

WASSERSTEIN GRADIENT FLOWS OF MMD FUNCTIONALS WITH DISTANCE KERNEL AND CAUCHY PROBLEMS ON QUANTILE FUNCTIONS

joint work with



Richard Duong, TU Berlin

Robert Beinert, TU Berlin

Johannes Hertrich, UCL

Gabriele Steidl, TU Berlin

University of South Carolina, Columbia, 30.08.2024.

Joint ACM and RTG data science seminar (Changhui Tan, Siming He, Wuchen Li).

Problem: Target measure $\nu \in \mathcal{P}_2(\mathbb{R}^d)$ is unknown, only **samples** are available.

Goal: Recover ν .

Solution: minimize the metric **kernel discrepancy**

$$\mathcal{F}_\nu := \text{MMD}_K(\cdot, \nu)^2: \mathcal{P}_2(\mathbb{R}^d) \rightarrow [0, \infty)$$

(which can be estimated using samples) to ν by finding curve of measures $(\gamma_t)_{t>0} \subset \mathcal{P}_2(\mathbb{R}^d)$ along which \mathcal{F}_ν decreases “the fastest”.

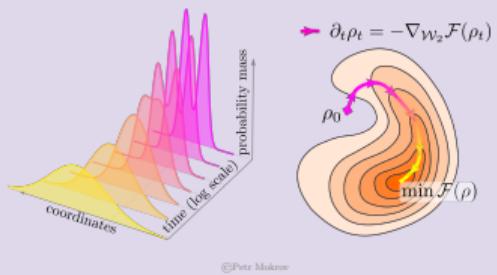
©Francis Bach

One way to construct $(\gamma_t)_{t>0}$: **Wasserstein gradient flows**.

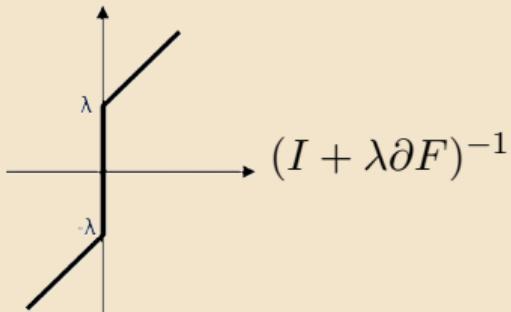
Different kernels $K \rightsquigarrow$ very different behavior of $(\gamma_t)_{t>0}$.

Aim of this preprint: study behavior of $(\gamma_t)_{t>0}$ for irregular kernel $K(x, y) := -|x - y|$ for $d = 1$.

