MMD-regularized  $f$ -divergence:

$${}^\lambda G_{f,\nu} \circ m = D_{f,\nu}^\lambda$$



- Dual formulation

$$D_{f,\nu}^\lambda(\mu) = \max \left\{ \mathbb{E}_\mu[p] - \mathbb{E}_\nu[f^* \circ p] - \frac{\lambda}{2} \|p\|_{\mathcal{H}_K}^2 : p \in \mathcal{H}_K, p \leq f'_\infty \right\}. \quad (2)$$

$\hat{p} \in \mathcal{H}_K$  maximizes (2)  $\iff \hat{g} = m(\mu) - \lambda \hat{p}$  is primal solution.

$$\frac{\lambda}{2} \|\hat{p}\|_{\mathcal{H}_K}^2 \leq D_{f,\nu}^\lambda(\mu) \leq \|\hat{p}\|_{\mathcal{H}_K} (\|m_\mu\|_{\mathcal{H}_K} + \|m_\nu\|_{\mathcal{H}_K}) \quad \text{and} \quad \|\hat{p}\|_{\mathcal{H}_K} \leq \frac{2}{\lambda} d_K(\mu, \nu).$$

- $D_{f,\nu}^\lambda$  is Fréchet differentiable on  $\mathcal{M}(\mathbb{R}^d)$  and its gradient is  $\lambda$ -Lipschitz with respect to  $d_K$ :

$$\nabla D_{f,\nu}^\lambda(\mu) = \operatorname{argmax} (2).$$

# THEOREM. (PROPERTIES OF $D_{f,\nu}^\lambda$ ) [NSSR24]

- **Asymptotic regimes:** Mosco resp. pointwise convergence (if  $0 \in \text{int}(\text{dom}(f^*))$  resp.  $f^*$  differentiable in 0)

$$D_{f,\nu}^\lambda \rightarrow D_{f,\nu} \quad \lambda \searrow 0 \quad \text{and} \quad (1 + \lambda)D_{f,\nu}^\lambda \rightarrow \frac{1}{2}d_K(\cdot, \nu)^2 \quad \lambda \rightarrow \infty$$

	MMD metric $d_K$	$f$ -divergence $D_f$
ingredients	characteristic kernel $K$	convex, lsc $f$ , $f(1) = 0$
examples	Gaussian, IMQ, Matérn	KL, JSD, Hellinger, $\chi^2$
positive definite	👍	👍
symmetric	👍	👎
triangle inequality	👍	👎
variational formulation	👍	👍
topology on measures	🚫	💪
geometry on measures	😐	😊

- **Divergence property:**  $D_{f,\nu}^\lambda(\mu) = 0 \iff \mu = \nu$ .
- If  $f^*$  is differentiable in 0, then  $(\mu, \nu) \mapsto D_{f,\nu}^\lambda(\mu)$  **metrizes weak convergence** on  $\mathcal{M}_+(\mathbb{R}^d)$ -balls.

1. RKHS & MMD

2. Moreau envelopes

3.  $f$ -divergences

4. MMD-Moreau envelopes  
of  $f$ -divergences

5. Wasserstein gradient flow

6. WGF of MMD-Moreau envelopes of  $f$ -divergences

$\mathcal{P}_2(\mathbb{R}^d) := \{\mu \in \mathcal{P}(\mathbb{R}^d) : \int_{\mathbb{R}^d} \|x\|_2^2 < \infty\}, \|\cdot\|_2$  Eucl. norm.

$$W_2(\mu, \nu)^2 = \min_{\pi \in \Gamma(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|_2^2 d\pi(x, y), \quad \mu, \nu \in \mathcal{P}_2(\mathbb{R}^d).$$

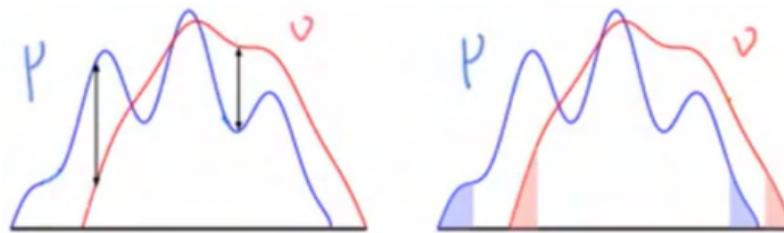


FIG. 2: Vertical ( $L_2$ ) vs. horizontal ( $W_2$ ) mass displacement.

