

Geometric Losses for Machine Learning - Sinkhorn Divergences for Unbalanced Optimal Transport

Thibault Séjourné

PGMO Seminar – 4th December, 2019

Joint work with Jean Feydy, Francois-Xavier Vialard, Alain Trouvé and Gabriel Peyré

Introduction

Csiszar divergences

Optimal Transport

Unbalanced Optimal Transport

Entropic Optimal Transport

Correcting the entropic bias - Sinkhorn divergence

Introduction

From discrete to continuous setting (and back)

Definitions

- Continuous function: $f \in \mathcal{C}(\mathcal{X})$
- Positive measure: Linear form $\alpha \in \mathcal{M}_+(\mathcal{X})$
- Dual product: $\langle \alpha, f \rangle = \int_{\mathcal{X}} f(x) d\alpha(x) = \mathbb{E}_{\alpha}[f]$.

Discretizing measures

When $\alpha = \sum_{i=1}^n \alpha_i \delta_{x_i}$ one implements α on a computer with $(\alpha_i) \in \mathbb{R}^n$ and $(x_i) \in \mathbb{R}^{n \times d}$. Then functions are vectors $(f_i) = (f(x_i))_i \in \mathbb{R}^n$ and $\langle \alpha, f \rangle = \sum \alpha_i f_i$.

Small Take home message

Some algorithms are better understood using a continuous formalism.

We require that the loss verifies at least the following axioms:

- Positivity: $\forall(\alpha, \beta), \mathcal{L}(\alpha, \beta) \geq 0$.
- Definiteness: $\forall(\alpha, \beta), \mathcal{L}(\alpha, \beta) = 0 \Leftrightarrow \alpha = \beta$.
- Metrizing weak* convergence (convergence in law):

$$\forall(\alpha, \beta), \mathcal{L}(\alpha, \beta) \rightarrow 0 \Leftrightarrow \alpha \rightarrow \beta,$$

where $\alpha \rightarrow \beta \Leftrightarrow \forall f \in \mathcal{C}(\mathcal{X}), \langle \alpha, f \rangle \rightarrow \langle \beta, f \rangle$.

- Differentiability (for backpropagation).

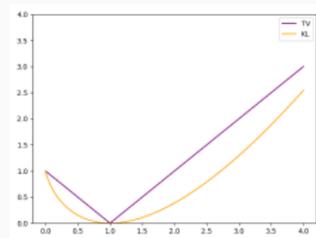
Csiszar divergences

Definitions [Csiszàr'67]

- Entropy φ : nonnegative, l.s.c., convex on \mathbb{R}_+ s.t. $\varphi(1) = 0$
 - Recession constant: $\varphi'^{\infty} = \lim_{x \rightarrow \infty} \varphi(x)/x$
 - Lebesgue decomposition: $\forall(\alpha, \beta), \alpha = \frac{d\alpha}{d\beta}\beta + \alpha^{\top}$
 - φ -divergence: $D_{\varphi}(\alpha, \beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha}{d\beta}\right)d\beta + \varphi'^{\infty} \int_{\mathcal{X}} d\alpha^{\top}$
- Discretized: $D_{\varphi}(\alpha, \beta) = \sum_{\beta_i \neq 0} \varphi\left(\frac{\alpha_i}{\beta_i}\right)\beta_i + \varphi'^{\infty} \sum_{\beta_i = 0} \alpha_i$

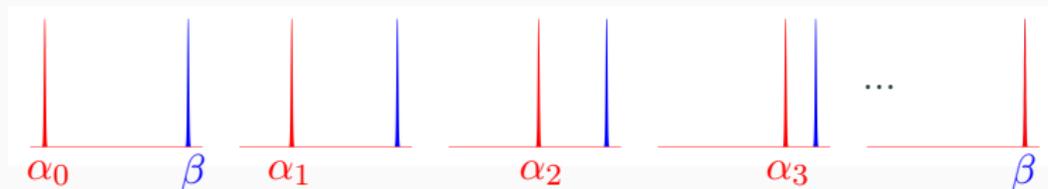
Examples:

- KL: $\varphi(x) = x \log x - x + 1, \varphi'^{\infty} = +\infty,$
- TV: $\varphi(x) = |x - 1|$ and $\varphi'^{\infty} = 1.$

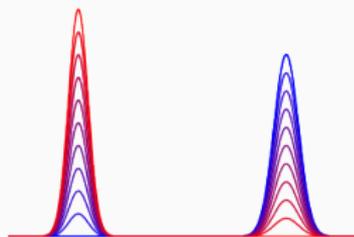


Properties of Csiszàr divergences

Consider the sequence $\alpha_n = \delta_{1/n}$ and $\beta = \delta_0$. One has $\alpha_n \rightarrow \beta$, but $\text{KL}(\alpha_n|\beta) = \infty$ and $\text{TV}(\alpha_n|\beta) = 2$.