1. Optimal transport and Wasserstein gradient flow



2. Maximal monotone inclusions on Hilbert spaces

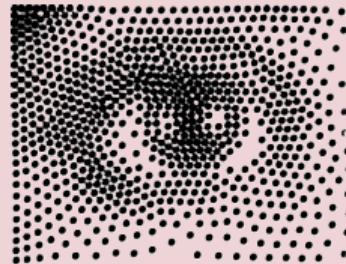


4. Wasserstein gradient flow of the MMD with negative distance kernel

5. Invariant subsets & smoothing properties

$$\gamma_0 = \delta_x,$$
$$\gamma_t \sim \mathcal{U}[a_t, b_t], \quad t > 0.$$

3. Negative distance kernel, Max. Mean Discrepancy



https://www.mis.mpg.de/research/DP_Boltzmann_detail

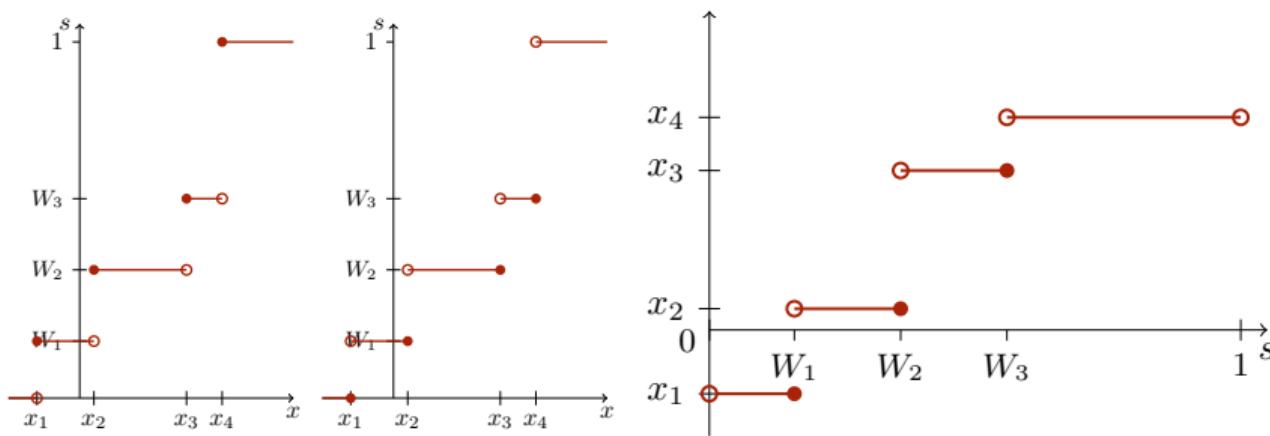
6. Numerical results

The **cumulative distribution functions** (CDFs) of $\nu \in \mathcal{P}(\mathbb{R})$ are given by

$$R_\nu^+ : \mathbb{R} \rightarrow [0, 1], \quad R_\nu^+(x) := \nu((-\infty, x]), \quad R_\nu^- : \mathbb{R} \rightarrow [0, 1], \quad R_\nu^-(x) := \nu((-\infty, x)).$$

and its (generalized inverse, the) **quantile function** by

$$Q_\nu : (0, 1) \rightarrow \mathbb{R}, \quad Q_\nu(s) := \min\{x \in \mathbb{R} : R_\nu^+(x) \geq s\}.$$



Left to right: R_ν^+, R_ν^-, Q_ν , for $\nu = \sum_{k=1}^4 w_k \delta_{x_k}$, $W_j := \sum_{k=1}^j w_k$.

WASSERSTEIN-2 SPACE ON THE REAL LINE

Consider the subset of probability measures $\mathcal{P}(\mathbb{R})$,

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

On $\mathcal{P}_2(\mathbb{R})$, the **Wasserstein-2 metric** is

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

$$\Gamma(\mu, \nu) := \{ \pi \in \mathcal{P}(\mathbb{R} \times \mathbb{R}) : (P_1)_\# \pi = \mu, (P_2)_\# \pi = \nu \}$$

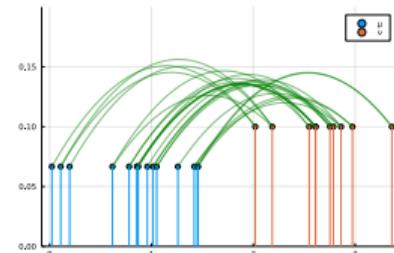
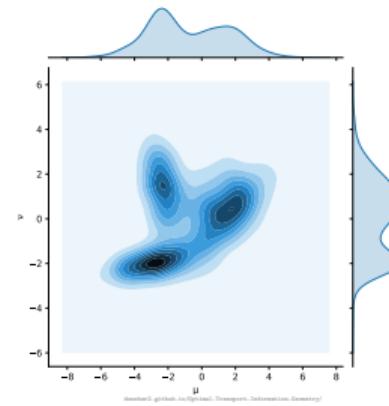
with the **projections** $P_i(x_1, x_2) := x_i$.

The **push-forward** $\#$ acts as $(f_\# \sigma)(A) := \sigma(f^{-1}(A))$ for all measurable sets $A \subset \mathbb{R}$.

The optimal plan $\hat{\pi}$ is unique (in 1D).

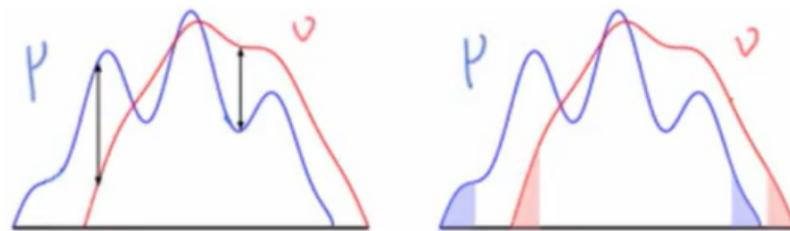
$$\begin{aligned} f &\text{--->} \bullet \\ \text{--->} \bullet &= \sum_i \delta_{x_i} \\ \mu &= \sum_i \delta_{f(x_i)} \\ f_\# \mu &\stackrel{\text{def}}{=} \sum_i \delta_{f(x_i)} \end{aligned}$$

Push-forward $\hat{\pi}$ realizes monotone rearrangement.



