

Stochastic Wasserstein Barycenters (Computational Optimal Transport Project)

Inès VATI ines.vati@eleves.enpc.fr
MVA, ENS Paris Saclay

January 17, 2024

Abstract

In this report, we studied the article [1] entitled "Stochastic Wasserstein Barycenters" published in 2018. Claici et al. proposed an innovative approach to compute the barycenter of a set of discrete or continuous probability measures under the Wasserstein metric from optimal transport theory in which the support of the barycenter is adjusted at each iteration. Their algorithm is straightforward parallelizable, easy to implement and is not restricted to input distributions supported on a fixed set of points.

We applied their method to both simple 1D and 2D cases. Additionally, we conduct a comparative analysis, evaluating the performance of this approach against existing methods in terms of both quality and computational speed.

The code is available at https://github.com/InesVATI/projectOT_stochastic_wasserstein_barycenters.

1 Introduction

Addressing the computation of barycenters for probability distributions stands as a fundamental problem in statistics and machine learning. It becomes compelling to summarize a collection of probability distributions with due consideration to the inherent geometric structure.

The theory of optimal transport (OT) provides a promising and theoretically-grounded avenue for averaging distributions over a geometric domain. The Wasserstein barycenter is a generalization of the notion of mean to probability measures. It is defined as the minimizer of the sum of the Wasserstein distances to the input measures. The Wasserstein distance itself serves as a metric on the space of probability measures that takes into account the underlying geometry of the space. It is defined as the minimum cost of transporting one measure onto the other.

The computation of barycenters finds diverse applications across various fields, including statistics, computer vision [2, 3], signal analysis [4, 5], and medical imaging [6]. For instance, Bruckstein et al. [2] propose a new approach to compute the average of discrete probability distributions as a barycenter over the Wasserstein space and apply their method to texture synthesis and texture mixing. Texture mixing problem consists in synthesizing a new texture from a collection of atoms, i.e. exemplars. In [4], the authors leverage these tools in time series modeling. Each segment in time series data can represent a state, such as running or walking in a human activity application. They propose an innovative Dynamical Wasserstein Barycenter model to estimate the distribution of pure states while improving state estimation for transition periods.

Related works. Several works have addressed the computation of Wasserstein barycenter. Cuturi et al. [7] presented a concise and efficient Sinkhorn algorithm, subsequently extended to barycenter problems through, for instance, iterative Bregman projection algorithms [8]. These algorithms introduce entropic regularization to the initial linear problem, demonstrating that the set of linear constraints can be split in an intersection of simpler constraints, allowing for closed-form projections.

Staib et al. [9] introduced a stochastic barycenter algorithm from samples, presenting a scalable and parallelized approach suitable for streaming data—continuously generated data from diverse sources. Their method is also robust to nonstationary input distributions. However, their method requires a finite, predetermined set of support points.

Several works have aimed at enhancing algorithmic speed. Dvurechenskii et al. [6] proposed a distributed algorithm for computing a discrete approximation of the regularized Wasserstein barycenter for a set of continuous probability distributions stored across a network, i.e. each agent constituting the network holds a private continuous probability measure. Grounded in an accelerated primal-dual stochastic gradient method for convex optimization with linear equality constraints, their approach seeks to expedite computation in a decentralized fashion.

Current state-of-the-art techniques typically define the barycenter on a fixed, refined grid and optimize the weights associated with each grid point. In this studied article [1], the optimization is extended to include the grid itself, and uniform weights are assigned to each grid point.

Problem statement. We consider of collection of J distributions $\{\mu_j\}_{j=1}^J$, either discrete or continuous, defined on a common domain $\mathcal{X} \subset \mathbb{R}^d$. The barycenter ν of these distributions is defined as the solution of the follow-

ing optimization problem

$$\min_{\nu} F_{\text{primal}} = \frac{1}{J} \sum_{j=1}^J W_2^2(\nu, \mu_j) \quad (1)$$

The squared 2-Wasserstein distance between two probability measures ν and μ is defined as

$$W_2^2(\mu, \nu) = \min_{\pi \in \mathbf{U}(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} d(x, y)^2 d\pi(x, y) \quad (2)$$

where $\mathbf{U}(\mu, \nu)$ is the set of all coupling π between μ and ν

$$\mathbf{U}(\mu, \nu) = \{\pi \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y}), (P_{\mathcal{X}})_{\#}\pi = \mu \text{ and } (P_{\mathcal{Y}})_{\#}\pi = \nu\}$$

$(P_{\mathcal{X}})_{\#}$ and $(P_{\mathcal{Y}})_{\#}$ are the push-forward by the projectors $P_{\mathcal{X}}$ and $P_{\mathcal{Y}}$ on the first and second marginals respectively.

They opt to a semidiscrete approximation of the barycenter. Let $\Sigma = \{x_i\}_{i=1}^M$ be a set of M points in \mathcal{X} . The barycenter ν is approximated by a discrete measure $\hat{\nu}$ with uniform weights on the points of $\Sigma \subset \mathcal{X}$, i.e.

$$\hat{\nu} = \frac{1}{M} \sum_{i=1}^M \delta_{x_i}$$

In this setting, the optimization problem reads

$$\min_{\forall i \leq M, x_i \in \mathcal{X}} \frac{1}{J} \sum_{j=1}^J W_2^2(\mu_j, \hat{\nu}) \quad (3)$$

The dual problem of each sub problem of (3) reads

$$\begin{aligned} \max_{\phi_j \in L^1(\mathcal{X})} F_{\text{dual}}(\phi_j, \{x_i\}_{i=1}^M) &= \int_{\mathcal{X}} \phi_j(x) d\nu(x) + \int_{\mathcal{X}} \bar{\phi}_j(y) d\mu_j(y) \\ &= \frac{1}{M} \sum_{i=1}^M \phi_j(x_i) + \int_{\mathcal{X}} \bar{\phi}_j(y) d\mu_j(y) \end{aligned}$$

where $\phi_j \in \mathbb{R}^M$ is the discrete Kantorovith potential associated with the OT problem (2) and $\bar{\phi}_j$ is the c-transform of ϕ_j defined by $\bar{\phi}(y) = \inf_{x \in \mathcal{X}} d(x, y) - \phi(x)$ [10].