- 😊 Simple and cheap to compute
- 😞 Ignores the geometry and do not metrize convergence in law



Optimal Transport

Optimal Transport (OT)

Optimal Transport Distance

$$\text{OT}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi \geq 0} \left\{ \langle \pi, C \rangle : \begin{array}{l} \pi \mathbf{1} = \alpha \\ \pi^\top \mathbf{1} = \beta \end{array} \right\}.$$

Called p-Wasserstein distance for $C = d^p$.

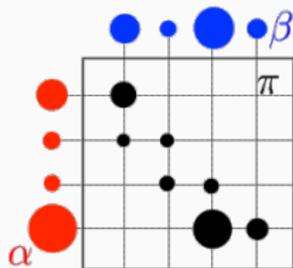
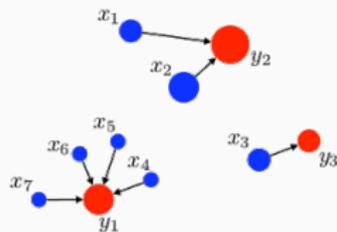
Discrete: $\langle \pi, C \rangle = \sum_{i,j} \pi_{ij} C_{ij}$

Intuition: Moving π_{ij} grams from x_i to y_j costs $\pi_{ij} \times C_{ij} = \pi_{ij} \times C(x_i, y_j)$.

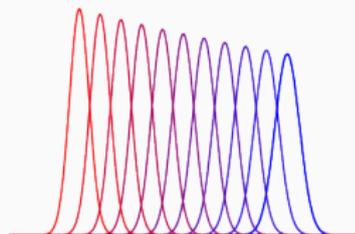
Choice of $C \rightarrow$ Choice of geometric prior.

\Rightarrow Learn it !

[Kantorovich'42]



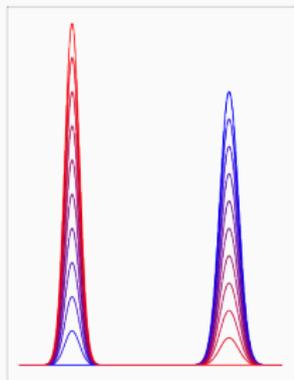
- One has $OT(\delta_x, \delta_y) = C(x, y)$
- $\Rightarrow OT(\delta_{1/n}, \delta_0) \xrightarrow{n \rightarrow \infty} 0$
- Metric on $\mathcal{X} \rightarrow$ metric on $\mathcal{M}_+(\mathcal{X})$



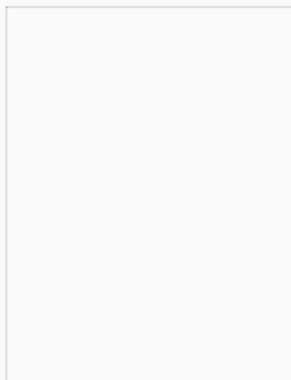
- 😊 Metrizes convergence in law
- ☹️ Computation complexity $\mathcal{O}(n^3 \log n)$, not differentiable
- ☹️ Only compares probabilities, i.e. normalized weighted point clouds

Unbalanced Optimal Transport

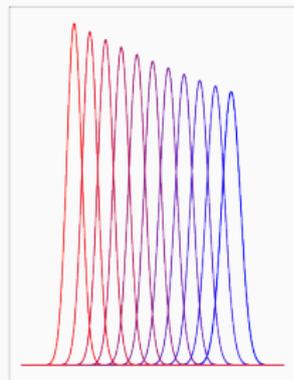
Hybridizing Vertical and Horizontal Geometries



Vertical



In between ?



Horizontal

Unbalanced optimal transport

Hybridizing \Rightarrow Soften the hard constraint $\pi_1 = \alpha \rightarrow \rho D_\varphi(\pi_1 | \alpha)$.