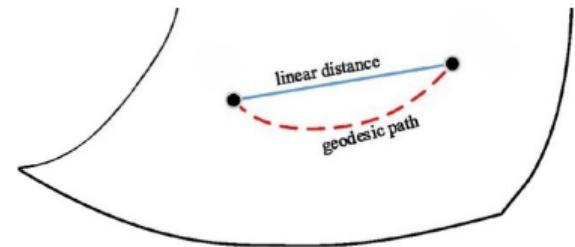
Geodesic (= shortest constant-speed path) between μ and ν is

$$\gamma_t := ((1-t)P_1 + tP_2)_\# \hat{\pi}, \quad t \in [0, 1], \quad (1)$$



Vertical (L_2 , $g_t := (1-t)f_\mu + tf_\nu$) vs. horizontal (W_2) mass displacement.

©Anne Krämer



Geodesic on a manifold.

arXiv.org/abs/2006.04600

DEFINITION (W_2 -GEODESIC CONVEXITY)

$\mathcal{F}: \mathcal{P}_2(\mathbb{R}) \rightarrow \mathbb{R}$ is **convex along geodesics** if

$$\mathcal{F}(\gamma_t) \leq (1-t)\mathcal{F}(\mu) + t\mathcal{F}(\nu), \quad \forall \mu, \nu \in \mathcal{P}_2(\mathbb{R}), \quad \gamma_t \text{ from (1).}$$

W_2 geodesic

©Léonard Claes

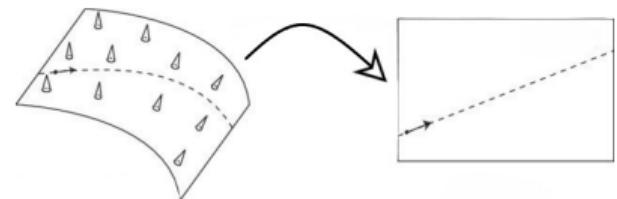
THEOREM (ISOMETRIC EMBEDDING OF $(\mathcal{P}_2(\mathbb{R}), W_2)$)

Let $\mathcal{C}(0, 1) := \{Q_\mu : \mu \in \mathcal{P}_2(\mathbb{R})\} \subset L_2(0, 1)$ be the set of quantile functions.

The map

$$\mathcal{P}_2(\mathbb{R}) \rightarrow \mathcal{C}(0, 1), \quad \mu \mapsto Q_\mu$$

is an **isometric isomorphism** with inverse $Q \mapsto Q \# \Lambda_{(0,1)}$, where Λ is the Lebesgue measure.



Isometric embedding.

<https://doi.org/10.1109/TC.2008.2211>

KEY IDEA: instead of working with $\mathcal{F}: \mathcal{P}_2(\mathbb{R}) \rightarrow (-\infty, \infty]$, find

$$F: L_2(0, 1) \rightarrow (-\infty, \infty] \quad \text{satisfying} \quad F(Q_\mu) = \mathcal{F}(\mu).$$

Given: Hilbert space H , energy functional $f: H \rightarrow (-\infty, \infty]$ (convex, bounded below).

Idea: minimize $f \rightsquigarrow$ find an *absolutely continuous* curve $t \mapsto x(t)$
such that $t \mapsto f(x(t))$ decreases “the fastest”.

$$\frac{d}{dt}x(t) = -\nabla f(x(t))$$

©Francis Bach

If f is not differentiable, there might be more than one descent direction.

