

Wasserstein Gradient Flows of Moreau Envelopes of f -Divergences in Reproducing Kernel Hilbert Spaces

Sebastian Neumayer, Viktor Stein, Gabriele Steidl, Nicolaj Rux

Goal: minimize **f -divergence loss** $D_{f,\nu}$ with target measure $\nu \in \mathcal{M}_+(\mathbb{R}^d)$ (e.g. generative adversarial networks, variational inference).

Often **only samples** are available \rightsquigarrow empirical measures.

BUT: D_f between empirical measures is $\infty \rightsquigarrow$ **regularize f -divergence**.

- **Contribution.** Prove identification of MMD-regularized f -divergence functional as Moreau envelope in RKHS. Existence and uniqueness of its Wasserstein gradient flow. Flow starting at empirical measure is particle flow.
- **Prior work.** Regularize MMD with f -divergence [5], MMD-Pasch-Hausdorff envelope of f -divergences [7], W_1 -Moreau envelope of f -divergences [8].
- **Method.** Euler forward discretize particle flow (= gradient descent on the positions).
- **Result.** We can simulate particle flows for divergences with finite and infinite recession constant f'_∞ . Tsallis- α divergence with moderately large α outperforms KL-divergence ($\alpha = 1$): faster target recovery and more stable.

Reproducing Kernel Hilbert Space, KME, Maximum Mean Discrepancy

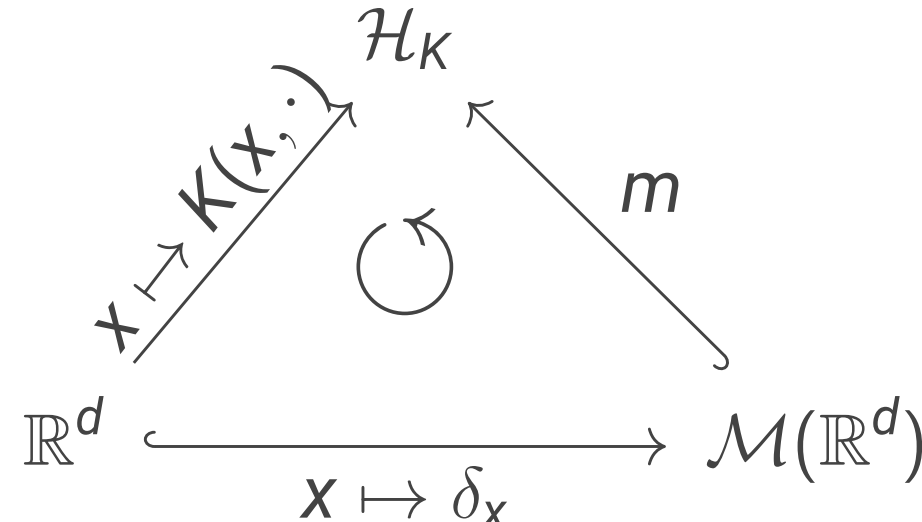
$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ **symmetric, positive definite, bounded** kernel with $K(x, \cdot) \in \mathcal{C}_0(\mathbb{R}^d)$. We focus on **radial** kernels $K(x, y) = \phi(\|x - y\|_2^2)$ with $\phi \in \mathcal{C}^2(\mathbb{R}^d)$ completely monotone.

Examples. Gaussian $\phi(r) = \exp(-\frac{1}{2s}r)$, IMQ $\phi(r) := (s + r)^{-\frac{1}{2}}$, spline $\phi(r) = (1 - \sqrt{r})_+^{q+2}$.

\rightsquigarrow **reproducing kernel Hilbert space (RKHS)** $\mathcal{H}_K := \overline{\text{span}}(\{K(x, \cdot) : x \in \mathbb{R}^d\})$.

