



Computing Wasserstein Barycenter via Operator Splitting: the Method of Averaged Marginals

Daniel Mimouni^{*+}

Joint work with:

Paul Malisani^{*}, Jiamin Zhu^{*}, Welington de Oliveira⁺

^{*}IFP Energies nouvelles

⁺Mines Paris, PSL



Table of contents

I. Motivations

I.1. Applications

II. Background on Discrete Optimal Transport

II.1. Wasserstein Distance

II.2. Wasserstein Barycenter (WB)

III. The Method of Averaged Marginals (MAM)

III.1. Reformulation of the LP

III.2. Douglas-Rachford (DR) theory

III.3. Algorithm

IV. Applications

IV.1. Qualitative comparison

IV.2. Quantitative comparison

IV.3. Influence of the support

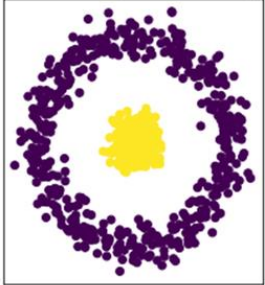
IV.4. Unbalanced Wasserstein Barycenter

I. Motivations

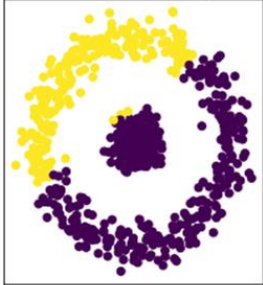
I.1. Applications

Clustering :

Wasserstein-Spectral clustering



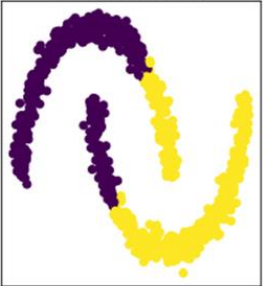
k-means clustering



Wasserstein-Spectral Clustering

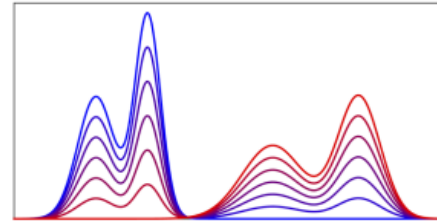
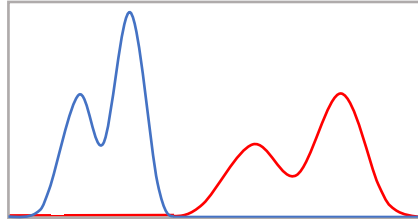


kmeans Clustering

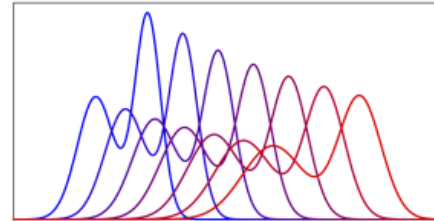


(a) Comparison of Wasserstein-Spectral clustering, spectral clustering, and k-means on Two-Circles dataset and Moons dataset [2]

Data preprocessing :



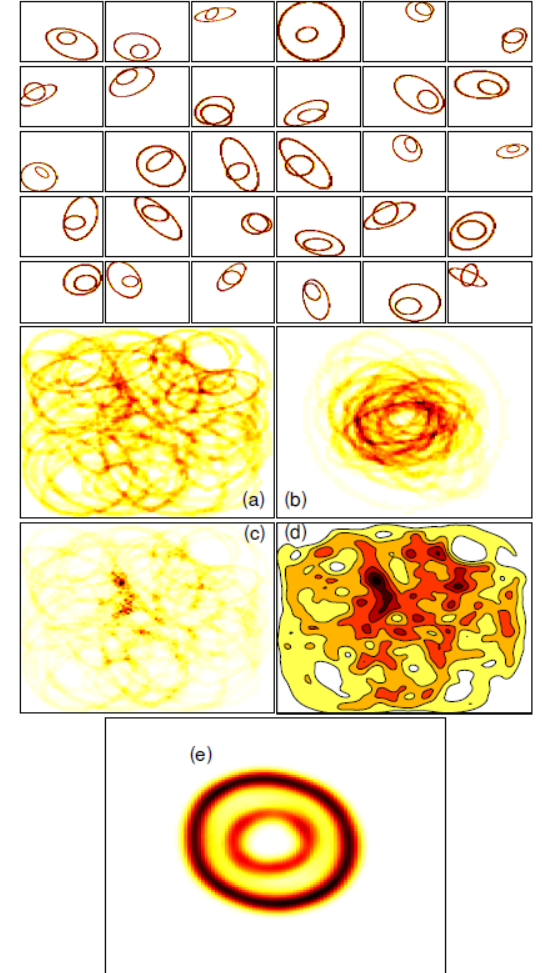
ℓ_2 interpolation



Wasserstein interpolation

Comparison between Euclidean (left) and Optimal Transport (right) barycenters between two densities, one being a translated and scaled version of the other. Colors encode the progression of the interpolation. The Euclidean interpolation results in mixtures of the two initial densities, while Optimal Transport results in a progressive translation and scaling [3]