DEFINITION

The **subdifferential** of a **convex** function $F: H \rightarrow \mathbb{R}$ is

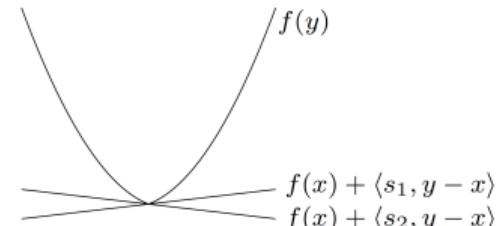
$$\partial F(u) := \{v \in H : F(w) \geq F(u) + \langle v, w - u \rangle \ \forall w \in H\}.$$

THEOREM (PROPERTIES OF THE RESOLVENT)

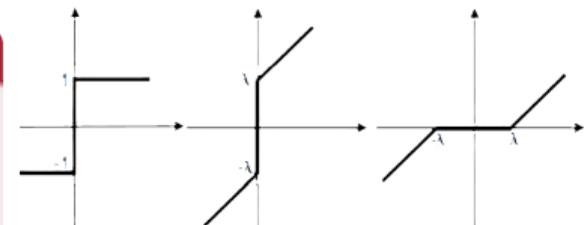
Let $F: H \rightarrow \mathbb{R}$ be convex and lower semicontinuous, $F \not\equiv \infty$. Then ∂F is maximal monotone, hence $\forall \varepsilon > 0$, the **resolvent**

$$J_\varepsilon^{\partial F} := (I + \varepsilon \partial F)^{-1} : H \rightarrow H$$

is single-valued.



©Florian Gohmann



∂F , $I + \lambda \partial F$, $(I + \lambda \partial F)^{-1}$, where $F := |\cdot|$.

This theorem tells us that gradient flows in Hilbert spaces exist and are unique.

THEOREM (EXISTENCE, REGULARITY OF STRONG SOLUTIONS TO CAUCHY PROBLEMS (BREZIS '67))

Let $f: H \rightarrow \mathbb{R}$ s.t. $\partial f: H \rightarrow 2^H$ is maximal monotone and $g_0 \in \text{dom}(\partial f)$. Then $\exists! g: [0, \infty) \rightarrow H$ s.t.

- $g(0) = g_0$ and $\frac{dg}{dt}(t) \in -(\partial f)(g(t))$ for almost all $t > 0$, and $g(t) \in \text{dom}(\partial f)$ for all $t > 0$
- g is Lipschitz continuous on $[0, \infty)$.
- g is given by the exponential formula

$$g(t) = \lim_{n \rightarrow \infty} \left(J_{\frac{t}{n}}^{\partial f} \right)^n (g_0) \quad \text{uniformly in } t \text{ on compact intervals.}$$

For $\mathcal{P}_2(\mathbb{R})$ instead of H : what are the analogs of the tangent vector $\frac{d}{dt}x(t)$ and the subdifferential ∂f ?

DEFINITION (FRÉCHET SUBDIFFERENTIAL IN WASSERSTEIN SPACE)

The (reduced) Fréchet subdifferential of $\mathcal{F}: \mathcal{P}_2(\mathbb{R}) \rightarrow \mathbb{R}$ at μ is

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

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

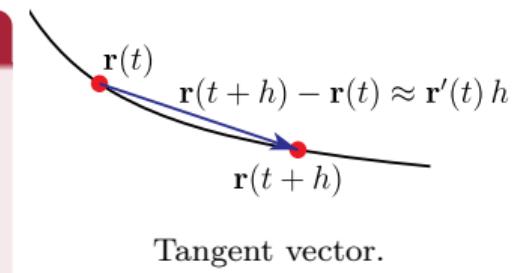
$$\partial_t \gamma_t + \nabla \cdot (v_t \gamma_t) = 0, \quad (t, x) \in (0, \infty) \times \mathbb{R}, \text{ weakly.} \quad (\text{Continuity Eq.})$$

DEFINITION (WASSERSTEIN GRADIENT FLOW)

A absolutely continuous curve $\gamma: (0, \infty) \rightarrow \mathcal{P}_2(\mathbb{R})$ with velocity field

$(v_t \in T_{\gamma_t} \mathcal{P}_2(\mathbb{R}))_{t>0}$ is a *Wasserstein gradient flow with respect to*
 $\mathcal{F}: \mathcal{P}_2(\mathbb{R}) \rightarrow \mathbb{R}$ if

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



THEOREM (AMBROSIO, GIGLI, SAVARÉ (2005))

Let $\mathcal{F}: \mathcal{P}_2(\mathbb{R}) \rightarrow \mathbb{R}$ be bounded from below, lower semicontinuous, **geodesically convex**.

Existence and uniqueness. Then $\exists!$ **Wasserstein gradient flow** $\gamma: (0, \infty) \rightarrow \mathcal{P}_2(\mathbb{R})$ with respect to \mathcal{F} starting at $\gamma(0+) = \mu_0 \in \mathcal{P}_2(\mathbb{R})$.

