



HAL
open science

A novel notion of barycenter for probability distributions based on optimal weak mass transport

Elsa Cazelles, Felipe Tobar, Joaquin Fontbona

► **To cite this version:**

Elsa Cazelles, Felipe Tobar, Joaquin Fontbona. A novel notion of barycenter for probability distributions based on optimal weak mass transport. Thirty-fifth Conference on Neural Information Processing Systems (NeurIPS 2021), NIPS: Neural Information Processing Systems Foundation, Dec 2021, Online, France. hal-03160475v1

HAL Id: hal-03160475

<https://hal.science/hal-03160475v1>

Submitted on 27 Oct 2021 (v1), last revised 10 Mar 2023 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A novel notion of barycenter for probability distributions based on optimal weak mass transport

Elsa Cazelles
IRIT, Université de Toulouse, CNRS
elsa.cazelles@irit.fr

Felipe Tobar
IDIA & CMM, Universidad de Chile
ftobar@dim.uchile.cl

Joaquin Fontbona
CMM, Universidad de Chile
fontbona@dim.uchile.cl

Abstract

We introduce weak barycenters of a family of probability distributions, based on the recently developed notion of optimal weak transport of mass [25], [9]. We provide a theoretical analysis of this object and discuss its interpretation in the light of convex ordering between probability measures. In particular, we show that, rather than averaging the input distributions in a geometric way (as the Wasserstein barycenter based on classic optimal transport does) weak barycenters extract common geometric information shared by all the input distributions, encoded as a latent random variable that underlies all of them. We also provide an iterative algorithm to compute a weak barycenter for a finite family of input distributions, and a stochastic algorithm that computes them for arbitrary populations of laws. The latter approach is particularly well suited for the *streaming setting*, i.e., when distributions are observed sequentially. The notion of weak barycenter and our approaches to compute it are illustrated on synthetic examples, validated on 2D real-world data and compared to standard Wasserstein barycenters.

1 Introduction

Optimal transport (OT) [40] has had a tremendous impact in the machine learning (ML) community recently, as it provides meaningful and implementable distances between probability distributions [33], thus advancing many aspects in the field, see e.g. [7, 42, 32]. The space of probability measures on \mathbb{R}^d with finite second moment can be *metrised* with the Wasserstein-2 distance, the computation of which amounts to finding a transport plan that minimises the quadratic average cost of transporting mass from a source probability measure onto a target one. In this context, a natural method for averaging a finite family of probability measures is to compute their Fréchet mean, with respect to the Wasserstein-2 distance, which corresponds to the *Wasserstein barycenter* introduced in [1].

The goal of the present work is to explore theoretical features and potential applications to ML of barycenters of probability measures analogously defined in terms of *optimal weak transport* (OWT, see [25]) or more precisely quadratic barycentric transport costs. In a nutshell, for a source measure μ and a target measure ν , the OWT problem aims to transport mass so that the conditional spatial mean of target support points y , given their source support points x , is close to x in average. This amounts to finding an intermediate measure η , possibly *more concentrated* than ν in the sense of convex ordering of probability measures, which is *close* to μ with respect to the Wasserstein-2 distance.

The main motivation of our work is to investigate the effect and meaning of combining a family of probability measures using OWT instead of OT. To that end, we will define the *weak barycenter* of

this family through an optimisation problem, and discuss some of its properties. Importantly, we will see that, rather than averaging the input distributions in a metric sense, solving a weak barycenter problem corresponds to finding probability measures that encode geometric or shape information *shared across* all of them. In fact, the weak barycenter problem will be interpreted as finding a latent random variable common to all the input distributions. Implications of this latent variable interpretation, in terms of robustness to outliers, will also be drawn in our work.

A second motivation for our work is to develop and implement computational methods for weak barycenters, capitalising on the fact that the optimal weak coupling between *any* pair of distributions, with finite second moments, is always realised by a unique optimal *map*. This property is in sharp contrast to standard OT, where the absolute continuity with respect to the Lebesgue measure of the source or target measure is typically needed to grant the existence and uniqueness of a map—the so-called Monge map—realising the optimal coupling between them. This map is often required in different ways to compute Wasserstein barycenters (see [5], [43] or [34]).

Similarly to the Wasserstein barycenter problem, we will develop a fixed-point formulation of the weak barycenter problem, based on OWT plans. This allows us to, following [5, 43], construct an iterative procedure to compute a weak barycenter for a finite family of distributions and analyse its convergence properties. We will also define and study the so-called *weak population barycenters*, that deal with a population of probability measures distributed according to a given law \mathbb{Q} supported on the Wasserstein-2 space, as in [29] for the OT case. Extending ideas from [34], we will then propose an iterative stochastic algorithm for online computation of the weak population barycenter, from a *stream* of probability measures sampled from \mathbb{Q} . We will then provide numerical simulations using this proposed method, in order to illustrate the geometric meaning of the weak barycenter, and we will compare it with related objects obtained with standard OT or its entropy-regularised counterpart.

Organisation of the paper. Sec. 2 documents the background on OWT and the assumptions underlying our work. Sec. 3 analyses the weak barycenter problem, interprets it in the light of convex ordering and a latent variable model and addresses the case of an infinite population of distributions. Sec. 4 introduces two algorithms for computing the weak barycenter in the finite or population settings. Sec. 5 and 6 present the experimental setting and validation of our proposal respectively. Lastly, Sec. 7 discusses our findings and future research questions. The Appendix contains all the proofs, additional details or our simulations and the code of our experiments.

2 Background : optimal weak transport and Wasserstein barycenter

The optimal transport (OT) problem [41] aims to find the lowest cost to transfer the mass from one probability measure onto another. Therefore, OT is a natural way to compare two probability distributions in terms of their geometric information. In particular, the Wasserstein- p distance W_p , associated with the Euclidean cost in \mathbb{R}^d , metrises the space $\mathcal{P}_p(\mathbb{R}^d)$ of probability measures on \mathbb{R}^d with finite p -moment. Precisely, for $\mu, \nu \in \mathcal{P}_p(\mathbb{R}^d)$,

$$W_p(\mu, \nu) = \left(\min_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^p d\pi(x, y) \right)^{1/p}, \quad (1)$$

where π is a *transport plan* between μ and ν , that is, an element of the set $\Pi(\mu, \nu)$ of probability measures on the product space $\mathbb{R}^d \times \mathbb{R}^d$ with marginals μ and ν . For $p = 2$ and μ absolutely continuous (*a.c.*), the unique optimal plan is concentrated on the graph of a measurable map called *Monge map* such that $\nu = T\#\mu$, see eq. (14) in Appendix A.1.

Optimal weak transport. We consider here the optimal weak transport (OWT) problem introduced in [25] and in particular the special case of barycentric transport costs. The OWT problem is then defined for $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ by

$$V(\mu|\nu) = \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d} \|x - \int_{\mathbb{R}^d} y d\pi_x(y)\|^2 d\mu(x), \quad (2)$$

where π_x is the *disintegration* of the transport plan π with respect to the first marginal μ , *i.e.* $\pi(dx dy) = \pi_x(dy)\mu(dx)$. As our work strongly leans on OWT theory, we recall in Appendix A.2, Th. 6, that V is continuous with respect to the Wasserstein metric [9]. Additionally, the two following results from [8] (stated for our specific setting) lay the ground for our proposed weak barycenters.

Theorem 1 ([8], Theorem 1.2). *The problem (2) admits a unique minimiser.*

This first result strongly differs from the classical OT setting, for which the uniqueness of an optimal transport plan is not guaranteed for arbitrary measures. The optimisation problem in Eq. (2) can also be reformulated thanks to the Brenier-Strassen theorem [24], [8], through the notion of convex ordering. We denote by $\eta \leq_c \nu$ the *convex order of measures*, meaning that $\int \phi d\eta \leq \int \phi d\nu$ for any convex function ϕ that is nonnegative or integrable with respect to $\eta + \nu$. By Strassen's theorem [39], two distributions are in convex order if and only if there exists a martingale coupling between them. The following theorem is a generalisation of the result originally proved in [24], Th. 1.2.

Theorem 2 ([8], Theorem 1.4). *Let $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ and $\nu \in \mathcal{P}_1(\mathbb{R}^d)$. There exists a unique $\eta^* \leq_c \nu$ such that*

$$W_2^2(\mu, \eta^*) = \inf_{\eta \leq_c \nu} W_2^2(\mu, \eta) = V(\mu|\nu). \quad (3)$$

Moreover, there exists a convex function $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$ of class C^1 with $\nabla \psi$ being 1-Lipschitz, such that $\nabla \psi \# \mu = \eta^$. Finally, the optimal coupling $\pi^{\mu, \nu} \in \Pi(\mu, \nu)$ verifies $\int y d\pi_x^{\mu, \nu}(y) = \nabla \psi(x)$ μ -a.s.*

The measurable map, or barycentric projection, $S_\mu^\nu(x) := \int_{\mathbb{R}^d} y d\pi_x^{\mu, \nu}(y)$ associated to the plan $\pi^{\mu, \nu}$ achieving the minimum in Eq. (2) is consequently uniquely defined and will be called *optimal barycentric projection*. From this notation, we can write the OWT cost in terms of an OT cost according to $V(\mu|\nu) = W_2^2(\mu, S_\mu^\nu \# \mu)$. We emphasise that S_μ^ν is directly related to the optimisation problem (2), whereas applied works such as [37, 35] make use of a barycentric projection constructed from a transport plan solving an OT problem between μ and ν (often regularised) as a substitute for the Monge map, which may not exist (more details on S_μ^ν are displayed in Appendix A.3).

Last, let us note that OWT is somehow also related to the martingale OT problem developed in the stochastic finance community [12, 2, 26], which puts the focus on the optimal transfer of mass between distributions assumed to be in convex order themselves.

Wasserstein barycenter. The classical Wasserstein barycenter problem for a set of probability measures $\nu_1, \dots, \nu_n \in \mathcal{P}_2(\mathbb{R}^d)$ with weights $\lambda_1, \dots, \lambda_n$ in the simplex (i.e. $\lambda_i \geq 0$ and $\sum_{i=1}^n \lambda_i = 1$) is defined [1] by

$$\arg \min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=1}^n \lambda_i W_2^2(\mu, \nu_i). \quad (4)$$

The Wasserstein barycenter has been extensively studied both theoretically and numerically [29, 43, 5, 14]. Regarding the numerical part, [38] focuses on the computation of Wasserstein barycenters for a fixed number of measures and a stream of observations per measure; additionally, [31] proposed an entropy-regularised alternative via stochastic optimisation for computing the Wasserstein barycenter of *a.c.* distributions only from observations. Constrained by their assumption of *a.c.*, [43] computes the Wasserstein barycenter by smoothing the observed empirical distributions. Furthermore, [20] compares the complexity of both the *sample* Wasserstein barycenter and a stochastic approximation to estimate a population barycenter (discrete measures and entropic regularisation). Finally, the authors of [3] recently proposed an algorithm to compute the barycenters in polynomial time.

3 Optimal weak transport barycenters and latent variable interpretation

3.1 Definition and basic properties

In a similar fashion, based on the weak transport cost in Eq. (2), we propose the following variant:

Definition 1. *The set of weak barycenters of a finite family of measures $\{\nu_i\}_{i=1, \dots, n} \in \mathcal{P}_2(\mathbb{R}^d)$ with weights $\{\lambda_i\}_{i=1, \dots, n}$ in the simplex is defined as*

$$\arg \min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=1}^n \lambda_i V(\mu|\nu_i). \quad (5)$$

Thus, a weak barycenter averages, with respect to the Wasserstein metric, an optimally chosen set of probability measures $\{\eta_1, \dots, \eta_n\}$ which are *more concentrated* than the corresponding ν_i , in the sense that $\eta_i \leq_c \nu_i$ for each $1 \leq i \leq n$. The existence of a solution is established as follows:

Proposition 1. *The weak barycenter problem in Eq. (5) admits a minimiser $\mu \in \mathcal{P}_2(\mathbb{R}^d)$.*

See Sec. B of the Appendix for the proof of the above Proposition (which relies on Prokhorov's theorem) and all the proofs for this Section. Uniqueness is in general not granted: we next show that the set of solutions is indeed an interval, with respect to the partial order of convex ordering of probability measures.