The **kernel mean embedding (KME)** of finite signed measures, $\mathcal{M}(\mathbb{R}^d)$, into \mathcal{H}_K is

$$m: \mathcal{M}(\mathbb{R}^d) \rightarrow \mathcal{H}_K, \quad \mu \mapsto m_\mu := \int_{\mathbb{R}^d} K(x, \cdot) d\mu(x). \quad (1)$$



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

Then the **maximum mean discrepancy (MMD)**

$$d_K: \mathcal{M}(\mathbb{R}^d) \times \mathcal{M}(\mathbb{R}^d) \rightarrow [0, \infty), \quad (\mu, \nu) \mapsto \|m_\mu - m_\nu\|_{\mathcal{H}_K}. \quad (2)$$

is an **incomplete** metric. We have for all $\mu, \nu \in \mathcal{M}(\mathbb{R}^d)$

$$d_K(\mu, \nu)^2 = \int_{\mathbb{R}^d \times \mathbb{R}^d} K(x, y) d(\mu - \nu)(x) d(\mu - \nu)(y). \quad (3)$$

Regularization in Convex Analysis

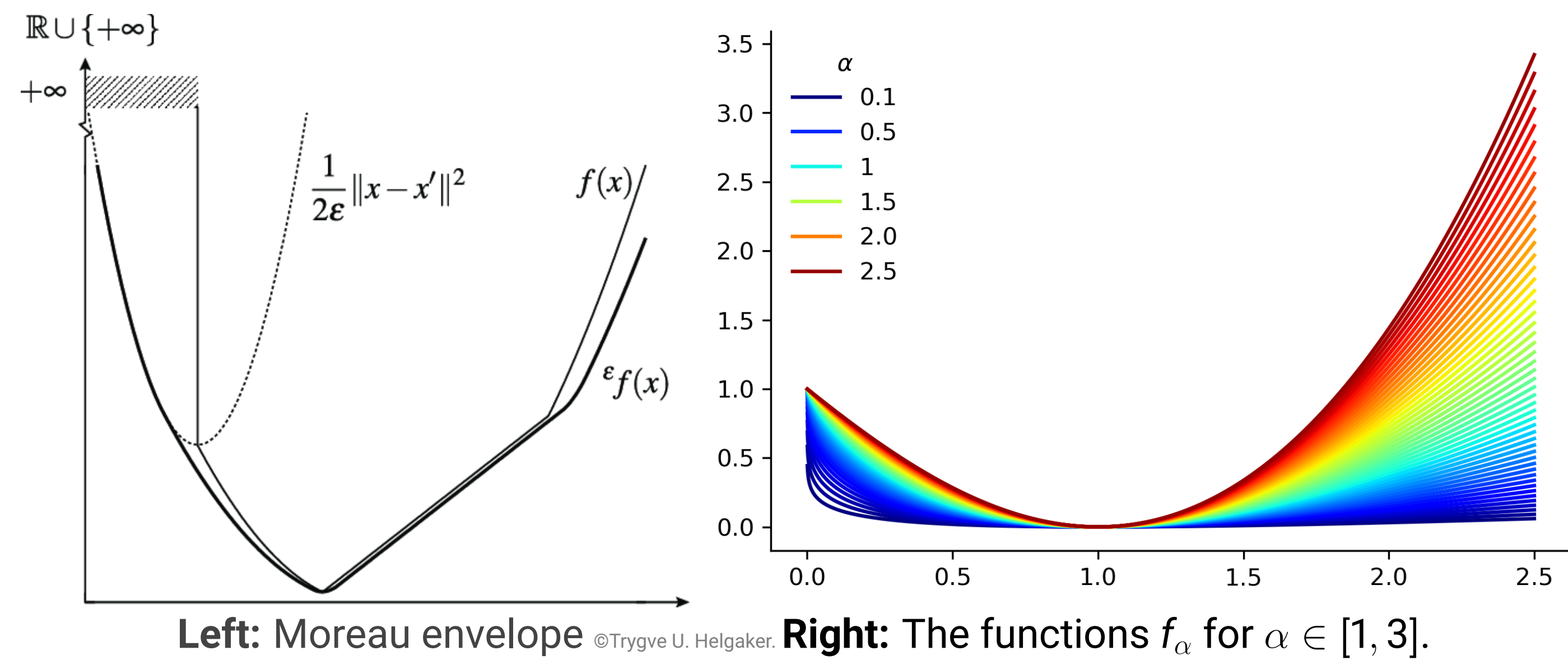
$(H, \langle \cdot, \cdot \rangle, \|\cdot\|)$ Hilbert space, $f: H \rightarrow (-\infty, \infty]$ **convex** lower semicontinuous (we write $f \in \Gamma_0(H)$) with $\text{dom}(f) := \{x \in H : f(x) < \infty\} \neq \emptyset$.