Allows for creation/destruction of mass locally.

Definition - Unbalanced OT [Liero, Mielke, Savaré '18]

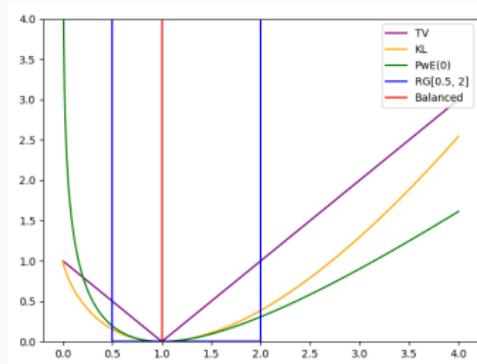
For any φ -divergence D_φ and any measures (α, β) one defines:

$$\text{OT}_\rho(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\pi \geq 0} \langle \pi, C \rangle + \rho D_\varphi(\pi_1, \alpha) + \rho D_\varphi(\pi_2, \beta).$$

- Add a parameter ρ : Transport radius. ($\text{OT}_\rho \xrightarrow{\rho \rightarrow +\infty} \text{OT}$).
- Choice of D_φ : prior on the mass variation dynamics
- Balanced OT is retrieved with $D_\varphi = \iota_{(=)}$

Examples of entropies

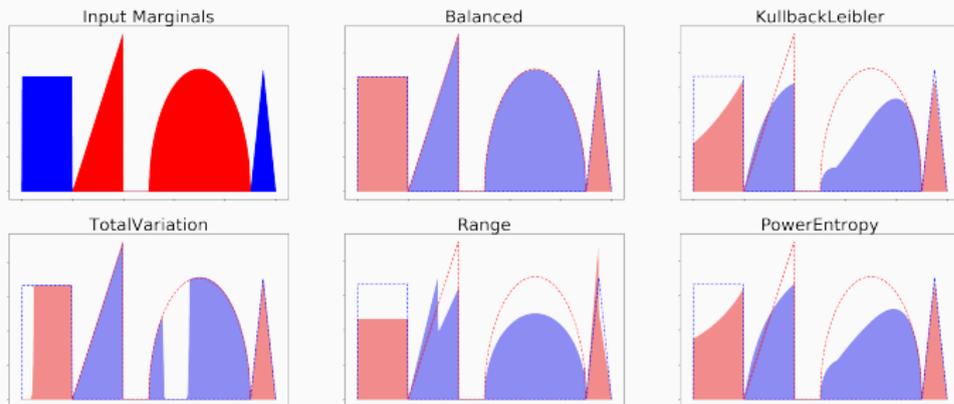
- **Balanced**: $\varphi(x) = \iota_{\{1\}}(x)$ with $D_\varphi(\pi_1, \alpha) = \iota_{(=)}(\pi_1, \alpha)$.
- **TV**: $\varphi(x) = |x - 1|$
- **KL**: $\varphi(x) = x \log x - x + 1$
- **Power entropy**: $\varphi(x) = \frac{1}{p(p-1)}(x^p - p(x-1) - 1)$, $p \in \mathbb{R}$.
- **Range**: $\varphi(x) = \iota_{[a,b]}(x)$ ($a \leq 1 \leq b$), i.e. $a\alpha \leq \pi_1 \leq b\alpha$.



Numerical examples

Reminder: Local mass creation and destruction is allowed

- Shows how α is matched onto β and vice versa through π .
- Plots $\pi_1 \approx \alpha$ and $\pi_2 \approx \beta$
- Input marginals are dashed.



Entropic Optimal Transport

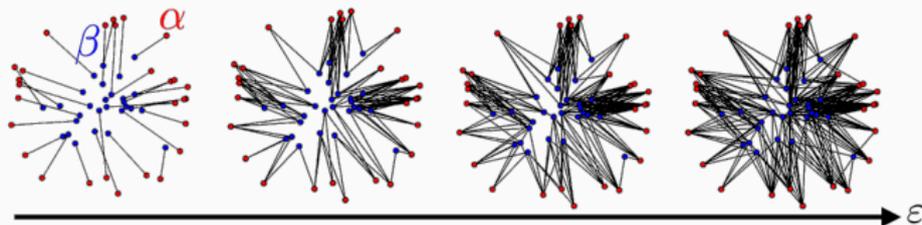
Regularization of OT

Reminder: OT is computationally expensive and non-smooth

Idea: Add an entropic penalty $\varepsilon \text{KL}(\pi, \alpha \otimes \beta)$

Entropic Unbalanced OT [Cuturi'13, Chizat'18]

$$\text{OT}_\varepsilon(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\pi \geq 0} \langle \pi, C \rangle + \rho D_\varphi(\pi_1, \alpha) + \rho D_\varphi(\pi_2, \beta) + \varepsilon \text{KL}(\pi, \alpha \otimes \beta)$$