Approximation scheme. The piecewise constant curves constructed from the iterates

$$\mu_{n+1} := \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \mathcal{F}(\mu) + \frac{1}{2\tau} W_2^2(\mu_n, \mu) \right\}, \quad \tau > 0, \quad (\text{Minimizing movement scheme})$$

i.e., γ_τ defined by $\gamma_\tau|_{(n\tau, (n+1)\tau]} := \mu_n$, $n \in \mathbb{N}$, converge locally uniformly to γ as $\tau \downarrow 0$.

Convergence speed. If $\bar{\mu}$ is a minimizer of \mathcal{F} , then

$$\mathcal{F}(\gamma_t) - \mathcal{F}(\bar{\mu}) \leq \frac{1}{2t} W_2^2(\mu_0, \bar{\mu}). \quad (\text{Sublinear convergence rate})$$

We establish a correspondence between $L_2(0, 1)$ -gradient flows of F and Wasserstein gradient flows of \mathcal{F} .

THEOREM (QUANTILE FUNCTION REFORMULATION [DSBHS24])

Let $F: L_2(0, 1) \rightarrow (-\infty, \infty]$ be convex and lsc, $F \not\equiv \infty$ with $F(Q_\mu) = \mathcal{F}(\mu)$ for all $\mu \in \mathcal{P}_2(\mathbb{R})$. Assume $J_\varepsilon^{\partial F}$ maps $\mathcal{C}(0, 1)$ into itself $\forall \varepsilon > 0$. Then, \forall initial datum $g_0 \in \mathcal{C}(0, 1) \cap \text{dom}(\partial F)$,

$$\begin{cases} \partial_t g(t) + \partial F(g(t)) \ni 0, & t \in (0, \infty), \\ g(0) = g_0, \end{cases} \quad (\text{Cauchy Problem})$$

has a unique strong solution g .

The curve $\gamma_t := (g(t))_{\#} \Lambda_{(0,1)}$ has quantile functions $Q_{\gamma_t} = g(t)$ and is a **Wasserstein gradient flow** of \mathcal{F} with $\gamma(0+) = (g_0)_{\#} \Lambda_{(0,1)}$.

Consider the **negative distance kernel** $K(x, y) := -|x - y|$.

K is only **conditionally positive definite**.

Motivation: electrostatic principles, interacting species, dithering.

DEFINITION (MMD)

The **maximum mean discrepancy** of with respect to K is

$$\mathcal{P}_2(\mathbb{R}) \times \mathcal{P}_2(\mathbb{R}) \rightarrow [0, \infty),$$

$$(\mu, \nu) \mapsto \frac{1}{2} \text{MMD}_K(\mu, \nu)^2 = \int_{\mathbb{R} \times \mathbb{R}} K(x, y) d(\mu - \nu)(x) d(\mu - \nu)(y).$$



Dithering

doi.org/10.1101/100790

$$\mathcal{F}_\nu(\mu) := \underbrace{-\frac{1}{2} \int_{\mathbb{R} \times \mathbb{R}} |x - y| d\mu(x) d\mu(y)}_{\text{interaction}} + \underbrace{\int_{\mathbb{R}^d \times \mathbb{R}^d} |x - y| d\mu(x) d\nu(y)}_{\text{potential}}. \quad (2)$$

We now apply this theorem to the MMD with the negative distance kernel, i.e. to

$$\mathcal{F}_\nu := \frac{1}{2} \text{MMD}_K(\cdot, \nu).$$

LEMMA (PROPERTIES OF F_ν)

The functional

$$F_\nu : L_2(0, 1) \rightarrow \mathbb{R}, \quad u \mapsto \int_0^1 \left((1 - 2s)(u(s) - Q_\nu(s)) + \int_0^1 |u(s) - Q_\nu(t)| dt \right) ds$$

is convex and continuous, we have $F_\nu(Q_\mu) = \mathcal{F}_\nu(\mu)$ $\forall \mu \in \mathcal{P}_2(\mathbb{R})$, and

$$\partial F_\nu(u) = \left\{ f \in L_2(0, 1) : f(s) \in 2[R_\nu^-(u(s)), R_\nu^+(u(s))] - 2s \text{ for a.e. } s \in (0, 1) \right\}, \quad u \in L_2(0, 1)$$

and $J_\varepsilon^{\partial F_\nu}$ maps $\mathcal{C}(0, 1)$ into itself $\forall \varepsilon > 0$.