Therefore, the objective function optimize in this work [1] is

$$F(\{\phi_j\}_{j=1}^J, \{x_i\}_{i=1}^M) = \frac{1}{J} \sum_{j=1}^J F_{\text{dual}}(\phi_j, \{x_i\}_{i=1}^M) \quad (4)$$

Contributions. The main contribution of the research conducted in [1] lies in the introduction of a stochastic approach for computing Wasserstein barycenters, thereby eliminating the need for a fixed grid implementation. Notably, the support of the estimated barycenter is optimized and empirically demonstrated to be contained within the support of the true barycenter. In a departure from conventional approaches, their problem formulation eschews regularization. The authors also provide a theoretical analysis elucidating the algorithm’s convergence dynamics.

In this report, we delve into the method proposed by [1], applying it to both simple 1D and 2D cases. Additionally, we conduct a comparative analysis, evaluating the performance of this approach against existing methods in terms of both quality and computational speed¹.

2 Presentation of the method

The proposed algorithm [1] boils down to two key steps :

1. **Stochastic Gradient Ascent for Potentials:** With $\{x_i\}$ fixed, a stochastic gradient ascent is employed to optimize the potentials $\{\phi_j^i\}$, leveraging the concavity of F (4) in ϕ_j .
2. **Fixed-Point Iteration for Positions:** With $\{\phi_j^i\}$ fixed, a single fixed point iteration is performed to update the barycenter grid $\{x_i\}$ using the Reynolds’ transport theorem [1]. This update is facilitated by a closed-form expression for the zeroed values of the gradient, i.e.,

$$\frac{\partial F}{\partial x_i} = 0$$

2.1 Stochastic Gradient Ascent for Potentials

For step 1, we note that F is concave in the potentials ϕ_j . Indeed, we have shown in the course [10] that the c -transform is concave for the euclidean cost $d(x, y)^2 = \|x - y\|_2^2$. To get the gradient, we need to compute the quantities

$$\begin{aligned} a_j^i &= \int_{V_{\phi_j}^i} d\mu_j(y) & b_j^i &= \int_{V_{\phi_j}^i} y d\mu_j(y) \\ &= \mathbf{E}_{y \sim \mu_j} \left[\mathbb{1}_{y \in V_{\phi_j}^i} \right] & &= \mathbf{E}_{y \sim \mu_j} \left[y \cdot \mathbb{1}_{y \in V_{\phi_j}^i} \right] \end{aligned}$$

¹Code is available at https://github.com/InesVATI/projectOT_stochastic_wasserstein_barycenters

where $\mathbb{1}$ indicates the indicator function of a set. In the derivative of (4), the *power cell* $V_{\phi_j}^i$ of point x_i comes into play and is defined as follows

$$V_{\phi}^i = \{x \in \mathcal{X}, d(x, x_i)^2 - \phi_i \leq d(x, x_{i'})^2 - \phi_{i'}, \forall i'\}$$

The computation of b_j^i is given in the algorithm 1, the computation of a_j^i is similar.

Algorithm 1 Estimate a_j^i and b_j^i by Monte Carlo approximation

Require: iid samples $\{Y\}_{k=1}^K \sim \mu_j$

```

1: for  $y$  in  $Y$  do
2:   if  $y \in V_{\phi_j}^i$  then
3:      $b_j^i \leftarrow b_j^i + y$ 
4:   end if
5: end for
6:  $b_j \leftarrow \frac{1}{K} b_j$ 
7: return  $b_j^i$ 
```

The well known gradient ascent step is given by

$$\phi_j^{(l+1)} = \phi_j^{(l)} + \alpha \frac{\partial F}{\partial \phi_j}(\phi_j^{(l)})$$

at iteration l , where λ is the step size. To improve performance, the authors apply the Nesterov acceleration. The *ascent step* is given in algorithm 2.

Algorithm 2 Ascent step

```

1: for  $j = 1, \dots, J$  do
2:    $z^{(0)} \leftarrow 0$ 
3:   while  $\left\| \widehat{\frac{\partial F}{\partial \phi_j}} \right\|_2^2 > \epsilon_{ascent}$  do
4:     Compute  $\hat{a}_j^i$  thanks to algorithm 1
5:      $z^{l+1} \leftarrow \beta z^{(l)} + \widehat{\partial_{\phi_j} F}(\phi_j)$  {Nesterov acceleration}
6:      $\phi_j^{(l+1)} \leftarrow \phi_j^{(l)} + \alpha z^{(l+1)}$ 
7:   end while
8: end for
```

In algorithm 2, $\widehat{\partial_{\phi_j} F}$ denote the stochastic approximation of the gradient of F with respect to ϕ_j .

2.2 Fixed Point Iteration for Positions

The update of the point $\{x_i\}$ is given by

$$\frac{\partial F}{\partial x_i} = 0 \implies x_i = \frac{\sum_{j=1}^J b_j^i}{\sum_{j=1}^J a_j^i}$$

This step boils down to a simple fixed point iteration and is denoted as the *snap step*. The algorithm 3 summarize the whole procedure.

Algorithm 3 Ascent and Snap algorithm for computing Stochastic Wasserstein Barycenters

```

1: for  $t = 1, \dots, T$  do
2:   Update  $\{\phi_j\}$  with algorithm 2 {Ascent step}
3:   Compute  $\hat{b}_j^i$  with algorithm 1
4:   for  $x_i \in \Sigma$  do
5:      $x_i \leftarrow \frac{\sum_{j=1}^J \hat{b}_j^i}{\sum_{j=1}^J \hat{a}_j^i}$  {Snap step}
6:   end for
7: end for
8: return Optimized barycenter support  $\Sigma^* = \{x_i\}_{i=1}^M$ 

```

2.3 Theoretical guarantees

Well behavior of the algorithm is ensured for the following assumptions.

1. At least one of the input distribution μ_j is absolutely continuous with respect to the Lebesgue measure.
2. The domain \mathcal{X} is a compact subset of the Euclidean space \mathbb{R}^d .

The first assumption ensures the existence of a unique Wasserstein barycenter [1].