Duality of regularized OT

Writing $\varphi^*(x) = \sup_{y \geq 0} xy - \varphi(y)$, the dual reads

$$\begin{aligned} \text{OT}_\varepsilon(\alpha, \beta) = \sup_{f, g \in \mathcal{C}(\mathcal{X})} & \langle \alpha, -(\rho\varphi)^*(-f) \rangle + \langle \beta, -(\rho\varphi)^*(-g) \rangle \\ & - \varepsilon \langle \alpha \otimes \beta, e^{\frac{f(x)+g(y)-C(x,y)}{\varepsilon}} - 1 \rangle \end{aligned}$$

The alternate dual ascent is straightforward to compute:

Alternate dual ascent

given any initialization $f_0 \in \mathcal{C}(\mathcal{X})$. At time t one has (f_t, g_t) .

Then

- Fix f_t and find optimal g in the dual $\rightarrow g_{t+1}$,
- Fix g_{t+1} and find optimal f in the dual $\rightarrow f_{t+1}$,
- Iterate until convergence.

Unbalanced Sinkhorn algorithm

Proposition - Unbalanced Sinkhorn algorithm

Define the following operators

- (Softmin / LogSumExp) $\text{Smin}_{\alpha}^{\varepsilon}(f) \stackrel{\text{def.}}{=} -\varepsilon \log \langle \alpha, e^{-f/\varepsilon} \rangle$
- (Anisotropic Prox) $\text{aprox}(p) = \arg \min_{q \in \mathbb{R}} \varepsilon e^{(p-q)/\varepsilon} + \varphi^*(q)$

The optimality condition defines the Sinkhorn algorithm

$$g_{t+1}(y) = -\text{aprox}(-\text{Smin}_{\alpha}^{\varepsilon}(C(\cdot, y) - f_t))$$

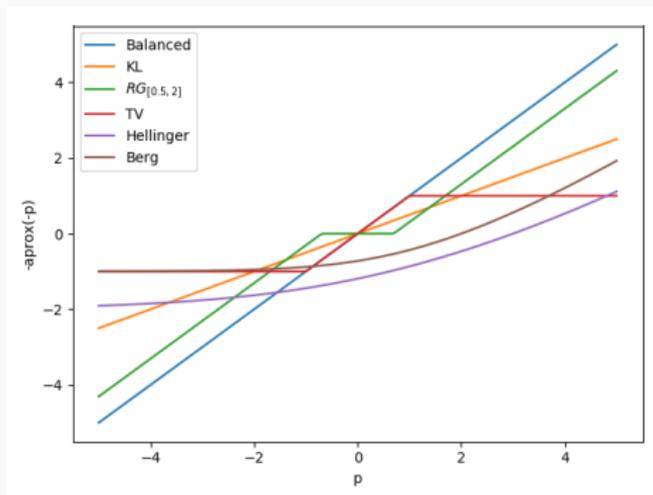
$$f_{t+1}(x) = -\text{aprox}(-\text{Smin}_{\beta}^{\varepsilon}(C(x, \cdot) - g_{t+1})).$$

Theorem [S., Feydy, Vialard, Trounev, Peyre '19]

The Sinkhorn algorithm converges towards the optimal (f, g) of $\text{OT}_{\varepsilon}(\alpha, \beta)$ when φ^* is strictly convex, but also for TV, Range and Balanced OT.