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

$${}^\varepsilon f: H \rightarrow \mathbb{R}, \quad x \mapsto \min \left\{ f(x') + \frac{1}{2\varepsilon} \|x - x'\|^2 : x' \in H \right\} \quad (4)$$

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

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



f -divergence

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

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

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

$$D_{f,\nu}(\rho\nu + \mu_s) := \int_{\mathbb{R}^d} f \circ \rho d\nu + f'_\infty \cdot \mu_s(\mathbb{R}^d) \quad (\infty \cdot 0 := 0) \quad (5)$$

$$= \sup_{h \in \mathcal{C}_b(\mathbb{R}^d, \text{dom}(f^*))} \mathbb{E}_\mu[h] - \mathbb{E}_\nu[f^* \circ h], \quad \mathbb{E}_\mu[h] := \int_{\mathbb{R}^d} h(x) d\mu(x) \quad (6)$$

$D_{f,\nu}$ is convex and weak* lower semicontinuous.

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- α divergence** T_α for $\alpha > 0$. We have $T_1 = \text{KL}$.

MMD-Regularized f -divergence

The **MMD-regularized f -divergence** functional is

$$D_{f,\nu}^\lambda(\mu) := \min \left\{ D_{f,\nu}(\sigma) + \frac{1}{2\lambda} d_K(\mu, \sigma)^2 : \sigma \in \mathcal{M}(\mathbb{R}^d) \right\}, \quad \mu \in \mathcal{M}(\mathbb{R}^d). \quad (7)$$

Generalizes the KALE-functional [4], which is recovered for $f = f_{\text{KL}}$.

Theorem. (Moreau envelope interpretation)

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

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

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

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

Theorem. (Properties of $D_{f,\nu}^\lambda$)

1. Dual formulation

$$D_{f,\nu}^\lambda(\mu) = \max \left\{ \mathbb{E}_\mu[h] - \mathbb{E}_\nu[f^* \circ h] - \frac{\lambda}{2} \|h\|_{\mathcal{H}_K}^2 : h \in \mathcal{H}_K, h \leq f'_\infty \right\}. \quad (8)$$

2. $D_{f,\nu}^\lambda$ is Fréchet diff'able with $\frac{1}{\lambda}$ -Lipschitz gradient with respect to d_K :

$$\nabla D_{f,\nu}^\lambda(\mu) = \arg\max(8).$$

3. Asymptotic regimes: Mosco resp. pointwise convergence (if $0 \in \text{int}(\text{dom}(f^*))$) resp. f^* diff'able in 0)

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

4. Divergence property: $D_{f,\nu}^\lambda(\mu) = 0 \iff \mu = \nu$.

5. If f^* is diff'able in 0, then $(\mu, \nu) \mapsto D_{f,\nu}^\lambda$ metrizes weak convergence on $\mathcal{M}_+(\mathbb{R}^d)$ -balls.

Wasserstein Gradient Flow with respect to $D_{f,\nu}^\lambda$

$D_{f,\nu}^\lambda$ is $(-M)$ -convex along generalized geodesics with $M := \frac{8}{\lambda} \sqrt{(d+2)\phi''(0)\phi(0)}$.

strong Fréchet **subdifferential**: $\partial D_{f,\nu}^\lambda(\mu) = \{\nabla \arg\max(8)\}$.

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 $\partial_t \gamma_t = \text{div}(\gamma_t \nabla(\partial D_{f,\nu}^\lambda(\gamma_t)))$, $\gamma_0 = \mu_0$.

If μ_0 is empirical, then so is μ_t for all $t > 0$ (particle flows are W_2 gradient flows).

Numerical Experiments - Particle Descent Algorithm

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

$$x_j^{(n+1)} = x_j^{(n)} - \tau \nabla \hat{p}_n(x_j^{(n)}), \quad \hat{p}_n = \arg\max \text{ in } D_{f,\nu}^\lambda(\mu_n) \quad j \in \{1, \dots, N\}, n \in \mathbb{N}.$$

Representer-type theorem. If $f'_\infty = \infty$ or if $\lambda > 2d_K(\mu_n, \nu) \sqrt{\phi(0)\frac{1}{f'_\infty}}$, then finding \hat{p}_n is a **finite-dimensional strongly convex** problem (we solve it with **L-BFGS-B**).