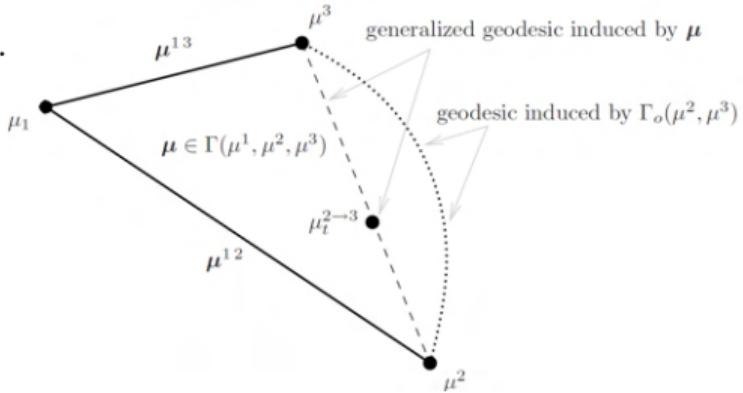


FIG. 3: Generalized geodesic from  $\mu_2$  to  $\mu_3$  with base  $\mu_1$  [AGS08].

## DEFINITION (GENERALIZED GEODESIC CONVEXITY)

A function  $\mathcal{F}: \mathcal{P}_2(\mathbb{R}^d) \rightarrow (-\infty, \infty]$  is *M-convex along generalized geodesics* if, for every  $\sigma, \mu, \nu \in \text{dom}(\mathcal{F})$ , there exists a  $\alpha \in \mathcal{P}_2(\mathbb{R}^{3d})$  with  $(P_{1,2})_\# \alpha \in \Gamma^{\text{opt}}(\sigma, \mu)$  and  $(P_{1,3})_\# \alpha \in \Gamma^{\text{opt}}(\sigma, \nu)$  such that

$$\mathcal{F}\left(\left((1-t)P_2 + tP_3\right)_\# \alpha\right) \leq (1-t)\mathcal{F}(\mu) + t\mathcal{F}(\nu) - \frac{M}{2}t(1-t) \int_{\mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^d} \|y - z\|_2^2 d\alpha(x, y, z), \quad \forall t \in [0, 1].$$

## DEFINITION (FRÉCHET SUBDIFFERENTIAL IN WASSERSTEIN SPACE)

The (reduced) Fréchet subdifferential of  $\mathcal{F}: \mathcal{P}_2(\mathbb{R}^d) \rightarrow (-\infty, \infty]$  at  $\mu \in \text{dom}(\mathcal{F})$  is

$$\partial \mathcal{F}(\mu) := \left\{ \xi \in L^2(\mathbb{R}^d; \mu) : \mathcal{F}(\nu) - \mathcal{F}(\mu) \geq \inf_{\pi \in \Gamma^{\text{opt}}(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \langle \xi(x_1), x_2 - x_1 \rangle d\pi(x, y) + o(W_2(\mu, \nu)) \right\}$$

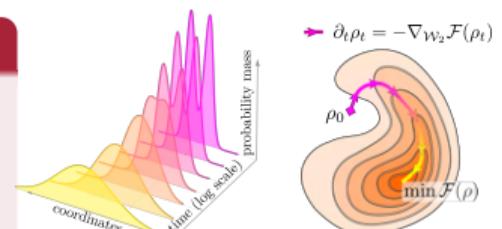
A curve  $\gamma: (0, \infty) \rightarrow \mathcal{P}_2(\mathbb{R}^d)$  is *absolutely continuous* if  $\exists$   $L^2$ -Borel velocity field  $v: \mathbb{R}^d \times (0, \infty) \rightarrow \mathbb{R}^d$  s.t.

$$\partial_t \gamma_t + \nabla \cdot (v_t \gamma_t) = 0, \quad (t, x) \in (0, \infty) \times \mathbb{R}^d, \text{ weakly.} \quad (\text{Continuity Eq.})$$

## DEFINITION (WASSERSTEIN GRADIENT FLOW)

A locally absolutely continuous curve  $\gamma: (0, \infty) \rightarrow \mathcal{P}_2(\mathbb{R}^d)$  with velocity field  $v_t \in T_{\gamma_t} \mathcal{P}_2(\mathbb{R}^d)$  is a *Wasserstein gradient flow with respect to*  $\mathcal{F}: \mathcal{P}_2(\mathbb{R}^d) \rightarrow (-\infty, \infty]$  if

$$v_t \in -\partial \mathcal{F}(\gamma_t), \quad \text{for a.e. } t > 0.$$



# WASSERSTEIN GRADIENT FLOW WITH RESPECT TO $D_{f,\nu}^\lambda$

## THEOREM (CONVEXITY AND GRADIENT OF $D_{f,\nu}^\lambda$ [NSSR24])

Since  $K$  being radial and smooth,  $D_{f,\nu}^\lambda$  is  **$M$ -convex along generalized geodesics** with  $M := -8\lambda^{-1}\sqrt{(d+2)\phi''(0)\phi(0)}$  and its (reduced) Fréchet **subdifferential** is  $\partial D_{f,\nu}^\lambda(\mu) = \{\nabla \operatorname{argmax}(2)\}$ .