Under these conditions, the authors prove that the Algorithm 3 converges to a local minimum of the functional F_{primal} (1) for a step size equal to the inverse of the Lipschitz constant of the gradient of F_{primal} given in [1].

However, empirical evidence presented by the authors indicates that their method converges even in cases where these assumptions are not met.

Moreover, Let ν_M^* be a uniform measure supported on M points that minimizes $\frac{1}{J} \sum_{j=1}^J W_2(\nu_M^*, \mu_j)$. Remind that ν denotes the true barycenter of the measures μ_j . Since the Wasserstein distance metrizes the convergence in law, they have that $\nu_M^* \rightarrow \nu$.

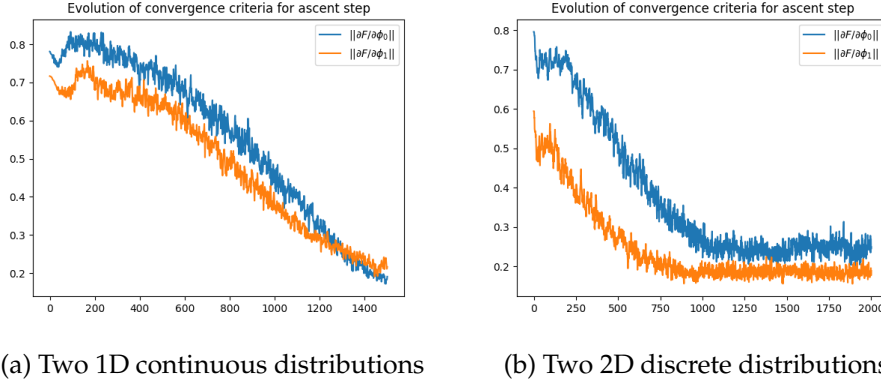


Figure 1: Evolution of $\left\| \frac{\partial F}{\partial \phi_j} \right\|$ during the ascent step (Algorithm 2) for each experiment. In the x-axis, is the number of iterations of the while loop.

While the authors do not provide guarantees that the algorithm converges toward ν_M^* , they show that the global minimizer of (1) efficiently approximate the true barycenter as $M \rightarrow \infty$.

Furthermore, they theoretically establish that their algorithm monotonically decreases the objective of the primal problem (1).

3 Experiences and Numerical Analysis

3.1 Analysis of the ascent step convergence

During the iterative procedure, we plot in Figure 1 the evolution of the norm of the gradient of F , for $j \in \llbracket 1, J \rrbracket$. The stochastic ascent gradient seems to converge as the norm of the gradient reaches a plateau. The convergence is faster for the 2D discrete distributions than for the 1D continuous distributions.

It is worth noting that the majority of the computation time is consumed by the ascent step (several hours) as shown in Table 1.

In [1], the authors use a step size of $\alpha = 10^{-3}$ and a stop criteria of $\epsilon_{ascent} = 10^{-6}$ in Algorithm 2. However, the while loop took an excessive amount of time to converge, exceeding 6 hours. To enhance computational efficiency, we opt for $\alpha = 0.05$ and $\epsilon_{ascent} = 10^{-4}$. Nesterov acceleration is set to $\beta = 0.99$ like in [1]. Additionally, we introduce a maximum number of iterations, limiting the while loop to $T_{max} = 1500$ for efficiency.

We employ a single iteration for both of our experiments, as the authors [1] assert that the entire algorithm converges within a single iteration.

Experiment	Ascent step	Snap step
1D skewnorm	8309.44 ± 0.33	5.28 ± 0.23
2D discrete	11938.70 ± 0.34	9.81 ± 0.34

Table 1: Computation time (in seconds) for each step of the Algorithm 3. The experiences are repeated 3 times.

3.2 Quantitative Analysis and Comparative Evaluation

3.2.1 Optimal Monge Map Approach to Barycenter Computation

In the case $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ and $d(x, y)^2 = \|x - y\|^2$ and if at least one of the two input measures has a density with respect to the Lebesgue measure, the Theorem 2.1 of [10] states that the optimal coupling π of (2) is unique and is given by $\pi = (Id, T)_\# \mu$ where $T : \mathcal{X} \rightarrow \mathcal{Y}$ denote the "optimal Monge map" with $T_\# \mu = \nu$.

In the 1D case [see 10, Remark 2.30], $d = 1$, the optimal Monge map between two distributions μ_1 and μ_2 writes

$$T = \mathcal{C}_{\mu_2}^{-1} \circ \mathcal{C}_{\mu_1} \quad (5)$$

where $\mathcal{C}_\mu : \mathbb{R} \rightarrow [0, 1]$ and $\mathcal{C}_\mu^{-1} : [0, 1] \rightarrow \mathbb{R}$ are respectively the cumulative distribution function and its pseudoinverse, also called the generalized quantile function of μ .

Using the Remark 7.1 [10], the McCann's interpolation [11] between two measures reads, for $t \in [0, 1]$

$$\begin{aligned} \mu_t &= (tT + (1 - t)Id)_\# \mu_1 \\ \mu_t &= (t\mathcal{C}_{\mu_2}^{-1} \circ \mathcal{C}_{\mu_1} + (1 - t)Id)_\# \mu_1 \end{aligned} \quad (6)$$

Figure 2 shows the displacement interpolation between 1D measures, using the cumulative distribution function as detailed in (6).

To obtain the interpolant, we generate K samples (X_k) independent and identically distributed (iid) according to μ_1 . Then, we compute the samples (Y_k) as follows

$$Y_k = tTX_k + (1 - t)X_k, \quad \forall k \in \llbracket 1, K \rrbracket$$

where T is given in (5). We have that (Y_k) is iid according to μ_t .

To compare with the algorithm 3, we use $t = 0.5$. The only hyperparameter to fix is the number of samples K to generate. We choose $K = 2000$. The inconvenient is that we can apply this computation only in the 1D case and for two input measures.