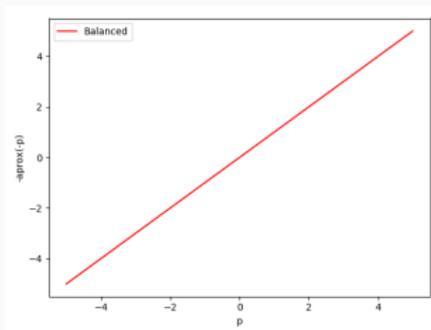
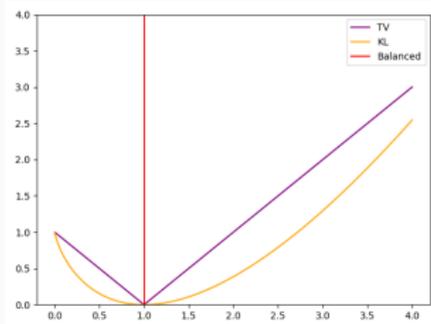
Examples of Anisotropic prox

- Balanced Sinkhorn = Softmin
 - Unbalanced Sinkhorn = $\mathbf{aprox} \circ \text{Softmin}$
- \Rightarrow Unbalanced Sinkhorn = Readjusting Balanced Sinkhorn with the operator \mathbf{aprox} .
- Sinkhorn algorithm is a (weakly) contractive algorithm



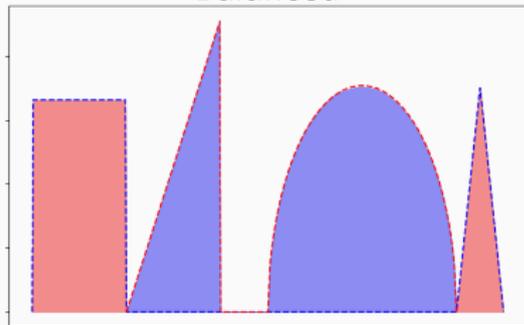
Examples of Anisotropic prox - Balanced

Entropy and Aprox



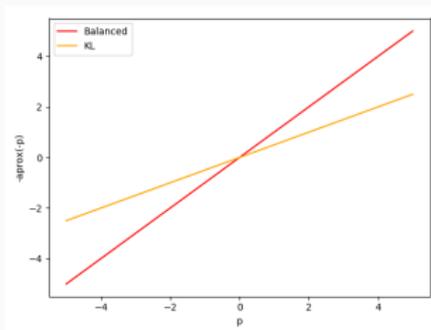
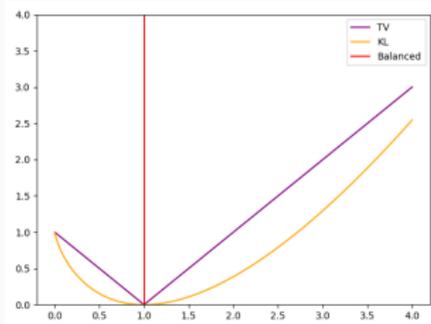
$$D_\varphi = \iota_{\{=\}}$$
$$\varphi(x) = \iota_{\{1\}}(x)$$
$$\text{aprox}(x) = x$$

Balanced



Examples of Anisotropic prox - KL

Entropy and Aprox

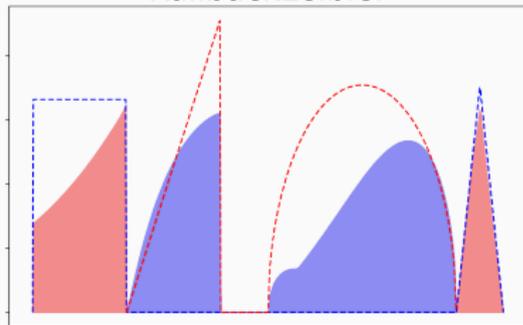


$$D_\varphi = \rho \text{KL}$$

$$\varphi(x) = \rho(x \log x - x + 1)$$

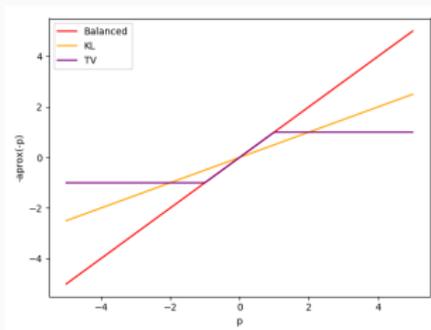
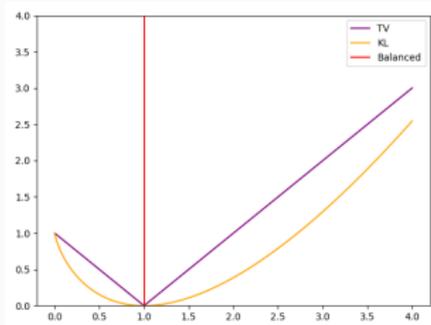
$$\text{aprox}(x) = \frac{\rho}{\rho + \varepsilon} x$$