*Remark.*  $M$  seems non-optimal, since for  $\lambda \rightarrow 0$ ,  $D_{f,\nu}^\lambda \rightarrow D_{f,\nu}$  and  $D_{f,\nu}$  is 0-convex, but  $M \rightarrow -\infty$ .

## COROLLARY

There **exists a unique Wasserstein gradient flow**  $(\gamma_t)_{t>0}$  of  $D_{f,\nu}^\lambda$  starting at  $\mu_0 \in \mathcal{P}_2(\mathbb{R}^d)$ , fulfilling the continuity equation  $\partial_t \gamma_t = \nabla \cdot (\gamma_t(\partial D_{f,\nu}^\lambda(\gamma_t)))$ ,  $\gamma_0 = \mu_0$ .

## LEMMA (PARTICLE FLOWS ARE $W_2$ GRADIENT FLOWS)

If  $\mu_0$  is empirical, then so is  $\gamma_t$  for all  $t > 0$ .

## NUMERICAL EXPERIMENTS - PARTICLE DESCENT ALGORITHM

Take i.i.d. samples  $(x_j^{(0)})_{j=1}^N \sim \mu_0$  and  $(y_j)_{j=1}^M \sim \nu$ . Forward Euler discretization in time with step size  $\tau > 0$  yields

$$\gamma_{n+1} := (\text{id} - \tau \nabla \hat{p}_n)_\# \gamma_n, \quad \hat{p}_n = \text{argmax in } D_{f,\nu}^\lambda(\gamma_n)$$

so  $(\gamma_n)_{n \in \mathbb{N}} = \frac{1}{N} \sum_{j=1}^N \delta_{x_j^{(n)}}$  with gradient step

$$x_j^{(n+1)} = x_j^{(n)} - \tau \nabla \hat{p}_n(x_j^{(n)}), \quad j \in \{1, \dots, N\}, n \in \mathbb{N}.$$

### THEOREM (REPRESENTER-TYPE THEOREM [NSSR24])

If  $f'_\infty = \infty$  or if  $\lambda > 2d_K(\gamma_n, \nu) \sqrt{\phi(0)} \frac{1}{f'_\infty}$ , then finding  $\hat{p}_n$  is a **finite-dimensional strongly convex** problem.

To find  $\hat{p}_n$ , we use **L-BFGS-B**, a quasi-Newton method. We use annealing strategy for  $\lambda$  if  $f'_\infty < \infty$ .

## NUMERICAL EXPERIMENTS

FIG. 4: IMQ kernel,  $\lambda = \frac{1}{100}$ ,  $\tau = \frac{1}{1000}$ , Top: Tsallis-3 divergence, Bottom: Tsallis- $\frac{1}{2}$  divergence, with annealing.

FIG. 5: Number of starting particles  $N$ , less than number of samples of target,  $M \rightsquigarrow$  quantization

- **Non-differentiable** (e.g. Laplace =  $\frac{1}{2}$ -Matérn) and unbounded (e.g. Riesz, Coulomb) kernels.
- **Convergence rates** in suitable metric.
- Prove consistency bounds [Leclerc, Mérigot, Santambrogio, Stra. 2020] and **better  $M$ -convexity estimates**.
- Convergence for annealing strategy?
- Different domains, e.g. compact subsets of  $\mathbb{R}^d$  (manifolds like sphere, torus), groups, infinite-dimensional spaces.
- Regularize other divergences, e.g. Rényi divergences, Bregman divergences.
- Gradient flow of  $D_{f,\nu}^\lambda$  with respect to other metrics, like Kantorovich-Hellinger (related to unbalanced OT), MMD, Fisher-Rao or Wasserstein- $p$  for  $p \in [1, \infty]$ .
- More elaborate time discretizations, variable step sizes.

- We created novel objective. Minimizing it allows sampling from a target measure of which only samples are known.
- Clear, rigorous interpretation using Convex Analysis and RKHS.
- Theory covers (almost) all  $f$ -divergences.
- Best of both worlds:  $D_{f,\nu}^\lambda$  interpolates between  $D_{f,\nu}$  and  $d_K(\cdot, \nu)^2$ .
- Effective algorithms due to (modified) representer theorem & GPU / PyTorch.