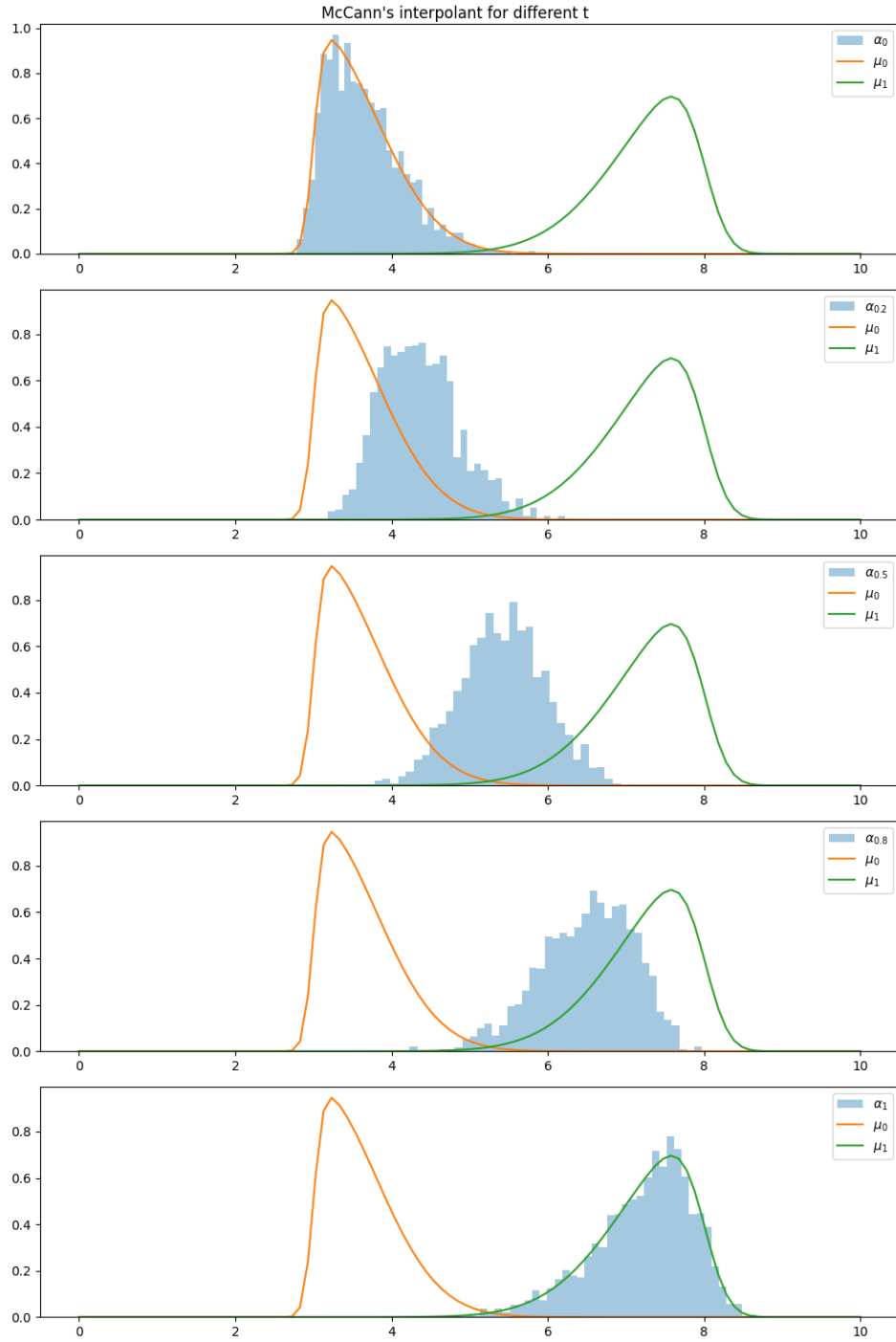


Figure 2: McCann's interpolant between two skew-normal distributions. From top to bottom, $t=0, 0.2, 0.5, 0.8, 1$. The barycenter that we will consider corresponds to $t = 0.5$, depicted in the middle plot.

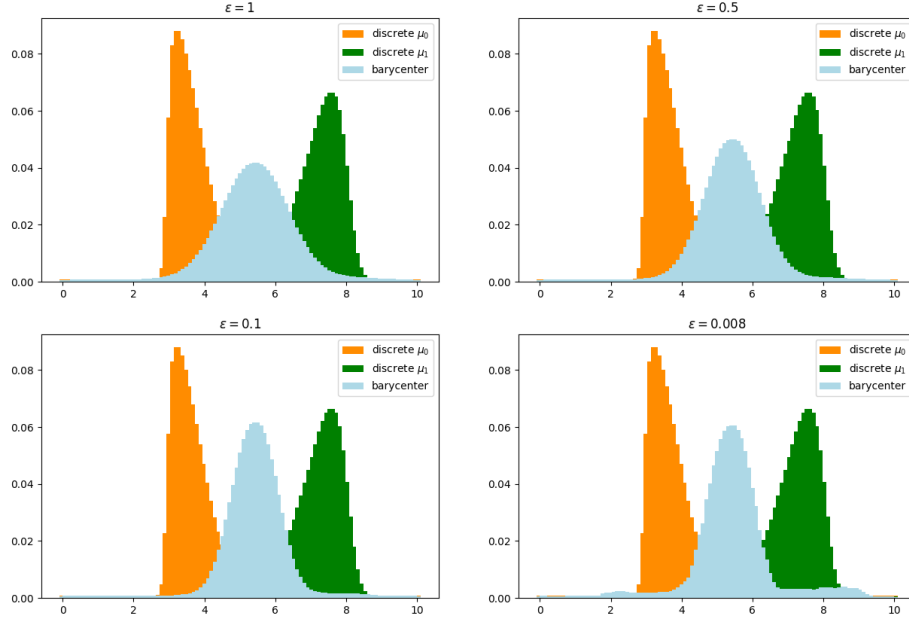


Figure 3: Histogram showing the input distributions (skew-normal distributions) and the barycenter computed using the Sinkhorn algorithm for different values of the regularized parameter ϵ .

3.2.2 Use of the Sinkhorn algorithm to compute the barycenter

We will also compare our algorithm 3 to the computation of the barycenter using the Sinkhorn algorithm [10]. This algorithm solved the discretized regularized optimal transport problem using the optimality condition that shows that the optimal coupling P_ϵ necessarily has the form

$$P_\epsilon = \text{diag}(u) K \text{diag}(v)$$

where the Gibbs kernel is defined as

$$K := e^{-\frac{C}{\epsilon}}.$$

where $C_{ij} = d(x_i, y_j)^2$ is the cost matrix. The vectors u and v are non negative vectors. ϵ is the regularization parameter.