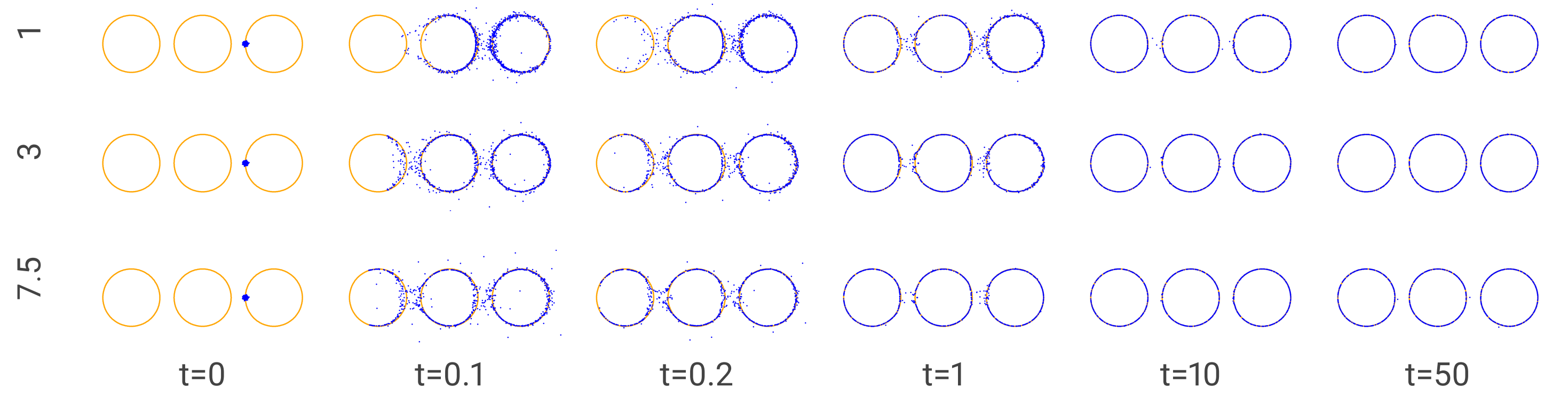


Figure 1. Wasserstein gradient flow of the regularized Tsallis- α divergence $D_{f,\nu}^\lambda$ for $\alpha \in \{1, 3, 7.5\}$, where ν are the three rings. Code: https://github.com/ViktorAJStein/Regularized_f_Divergence_Particle_Flows

Further work. Non-differentiable (e.g. Laplace = $\frac{1}{2}$ -Matérn) and unbounded (e.g. Riesz, Coulomb) kernels. Other divergences, e.g. Rényi. Different time discretizations. Prove consistency bounds [1] and convergence rates.

- [1] H. Leclerc, Q. Mérigot, F. Santambrogio and F. Stra (2020) Lagrangian discretization of crowd motion and linear diffusion. SIAM J. Numer. Anal. **58** (4).
- [2] L. Ambrosio, N. Gigli and G. Savaré (2008). Gradient flows in metric spaces and in the space of probability measures, 2nd edition. Lectures in Mathematics ETH Zürich. Birkhäuser.
- [3] J. Hertrich, C. Wald, F. Altekürger and P. Hagemann (2024). Generative sliced MMD flows with Riesz kernels. ICLR'24.
- [4] P. Glaser, M. Arbel and A. Gretton (2021) KALE flow: A relaxed KL gradient flow for probabilities with disjoint support. NeurIPS'21.
- [5] H. Kremer, Y. Nimmour, B. Schölkopf and J.-J. Zhu (2023) Estimation beyond data reweighting: kernel methods of moments. ICML'23.
- [6] M. Liero, A. Mielke and G. Savaré (2017) Optimal entropy-transport problems and a new Hellinger–Kantorovich distance between positive measures. Invent. Math. 211 3.
- [7] J. Birrell et. al. (2022) (f, Γ) -Divergences: Interpolating between f -Divergences and Integral Probability Metrics. J. Mach. Learn. Res. 23 **39**.
- [8] D. Terjék (2021) Moreau-Yosida f -divergences. ICML'21.