In the following, we denote by X and Y_i random variables with respective laws μ and ν_i , for $1 \leq i \leq n$, and δ_a the Dirac measure supported on $a \in \mathbb{R}^d$.

Lemma 1. *If μ is a weak barycenter of $\{\nu_i\}_{i=1,\dots,n}$ and $\mu' \leq_c \mu$, then μ' also is a weak barycenter. In particular, the Dirac measure supported on $\mathbb{E}_\mu(X)$ is always a weak barycenter. Moreover, a Dirac distribution $\delta_{\bar{\omega}}$ is a weak barycenter if and only if $\bar{\omega} = \sum_{i=1}^n \lambda_i \mathbb{E}_{\nu_i}(Y_i)$.*

A consequence of the above lemma is that for any weak barycenter μ ,

$$\mathbb{E}_\mu(X) = \sum_{i=1}^n \lambda_i \mathbb{E}_{\nu_i}(Y_i), \quad (6)$$

and the value of the weak barycenter problem is given by

$$\inf_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=1}^n \lambda_i V(\mu | \nu_i) = \sum_{i=1}^n \lambda_i \|\mathbb{E}(Y_i)\|^2 - \left\| \sum_{i=1}^n \lambda_i \mathbb{E}(Y_i) \right\|^2. \quad (7)$$

We can also derive the following characterisation on the set of weak barycenters:

Proposition 2. *A measure $\mu \in \mathcal{P}(\mathbb{R}^d)$ is a weak barycenter of $\{\nu_i\}_{i=1,\dots,n}$ if and only if its mean satisfies (6) and $\hat{\mu} \leq_c \hat{\nu}_i$ holds for all $1 \leq i \leq n$, where $\hat{\nu}$ denotes the centered version of a law ν .*

For instance, in the case of one dimensional Gaussian distributions $\nu_i = \mathcal{N}(m, \sigma_i^2)$, the set of weak barycenters includes $\{\mu = \mathcal{N}(m, \sigma^2) \mid 0 \leq \sigma^2 \leq \min_{1 \leq i \leq n} \sigma_i^2\}$.

A natural question is whether a "maximal" weak barycenter exists, in the sense of convex ordering (up to translation by the mean). For $d = 1$, the answer is affirmative. When the means $\mathbb{E}(Y_i)$ are equal, this follows from the complete lattice property of the set of probability measures with respect to the convex ordering (see [28]); the general case can then be reduced to the latter using Proposition 2. For $d \geq 2$, this property is in general not true and the answer depends on the family $\{\nu_i\}_{i=1,\dots,n}$.

In the particular case of *a.c.* input measures, we can bound the distance between the Wasserstein and weak barycenters by the variances of the distributions $(\nu_i)_{1 \leq i \leq n}$. The barycenters are then closer the more concentrated each ν_i is.

Lemma 2. *Let $\nu_1, \dots, \nu_n \in \mathcal{P}_2(\mathbb{R}^d)$ be *a.c.*, at least one of them with bounded density. Let $\bar{\mu}$ and $\tilde{\mu}$ respectively denote the weak and the Wasserstein barycenters. Then*

$$W_2^2(\bar{\mu}, \tilde{\mu}) \leq 2 \sum_{i=1}^n \lambda_i (\mathbb{E}\|Y_i\|^2 - \|\mathbb{E}Y_i\|^2).$$

3.2 Weak barycenters as latent variables

The weak barycenter encodes common geometric information present in all the input measures considered, therefore, it can be intuitively and rigorously interpreted as being the distribution of a latent variable underlying the realisations of random variables of laws ν_i for all $1 \leq i \leq n$.

Theorem 3. *Let μ be a weak barycenter of $\{\nu_i\}_{i=1,\dots,n}$. Then, for each $1 \leq i \leq n$, a random variable $Y_i \sim \nu_i$ can be realised as*

$$Y_i = X + (\mathbb{E}Y_i - \mathbb{E}X) + \bar{Y}_i,$$

where $X \sim \mu$ and $\bar{Y}_i = Y_i - \mathbb{E}(Y_i|X)$ is centered conditionally on X . Moreover, one has $S_\mu^\nu(X) = X + (\mathbb{E}Y_i - \mathbb{E}X)$ for all $i = 1, \dots, n$. Finally, we have $\mathbb{E}(Y_i - \mathbb{E}Y_i | X - \mathbb{E}X) = X - \mathbb{E}X$ or, equivalently, $\hat{\mu} \leq_c \hat{\nu}_i$, with $\hat{\mu}$ and $\hat{\nu}_i$ the laws of $X - \mathbb{E}X$ and $Y_i - \mathbb{E}Y_i$ respectively.

That is to say, each $Y_i \sim \nu_i$ can be realised by sampling a random variable X common to all $i = 1, \dots, n$ and distributed according to the weak barycenter μ , translating that value by $\mathbb{E}Y_i - \mathbb{E}X$ and adding a *cluster-specific* component \bar{Y}_i or idiosyncratic noise, centered conditionally on X .

Remark 1. *The observations of each class (i.e. input measure) can be interpreted as outliers with respect to the (translated) law of the weak barycenter, which are statistically different and are thus left aside of its support. This way, the weak barycenter is robust to outliers, as it tends to discard them, by construction. Furthermore, this "robustness" property results in the stability of weak barycenter upon perturbation of a class with larger noise (or more scattered, outlying values). More precisely, if a class is corrupted in such a way that their observations result in a stochastically larger distribution than the original one, a weak barycenter computed in terms of the original (stochastically smaller) class will still be a weak barycenter in the new corrupted setting. An intuitive and simple way to illustrate this point follows by considering a weak barycenter μ of a one dimensional and centered family of input distributions $\{\nu_i\}_{i=1,\dots,n}$. By Proposition 2, μ must verify $\mu \leq_c \nu_i$ for all $i = 1, \dots, n$. In particular, from Theorem 3.A.1. in [36], we have that $\int_x^\infty \mathbb{P}(X > u)du \leq \int_x^\infty \mathbb{P}(Y_i > u)du$ for all $x \in \mathbb{R}$, where $X \sim \mu$ and $Y_i \sim \nu_i$. Therefore, μ is likely to avoid outliers. Another supportive intuition in terms of robustness is that a maximal weak barycenter would be one that includes the most possible points of all classes (or distributions) in its support (all this, after re-centering) and leaves out only "outliers". A non-maximal weak barycenter is then more conservative, meaning that it counts on fewer points and leaves out more possible outliers.*

3.3 Extension for the population barycenter

The population Wasserstein barycenter introduced in [29] and [4] extends the definition of Wasserstein barycenter for an infinite number of measures. This formulation is particularly relevant for the construction of an iterative algorithm to compute the barycenter for the streaming case, that is, when the measures are received *online*. The proofs are reported in Section C of the Appendix.

Let us consider a probability measure $\mathbb{Q} \in \mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$, meaning that \mathbb{Q} is supported on a set of measures with finite moments of order 2, such that for some (and thus all) $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we have that $\int_{\mathcal{P}_2(\mathbb{R}^d)} W_2^2(\mu, \nu) d\mathbb{Q}(\nu) < \infty$.

Definition 2. *We define the set of weak population barycenters of a distribution $\mathbb{Q} \in \mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$ as*

$$\arg \min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \int_{\mathcal{P}_2(\mathbb{R}^d)} V(\mu|\nu) d\mathbb{Q}(\nu). \quad (8)$$

The following lemma guarantees that the map $(x, \nu) \mapsto S_\mu^\nu(x)$ appearing in Eq. (8) through $V(\mu|\nu) = \int \|x - S_\mu^\nu(x)\|^2 d\mu(x)$ is well defined.

Lemma 3. *The function $(\mu, \nu) \in (\mathcal{P}_2(\mathbb{R}^d))^2 \mapsto \pi^{\mu, \nu} \in \mathcal{P}_2(\mathbb{R}^{2d})$ mapping (μ, ν) to the unique optimal plan $\pi^{\mu, \nu}$ realising $V(\mu|\nu)$ in Eq. (2) is continuous. As a consequence, for each $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ the function $(x, \nu) \in \mathbb{R}^d \times \mathcal{P}_2(\mathbb{R}^d) \mapsto S_\mu^\nu(x)$ is measurable.*

Using similar arguments as those of Proposition 1 and the fact that any probability measure can be approximated by a sequence of probability measures with finite support, the following proposition confirms that the weak population barycenter problem is also well defined.

Proposition 3. *The minimisation problem in Eq. (8) admits a solution.*

4 Algorithms via fixed-point representations

4.1 Weak barycenter

For the Wasserstein barycenter problem in Eq. (4), the authors in [1] proved that if at least one of the measures ν_1, \dots, ν_n is *a.c.*, the Wasserstein barycenter is unique. Furthermore, if all the ν_i 's are *a.c.*, and at least one of them has a bounded density, then the unique Wasserstein barycenter is also *a.c.* and verifies a fixed-point equation. This last property has been thoroughly studied by [5] and [43] and leveraged to compute an approximation of the barycenter via an iterative algorithm based on Monge maps, whose existence and uniqueness are guaranteed by the *a.c.* of the measures involved.

Akin to the fixed-point methodology in the classical Wasserstein scenario, we define an iterative procedure based on the barycentric projection computed in the optimal weak transport problem in Eq. (2), that is valid for arbitrary distributions. Therefore, we consider the following iterative rule for

probability measures $\nu_1, \dots, \nu_n \in \mathcal{P}_2(\mathbb{R}^d)$:

$$\mu_{k+1} = G(\mu_k), \text{ with } G(\mu) = \left(\sum_{i=1}^n \lambda_i S_\mu^{\nu_i} \right) \# \mu, \quad (9)$$

where for each $i = 1, \dots, n$, the optimal barycentric projection is given by $S_\mu^{\nu_i} : x \mapsto \int y d\pi_x^{\mu, \nu_i}(y)$, for $\pi_x^{\mu, \nu_i} \in \Pi(\mu, \nu_i)$ achieving the minimum in the OTW problem in Eq. (2). The proposed iterative procedure is presented in Algorithm 1.

A fundamental difference between the fixed-point computation of the Wasserstein barycenter [5] and a weak barycenter is that the optimal Monge map T_μ^ν in the OT problem verifies $T_\mu^\nu \# \mu = \nu$, whereas the pushforward measure $S_\mu^\nu \# \mu$ in the OTW setting still depends on μ . We will then prove that the iterative algorithm in Eq. (9), based on the maps $S_\mu^{\nu_i}$, admits converging subsequences. A convenient result is the continuity of the functional G in Eq. (9), which can be proven using Arzela-Ascoli theorem on a set of barycentric projections as well as the Skorohod's representation theorem.

Theorem 4. *The function $\mu \mapsto G(\mu)$ defined in Eq. (9) is W_2 -continuous from $\mathcal{P}_2(\mathbb{R}^d)$ to $\mathcal{P}_2(\mathbb{R}^d)$.*

Using an approach similar to [5] for the Wasserstein barycenter, we can state the following results for the proposed fixed-point procedure.

Proposition 4. *If μ is a weak-barycenter, that is a solution of problem (5), then $G(\mu) = \mu$ i.e. $x = \sum_{i=1}^n \lambda_i S_\mu^{\nu_i}(x), \mu(x)$ -a.s.*

The inverse implication of Proposition 4 is not necessarily true, that is, some fixed points may not be weak barycenters. However, a Dirac delta $\delta_\omega, \omega \in \mathbb{R}^d$, that meets the fixed-point condition $\delta_\omega = G(\delta_\omega)$, is a weak barycenter (see Lemma 1).