Proof. By elementary means, nothing fancy happens here.

The lower Lipschitz constant of $g \in \mathcal{C}(0, 1)$ is

$$L_{\text{low}}(g) := \max \left\{ L \geq 0 : \frac{g(s_1) - g(s_2)}{s_1 - s_2} \geq L \quad \forall s_1, s_2 \in (0, 1) \right\} \geq 0,$$

If $\mu = \delta_x$, then $Q_\mu \equiv x$ and $L_{\text{low}}(Q_\mu) = 0$.

THEOREM (TIME EVOLUTION OF CDF'S LOWER LIPSCHITZ CONSTANTS)

Let $\nu \in \mathcal{P}_2(\mathbb{R})$ with $L_{\text{low}}(Q_\nu) > 0$, and $g_0 = Q_{\mu_0} \in \mathcal{C}(0, 1)$. We have

$$\text{Lip}(R_{\gamma_t}) \leq \left(L_{\text{low}}(g_0) \cdot e^{-\frac{2t}{L_{\text{low}}(Q_\nu)}} + L_{\text{low}}(Q_\nu) \cdot (1 - e^{-\frac{2t}{L_{\text{low}}(Q_\nu)}}) \right)^{-1} < \infty.$$

THEOREM (CONTINUITY IS PRESERVED & MONOTONICITY OF THE SUPPORT)

Let $\nu \in \mathcal{P}_2(\mathbb{R})$, g_0 be continuous and g the solution of the Cauchy problem starting in g_0 .

- $g(t)$ is continuous for all $t \geq 0$.
- The ranges fulfill $\overline{g(t_1)(0, 1)} \subseteq \overline{g(t_2)(0, 1)}$ for all $0 \leq t_1 \leq t_2$.

COROLLARY (POINT MEASURE TARGET)

Let $\nu := \sum_{j=1}^n w_j \delta_{x_j}$. Then F_ν flow is given by

$$[g(t)](s) := \begin{cases} Q_{\mu_0}(s) + 2(s - R_{s,0})t, & t \in [t_{s,0}, t_{s,1}), \\ x_{s,j} + 2(s - R_{s,j})(t - t_{s,j}), & t \in [t_{s,j}, t_{s,j+1}), \\ Q_\nu(s), & t \geq t_{s,|\ell_s - k_s|}, \end{cases}$$

where

$$t_{s,0} := 0, \quad t_{s,1} := \frac{x_{s,1} - Q_{\mu_0}(s)}{2(s - R_{s,0})}, \quad t_{s,j+1} := t_{s,j} + \frac{x_{s,j+1} - x_{s,j}}{2(s - R_{s,j})},$$

$$Q_{\mu_0}(s) \leq Q_\nu(s)$$

$$Q_{\mu_0}(s) \geq Q_\nu(s)$$

$$\ell_s \quad W_{\ell_s-1} < s < W_{\ell_s} \quad W_{\ell_s-1} < s < W_{\ell_s}$$

$$k_s \quad x_{k_s} \leq Q_{\mu_0}(s) < x_{k_s+1} \quad x_{k_s-1} < Q_{\mu_0}(s) \leq x_{k_s}$$

$$x_{s,j} \quad x_{k_s+j} \quad x_{k_s-j} \quad j \leq |\ell_s - k_s|$$

$$R_{s,j} \quad W_{k_s+j} \quad W_{k_s-j-1} \quad j \leq |\ell_s - k_s| - 1$$

NUMERICAL EXPERIMENTS - IMPLICIT EULER (BACKWARD) SCHEME

Let $\tau > 0$. The minimizing movement (or JKO) scheme,

$$\mu_{n+1} := \operatorname{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \left\{ \mathcal{F}_\nu(\mu) + \frac{1}{2\tau} W_2^2(\mu_n, \mu) \right\},$$

can be rewritten using the isometry $\mathcal{P}_2(\mathbb{R}) \rightarrow \mathcal{C}(0, 1)$ as

$$g_{n+1} = \operatorname{argmin}_{g \in \mathcal{C}(0, 1)} \left\{ F_\nu(g) + \frac{1}{2\tau} \int_0^1 |g - g_n|^2 \, ds \right\}$$