Thank you for your attention!

I am happy to take any questions.

Paper link: [arxiv.org/abs/2402.04613](https://arxiv.org/abs/2402.04613)

My website: [viktorajstein.github.io](https://viktorajstein.github.io)

## REFERENCES I

- [AGS08] Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré, *Gradient flows: in metric spaces and in the space of probability measures*, 2 ed., Springer Science & Business Media, 2008.
- [BDK<sup>+</sup>22] Jeremiah Birrell, Paul Dupuis, Markos A. Katsoulakis, Yannis Pantazis, and Luc Rey-Bellet, *( $f, \Gamma$ )-divergences: Interpolating between  $f$ -divergences and integral probability metrics*, J. Mach. Learn. Res. **23** (2022), no. 39, 1–70.
- [GAG21] Pierre Glaser, Michael Arbel, and Arthur Gretton, *KALE flow: A relaxed KL gradient flow for probabilities with disjoint support*, Advances in Neural Information Processing Systems (Virtual event), vol. 34, 6–14 Dec 2021, pp. 8018–8031.
- [HWAH24] J. Hertrich, C. Wald, F. Altekrüger, and P. Hagemann, *Generative sliced MMD flows with Riesz kernels*, International Conference on Learning Representations (ICLR) (Vienna, Austria), 7 – 11 May 2024.

## REFERENCES II

- [KYSZ23] H. Kremer, Nemmour Y., B. Schölkopf, and J.-J. Zhu, *Estimation beyond data reweighting: kernel methods of moments*, ICML'23: Proceedings of the 40th International Conference on Machine Learning (Honolulu, Hawaii, USA), vol. 202, July 23 - 29 2023, p. 17745–17783.
- [LMS17] Matthias Liero, Alexander Mielke, and Giuseppe Savaré, *Optimal entropy-transport problems and a new Hellinger–Kantorovich distance between positive measures*, Invent. Math. **211** (2017), no. 3, 969–1117.
- [LMSS20] Hugo Leclerc, Quentin Mérigot, Filippo Santambrogio, and Federico Stra, *Lagrangian discretization of crowd motion and linear diffusion*, SIAM J. Numer. Anal. **58** (2020), no. 4, 2093–2118. MR 4123686
- [Ter21] Dávid Terjék, *Moreau-Yosida  $f$ -divergences*, International Conference on Machine Learning (ICML) (Virtual event), PMLR, Jul 18–24 2021, pp. 10214–10224.

## Interpolating between OT and KL regularized OT using Rényi Divergences

Rényi divergence  $\notin \{f\text{-div.}, \text{Bregman div.}\}$ ,  $\alpha \in (0, 1)$

$$R_\alpha(\mu \mid \nu) := \frac{1}{\alpha - 1} \ln \left[ \int_X \left( \frac{d\mu}{d\tau} \right)^\alpha \left( \frac{d\nu}{d\tau} \right)^{1-\alpha} d\tau \right],$$

$$\text{OT}_{\varepsilon, \alpha}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \langle c, \pi \rangle + \varepsilon R_\alpha(\pi \mid \mu \otimes \nu)$$

is a metric, where  $\varepsilon > 0$ ,  $\mu, \nu \in \mathcal{P}(X)$ ,  $X$  compact.

$$\text{OT}(\mu, \nu) \xleftarrow[\text{or } \varepsilon \rightarrow 0]{\alpha \searrow 0} \text{OT}_{\varepsilon, \alpha}(\mu, \nu) \xrightarrow{\alpha \nearrow 1} \text{OT}_\varepsilon^{\text{KL}}(\mu, \nu).$$

In the works: **debiased** Rényi-Sinkhorn divergence

$$\text{OT}_{\varepsilon, \alpha}(\mu, \nu) - \frac{1}{2} \text{OT}_{\varepsilon, \alpha}(\mu, \mu) - \frac{1}{2} \text{OT}_{\varepsilon, \alpha}(\nu, \nu).$$

$W_2$  gradient flows of  $d_K(\cdot, \nu)^2$  with  
 $K(x, y) := -|x - y|$  in 1D.

Reformulation as **maximal monotone** inclusion Cauchy problem in  $L_2(0, 1)$  via **quantile functions**.

Comprehensive description of solutions' behavior, **instantaneous measure-to- $L^\infty$  regularization**, implicit Euler is simple.