Proposition 5. *Let $(\mu_k)_k$ be the sequence defined by the iterative procedure $\mu_{k+1} = G(\mu_k)$ and starting from $\mu_0 \in \mathcal{P}_2(\mathbb{R}^d)$. Then $(\mu_k)_k$ is tight and every converging subsequence must converge to a fixed point of G .*

We observe that these results also hold for the classical Wasserstein barycenter of *a.c.* measures $\{\nu_i\}_{i=1, \dots, n}$ such that at least one of them has a bounded density. Moreover, the inverse implication, namely if μ is a fixed-point then it is a barycenter, is not straightforward even in the Wasserstein barycenter case, for which one considers the fixed-point equation given by $\mu = (\sum_{i=1}^n \lambda_i T_\mu^{\nu_i}) \# \mu$, with $T_\mu^{\nu_i}$ the Monge map verifying $\nu_i = T_\mu^{\nu_i} \# \mu$. Indeed, [1] prove that if μ checks $x = \sum_{i=1}^n \lambda_i T_\mu^{\nu_i}(x)$ for every $x \in \mathbb{R}^d$, not only μ -almost everywhere, then μ is a Wasserstein barycenter. Also, [43, Theorem 2] provide additional conditions for this to be true by essentially invoking more smoothness on the distributions $\{\nu_i\}_{i=1, \dots, n}$. Additionally, they only conjecture that under the same assumptions, the fixed-point is unique. Our method, however, includes arbitrary probability measures. Therefore, we do not expect to obtain similar results as in the Wasserstein barycenter case, for which smoothness is required.

4.2 Weak population barycenter

Based on [34], we construct a stochastic iterative algorithm for computing the weak population barycenter in Eq. (8). We clarify that [34] is constrained to probability measures \mathbb{Q} supported on distributions that are *a.c.*, whereas in our setting these distributions only need to belong to $\mathcal{P}_2(\mathbb{R}^d)$. Let us notice that our algorithms can be interpreted as geodesic gradient descent as in [34] and [16], however, OTW is not a metric and its potential geodesic structure is so far unknown. Therefore, the proposed algorithm only aims to mimic Riemmanian gradient descent. Our fixed-point result for the weak population barycenter problem is stated in the following Lemma:

Lemma 4. *If μ is a weak population barycenter of \mathbb{Q} , then $x = \int S_\mu^\nu(x) d\mathbb{Q}(\nu), \mu(x)$ -a.s.*

As in the finite case, the inverse implication is difficult to obtain. In particular, this has not been proven for the classical population Wasserstein barycenter in [34], where it boils down to prove the uniqueness of an absolutely continuous fixed point of $\mu \mapsto (\int T_\mu^\nu d\mathbb{Q}(\nu)) \# \mu$, where T_μ^ν is the Monge map between μ and ν . As explained in [34], the uniqueness of such fixed points has also been studied under some strong assumptions in [14] by considering parametric classes of random probability measures with compact support. Note that this result is expected to be true by again invoking more smoothness on the distributions at hand. As our method focuses (in particular) on

discrete probability measures, the conditions under which the inverse equality holds are beyond the scope of our work. However, from the experimental results in Section 6, we believe our method presents practical advantages.

We next develop an iterative procedure that converges towards a distribution μ verifying the fixed-point equation $x = \int S_{\mu}^{\nu}(x)d\mathbb{Q}(\nu), \mu(x)$ -a.s. Our method is presented in the following definition and illustrated in Algorithm 2.

Definition 3. Let $\mu_0 \in \mathcal{P}_2(\mathbb{R}^d), \nu^k \stackrel{i.i.d.}{\sim} \mathbb{Q}$ and $\gamma_k > 0$. We define the following iterative procedure for $k \geq 0$:

$$\mu_{k+1} = \left[(1 - \gamma_k)\text{id} + \gamma_k S_{\mu_k}^{\nu^k} \right] \# \mu_k, \quad (10)$$

where $S_{\mu_k}^{\nu^k}$ is the optimal barycentric projection between μ_k and ν^k and id is the identity operator.

Assuming the following conditions on the steps γ_k :

$$\sum_{k=1}^{\infty} \gamma_k^2 < \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \gamma_k = \infty, \quad (11)$$

we are able to prove the convergence of the iterative scheme in a similar fashion as Theorem 4.7 in [34]. The proof is to a large extent very similar to that of the classical Wasserstein case.

Theorem 5. Under the conditions in Eq. (11), the sequence $(\mu_k)_k$ defined in Eq. (10) is a.s. relatively compact in W_q for all $q < 2$ (in particular it is tight). Moreover, a limit point μ verifies $x = \int S_{\mu}^{\nu}(x)d\mathbb{Q}(\nu), \mu(x)$ -a.s.

Note that we prove in Prop. 7 in Appendix that the function $\mu \mapsto \left\| \int S_{\mu}^{\nu}d\mathbb{Q} - \text{id} \right\|_{\mathbb{L}^2(\mu)}^2$ is continuous w.r.t. W_2 . This result, however, was obtained under additional constraints for the OT case in [34].

Algorithm 1: Weak barycenter

Input: distributions ν_1, \dots, ν_n , # steps K ;
initialisation: $\mu_0 = \nu_1$;
for $k = 0, 1, \dots, K$ **do**
 for $i = 1, 2, \dots, n$ **do**
 Solve the OWT problem between μ_k
 and ν_i to obtain π^{μ_k, ν_i} ;
 $S_i = \int y d\pi_x^{\mu_k, \nu_i}(y)$
 end
 $\mu_{k+1} = (\sum_{i=1}^n \lambda_i S_i) \# \mu_k$
end

Algorithm 2: Weak population barycenter

Input: number of steps K ;
initialise distribution $\mu_0 \sim \mathbb{Q}$;
for $k = 0, 1, \dots, K$ **do**
 Sample $\nu^k \sim \mathbb{Q}$;
 Update γ_k ;
 Solve the OWT problem to obtain π^{μ_k, ν^k}
 $S_k = \int y d\pi_x^{\mu_k, \nu^k}(y)$;
 $\mu_{k+1} = [(1 - \gamma_k)\text{id} + \gamma_k S_k] \# \mu_k$;
end

5 Computational aspects

Setting and computation of OWTs. Both Algorithms 1 and 2 require the computation of the optimal barycentric projection associated to the OWT problem in Eq. (2). For two discrete measures $\mu = \sum_{i=1}^r a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^m b_j \delta_{y_j}$, this boils down to solving the following quadratic programming problem

$$\min_{\pi \in \mathbb{R}^{r \times m}} \left\{ \sum_{i=1}^r a_i \left\| x_i - \left(\frac{\pi \mathbf{y}}{\mathbf{a}} \right)_i \right\|^2, \pi_{ij} \geq 0, \pi \mathbf{1} = \mathbf{a}, \pi^T \mathbf{1} = \mathbf{b} \right\}, \quad (12)$$

which can be solved using a solver such as **cvxpy**. We also propose to solve the OWT problem in Eq. (12) with a proximal algorithm. The optimal barycentric projection is then constructed as $\frac{\pi \mathbf{y}}{\mathbf{a}}$. The details and examples are presented in Appendix E.1.

Comparison setting. In the next section, we compare our proposed computation for weak barycenters in Definition 2 (Algorithm 2) to the classic Wasserstein barycenter in particular for a stream of a.c. measures. Namely, we will run Algorithm 2 by, following [18, 35], replacing optimal barycentric

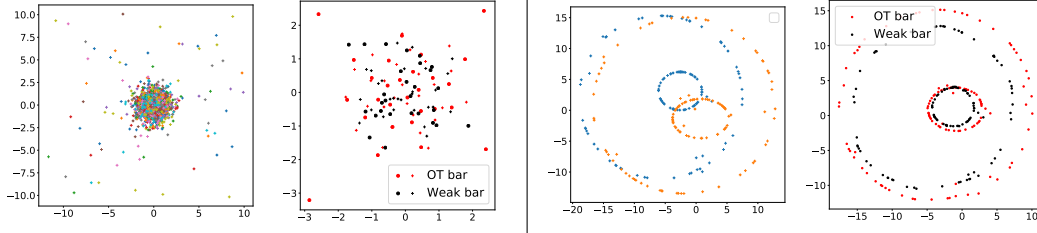


Figure 1: (left) Empirical Gaussian distributions and their OWT (black) and OT (red) barycenters for Gaussian observations (crosses) and corrupted observations (dots). (right) Empirical distributions supported on two ellipses and their OWT (black) and OT (red) barycenters.

projections by the barycentric projections associated either to i) an optimal plan in the Kantorovich problem (1), or ii) the optimal Sinkhorn plan in the entropy regularised OT problem [17] given by

$$\arg \min_{\pi \in \Pi(\mu, \nu)} \int \|x - y\|^2 d\pi(x, y) + \varepsilon KL(\pi | \mu \otimes \nu), \quad (13)$$

where KL denotes the Kullback-Leibler divergence. The associated barycenters will be referred to as *OT barycenter* and *OT Sinkhorn barycenter* respectively. The optimal plans for OT and regularised OT problem were computed using *POT toolbox* [22]. Notice that what we call OT barycenter (resp. OT Sinkhorn barycenter) is not solving a Wasserstein barycenter problem (resp. a regularised Wasserstein barycenter problem). Therefore, our method for barycentric computation differs from previous ones in the literature (see Section 2) in that it i) can process a *stream* of an unknown number of measures, ii) does not require the measures to be a.c., and iii) does not appeal to additional regularisation of the measures or the Wasserstein metric.

6 Experimental results

This section is devoted to the empirical validation of our proposal on both synthetic and real-world data. We first focused on Algorithm 2 since multiple algorithms to compute a Wasserstein barycenter for a fixed number of distributions are already available [18, 38]. We present two robustness to outliers experiments, then we validate our OWT barycenter on synthetic dataset and real-world ones. The overall conclusion of our experiments is that the weak barycenter is more likely to maintain the common (or shared) geometric features of the measures involved, as expected from Theorem 3. Additional experiments are presented in Appendix E.2, including the comparison of the energy for the computed weak barycenter in Algorithm 1 against the approximated optimal energy (using Eq.(7) and the plug-in estimator).

6.1 Robustness to outliers

OT’s sensitivity to outliers is a well-known problem that can be addressed *e.g.* with unbalanced OT [10]. We observed that OWT also allows to deal with outliers, which is coherent with the latent variable interpretation (see Remark 1). We illustrate this with two experiments. In Fig. 1 (left), we consider 50 sets of 20 – 30 observations from different 2D Gaussian measures, where each observation may be corrupted by random translations (Bernoulli $p = 0.05$) thus producing outliers. We show the resulting barycenters (dots), and barycenters without outliers (crosses) for Wasserstein barycenter (red) and weak barycenter (black), which shows robustness to outliers. In Fig. 1 (right), we consider two distributions supported on pair-of-ellipses, and 120 observations per distribution. Again, each observation may be corrupted by random translations (Bernoulli $p = 0.05$). The weak barycenter (black) shows a better preservation of the shapes than the Wasserstein barycenter (red), in particular, the red dots are more often located outside the ellipses.

6.2 Synthetic distributions

We implemented the proposed sequential computation of weak barycenters (Algorithm 2) on two examples of synthetic distributions: Gaussians and spirals. In each case, we sampled r observations from a random distribution at each step, and considered K steps (and thus K measures for each case).