$$F_\nu \in \Gamma_0(L_2(0, 1)) \quad (I + \tau \partial F_\nu)^{-1}(g_n),$$

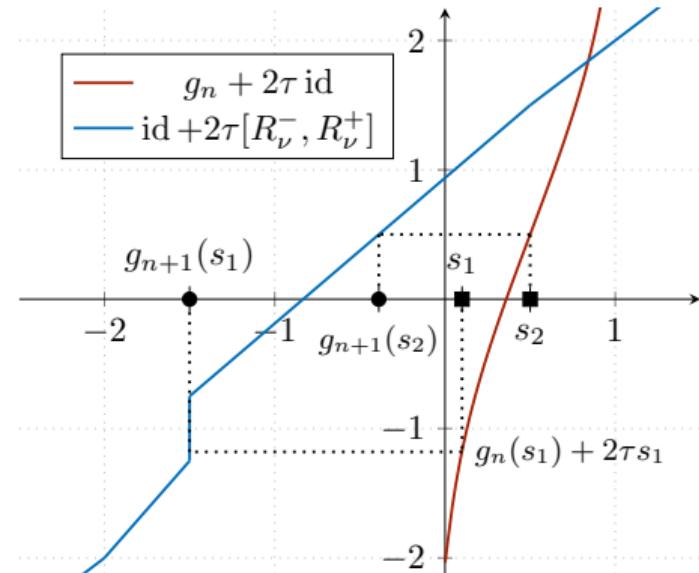
which is equivalent to

$$g_n(s) + 2\tau s \in g_{n+1}(s) + 2\tau [R_\nu^-(g_{n+1}(s)), R_\nu^+(g_{n+1}(s))] \quad (3)$$

for all $s \in (0, 1)$.

$(g_n)_\# \Lambda_{(0,1)}$ approximates WGF for $t \in (n\tau, (n+1)\tau]$.

$(g_n)_{n \in \mathbb{N}} \rightarrow Q_\nu$ weakly in $L_2(0, 1)$ and $\mu_n := (g_n)_\# \Lambda_{(0,1)} \rightarrow \nu$ narrowly for fixed $\tau > 0$.



Implicit Euler step visualized.

If $R_\nu^+ = R_\nu^- =: R_\nu$, we can also use explicit Euler discretization

$$g_{n+1} = g_n - \tau \nabla F_\nu(g_n) = g_n - 2\tau(R_\nu \circ g_n - \text{id}),$$

Advantage: we don't have to solve an inclusion in each step.

Disadvantage: weaker convergence guarantees, might not preserve $\mathcal{C}(0, 1)$
(iterates not monotone).

$$\mu_0 = \mathcal{N}(-5, 1), \nu = \mathcal{N}(5, 1).$$

- Extension to higher dimensions. Lagrangian reformulation now involves flow map.
- Deterministic particle approximations and mean field limits.
- Convergence of explicit scheme, distance of implicit and explicit iterates.
- Generalize Lipschitz properties to Hölder properties.
- Applying similar techniques to non-convex functional and then regularize.

- Reformulation as **maximal monotone** inclusion Cauchy problem in $L_2(0, 1)$ via **quantile functions**.
- Comprehensive description of solutions' behavior, qualitative description of **instantaneous measure-to- L^∞ regularization**.
- Implicit Euler is simple.

Thank you for your attention!

I am happy to take any questions.

Paper link: <https://arxiv.org/abs/2408.07498>

My website: viktorajstein.github.io

- [AGS08] L. Ambrosio, N. Gigli, and G. Savare, *Gradient flows*, 2nd ed., Lectures in Mathematics ETH Zürich, Birkhäuser, Basel, 2008.
- [Bre73] Haim Brezis, *Operateurs maximaux monotones*, North-Holland Mathematics Studies, 1973 (French).
- [CL71] M. G. Crandall and T. M. Liggett, *Generation of semi-groups of nonlinear transformations on general banach spaces*, American Journal of Mathematics **93** (1971), no. 2, 265–298.
- [JKO98] Richard Jordan, David Kinderlehrer, and Felix Otto, *The variational formulation of the Fokker–Planck equation*, SIAM Journal on Mathematical Analysis **29** (1998), no. 1, 1–17.

Interpolating between OT and KL regularized OT using Rényi Divergences

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

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

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

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

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

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

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