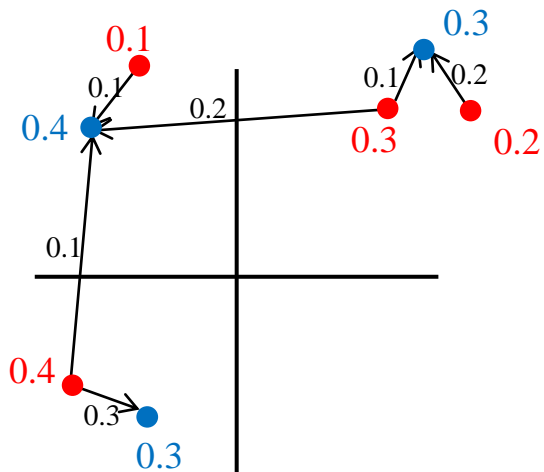
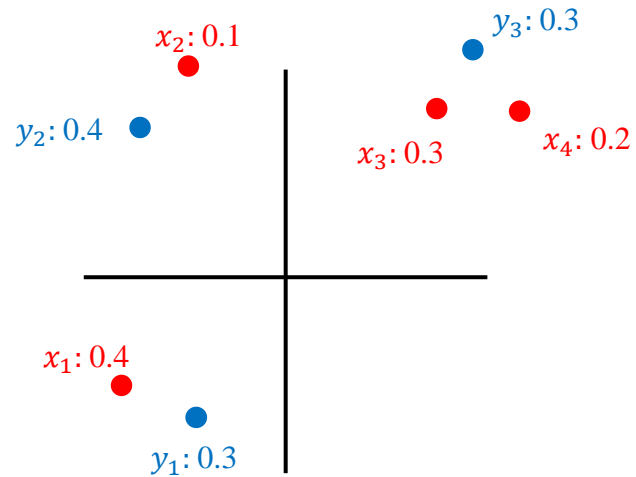
Visualization :



(Top) 30 artificial images of two nested random ellipses. Mean measures using the (a) Euclidean distance (b) Euclidean after re-centering images (c) Jeffrey centroid (Nielsen, 2013) (d) RKHS distance (Gaussian kernel, $\sigma= 0.002$) (e) 2-Wasserstein distance. [1]

II. Background on Discrete Optimal Transport

II.1. Wasserstein distance



the ι -Wasserstein distance $W_\iota(\mu, \nu)$

$$\text{OT}_{\Xi, Z}(p, q) := \begin{cases} \min_{\pi \geq 0} & \sum_{r=1}^R \sum_{s=1}^S d(\xi_r, \zeta_s)^\iota \pi_{rs} \\ \text{s.t.} & \sum_{r=1}^R \pi_{rs} = q_s, \quad s = 1, \dots, S \\ & \sum_{s=1}^S \pi_{rs} = p_r, \quad r = 1, \dots, R \end{cases}$$

π	x_1	x_2	x_3	x_4
y_1	0.3	0	0	0
y_2	0.4	0.1	0.2	0
y_3	0.3	0	0.1	0.2

II. Background on Discrete Optimal Transport

II.1. Wasserstein distance

Recall:

$$\Delta_n(\tau) := \left\{ u \in \mathbb{R}_+^n : \sum_{i=1}^n u_i = \tau \right\}$$

the ι -Wasserstein distance $W_\iota(\mu, \nu)$

finitely many R scenarios $\Xi := \{\xi_1, \dots, \xi_R\}$ for ξ and $S^{(m)}$ scenarios $Z^{(m)} := \{\zeta_1^{(m)}, \dots, \zeta_{S^{(m)}}^{(m)}\}$ for $\zeta^{(m)}$, $m = 1, \dots, M$, i.e., measures of the form

$$\mu = \sum_{r=1}^R p_r \delta_{\xi_r} \quad \text{and} \quad \nu^{(m)} = \sum_{s=1}^{S^{(m)}} q_s^{(m)} \delta_{\zeta_s^{(m)}}, \quad m = 1, \dots, M,$$

with δ_u the Dirac unit mass on $u \in \Omega$, $p \in \Delta_R$, and $q^{(m)} \in \Delta_{S^{(m)}}$, $m = 1, \dots, M$.