2D Gaussians ($r = 100$ & $K = 15$). We considered distributions $\mathcal{N}(m, I)$, with m uniformly distributed on $(-3, 3) \times (-5, 5)$ and I the identity matrix. Fig. 2 (left) shows the empirical distributions together with the OWT and OT barycenters, the weak barycenter being the less spread out as expected. The three remaining plots illustrate the behaviour of the barycenters constructed as stated in Sec. 5. For a small regularisation parameter ε in Eq. (13), the OT and OT Sinkhorn barycenters are similar, however, as ε increases the OT Sinkhorn (OTS) barycenter becomes closer to the weak barycenter and thus even more concentrated, meaning that its samples tend to be closer to each other. Critically, for a very large ε , as the entropy tends to spread the mass in the regularised optimal plan, the associated barycentric projection will roughly move the mass to the spatial mean of the target distribution’s support.

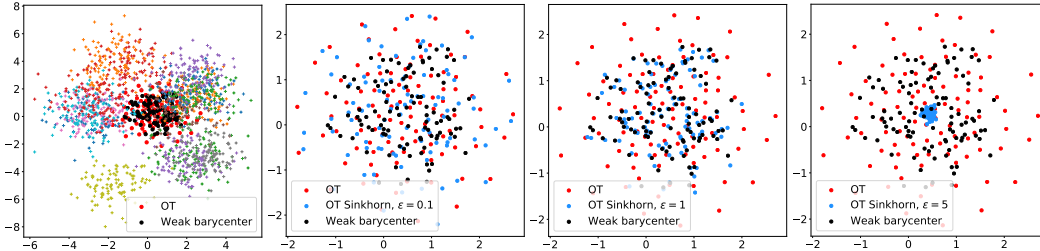


Figure 2: (left) Empirical Gaussian distributions and their OWT (black) and OT (red) barycenters computed with Algorithm 2. Illustration of the weak (black), OT (red) and OT Sinkhorn (blue) barycenters for different values of $\varepsilon = 0.1, 1, 5$.

Spiral distributions. ($r \in (200, 225)$ & $K = 10$). In this experiment, we considered distributions supported on a spiral—see Fig. 3 (left), with random ratio in $(0, 3)$. The OT and OWT barycenters are presented in Fig. 3 (right). Again, the weak barycenter seems to better preserve the shape of the spiral than the OT barycenter.

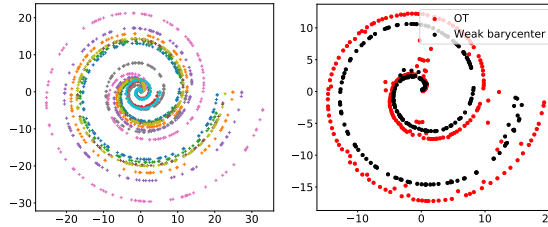


Figure 3: (left) Distributions supported on spiral. (right) OWT (black) and OT (red) barycenters computed with Algorithm 2.

6.3 Real-world dataset

MNIST dataset. We considered the well-known MNIST dataset [30] of grayscale images of handwritten digits. The images, of size 28×28 pixels, can be normalised and thus be interpreted as discrete probability measures supported on a two-dimensional grid of size 28×28 . We computed the barycenters with 30 steps of Algorithm 1 between two digits "8", that are noisy versions of the same digit with the aim to produce a more stable barycenter. To produce noisy data, we randomly (Bernoulli $p = 0.1$) move pixels of the prototype digit displayed in Fig. 4 (left). Fig. 4 (right) shows the barycenters using the OWT, OT, and entropic-OT (for $\varepsilon = 1$). This example illustrates how OWT reduces dispersion, so that weak barycenter provides the best uniformly spread results among the barycenters considered, with the two loops of the "8" well shaped.

Cytometry dataset. In biotechnology, *flow cytometry* is measured through intracellular markers of single cells in a biological sample with the objective of recognising common features across patients. However, these measurements are often disrupted by acquisition, rather than biological artefacts [27], thus hindering the identification of common features. To address this challenge, we compute the weak barycenter for the forward-scattered light (FSC) and side-scattered light (SSC) cell’s markers (using the flowStats package of Bioconductor [23]). We considered $K = 15$ patients and a variable number of cells per patient between 88 and 2185. Fig. 5 shows the 15 distributions (left) and the computed barycenters (right), thus confirming the ability of the weak barycenter to resolve the alignment of the dataset, while maintaining the expected diamond-shape. Moreover, the advantage of our proposed streaming procedure is fully exploited in this setting, since data from one or several patient can arrive

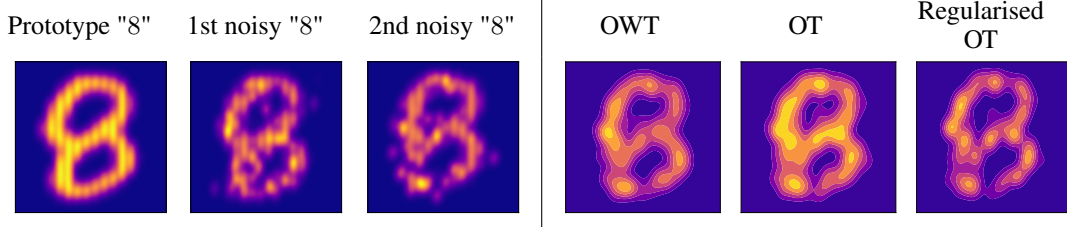


Figure 4: Digit "8" (MNIST). From left to right: Prototype "8", first and second noisy versions of the prototype by randomly (Bernoulli $p = 0.1$) moving pixels, three barycenters constructed with Algorithm 1 associated to the OWT plan, an OT plan and the entropy regularised OT plan for $\varepsilon = 1$.

sequentially. Though this setting has been addressed with the Wasserstein barycenter in [13], also in Fig. 5, such method required a fixed grid to compute the barycenter unlike our method, thus revealing the computational simplicity of the weak barycenter.

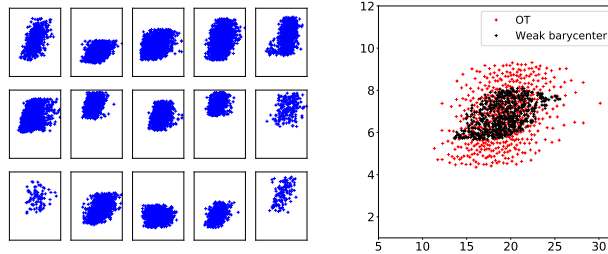


Figure 5: (left) Cytometry dataset for $n = 15$ patients and FSC vs. SSC cell's marker. (right) The weak-barycenter (black) computed with Algorithm 2 and the OT barycenter (red). The data are represented with the same axis as the figure of barycenters.

7 Discussion

We have introduced the weak barycenter, which extracts common geometric information of probability measures on \mathbb{R}^d based on optimal weak transport, and showed that it can be interpreted as a latent variable model. From the fixed-point formulation defined in terms of optimal weak transport maps, irrespective of the regularity assumptions on the measures involved, we developed practical computation via an iterative algorithm with guaranteed convergence. In particular, the proposed algorithms do not require a common grid on the sample space, when processing either observed data or samples from distributions. We have also proposed weak barycenters of a possibly infinite population of measures and developed a stochastic procedure for computing it in the streaming data regime where distributions are processes into the weak barycenter as they arrive. This has critical implications for continual-learning methods in the ML community.

Additional studies will focus on deepen the latent variable interpretation of weak barycenters, and its relationship to the aggregate information represented by the Wasserstein barycenter. Also, we identify two relevant theoretical aspects for further research: i) to exhibit general conditions on the family of input measures (or on the law of the population) for the existence of weak barycenters that are not Dirac masses; and ii) to provide conditions on those input measures for a "maximal" weak barycenter (in terms of convex ordering) to exist when $d \geq 2$, among all the solutions of the weak barycenter problem (and, if possible, a way of constructing it by regularisation most probably). The statistical behaviour of the weak barycenter can also be investigated, in particular when constructed from large empirical random samples of given distributions. Lastly, the weak population barycenter could also be used to construct a predictive posterior in the context of Bayesian learning, as was done for Wasserstein barycenters in [34].

Acknowledgments. We thank Julio Backhoff-Veraguas for his valuable insight during the writing of this paper. This work was funded by ANID grants: AFB170001 & ACE210010 (CMM), FB0008 (AC3E), Fondecyt-Postdoctorado #3190926 and Fondecyt-Iniciación #1210606.

References

- [1] M. Agueh and G. Carlier. Barycenters in the Wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011.
- [2] A. Alfonsi, J. Corbetta, and B. Jourdain. Sampling of probability measures in the convex order and approximation of martingale optimal transport problems. *Available at SSRN 3072356*, 2017.
- [3] Jason M Altschuler and Enric Boix-Adsera. Wasserstein barycenters can be computed in polynomial time in fixed dimension. *J. Mach. Learn. Res.*, 22:44–1, 2021.
- [4] P.C. Alvarez-Esteban, E Del Barrio, J.A. Cuesta-Albertos, and C. Matrán. Wide consensus for parallelized inference. *arXiv: 1511.05350*, 2015.
- [5] P.C. Álvarez-Esteban, E. Del Barrio, J.A. Cuesta-Albertos, and C. Matrán. A fixed-point approach to barycenters in Wasserstein space. *Journal of Mathematical Analysis and Applications*, 441(2):744–762, 2016.
- [6] L. Ambrosio, N. Gigli, and G. Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2008.
- [7] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 214–223. PMLR, 2017.
- [8] J. Backhoff-Veraguas, M. Beiglböck, and G. Pammer. Existence, duality, and cyclical monotonicity for weak transport costs. *Calculus of Variations and Partial Differential Equations*, 58(6):203, 2019.
- [9] J. Backhoff-Veraguas, M. Beiglböck, and G. Pammer. Weak monotone rearrangement on the line. *Electronic Communications in Probability*, 25, 2020.
- [10] Y. Balaji, R. Chellappa, and S. Feizi. Robust optimal transport with applications in generative modeling and domain adaptation. *arXiv:2010.05862*, 2020.
- [11] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [12] M. Beiglböck, P. Henry-Labordère, and F. Penkner. Model-independent bounds for option prices—a mass transport approach. *Finance and Stochastics*, 17(3):477–501, 2013.
- [13] J. Bigot, E. Cazelles, and N. Papadakis. Data-driven regularization of Wasserstein barycenters with an application to multivariate density registration. *Information and Inference: A Journal of the IMA*, 8(4):719–755, 2019.
- [14] J. Bigot and T. Klein. Characterization of barycenters in the Wasserstein space by averaging optimal transport maps. *ESAIM: Probability and Statistics*, 22:35–57, 2018.
- [15] M.W. Botsko. An elementary proof of Lebesgue’s differentiation theorem. *The American Mathematical Monthly*, 110(9):834–838, 2003.
- [16] S. Chewi, T. Maunu, P. Rigollet, and A.J. Stromme. Gradient descent algorithms for Bures-Wasserstein barycenters. *arXiv: 2001.01700*, 2020.
- [17] M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems 26*, pages 2292–2300. 2013.
- [18] M. Cuturi and A. Doucet. Fast computation of Wasserstein barycenters. *Proceedings of the 31st International Conference on Machine Learning*, 32(2):685–693, 2014.
- [19] A. Dessein, N. Papadakis, and J.-L. Rouas. Regularized optimal transport and the ROT mover’s distance. *The Journal of Machine Learning Research*, 19(1):590–642, 2018.
- [20] D. Dvinskikh. Stochastic approximation versus sample average approximation for population Wasserstein barycenters. *arXiv e-prints*, pages arXiv–2001, 2020.
- [21] E.A. Feinberg, P.O. Kasyanov, and Y. Liang. Fatou’s lemma in its classical form and Lebesgue’s convergence theorems for varying measures with applications to Markov decision processes. *Theory of Probability & Its Applications*, 65(2):270–291, 2020.