KullbackLeibler



Examples of Anisotropic prox - TV

Entropy and Aprox

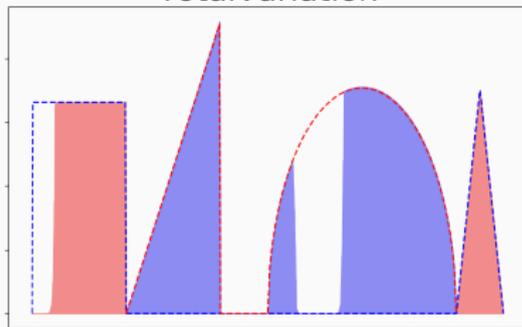


$$D_\varphi = \rho \text{TV}$$

$$\varphi(x) = \rho |x - 1|$$

$$\text{aprox}(x) = x \text{ if } x \in [-\rho, \rho], \rho \text{ if } x \geq \rho \text{ and } -\rho \text{ if } x \leq -\rho$$

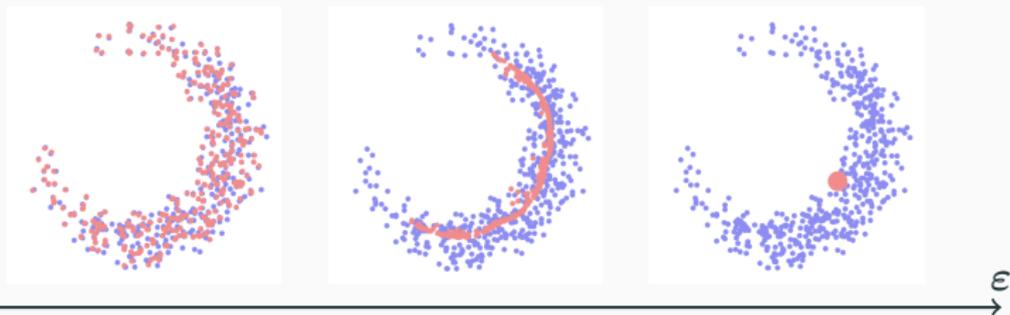
TotalVariation



Correcting the entropic bias - Sinkhorn divergence

Problem: OT_ε does not metrize weak* convergence for $\varepsilon > 0$. ☹️

$$\begin{aligned} \exists \alpha \in \mathcal{M}_1^+(\mathcal{X}), \text{OT}_\varepsilon(\alpha, \beta) < \text{OT}_\varepsilon(\beta, \beta). \\ \text{OT}_0(\alpha, \beta) \xleftarrow{0 \leftarrow \varepsilon} \text{OT}_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \rightarrow \infty} \alpha^\top C \beta. \end{aligned}$$



Unbalanced Sinkhorn Divergence

Setting $m(\mu)$ to be the total mass of the measure μ , we define

$$S_{\varepsilon, \rho}(\alpha, \beta) \stackrel{\text{def.}}{=} \text{OT}_{\varepsilon, \rho}(\alpha, \beta) - \frac{1}{2} \text{OT}_{\varepsilon, \rho}(\alpha, \alpha) - \frac{1}{2} \text{OT}_{\varepsilon, \rho}(\beta, \beta) + \frac{\varepsilon}{2} (m(\alpha) - m(\beta))^2.$$

It extends the balanced case from [Ramdas '15][Genevay '18].