Figure 3 depicts the barycenter obtained for various ϵ . It's worth noting that the algorithm encounters instability issues for ϵ values below 0.005. As illustrated, the lower ϵ , the sharper the density of the barycenter.

In our experiments, we use 500 iterations and we fix $\epsilon = 0.01$ for the 1D case and $\epsilon = 0.04^2$ for the 2D case.

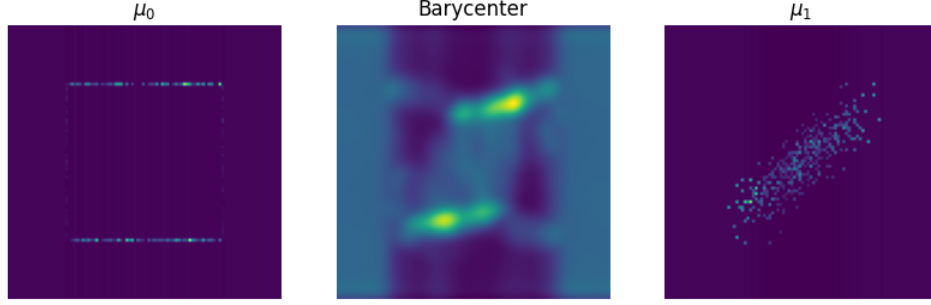


Figure 4: Results of the Sinkhorn algorithm for the 2D discrete case. The regularization is $\epsilon = 0.04^2$.

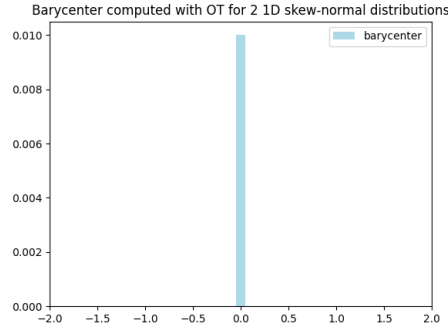


Figure 5: Results of the POT toolbox for the 1D continuous case.

3.2.3 A free grid Sinkhorn algorithm approach from POT toolbox

The POT toolbox (Python Optimal Transport) [12] contains implementations of a number of founding works of OT for machine learning such as Sinkhorn algorithm and Wasserstein barycenters. In particular, it provides an algorithm² based on [13] (Algorithm 2) that solves the free support regularized Wasserstein barycenter problem. Like the method §2, we optimize over the locations of the barycenter but not over its weights.

The figures 8, 9 and 10 in the Appendices shows the influence of the regularization ϵ on the free support barycenters for the 2D dataset.

3.2.4 Experience 1 : 1D skewed normal distributions

In this experience, we compute the barycenter between two 1D skewed normal distributions whose skew parameters are $s_1 = 8$ and $s_2 = 4$, means

²https://pythonot.github.io/gen_modules/ot.bregman.html#ot.bregman.free_support_sinkhorn_barycenter

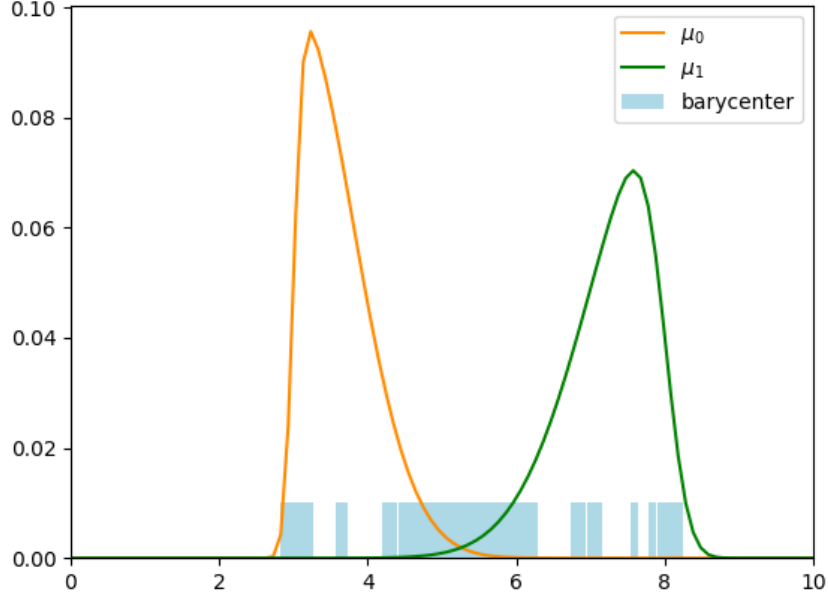


Figure 6: Results of the ascent/snap Algorithm 3 for the 1D continuous case. We use 16000 samples from the input distributions, and 1 iterations. See Section 3.1 for the hyperparameters.

are $m_1 = 3$ and $m = 8$ and standard deviations are $\sigma_1 = 0.8$ and $\sigma_2 = 1$, respectively.

The figures 2, 5, 6 and 3 depict the results for each methods. The support of the barycenter computed by the proposed method is larger than the support of the McCann’s interpolation and the one obtained with the Sinkhorn algorithm.

Moreover, uniform weights does not seem to be a good way to approximate barycenters of non-uniform distributions.

Running the free support sinkhorn algorithm from the POT toolbox did not yield any results as can be see on the Figure 5.

Regarding the computation time, the ascent/snap algorithm is by far the slowest algorithm. The table 2 shows the computation time for each method in each test case.

3.2.5 Experience 2 : 2D uniform distributions

In the second experiment, we generate two 2D uniform distributions. We generate 160 points from each distribution with a different geometrical structures. The first distribution is a square and the points from the second

Experiment	McCann	Sinkhorn	Free support Sinkhorn	Ascent Snap
1D skewnorm	0.02 ± 0.01	0.22 ± 0.02	0.74 ± 0.28	7814.27 ± 0.37
2D discrete	-	2.18 ± 0.34	0.38 ± 0.30	16351.66 ± 0.23

Table 2: Computation time (in seconds) for the different methods. The experiences were repeated 3 times for the Ascent snap algorithm and 10 times for the others.

distribution are generated from a multivariate normal distribution with mean $(0, 0)$ and covariance matrix $\Sigma = \begin{pmatrix} 0.6 & 0.5 \\ 0.5 & 0.6 \end{pmatrix}$.