- [22] R. Flamary, N. Courty, A. Gramfort, M.Z. Alaya, A. Boisbunon, 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, and T. Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):1–8, 2021.
- [23] R.C. Gentleman, V.J. Carey, D.M. Bates, B. Bolstad, M. Dettling, S. Dudoit, B. Ellis, L. Gautier, Y. Ge, J. Gentry, et al. Bioconductor: open software development for computational biology and bioinformatics. *Genome biology*, 5(10):R80, 2004.
- [24] N. Gozlan and N. Juillet. On a mixture of Brenier and Strassen theorems. *Proceedings of the London Mathematical Society*, 120(3):434–463, 2020.
- [25] N. Gozlan, C. Roberto, P.-M. Samson, and P. Tetali. Kantorovich duality for general transport costs and applications. *Journal of Functional Analysis*, 273(11):3327–3405, 2017.
- [26] J. Guyon and P. Henry-Labordere. Nonlinear option pricing. *CRC Press*, 2013.
- [27] F. Hahne, A.H. Khodabakhshi, A. Bashashati, C.-J. Wong, R.D. Gascoyne, A.P. Weng, V. Seyfert-Margolis, K. Bourcier, A. Asare, T. Lumley, et al. Per-channel basis normalization methods for flow cytometry data. *Cytometry Part A: The Journal of the International Society for Advancement of Cytometry*, 77(2):121–131, 2010.
- [28] R.P. Kertz and U. Rösler. Complete lattices of probability measures with applications to martingale theory. *Lecture Notes-Monograph Series*, pages 153–177, 2000.
- [29] T. Le Gouic and J.-M. Loubes. Existence and consistency of Wasserstein barycenters. *Probability Theory and Related Fields*, 168(3-4):901–917, 2017.
- [30] Y. LeCun. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [31] L. Li, A. Genevay, M. Yurochkin, and J. Solomon. Continuous regularized wasserstein barycenters. *arXiv preprint arXiv:2008.12534*, 2020.
- [32] X. Nguyen. Convergence of latent mixing measures in finite and infinite mixture models. *The Annals of Statistics*, 41(1):370–400, 2013.
- [33] G. Peyré and M. Cuturi. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- [34] G. Rios, J. Backhoff-Veraguas, J. Fontbona, and F. Tobar. Bayesian learning with Wasserstein barycenters. *arXiv preprint arXiv:1805.10833v3*, 2018.
- [35] V. Seguy, B.B. Damodaran, R. Flamary, N. Courty, A. Rolet, and M. Blondel. Large-scale optimal transport and mapping estimation. In *International Conference on Learning Representations (ICLR)*, 2018.
- [36] Moshe Shaked and J George Shanthikumar. *Stochastic orders*. Springer Science & Business Media, 2007.
- [37] J. Solomon, F. De Goes, G. Peyré, M. Cuturi, A. Butscher, A. Nguyen, T. Du, and L. Guibas. Convolutional Wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Transactions on Graphics (TOG)*, 34(4):1–11, 2015.
- [38] M. Staib, S. Claiici, J. Solomon, and S. Jegelka. Parallel streaming wasserstein barycenters. In *Advances in Neural Information Processing Systems*, pages 2647–2658, 2017.
- [39] V. Strassen. The existence of probability measures with given marginals. *Annals of Mathematical Statistics*, 36(2):423–439, 1965.
- [40] C. Villani. *Topics in Optimal Transportation*. Number 58. American Mathematical Soc., 2003.
- [41] C. Villani. *Optimal Transport: Old and New*. Springer Berlin Heidelberg, 2008.
- [42] J. Ye, P. Wu, J. Z. Wang, and J. Li. Fast discrete distribution clustering using Wasserstein barycenter with sparse support. *IEEE Transactions on Signal Processing*, 65(9):2317–2332, 2017.
- [43] Y. Zemel and V. Panaretos. Fréchet means and Procrustes analysis in Wasserstein space. *Bernoulli*, 25(2):932–976, 2019.

A Additional mathematical background

A.1 p -Wasserstein distance

For $p = 2$ and $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ such that μ is absolutely continuous (*a.c.*) with respect to Lebesgue measure, the unique optimal plan is concentrated on the graph of a measurable map and Eq. (1) boils down to Monge's problem:

$$W_2(\mu, \nu) = \left(\min_{T \in \mathbb{T}(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - T(x)\|^2 d\mu(x) \right)^{1/2}, \quad (14)$$

where $\mathbb{T}(\mu, \nu)$ is the set of measurable functions $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that $\nu = T\#\mu$. The *pushforward* operator $\#$ is defined such that for any measurable set $B \subset \mathbb{R}^d$, we have $\nu(B) = \mu(T^{-1}(B))$. In such a case, the optimal measurable map T in Eq. (14) is uniquely defined (see *e.g.* Th. 9.4 in [41]) and called *Monge map*.

A.2 Continuity of V

Theorem 6 ([9], Theorem 1.5). *Let $(\mu_n)_n \subset \mathcal{P}_2(\mathbb{R}^d)$ and $(\nu_n)_n \subset \mathcal{P}_1(\mathbb{R}^d)$. Then*

$$\begin{cases} \mu_n \rightarrow \mu & \text{in } W_2 \\ \nu_n \rightarrow \nu & \text{in } W_1 \end{cases} \implies \lim_n V(\mu_n | \nu_n) = V(\mu | \nu).$$

A.3 On the barycentric projection

For a given transport plan $\pi \in \Pi(\mu, \nu)$, with $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$, the associated barycentric projection is given by

$$S : x \mapsto \int_{\mathbb{R}^d} y d\pi_x(y).$$

First, for each $x \in \mathbb{R}^d$, $S(x)$ realises $\min_z \mathbb{E}_{Y \sim \pi_x} (\|z - Y\|^2)$. Second, this barycentric map S is actually optimal for the Monge's problem Eq. (14) between μ and $S\#\mu$, by Theorem 2.

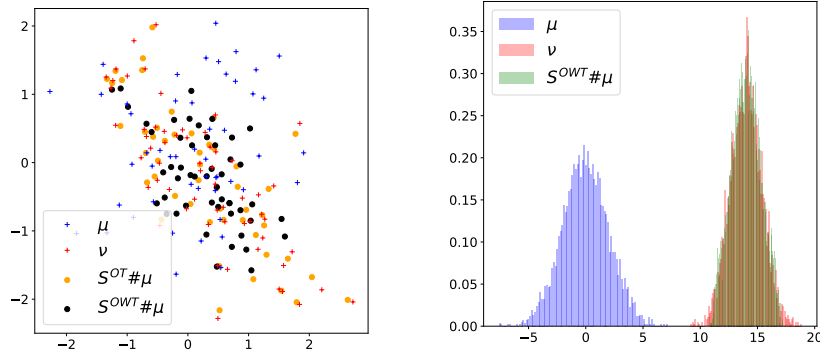


Figure 6: Example of pushforward measures constructed from barycentric projections for two measures μ and ν in two dimensions (left) and one dimension (right).

We next illustrate the differences between the optimal barycentric map and a barycentric map constructed from an OT plan in the classical Kantorovich formulation in Eq. (1). We sampled $r = 50$ observations X_i and $m = 60$ observations Y_i , each sets from a 2D Gaussian. We then defined the source and target distributions as $\mu = \frac{1}{r} \sum \delta_{X_i}$ and $\nu = \frac{1}{m} \sum \delta_{Y_i}$ respectively. Figure 6(left), shows these discrete distributions together with the pushforward measures $S^{OWT}\#\mu$ and $S^{OT}\#\mu$ constructed from the optimal weak plan π^{OWT} and an optimal plan π^{OT} respectively. The measure $S^{OT}\#\mu$ reasonably fits the target distribution ν , since when μ is *a.c.*, $S^{OT}\#\mu = \nu$. In particular, if μ and ν had the same number of points, $S^{OT}\#\mu$ would have matched ν . Regarding the measure $S^{OWT}\#\mu$, recall that $V(\mu | \nu) = \inf_{\eta \leq_c \nu} W_2^2(\mu, \eta) = W_2^2(\mu, S^{OWT}\#\mu)$, and therefore $S^{OWT}\#\mu \leq_c \nu$. Lastly, we have $W_2^2(\mu, \nu) = 0.85$, and $V(\mu | \nu) = 0.52 \leq W_2^2(\mu, S^{OT}\#\mu) = 0.81$ as expected.

In Figure 6(right), we present an example in one dimension, where we sample 4000 observations from $\mathcal{N}(0, 2)$ (resp. $\mathcal{N}(14, 1.4)$) to construct the empirical source measure μ (resp. empirical target measure ν). The distributions μ and ν are presented in the form of histograms. The distribution resulting from the optimal weak transport map $S_\mu^\nu \# \mu$ is in convex order with ν .

B Proofs of Section 3

Proof of Proposition 1. Let $(\mu_m)_m \subset \mathcal{P}_2(\mathbb{R}^d)$ be a minimising sequence of $F(\mu) := \sum_{i=1}^n \lambda_i V(\mu|\nu_i)$ and let $M < \infty$ be such that $F(\mu_m) \leq M$ for all m . Then $(\mu_m)_m$ is tight. Indeed,

$$\begin{aligned} \int \|x\|^2 d\mu_m(x) &\leq 2 \sum_{i=1}^n \lambda_i \inf_{\pi \in \Pi(\mu_m, \nu_i)} \left[\int \|x - \int y d\pi_x(y)\|^2 d\mu_m(x) + \int \left\| \int y d\pi_x(y) \right\|^2 d\mu_m(x) \right] \\ &\leq 2M + 2 \iint \|y\|^2 d\pi_x(y) d\mu_m(x) \leq 2M + 2 \sum_{i=1}^n \lambda_i \int \|y\|^2 d\nu_i(y), \end{aligned}$$

where the second inequality comes from Jensen's inequality. By Prokhorov's theorem, there exists a subsequence still denoted $(\mu_m)_m$ that weakly converges toward a probability measure μ^* . Recall that μ belongs to $\mathcal{P}_2(\mathbb{R}^d)$ since $\|\cdot\|^2$ is a l.s.c. function bounded from below and therefore $\int \|x\|^2 d\mu(x) \leq \liminf_m \int \|x\|^2 d\mu_m(x) < \infty$. By Theorem 6, we have that

$$F(\mu^*) = \sum_{i=1}^n \lambda_i \lim_m V(\mu_m|\nu_i) = \lim_m F(\mu_m) = \min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} F(\mu),$$

thus F admits at least a minimiser. \square

Proof of Lemma 1. By Strassen's theorem, we can build $X' \sim \mu'$ and $X \sim \mu$ in the same probability space, in such a way that $\mathbb{E}(X|X') = X'$. Denote by η^* the law η attaining $\inf_{\eta \leq_c \nu} W_2^2(\mu, \eta)$, and let $(X, Z) = (X, S_\mu^\nu(X))$ be the realisation of the optimal coupling for W_2 of μ and η^* , which can also be constructed in the same probability space due to its specific form. Then, by (the conditional version of) Jensen's inequality we have

$$V(\mu|\nu) = W_2^2(\mu, \eta^*) = \mathbb{E} \left[\mathbb{E}(\|X - S_\mu^\nu(X)\|^2 | X') \right] \geq \mathbb{E} \|X' - \mathbb{E}(S_\mu^\nu(X) | X')\|^2.$$

Recall now that $S_\mu^\nu(X) = \mathbb{E}(Y|X)$, where the conditional expectation is a measurable function only of X , constructed from the joint law $\pi^{\mu, \nu}$. Thus, for every nonnegative convex function ϕ , by applying twice Jensen's inequality we get

$$\mathbb{E}\phi(\mathbb{E}(S_\mu^\nu(X) | X')) \leq \mathbb{E}\phi(S_\mu^\nu(X)) = \mathbb{E}\phi(\mathbb{E}(Y|X)) \leq \mathbb{E}\phi(Y),$$

where $Y \sim \nu$. That is to say, the law η of the r.v. $\mathbb{E}(S_\mu^\nu(X) | X')$ satisfies $\eta \leq_c \nu$. It follows that

$$V(\mu|\nu) \geq W_2^2(\mu', \eta) \geq \inf_{\tilde{\eta} \leq_c \nu} W_2^2(\mu', \tilde{\eta}) = V(\mu'|\nu).$$