$$\text{OT}_{\Xi, Z}(p, q) := \begin{cases} \min_{\pi \geq 0} & \sum_{r=1}^R \sum_{s=1}^S \mathbf{d}(\xi_r, \zeta_s)^\iota \pi_{rs} \\ \text{s.t.} & \sum_{r=1}^R \pi_{rs} = q_s, & s = 1, \dots, S \\ & \sum_{s=1}^S \pi_{rs} = p_r, & r = 1, \dots, R \end{cases}$$

II. Background on Discrete Optimal Transport

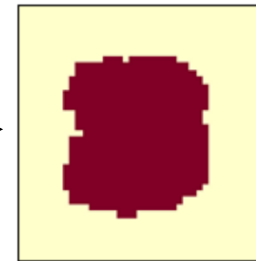
II.2. Wasserstein Barycenter (WB)

General formulation of the WB problem:

$$\min_{\Xi, p \in \Delta_R} \sum_{m=1}^M \alpha_m \text{OT}_{\Xi, Z^{(m)}}(p, q^{(m)})$$



Support optimization



Probability optimization



II. Background on Discrete Optimal Transport

II.2. Wasserstein Barycenter (WB)

General formulation of the WB problem:

$$\min_{\Xi, p \in \Delta_R} \sum_{m=1}^M \alpha_m \text{OT}_{\Xi, Z^{(m)}}(p, q^{(m)})$$

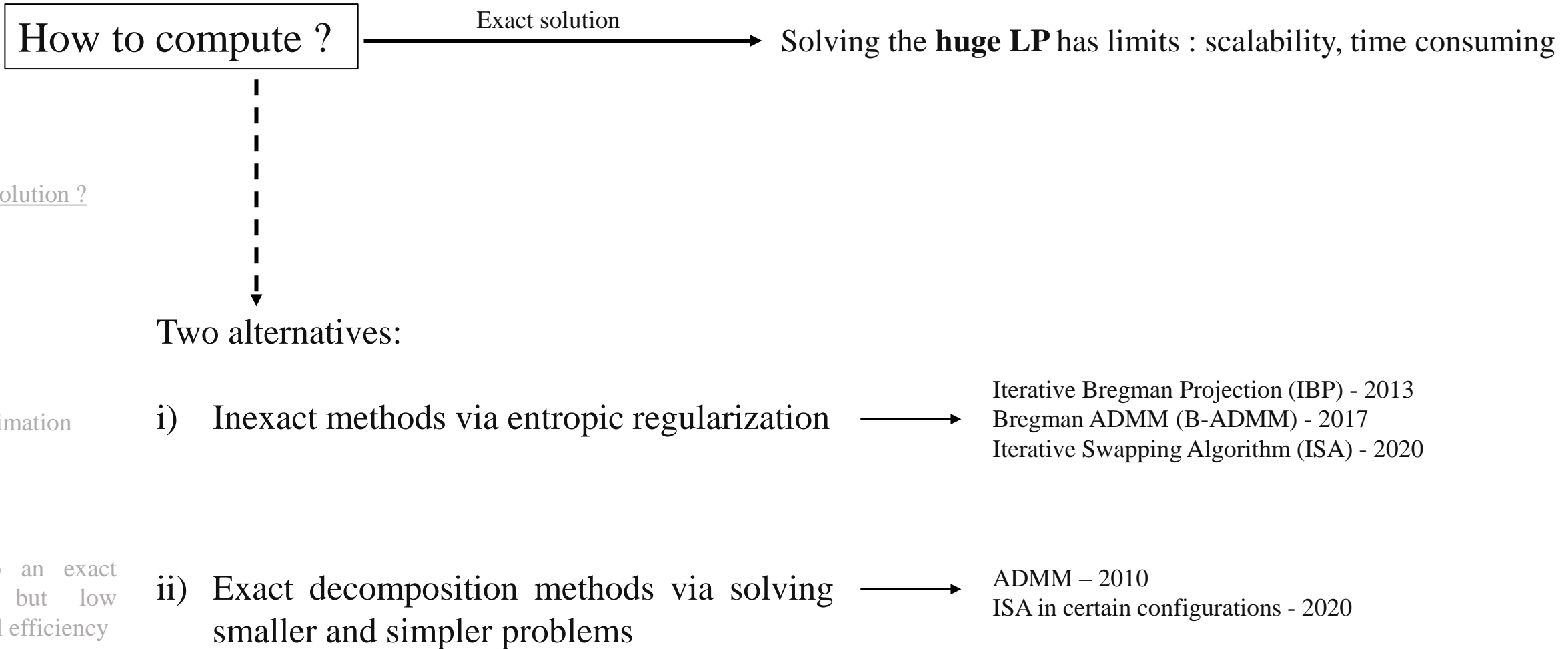
Block coordinate optimization:

- Step 1: support optimization \longrightarrow Straightforward solution exists if $\iota = 2$ (Euclidean norm) $\min_{\Xi} \sum_{m=1}^M \alpha_m \text{OT}_{\Xi, Z^{(m)}}(p^k, q^{(m)})$
- Step 2: probability optimization \longrightarrow $\min_{p \in \Delta_R} \sum_{m=1}^M \alpha_m \text{OT}_{\Xi, Z^{(m)}}(p, q^{(m)})$

Repeat until convergence.

II. Background on Discrete Optimal Transport

II.2. Wasserstein Barycenter (WB) - Previous works



III. The Method of Averaged Marginals (MAM)

III.1. Reformulation of the LP

Reformulation of the problem:

$$1. \quad f^{(m)}(\pi^{(m)}) := \sum_{r=1}^R \sum_{s=1}^{S^{(m)}} d_{rs}^{(m)} \pi_{rs}^{(m)} + \mathbf{i}_{\Pi^{(m)}}(\pi^{(m)})$$

$$2. \quad f(\pi) := \sum_{m=1}^M f^{(m)}(\pi^{(m)}) \quad \text{and} \quad g(x) := \mathbf{i}_{\mathcal{B}}(\pi)$$

$$3. \quad \min_{\pi} f(\pi) + g(\pi)$$

find π such that $0 \in \partial f(\pi) + \partial g(\pi)$

New problem : finding the zero of the sum of two maximal monotone operators

→ Several methods exist

→ Douglas-Rachford operator splitting is the most popular one (see ADMM or progressive hedging methods)

III. The Method of Averaged Marginals (MAM)

III.2. Douglas-Rachford (DR) theory

Douglas-Rachford steps:

given initial point $\theta^0 = (\theta^{(1),0}, \dots, \theta^{(M),0})$ and prox-parameter $\rho > 0$:

$$\left\{ \begin{array}{l} \pi^{k+1} = \text{prox}_{g/\rho}(\theta^k) \\ \hat{\pi}^{k+1} = \text{prox}_{f/\rho}(2\pi^{k+1} - \theta^k) \\ \theta^{k+1} = \theta^k + \hat{\pi}^{k+1} - \pi^{k+1} \end{array} \right. \begin{array}{l} \xrightarrow{\text{Projection onto } \mathcal{B} \text{ is explicit}} \\ \xrightarrow{\text{Projections onto the simplex}} \\ \Delta_n(\tau) := \left\{ u \in \mathbb{R}_+^n : \sum_{i=1}^n u_i = \tau \right\} \end{array} \begin{array}{l} \pi^{k+1} = \text{Proj}_{\mathcal{B}}(\theta^k) \\ \begin{pmatrix} \hat{\pi}_{1s}^{(m)} \\ \vdots \\ \hat{\pi}_{Rs}^{(m)} \end{pmatrix} = \text{Proj}_{\Delta_R(q_s^{(m)})} \begin{pmatrix} y_{1s} - \frac{1}{\rho} d_{1s}^{(m)} \\ \vdots \\ y_{Rs} - \frac{1}{\rho} d_{Rs}^{(m)} \end{pmatrix}, \quad s = 1, \dots, S^{(m)} \end{array}$$

III. The Method of Averaged Marginals (MAM)

III.3. Algorithm - Main steps

Step 1: Given a multi-transportation plan θ^k

- Marginals $p^{(m),k} = \theta^{(m),k} \mathbb{1}$, $m = 1, \dots, M$
- p^k is a weighted average of $\{p^{(1),k}, \dots, p^{(M),k}\}$

Step 2: Given θ^k , p^k and distance matrices

- Compute a multi-transportation plan π^k by performing $\sum_{m=1}^M S^{(m)}$ independent projections onto the simplex Δ_R