We compare the barycenter of the proposed method showed in Figure 7 with the barycenter plotted Figure 9 with a regularization of $\epsilon = 0.005$. The barycenter computed by the proposed method is slightly different from the one obtained with the free support Sinkhorn algorithm from POT §3.2.3. Our algorithm gives point inside the warped square whereas the POT barycenter only warped the segments. However, their shapes are similar.

As illustrated in The Figure 4, we had to construct an image to apply the Sinkhorn algorithm to 2D distributions corresponding to this scheme. However, this algorithm do not yield qualitative results as the Gibbs kernel tend to smooth the image.

4 Conclusion and Perspectives

We have applied the proposed method on both 1D and 2D cases. Even though the algorithm takes a lot of time to run, the results remains promising. This method is slower but more robust and versatile than the other methods tested. We also have shown that the support of the barycenter is not necessarily include in the support of the true barycenter as shown the comparaisn between the McCann’s interpolant (Fig. 2) and the barycenter computed with the proposed method (Fig 6) in the 1D case. A possible improvement of our result could be a better choice of the initialization of the positions x_i .

For computation time sake, we have computed the barycenter between two input distributions but it would be interesting to compare the results with three or more input measures. This extension would allow us to assess the robustness and scalability of the proposed method across a broader

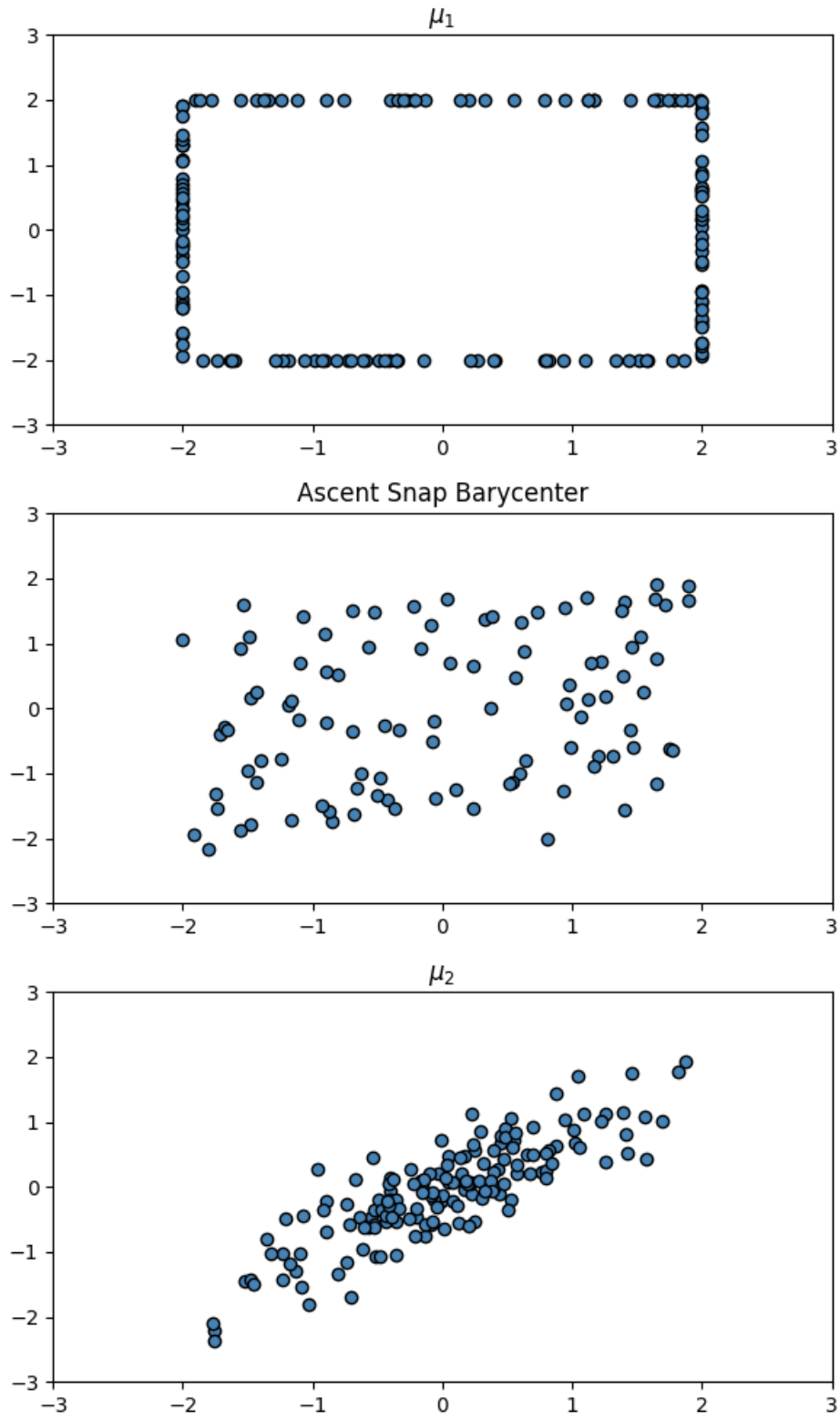


Figure 7: Results of the ascent/snap Algorithm 3 for the 2D discrete case. We use 16000 samples from the input distributions, and 1 iteration.

range of input scenarios.

Futhermore, comparing the proposed method with the log-domain Sinkhorn algorithm would also be interesting. The log-domain Sinkhorn algorithm is known to be more stable than the Sinkhorn algorithm for low regularization parameters but also slower.

5 Connection with the course

This study directly aligns with the course material. The article [1] is based on the notion of Wasserstein distance introduced during the course.

Additionally, we have conducted a comparative analysis of the proposed method with two approaches presented in the course.

We employed the Sinkhorn algorithm to calculate the barycenter of two measures in both 1-D and 2-D. In the case of 1-D, we have also computed the McCann's interpolation between two continuous measures, as depicted in Figure 3.