This immediately implies that μ' is a weak barycenter whenever μ is. In particular, if μ is a weak barycenter, then so is the Dirac mass supported on its mean. We then deduce that the set of minimisers of $\sum_i \lambda_i V(\mu|\nu_i)$ admits at least a Dirac mass δ_ω and

$$V(\delta_\omega|\nu_i) = \int \|x - \int y d\pi_x(y)\|^2 d\delta_\omega(x) = \|\omega - \mathbb{E}Y_i\|^2.$$

This implies that $\inf_\omega \sum \lambda_i V(\delta_\omega|\nu_i)$ is uniquely attained for $\bar{\omega} = \sum \lambda_i \mathbb{E}Y_i$. \square

Proof of Proposition 2. A probability measure μ is a weak barycenter if and only if $\sum_{i=1}^n \lambda_i V(\mu|\nu_i)$ is equal to the r.h.s. of (7). Let us suppose first that $\mathbb{E}Y_i = m$ for all $1 \leq i \leq n$, in which case the infimum in (5) is equal to 0. Then, μ is a weak barycenter if and only if $\mu \leq_c \nu_i$ for all $1 \leq i \leq n$ by definition of weak optimal transport (2), since in this case $V(\mu|\nu_i) = 0$. The general case can be reduced to the previous one, noting that

$$\begin{aligned} V(\mu|\nu_i) &= \inf_{\eta \leq_c \nu_i} W_2^2(\mu, \eta) \\ &= \inf_{\eta \leq_c \hat{\nu}_i} W_2^2(\hat{\mu}, \eta) + \|\mathbb{E}_\mu(X) - \mathbb{E}_{\nu_i}(Y_i)\|^2 \\ &= V(\hat{\mu}|\hat{\nu}_i) + \|\mathbb{E}_\mu(X) - \mathbb{E}_{\nu_i}(Y_i)\|^2, \end{aligned}$$

so that minimising $\sum_{i=1}^n \lambda_i V(\mu|\nu_i)$ over $\mu \in \mathcal{P}(\mathbb{R}^d)$ is equivalent to minimising $\sum_{i=1}^n \lambda_i V(\mu'|\hat{\nu}_i) + \sum_{i=1}^n \lambda_i \|\omega - \mathbb{E}_{\nu_i} Y_i\|^2$ over the two independent parameters (ω, μ') , with $\omega \in \mathbb{R}^d$ and $\mu' \in \mathcal{P}(\mathbb{R}^d)$ centered, taking μ as the law of $X = X' + \omega$ with $X' \sim \mu'$. \square

Proof of Lemma 2. Thanks to Prop. 3.3 in [5], we have that the (unique) Wasserstein barycenter verifies $\tilde{\mu} = \left(\sum_{i=1}^n \lambda_i T_{\tilde{\mu}}^{\nu_i}\right) \# \tilde{\mu}$ where $T_{\tilde{\mu}}^{\nu_i}$ is the optimal Monge map between $\tilde{\mu}$ and ν_i (see (14)). Moreover, from Proposition 4, a weak barycenter $\bar{\mu}$ also checks $\bar{\mu} = \left(\sum_{i=1}^n \lambda_i S_{\bar{\mu}}^{\nu_i}\right) \# \bar{\mu}$, where $S_{\bar{\mu}}^{\nu_i}$ is the optimal barycentric projection associated to $\bar{\pi}^i$ for $V(\bar{\mu}|\nu_i)$. Therefore, by Jensen's inequality applied twice,

$$\begin{aligned} W_2^2(\bar{\mu}, \tilde{\mu}) &\leq \iint \|x - y\|^2 d\bar{\mu}(x) d\tilde{\mu}(y) = \iint \left\| \sum_{i=1}^n \lambda_i S_{\bar{\mu}}^i(x) - \sum_{i=1}^n \lambda_i T_{\tilde{\mu}}^i(y) \right\|^2 d\bar{\mu}(x) d\tilde{\mu}(y) \\ &\leq \sum_{i=1}^n \lambda_i \iint \|T_{\bar{\mu}}^i(y) - z\|^2 d\bar{\pi}_x^i(z) d\bar{\mu}(x) d\tilde{\mu}(y) = \sum_{i=1}^n \lambda_i \iint \|T_{\bar{\mu}}^i(y) - z\|^2 d\tilde{\mu}(y) d\bar{\pi}^i(x, z) \\ &= \sum_{i=1}^n \lambda_i \iint \|y - z\|^2 d\nu_i(y) d\nu_i(z) = 2 \sum_{i=1}^n \lambda_i \iint (\mathbb{E}\|Y_i\|^2 - \|\mathbb{E}Y_i\|^2). \end{aligned}$$

\square

Proof of Theorem 3. Observe first that, by Theorem 2 and Strassen's theorem, solving the OWT problem (2) provides a unique (in law) coupling of three random variables (X, Y, Z) such that:

- i) (X, Y) has joint law $\pi^{\mu, \nu}$; in particular X and Y have the laws μ and ν respectively,
- ii) $Z = S_{\mu}^{\nu}(X) = \mathbb{E}(Y|X)$ a.s., it has law η^* and it is optimally coupled to X in the sense of the optimal transport problem (1),
- iii) (Z, Y) is a martingale, that is $\mathbb{E}(Y|Z) = Z$ a.s..

Bringing all together we get the decomposition:

$$Y = Z + Y - Z = S_{\mu}^{\nu}(X) + Y - \mathbb{E}(Y|X). \quad (15)$$

Now, by Lemma 1, if $X \sim \mu$ then the Dirac mass $\delta_{\mathbb{E}X}$ is a weak barycenter too. Thus we have on one hand:

$$\sum_{i=1}^n \lambda_i V(\mu|\nu_i) = \sum_{i=1}^n \lambda_i V(\delta_{\mathbb{E}X}|\nu_i). \quad (16)$$

Using Jensen's inequality, we see on the other hand that

$$\begin{aligned} V(\mu|\nu_i) &= \inf_{\eta \leq c\nu_i} W_2^2(\mu, \eta) \\ &= \inf_{\eta \leq c\nu_i} \mathbb{E}\|X - Z\|^2, \text{ with } (X, Z) \text{ an optimal coupling for } W_2^2 \text{ of } \mu \text{ and } \eta \\ &\geq \inf_{\eta \leq c\nu_i} \|\mathbb{E}X - \mathbb{E}Z\|^2 = \inf_{\eta \leq c\nu_i} W_2^2(\delta_{\mathbb{E}X}, \delta_{\mathbb{E}Z}) \\ &\geq \inf_{\tilde{\eta} \leq c\nu_i} W_2^2(\delta_{\mathbb{E}X}, \tilde{\eta}) = V(\delta_{\mathbb{E}X}|\nu_i). \end{aligned}$$

Identity (16) thus implies $V(\mu|\nu_i) = V(\delta_{\mathbb{E}X}|\nu_i)$ for all i . Denoting by η_i the law η attaining $\inf_{\eta \leq c\nu_i} W_2^2(\mu, \eta)$, and by (X, Z_i) the optimal coupling for W_2 of μ and η_i , we see that the latter can only occur if the equality case $\mathbb{E}\|X - Z_i\|^2 = \|\mathbb{E}X - \mathbb{E}Z_i\|^2$ in Jensen's inequality holds. This implies that $X - Z_i$ is deterministic for each i . Since $\mathbb{E}Z_i = \mathbb{E}Y_i$, we thus must have $X - Z_i = \mathbb{E}X - \mathbb{E}Y_i$. Taking $Z = Z_i$ and $Y = Y_i$ in Eq. (15), and noting that $S_{\mu}^{\nu_i}(X) = Z_i = X - (\mathbb{E}X - \mathbb{E}Y_i)$ the statement follows. \square

C Proofs of Section 3.3

Proof of Lemma 3. Let (μ_m) and (ν_m) , $m \in \mathbb{N}$, be two sequences in $\mathcal{P}_2(\mathbb{R}^d)$ respectively converging to μ and ν w.r.t. W_2 . Then, (μ_m) and (ν_m) are tight and thus the sequence $(\pi^m) := (\pi^{\mu_m, \nu_m})$ is tight too. Let π^{m_k} be a weakly convergent subsequence and π its limit. By Proposition 2.8 in [8] we have

$$\liminf_m V(\mu_m | \nu_m) = \liminf_m \int \|x - \int y d\pi_x^{m_k}\|^2 d\mu_{m_k}(x) \geq \int \|x - \int y d\pi_x\|^2 d\mu(x) \geq V(\mu | \nu).$$

However, we have $\lim_m V(\mu_m | \nu_m) = V(\mu | \nu)$ thanks to Theorem 6, hence $\int \|x - \int y d\pi_x\|^2 d\mu(x) = V(\mu | \nu)$. By uniqueness of the optimum for problem (2) we deduce that $\pi = \pi^{\mu, \nu}$. Since the same holds true for any weak limiting point of (π^m) , it follows that π^m weakly converges to $\pi^{\mu, \nu}$. Last, since $\int \|x\|^2 + \|y\|^2 d\pi^m(x, y) = \int \|x\|^2 d\mu_m(x) + \int \|y\|^2 d\nu_m(y)$, this quantity converges to $\int \|x\|^2 d\mu(x) + \int \|y\|^2 d\nu(y) = \int \|x\|^2 + \|y\|^2 d\pi(x, y)$, whence $W_2(\pi^m, \pi) \rightarrow 0$, and $(\mu, \nu) \in (\mathcal{P}_2(\mathbb{R}^d))^2 \mapsto \pi^{\mu, \nu} \in \mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d)$ is continuous, as required (hence measurable).

We now establish the joint measurability of $(x, \nu) \in \mathbb{R}^d \times \mathcal{P}_2(\mathbb{R}^d) \mapsto S_\mu^\nu(x)$ for fixed μ . Notice this is a stronger statement than just measurability in the x variable, for each (μ, ν) . Write $\bar{B}(x, r)$ for the closed ball of radius $r > 0$ centered at x . One easily checks that the function

$$(x, \pi) \mapsto \Psi_r(x, \pi) := \frac{\int y \mathbf{1}_{\{(y, z): z \in \bar{B}(x, r)\}} d\pi(z, y)}{\int \mathbf{1}_{\{(y, z): z \in \bar{B}(x, r)\}} d\pi(z, y)}$$

is measurable w.r.t. the pair (x, π) , the two integrals being limits of integrals with respect to $d\pi(z, y)$, of some bounded continuous functions of (x, y, z) . Thus, $\limsup_{r \rightarrow 0} \Psi_r(x, \pi)$, $\liminf_{r \rightarrow 0} \Psi_r(x, \pi)$ and the function $\Phi(x, \pi) := \limsup_{r \rightarrow 0} \Psi_r(x, \pi) \mathbf{1}_{\{\limsup_{r \rightarrow 0} \Psi_r(x, \pi) = \liminf_{r \rightarrow 0} \Psi_r(x, \pi)\}}$ depend in a measurable way on (x, π) . It follows that $(x, \mu, \nu) \mapsto \Phi(x, \pi^{\mu, \nu})$ is measurable as the composition of two measurable functions. But notice that for each fixed $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ one has $\Psi_r(x, \pi^{\mu, \nu}) = \frac{\int_{\bar{B}(x, r)} [\int y d\pi_z(y)] d\mu(z)}{\mu(\bar{B}(x, r))}$ which, by the Lebesgue derivation theorem for Radon measures (*see e.g.* [15]), converges $d\mu(x)$ a.s. in x , to $\int y d\pi_x(y) = S_\mu^\nu(x)$. Thus, for each $\mu \in \mathcal{P}_2(\mathbb{R}^d)$,

$$S_\mu^\nu(x) = \Phi(x, \pi^{\mu, \nu}) \quad \text{for all } \nu \in \mathcal{P}_2(\mathbb{R}^d) \text{ and } d\mu(x) \text{ a.e. } x,$$

with $(x, \nu) \mapsto \Phi(x, \pi^{\mu, \nu})$ a measurable function. The conclusion follows. \square