Step 3: Given θ^k , p^k and π^k

- Compute θ^{k+1} by a straightforward operation
- Set $k = k + 1$ and repeat

III. The Method of Averaged Marginals (MAM)

III.3. Algorithm - Feelings and philosophy

Algorithm 5.1 METHOD OF AVERAGED MARGINALS - MAM

1: Given $\rho > 0$, the distance matrix and initial point $d, \theta^0 \in \mathbb{R}^{R \times \sum_{m=1}^M S^{(m)}}$, and $a \in \Delta_M$ as in (5.3a), set $k \leftarrow 0$ and $p_r^{(m)} \leftarrow \sum_{s=1}^{S^{(m)}} \theta_{rs}^{(m),0}$, $r = 1, \dots, R$, $m = 1, \dots, M$ ▷ Step 0: input
 2: Set $\gamma \leftarrow \infty$ if $q^{(m)} \in \mathbb{R}_+^{S^{(m)}}$, $m = 1, \dots, M$, are balanced; otherwise, choose $\gamma \in (0, \infty)$
 3: **while** not converged **do** ▷ Step 1: average the marginals
 4: Compute $p^k \leftarrow \sum_{m=1}^M a_m p^{(m)}$
 5: Set $t^k = 1$ if $\rho \sqrt{\sum_{m=1}^M \frac{\|p^k - p^{(m)}\|^2}{S^{(m)}}} \leq \gamma$; otherwise, $t^k \leftarrow \gamma / \left(\rho \sqrt{\sum_{m=1}^M \frac{\|p^k - p^{(m)}\|^2}{S^{(m)}}} \right)$
 6: Choose an index set $\emptyset \neq \mathcal{M}^k \subseteq \{1, \dots, M\}$
 7: **for** $m \in \mathcal{M}^k$ **do** ▷ Step 2: update the m^{th} plan
 8: **for** $s = 1, \dots, S^{(m)}$ **do**
 9: Define $w_r \leftarrow \theta_{rs}^{(m),k} + 2t^k \frac{p_r^k - p_r^{(m)}}{S^{(m)}} - \frac{1}{\rho} d_{rs}^{(m)}$, $r = 1, \dots, R$
 10: Compute $(\hat{\pi}_{1s}^{(m)}, \dots, \hat{\pi}_{Rs}^{(m)}) \leftarrow \text{Proj}_{\Delta_R(q_s^{(m)})}(w)$ → Projection onto the simplex performed exactly by using efficient methods
 11: Update $\theta_{rs}^{(m),k+1} \leftarrow \hat{\pi}_{rs}^{(m)} - t^k \frac{p_r^k - p_r^{(m)}}{S^{(m)}}$, $r = 1, \dots, R$ → Projection onto \mathcal{B}
 12: **end for**
 13: Update $p_r^{(m)} \leftarrow \sum_{s=1}^{S^{(m)}} \theta_{rs}^{(m),k+1}$, $r = 1, \dots, R$ ▷ Step 3: update the m^{th} marginal
 14: **end for**
 15: **end while**
 16: Return $\bar{p} \leftarrow p^k$