Regarding the results, we observed the same influence of the regularization parameter ϵ on the Sinkhorn algorithm, as illustrated in Figure 3.

6 Conclusion and Perspectives

We have applied the proposed method on both 1D and 2D cases. Even though the algorithm takes a lot of time to run, the results remains promising. This method is slower but more robust and versatile than the other methods tested. We also have shown that the support of the barycenter is not necessarily include in the support of the true barycenter as shown the comparaisn between the McCann's interpolant (Fig. 2) and the barycenter computed with the proposed method (Fig 6) in the 1D case. A possible improvement of our result could be a better choice of the initialization of the positions x_i .

For computation time sake, we have computed the barycenter between two input distributions but it would be interesting to compare the results with three or more input measures. This extension would allow us to assess the robustness and scalability of the proposed method across a broader range of input scenarios.

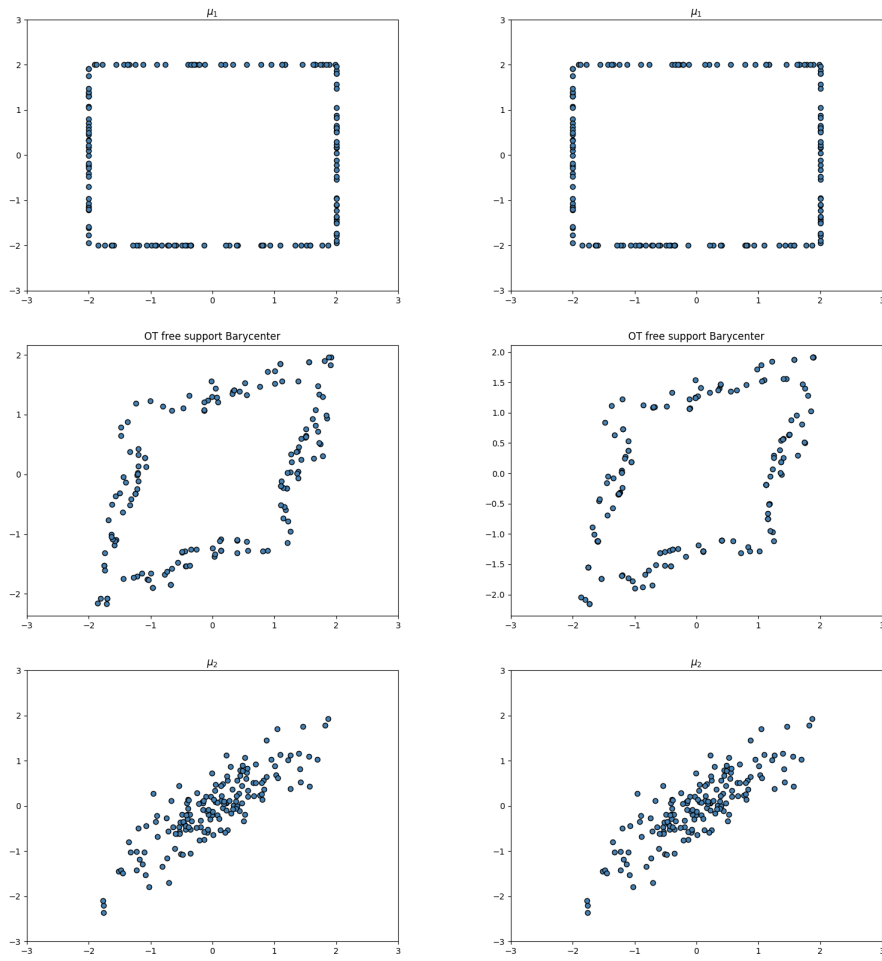
Futhermore, comparing the proposed method with the log-domain Sinkhorn algorithm would also be interesting. The log-domain Sinkhorn algorithm is known to be more stable than the Sinkhorn algorithm for low regularization parameters but also slower.

References

- [1] S. Clatici, E. Chien, J. Solomon, Stochastic Wasserstein Barycenters, **2018**.
- [2] J. Rabin, G. Peyré, J. Delon, M. Bernot in *Scale Space and Variational Methods in Computer Vision*, Vol. 6667, (Eds.: A. M. Bruckstein, B. M. Ter Haar Romeny, A. M. Bronstein, M. M. Bronstein), Series Title: Lecture Notes in Computer Science, Springer Berlin Heidelberg, Berlin, Heidelberg, **2012**, pp. 435–446.
- [3] A. Korotin, V. Egiazarian, L. Li, E. Burnaev.
- [4] K. Cheng, S. Aeron, M. C. Hughes, E. L. Miller in Advances in Neural Information Processing Systems, Vol. 34, Curran Associates, Inc., **2021**, pp. 27991–28003.
- [5] K. C. Cheng, S. Aeron, M. C. Hughes, E. L. Miller, Nonparametric and Regularized Dynamical Wasserstein Barycenters for Sequential Observations, **2023**.
- [6] P. Dvurechenskii, D. Dvinskikh, A. Gasnikov, C. Uribe, A. Nedich in Advances in Neural Information Processing Systems, Vol. 31, Curran Associates, Inc., **2018**.
- [7] M. Cuturi, Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances, **2013**.
- [8] J.-D. Benamou, G. Carlier, M. Cuturi, L. Nenna, G. Peyré, Iterative Bregman Projections for Regularized Transportation Problems, **2014**.
- [9] M. Staib, S. Clatici, J. Solomon, S. Jegelka, Parallel Streaming Wasserstein Barycenters, **2017**.
- [10] G. Peyré, M. Cuturi, Computational Optimal Transport, **2020**.
- [11] R. J. McCann, *Advances in Mathematics* **1997**, 128, 153–179.
- [12] R. Flamary, N. Courty, A. Gramfort, M. Z. Alaya, A. Boissunon, S. Chambon, L. Chapel, A. Corenflos, K. Fatras, N. Fournier, L. Gautheron, N. T. H. Gayraud, H. Janati, A. Rakotomamonjy, I. Redko, A. Rolet, A. Schutz, V. Seguy, D. J. Sutherland, R. Tavenard, A. Tong, T. Vayer, *Journal of Machine Learning Research* **2021**, 22, 1–8.
- [13] M. Cuturi, A. Doucet, Fast Computation of Wasserstein Barycenters, **2014**.