Proof of Proposition 3. By Theorem 6.16 in [41], we know that there exists a sequence of discretely supported distributions $(\mathbb{Q}_n)_n \subset \mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$ of the form $\mathbb{Q}_n = \sum_{i=1}^n \lambda_i \delta_{\nu_i}$, with $(\lambda_i)_{1 \leq i \leq n}$ in the simplex, and such that $W_2^2(\mathbb{Q}, \mathbb{Q}_n) := \inf_{\pi \in \Pi(\mathbb{Q}, \mathbb{Q}_n)} \int W_2^2(\nu, \tilde{\nu}) d\pi(\nu, \tilde{\nu}) \rightarrow 0$. We set

$$L_n(\mu) := \int_{\mathcal{P}_2(\mathbb{R}^d)} V(\mu | \nu) d\mathbb{Q}_n(\nu) = \sum_{i=1}^n \lambda_i V(\mu | \nu_i).$$

We denote $\mu^n \in \mathcal{P}_2(\mathbb{R}^d)$ the minimiser of L_n . Let us prove that $(\mu^n)_n$ is tight. First, μ^n admits moments of order 2 thanks to Jensen's inequality:

$$\begin{aligned} \int \|x\|^2 d\mu^n(x) &\leq \sum_{i=1}^n \lambda_i \left[\int \|x - S_{\mu^n}^{\nu_i}(x)\|^2 d\mu^n(x) + \int \|S_{\mu^n}^{\nu_i}(x)\|^2 d\mu^n(x) \right] \\ &\leq \sum_{i=1}^n \lambda_i V(\mu^n | \nu_i) + \sum_{i=1}^n \lambda_i \int \|y\|^2 d\nu_i(y) \\ &\leq \sum_{i=1}^n \lambda_i V(\mu | \nu_i) + \sum_{i=1}^n \lambda_i \int \|y\|^2 d\nu_i(y) \quad \text{for some } \mu \in \mathcal{P}_2(\mathbb{R}^d) \text{ since } \mu^n \text{ minimises } L_n \\ &\leq 2 \int \|x\|^2 d\mu(x) + 3 \sum_{i=1}^n \lambda_i \int \|y\|^2 d\nu_i(y), \end{aligned}$$

where the last inequality comes from $V(\mu|\nu_i) = \int \|x - S_{\mu}^{\nu_i}(x)\|^2 d\mu(x) \leq 2 \int \|x\|^2 d\mu(x) + 2 \int \|S_{\mu}^{\nu_i}(x)\|^2 d\mu(x)$. Moreover, since $W_2^2(\mathbb{Q}, \mathbb{Q}_n) \rightarrow 0$, we have (Lemma 5.1.7 in [6]) that $\int \psi(\nu) d\mathbb{Q}_n(\nu) \rightarrow \int \psi(\nu) d\mathbb{Q}(\nu)$ for any function ψ such that $|\psi(\nu)| \leq a + bW_2^2(\nu, \nu_0)$, $a, b \geq 0$. In particular, choosing $\psi(\nu) = W_2^2(\nu, \delta_0) = \int \|y\|^2 d\nu(y)$, it implies that $\sum_{i=1}^n \lambda_i \int \|y\|^2 d\nu_i(y) \rightarrow \int \int \|y\|^2 d\nu(y) d\mathbb{Q}(\nu) < \infty$. Therefore $(\sum_{i=1}^n \lambda_i \int \|y\|^2 d\nu_i(y))_n$ is bounded and $(\mu^n)_n$ is tight. Thus by Prokhorov's theorem, there exists a subsequence, still denoted $(\mu^n)_n$, that converges towards $\bar{\mu}$.

Let us now prove that this particular $\bar{\mu}$ minimises the function $L : \mu \mapsto \int_{\mathcal{P}_2(\mathbb{R}^d)} V(\mu|\nu) d\mathbb{Q}(\nu)$. First, let $\eta \in \mathcal{P}_2(\mathbb{R}^d)$, still by Lemma 5.1.7 in [6] and since $V(\eta|\nu) \leq W_2^2(\eta, \nu)$, we get that $L(\eta) = \int V(\eta|\nu) d\mathbb{Q}(\nu) \geq \liminf_{n \rightarrow \infty} \int V(\eta|\nu) d\mathbb{Q}_n(\nu)$. Since for each n , the distribution μ^n minimises L_n , we have

$$\liminf_{n \rightarrow \infty} \int V(\eta|\nu) d\mathbb{Q}_n(\nu) \geq \liminf_{n \rightarrow \infty} \int V(\mu^n|\nu) d\mathbb{Q}_n(\nu). \quad (17)$$

Thanks to Fatou's Lemma for sequences of measures $(\mathbb{Q})_n$ (see [21]), we have that

$$\liminf_{n \rightarrow \infty} \int V(\mu^n|\nu) d\mathbb{Q}_n(\nu) \geq \int \liminf_{n \rightarrow \infty} V(\mu^n|\nu) d\mathbb{Q}(\nu) = \int V(\bar{\mu}|\nu) d\mathbb{Q}(\nu),$$

where the last equality comes from the lower semi-continuity of V (Theorem 2.9 in [8]). This proves that $\bar{\mu}$ minimises L . \square

D Proofs of Section 4

The proof of Theorem 4, on the continuity of $G : \mu \mapsto (\sum_{i=1}^n \lambda_i S_{\mu}^{\nu_i}) \# \mu$, leans on the two following technical lemmas.

Lemma 5. *Let $(\rho_m)_m$ be a given sequence and ν be a fixed law in $\mathcal{P}_2(\mathbb{R}^d)$. For each m , let $S_m := S_{\rho_m}^{\nu}$ denote the barycenter map associated with the optimal coupling $\pi^{\rho_m, \nu}$ for (2). Then, the sequence of laws $(S_m \# \rho_m)_m$ has uniformly integrable second moments.*

Proof of Lemma 5. Let (X_m, Y_m) be a pair of random variables (r.v.) with joint law $\pi^{\rho_m, \nu}$, defined on some probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Notice that $S_m \# \rho_m$ is the law of the r.v. $\mathbb{E}(Y_m|X_m)$. Then, for each $M, K \geq 0$ we have

$$\begin{aligned} \int_{\{\|x\|^2 \geq M\}} \|x\|^2 dS_m \# \rho_m(x) &= \mathbb{E}(\|\mathbb{E}(Y_m|X_m)\|^2 \mathbf{1}_{\{\|\mathbb{E}(Y_m|X_m)\|^2 \geq M\}}) \\ &\leq \mathbb{E}(\mathbb{E}(\|Y_m\|^2|X_m) \mathbf{1}_{\{\|\mathbb{E}(Y_m|X_m)\|^2 \geq M\}}) \\ &= \mathbb{E}(\|Y_m\|^2 \mathbf{1}_{\{\|\mathbb{E}(Y_m|X_m)\|^2 \geq M, \|Y_m\|^2 \geq K\}}) \\ &\quad + \mathbb{E}(\|Y_m\|^2 \mathbf{1}_{\{\|\mathbb{E}(Y_m|X_m)\|^2 \geq M, \|Y_m\|^2 < K\}}) \\ &\leq \mathbb{E}(\|Y_m\|^2 \mathbf{1}_{\|Y_m\|^2 \geq K}) + \frac{K}{M} \mathbb{E}(\|\mathbb{E}(Y_m|X_m)\|^2), \end{aligned}$$

where we have used Jensen's inequality and the fact that $\mathbb{E}(Y_m|X_m)$ is measurable w.r.t. the σ -field generated by X_m . Applying Jensen's inequality to the last term again, and recalling that Y_m has law $\nu \in \mathcal{P}_2(\mathbb{R}^d)$, we deduce that

$$\sup_m \int_{\{\|x\|^2 \geq M\}} \|x\|^2 dS_m \# \rho_m(x) \leq \int_{\{\|x\|^2 \geq K\}} \|x\|^2 d\nu(x) + \frac{K}{M} \int \|x\|^2 d\nu(x), \quad (18)$$

which is smaller than a given $\varepsilon > 0$, by choosing $K > 0$ and then $M > 0$ large enough. \square

Lemma 6. *Let $(\rho_m)_n, \rho$ in $\mathcal{P}_2(\mathbb{R}^d)$ be such that $W_2(\rho_m, \rho) \rightarrow 0$. We have:*

- i) *For each $\nu \in \mathcal{P}_2(\mathbb{R}^d)$ the sequence of laws $((id, S_{\rho_m}^{\nu}) \# \rho_m)_m$ converges w.r.t. W_2 in $\mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d)$ to $(id, S_{\rho}^{\nu}) \# \rho$.*

ii) There exists in some probability space $(\Omega, \mathcal{F}, \mathbb{P})$, a sequence of r.v. $(X_m)_m$ of laws $(\rho_m)_m$ and a r.v. X of law ρ such that, for each $\nu \in \mathcal{P}_2(\mathbb{R}^d)$, the sequence $(X_m, S_{\rho_m}^\nu(X_m))_m$ (with laws $((\text{id}, S_{\rho_m}^\nu) \# \rho_m)_m$) converges in $\mathbb{L}^2(\Omega, \mathcal{F}, \mathbb{P})$ to $(X, S_\rho^\nu(X))$ (with law $(\text{id}, S_\rho^\nu) \# \rho$).

Proof of Lemma 6. For the entire proof, we fix a $\nu \in \mathcal{P}_2(\mathbb{R}^d)$ and we write $S_m := S_{\rho_m}^\nu$ and $S := S_\rho^\nu$ for simplicity.

i) By Theorem 2 and Part 1. of Theorem 1.5 in [9], $(S_m \# \rho_m)_m$ converges to $S \# \rho$ w.r.t. W_1 and, by Lemma 5, also with respect to W_2 . In particular, the sequence $((\text{id}, S_m) \# \rho_m)_m$ has tight marginals, and therefore it is tight too.

Let us identify its weak limiting points. For simplicity we rename $((\text{id}, S_m) \# \rho_m)_m$ a weakly convergent subsequence. By the previous discussion, its weak limit $d\hat{\rho}(x, z)$ clearly has first and second marginal laws equal to $d\rho(x)$ and $dS \# \rho(z)$ respectively. Moreover, $\int \|x\|^2 d\rho_m(x) + \int \|z\|^2 dS_m \# \rho_m(z) \rightarrow \int \|x\|^2 + \|z\|^2 d\hat{\rho}(x, z)$, hence $((\text{id}, S_m) \# \rho_m)_m$ converges to some $\hat{\pi}$ with respect to W_2 in $\mathcal{P}_2(\mathbb{R}^d \times \mathbb{R}^d)$.

Now, by the characterisation of optimisers in Theorem 2, we have $V(\rho_m | \nu) = W_2^2(\rho_m, S_m \# \rho_m) = \int \|x - S_m(x)\|^2 d\rho_m(x)$. Taking $m \rightarrow \infty$, and thanks to Theorem 6, we finally obtain

$$V(\rho | \nu) = W_2^2(\rho, S \# \rho) = \int \|x - z\|^2 d\hat{\pi}(x, z).$$

In particular, using again Theorem 2, we conclude that $d\hat{\pi}(x, z)$ must be of the form $(\text{id}, S) \# \rho$.

ii) By Skorohod's representation theorem, one can construct simultaneously in some probability space $(\Omega, \mathcal{F}, \mathbb{P})$, a sequence of r.v. $(X_m)_m$ of laws $(\rho_m)_m$ and a r.v. X of law ρ such that $(X_m)_m$ converges \mathbb{P} -a.s. to X . Moreover, since the sequence $(\rho_m)_m$ converges w.r.t. W_2 in $\mathcal{P}_2(\mathbb{R}^d)$, it has uniformly integrable second order moments. It follows that the sequence of r.v. $(|X_m|^2)_m$ is uniformly integrable and, by the Vitali convergence theorem, that $(X_m)_m$ also converges to X in $\mathbb{L}^2(\Omega, \mathcal{F}, \mathbb{P})$.