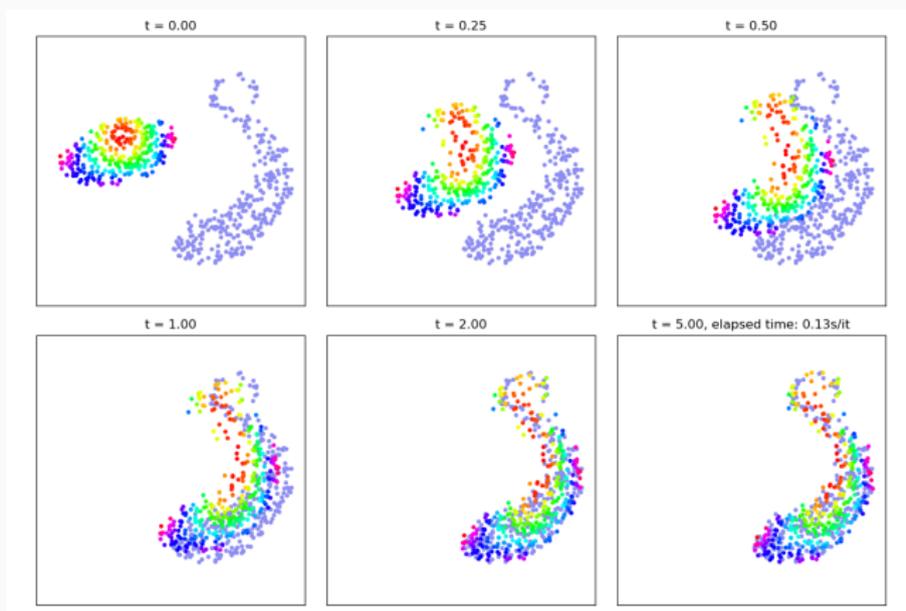
Theorem [S., Feydy, Vialard, Trounev, Peyre '19]

For any Lipschitz cost C s.t. $k_{\varepsilon} \stackrel{\text{def.}}{=} e^{-\frac{C}{\varepsilon}}$ is a positive universal kernel, for any $\varepsilon > 0$ and strictly convex φ^*

- $S_{\varepsilon, \rho}$ is convex, positive, definite.
- It is (weakly) differentiable.
- $S_{\varepsilon, \rho}(\alpha, \beta) \rightarrow 0 \Leftrightarrow \alpha \rightarrow \beta$.

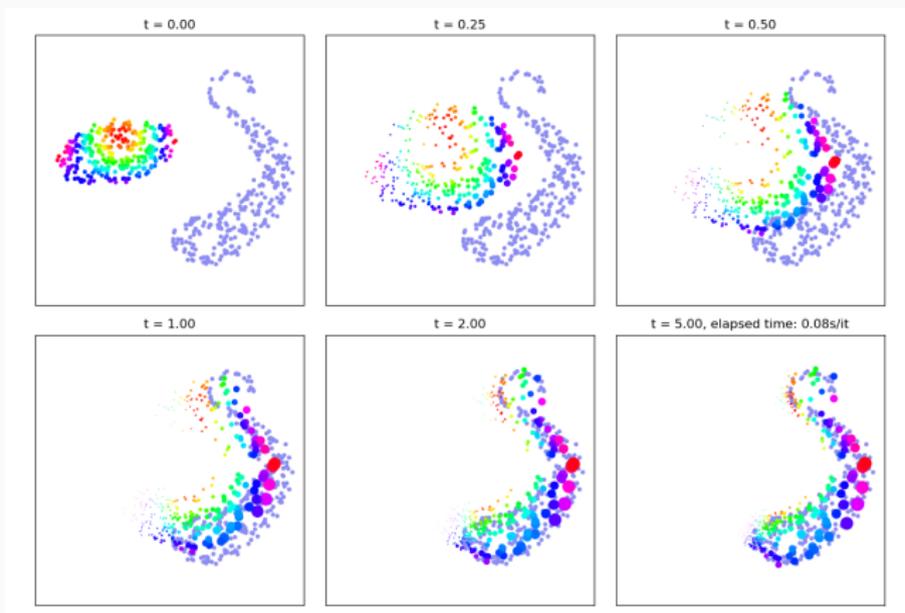
Numerics - Balanced Gradient Flow

- Model: $\alpha_\theta = \sum_{i=1}^n \alpha_i \delta_{x_i}$ with $\theta = (x_i)$
- Loss: S_ε with balanced OT and $\varepsilon = 0.01$