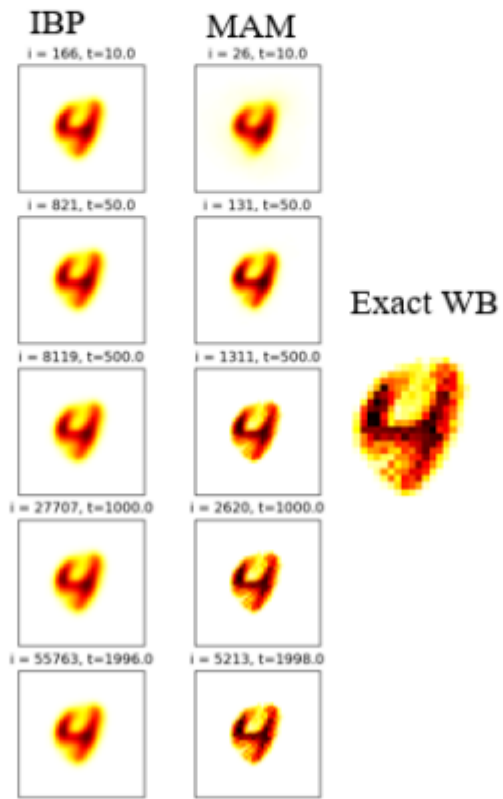
Unbalanced formulation

Can be executed in parallel or randomized

IV. Applications

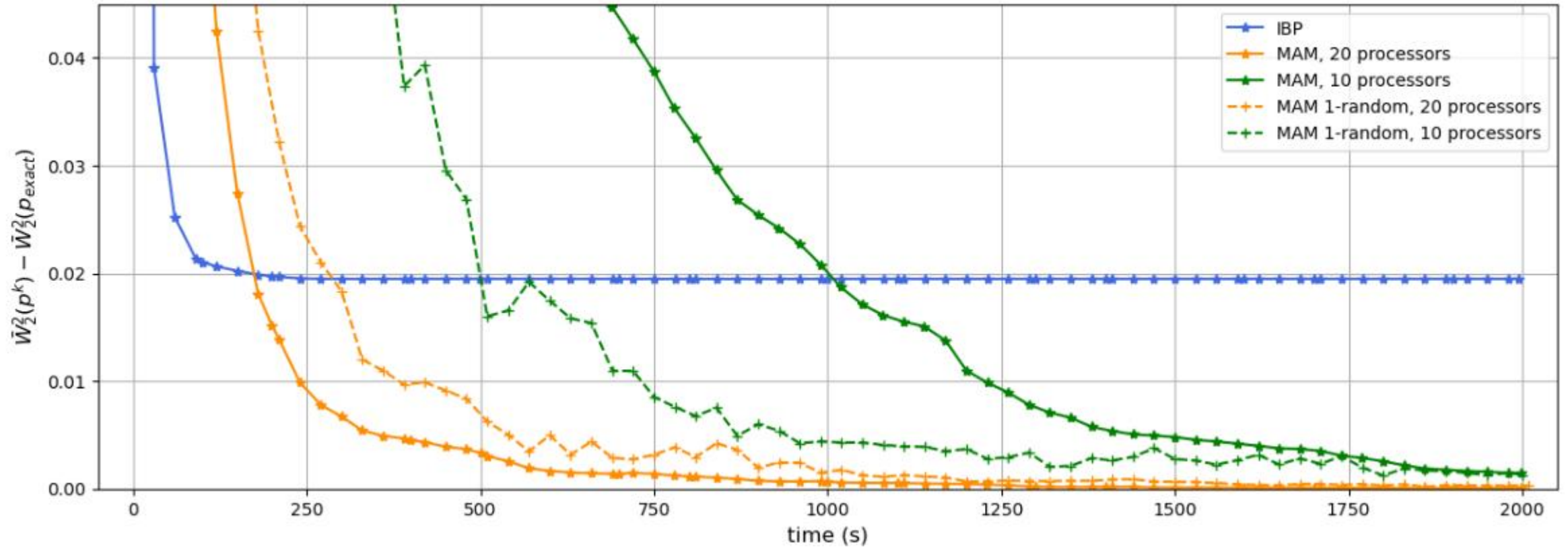
IV.1. Qualitative comparison

MAM	IBP
Exact algorithm	Iterative Bregman Projection <ul style="list-style-type: none"> state-of-the-art algorithm for WB based on an entropic regularization of the problem thus computes inexact WB