Appendices

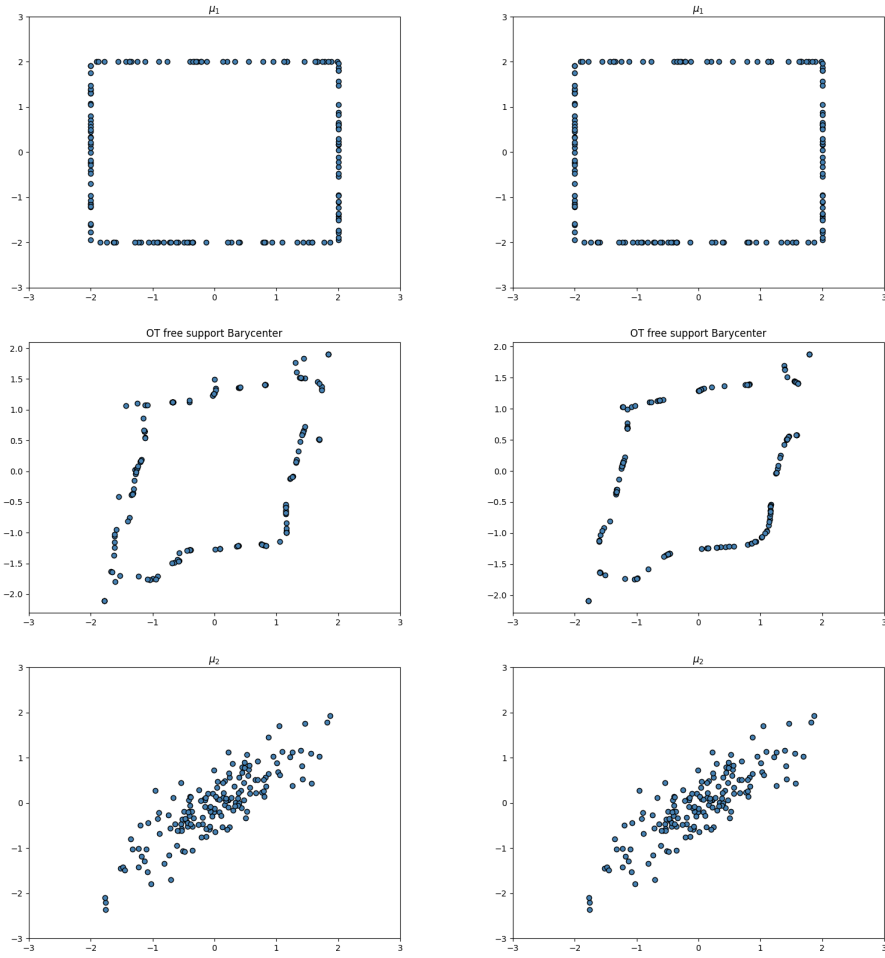
Barycenters computed with POT toolbox



(a) $\epsilon = 0.005$

(b) $\epsilon = 0.01$

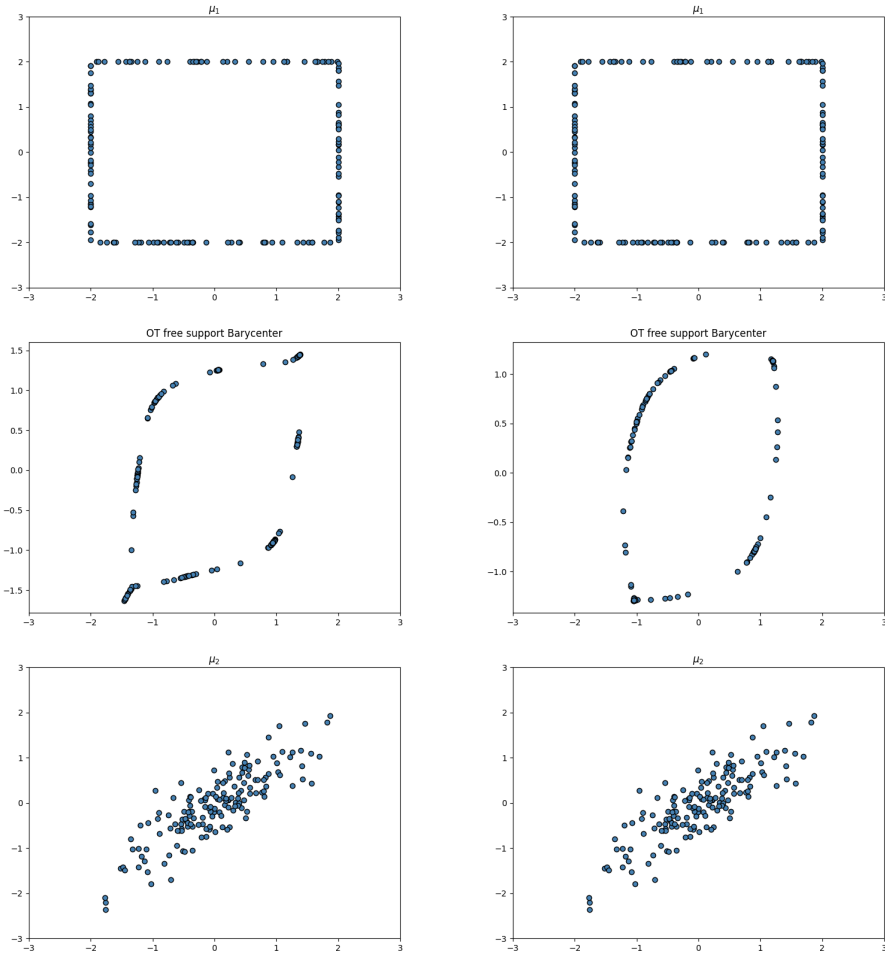
Figure 8: Free support barycenters computed with POT toolbox



(a) $\epsilon = 0.05$

(b) $\epsilon = 0.1$

Figure 9: Free support barycenters computed with POT toolbox



(a) $\epsilon = 0.5$

(b) $\epsilon = 1$

Figure 10: Free support barycenters computed with POT toolbox