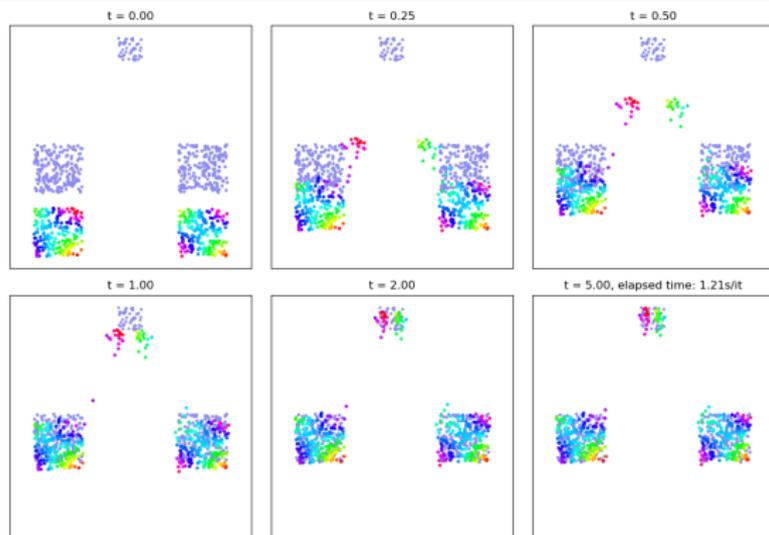
Numerics - Unbalanced Gradient Flow

- Model: $\alpha_\theta = \sum_{i=1}^n \alpha_i \delta_{x_i}$ with $\theta = (x_i, \alpha_i)$
- Loss: S_ε with KL UOT and $(\varepsilon, \rho) = (0.01, 0.3)$



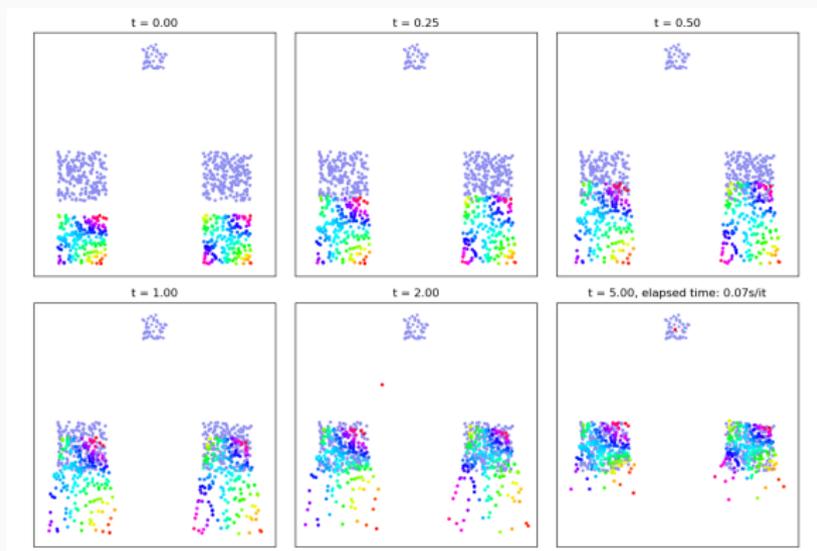
Numerics - Avoiding overfitting

- Model: $\alpha_\theta = \sum_{i=1}^n \alpha_i \delta_{x_i}$ with $\theta = (x_i)$
- Loss: S_ε with balanced OT and $\varepsilon = 0.01$



Numerics - Avoiding overfitting

- Model: $\alpha_\theta = \sum_{i=1}^n \alpha_i \delta_{x_i}$ with $\theta = (x_i)$
- Loss: S_ε with KL UOT and $(\varepsilon, \rho) = (0.01, 0.3)$



Unbalanced OT allows to avoid overfitting of outliers !

Implementation

- Sinkhorn divergences can be fastly computed via GPU-friendly routines
- + Efficient optimization heuristics (annealing + subsampling)
- Available losses (Balanced + KL) on Jean Feydy's repository:
<http://www.kernel-operations.io/geomloss/>
- Two modes:
 - Keeps backend for huge measures without overflow (~ 1 million points)
 - Mini-batch mode for machine learning.
- Implementation of other unbalanced divergences at:
<https://github.com/thibsej/unbalanced-ot-functionals>

- Family of parametric losses with appealing properties (convexity, differentiability, positivity...)
- Algorithm with linear convergence
- Consistent behavior which allows to crossvalidate w.r.t. ε
- Improvement of the statistical complexity (Not detailed here)

It remains to experiment new ML applications!

<http://www.kernel-operations.io/geomloss/>
<https://github.com/thibsej/unbalanced-ot-functionals>