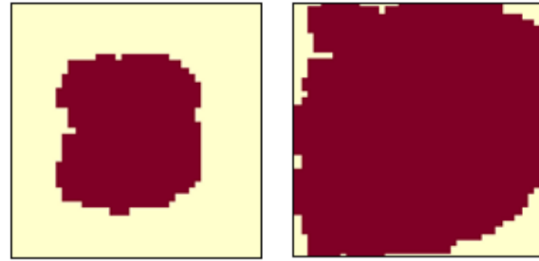
IV. Applications

IV.2. Quantitative comparison

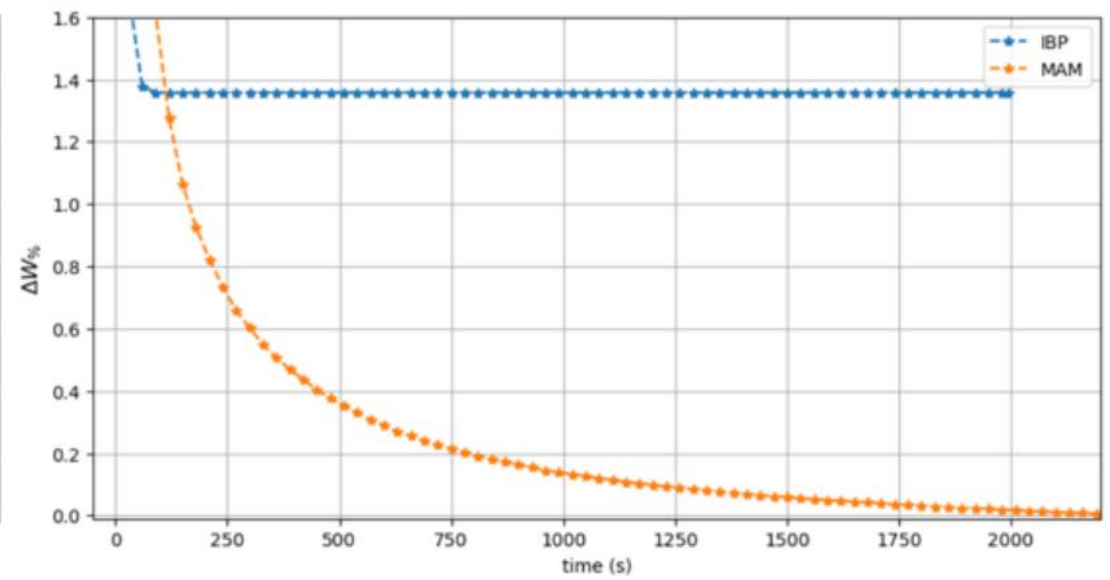
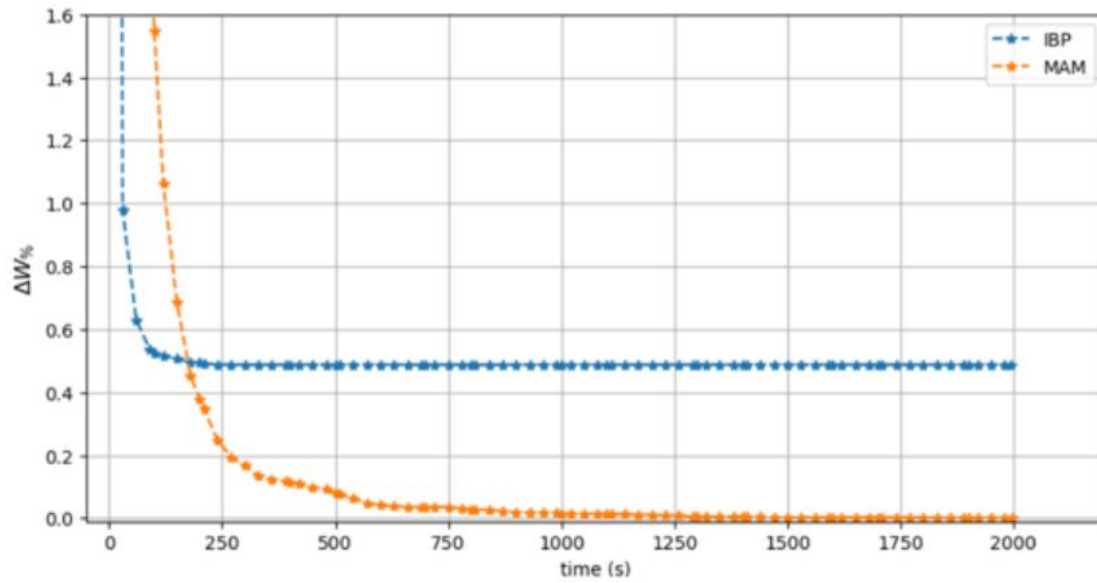


IV. Applications

IV.3. Influence of the support



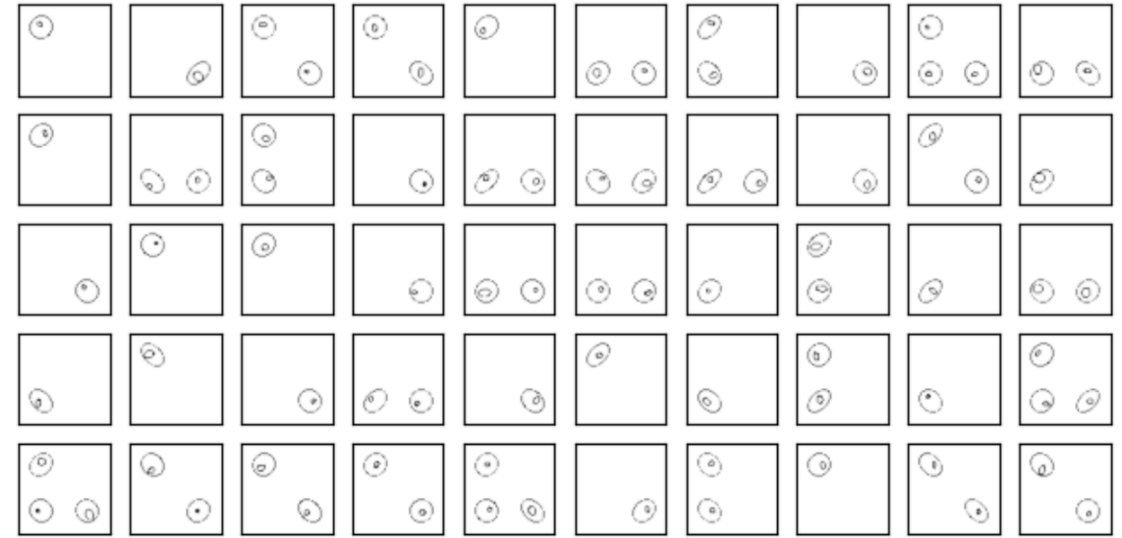
Union of the dataset support



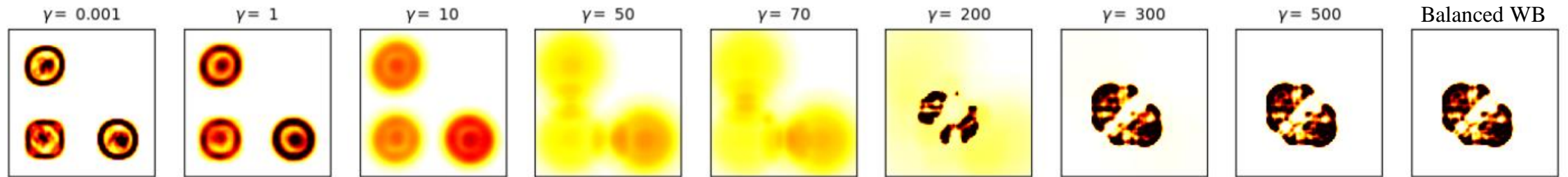
IV. Applications

IV.4. Unbalanced Wasserstein Barycenter

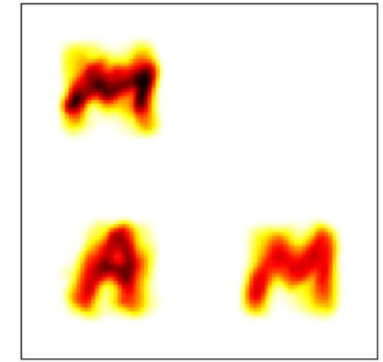
Dataset composed by 50 pictures with nested ellipses randomly positioned in the top left, bottom right and left corners :



The standard (balanced) WB is not always the best tool for summarizing:



Conclusion and future works



What have been introduced?

A novel approach for computing Wasserstein barycenters of discrete measures, that:

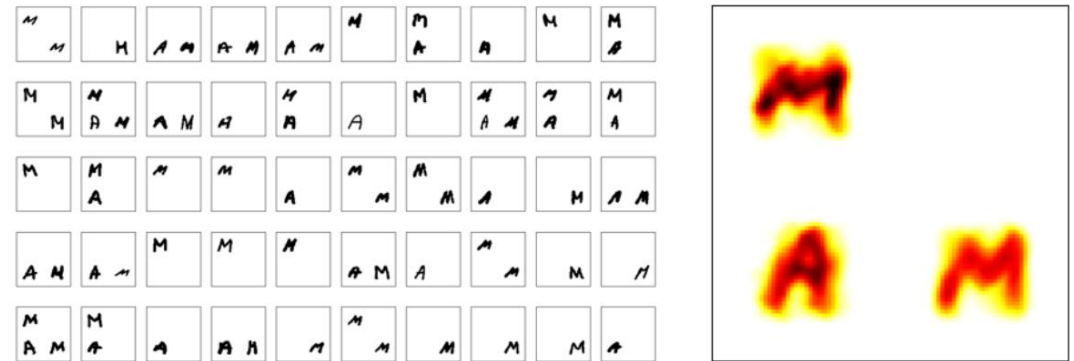
- Asymptotically exact
- Embarrassingly parallelizable and can be used on a randomized manner (almost surely convergence)
- Can tackle both the balanced and unbalanced case!

What can be done now?

- Adapt MAM to tackle scenario trees reduction problem in stochastic optimization
- Real life examples



References



- [1] Cuturi, M., & Doucet, A. (2014, June). Fast computation of Wasserstein barycenters. In *International conference on machine learning* (pp. 685-693). PMLR
- [2] Bonneel, N., Peyré, G., & Cuturi, M. (2016). Wasserstein barycentric coordinates: histogram regression using optimal transport. *ACM Trans. Graph.*, 35(4), 71-1.
- [3] El Hamri, M., Bennani, Y., & Falih, I. (2022). Hierarchical optimal transport for unsupervised domain adaptation. *Machine Learning*, 111(11), 4159-4182.
- [4] Mimouni, D., Malisani, P., Zhu, J., & de Oliveira, W. (2023). Computing Wasserstein Barycenter via operator splitting: the method of averaged marginals. *arXiv preprint arXiv:2309.05315*.