Now, by Lemma 5, the sequence of r.v. $(|S_m(X_m)|^2)_n$ is uniformly integrable too. Thus, by the Vitali convergence theorem, the statement will follow by proving that $S_m(X_m)$ converges in \mathbb{P} -probability to $S(X)$.

For each $N \in \mathbb{N}$, let $y \mapsto (y)^N$ denote the truncation of a vector $y \in \mathbb{R}^d$ obtained by projecting it onto the centered ball of radius N , $(y)^N := (1 \wedge \frac{N}{|y|})y$, which is a 1-Lipschitz function bounded by N . By Theorem 2, the functions $S_m^N := (S_m)^N$ are then 1-Lipschitz and bounded uniformly in $m \in \mathbb{N}$. Therefore, by the Arzela-Ascoli theorem, their restrictions to each compact cylinder set R of \mathbb{R}^d defines a relatively compact set of functions, with respect to the uniform topology in $C(R, \mathbb{R}^d)$. It follows by a diagonal argument that some subsequence $(S_{m_k}^N)_k$ converges, uniformly on compact sets, to some continuous function \tilde{S} on \mathbb{R}^d . Since X_n converges a.s. to the finite value X , we deduce that \mathbb{P} -a.s. as $k \rightarrow \infty$,

$$(X_{m_k}, S_{m_k}^N(X_{m_k})) \rightarrow (X, \tilde{S}(X)).$$

Notice now that $(X_{m_k}, S_{m_k}^N(X_{m_k}))$ has the law $(\text{id}, (\cdot)^N \circ S_{m_k}) \# \rho_{m_k}$ for each k and thus, by part a) and continuity of the mapping $(x, y) \mapsto (x, (y)^N)$, the r.v. $(X, \tilde{S}(X))$ has the law $(\text{id}, (\cdot)^N \circ S) \# \rho$. Hence we deduce that

$$(X, \tilde{S}(X)) = (X, (S(X))^N)$$

\mathbb{P} -almost surely. The previous arguments can be applied not just to $(X_m)_m$ but to any subsequence of it. That is, we can similarly prove that any subsequence of $(X_m, (S_m(X_m))^N)_m$ has a subsequence that a.s. converges to $(X, (S(X))^N)$. This means that, for each $N \in \mathbb{N}$

$$(X_m, (S_m(X_m))^N) \rightarrow (X, (S(X))^N)$$

in \mathbb{P} -probability when $n \rightarrow \infty$. To conclude, by tightness we can find for each $\eta > 0$ some $N \in \mathbb{N}$ large enough so that $\mathbb{P}(|S(X)| \geq N) \leq \eta$ and $\mathbb{P}(|S_m(X_m)| \geq N) \leq \eta$ for all $m \in \mathbb{N}$, which yields for each $\varepsilon > 0$,

$$\mathbb{P}(|S_m(X_m) - S(X)| \geq \varepsilon) \leq 2\eta + \mathbb{P}(|(S_m(X_m))^N - (S(X))^N| \geq \varepsilon).$$

Thus $\limsup_m \mathbb{P}(|S_m(X_m) - S(X)| \geq \varepsilon) \leq 2\eta$ for arbitrary $\eta > 0$ or, equivalently, $\mathbb{P}(|S_m(X_m) - S(X)| \geq \varepsilon) \rightarrow 0$ as $m \rightarrow \infty$, which concludes the proof of b). \square

We can now proceed to the proof of continuity of $G : \mu \mapsto (\sum_{i=1}^n \lambda_i S_\mu^{\nu_i}) \# \mu$.

Proof of Theorem 4. Let $(\rho_m)_m, \rho$ in $\mathcal{P}_2(\mathbb{R}^d)$ such that $W_2(\rho_m, \rho) \rightarrow 0$. We need to prove that $W_2^2(G(\rho_m), G(\rho)) \rightarrow 0$. For each m , we write $S_m^i := S_{\rho_m}^{\nu_i}$ and $S^i := S_\rho^{\nu_i}$.

By Lemma 6.ii), there exists in some probability space a sequence $(X_m)_m$ of laws $(\rho_m)_m$ and a r.v. X of law ρ such that

$$(S_m^1(X_m), \dots, S_m^n(X_m)) \rightarrow (S^1(X), \dots, S^n(X)) \quad \text{in } \mathbb{L}^2(\mathbb{P}).$$

Therefore, $\sum_{i=1}^n \lambda_i S_m^i(X_m)$ converges to $\sum_{i=1}^n \lambda_i S^i(X)$ in $\mathbb{L}^2(\mathbb{P})$. Since $\sum_{i=1}^n \lambda_i S_m^i(X_m)$ has law $G(\rho_m)$ and $\sum_{i=1}^n \lambda_i S^i(X)$ has law $G(\rho)$, the proof is complete. \square

Proof of Proposition 4. As in [5], we easily see that

$$\sum_{i=1}^n \lambda_i \int \|x - S_\mu^{\nu_i}(x)\|^2 d\mu(x) = \sum_{i=1}^n \lambda_i \int \|\bar{S}(x) - S_\mu^{\nu_i}(x)\|^2 d\mu(x) + \int \|x - \bar{S}(x)\|^2 d\mu(x).$$

But $\int \|x - S_\mu^{\nu_i}(x)\|^2 d\mu(x) = W_2^2(\mu, S_\mu^{\nu_i} \# \mu)$ since from Thm 1.4 in [8] the barycentric map $S_\mu^{\nu_i}$ is an optimal map for the Monge problem between μ and $S_\mu^{\nu_i} \# \mu$. Moreover, by definition $G(\mu) = \bar{S} \# \mu$, therefore $\int \|x - \bar{S}(x)\|^2 d\mu(x) \geq W_2^2(\mu, G(\mu))$. Finally, since $S_\mu^{\nu_i} \# \mu \leq_c \nu_i$, we have that $\int \|\bar{S}(x) - S_\mu^{\nu_i}(x)\|^2 d\mu(x) \geq V(G(\mu)|\nu_i)$. This, recalling that $V(\mu|\nu_i) = W_2^2(\mu, S_\mu^{\nu_i} \# \mu)$, yields

$$\sum_{i=1}^n \lambda_i V(\mu|\nu_i) \geq \sum_{i=1}^n \lambda_i V(G(\mu)|\nu_i) + W_2^2(\mu, G(\mu)). \quad (19)$$

Therefore, if μ is a weak barycenter, we readily get that $\mu = G(\mu)$. \square

Proof of Proposition 5. As in the proof of Theorem 4, we denote S_k^i the optimal barycentric projection associated to $\pi^{k,i} \in \Pi(\mu_k, \nu_i)$. First, we easily have that $\mu_{k+1} \in \mathcal{P}_2(\mathbb{R}^d)$, indeed by Jensen's inequality

$$\int \|x\|^2 d\mu_{k+1}(x) = \iint \left\| \sum_{i=1}^n \lambda_i S_k^i(x) \right\|^2 d\mu_k(x) \leq \sum_{i=1}^n \lambda_i \int \|y\|^2 d\nu_i(y) < \infty.$$

Then $(\mu_k)_k$ is tight, with uniformly integrable 2-moments by Lemma 5. Therefore $(\mu_k)_k$ admits a convergent subsequence in W_2 . Let $\tilde{\mu}$ be a weak limit of a subsequence $(\mu_{k_j})_j$, then we have $W_2(\mu_{k_j}, \tilde{\mu}) \xrightarrow{j \rightarrow \infty} 0$. By continuity of G in Theorem 4, we get $W_2(\mu_{k_j+1}, G(\tilde{\mu})) \xrightarrow{j \rightarrow \infty} 0$.

Moreover, by Theorem 6 we have for $F(\mu) := \sum_{i=1}^n \lambda_i V(\mu|\nu_i)$ that $F(\mu_{k_j}) \rightarrow F(\tilde{\mu})$ and $F(\mu_{k_j+1}) \rightarrow F(G(\tilde{\mu}))$ as $j \rightarrow \infty$. Let us prove that these two limits coincide. From (19), we have

$$F(\mu_{k_j}) \geq \sum_{i=1}^n \lambda_i V(G(\mu_{k_j})|\nu_i) = \sum_{i=1}^n \lambda_i V(\mu_{k_j+1}|\nu_i) = F(\mu_{k_j+1}).$$

Iterating this inequality leads to $F(\mu_{k_j}) \geq F(\mu_{k_j+1}) \geq F(\mu_{k_j+2})$ which yields $F(\tilde{\mu}) = F(G(\tilde{\mu}))$ and then $\tilde{\mu} = G(\tilde{\mu})$, using inequality (19). Thus $(\mu_{k_j})_j$ converges w.r.t. W_2 to a probability distribution $\tilde{\mu}$ which is a fixed point of G . \square

Proof of Lemma 4. The proof is similar to that of [34, Lemma 3.8]. For the sake of clarity, we rewrite it in our setting. We assume that $x = \int S_\mu^\nu(x) d\mathbb{Q}(\nu)$, $\bar{\mu}(x)$ -a.s. is not true, then

$$\begin{aligned} 0 &< \int \|x - \int S_\mu^\nu(x) d\mathbb{Q}(\nu)\|^2 d\bar{\mu}(x) \\ &= \int \|x\|^2 d\bar{\mu}(x) - 2 \iint \langle x, S_\mu^\nu(x) \rangle d\mathbb{Q}(\nu) d\bar{\mu}(x) + \int \left\| \int S_\mu^\nu(x) d\mathbb{Q}(\nu) \right\|^2 d\bar{\mu}(x). \end{aligned}$$

Moreover, $S_{\bar{\mu}}^{\mu} \# \bar{\mu} \leq_c \mu$, therefore by Theorem 1.4 in [8], we get

$$\begin{aligned} \int V \left(\left[\int S_{\bar{\mu}}^{\nu} d\mathbb{Q}(\nu) \right] \# \bar{\mu} | \mu \right) d\mathbb{Q}(\mu) &\leq \int \int \| S_{\bar{\mu}}^{\nu} d\mathbb{Q}(\nu) - S_{\bar{\mu}}^{\mu} \|_{\mathbb{L}^2(\bar{\mu})}^2 d\mathbb{Q}(\mu) \\ &= \int \int \| S_{\bar{\mu}}^{\nu}(x) \|^2 d\bar{\mu}(x) d\mathbb{Q}(\nu) - \int \int \| S_{\bar{\mu}}^{\nu} d\mathbb{Q}(\nu) \|^2 d\bar{\mu}(x). \end{aligned}$$

Finally, noticing that $\int \int \|x - S_{\bar{\mu}}^{\nu}(x)\|^2 d\bar{\mu}(x) d\mathbb{Q}(\nu) = \int V(\bar{\mu} | \nu) d\mathbb{Q}(\nu)$, we hence get

$$\int V \left(\left[\int S_{\bar{\mu}}^{\nu} d\mathbb{Q}(\nu) \right] \# \bar{\mu} | \mu \right) d\mathbb{Q}(\mu) < \int V(\bar{\mu} | \nu) d\mathbb{Q}(\nu),$$

which is in contradiction with $\bar{\mu}$ weak barycenter of \mathbb{Q} . \square

In order to study the convergence of the iterative scheme in (10), we define the following objects:

$$L(\mu) := \frac{1}{2} \int V(\mu | \nu) d\mathbb{Q}(\nu) \quad (20)$$

$$H(\mu)(x) := - \int (S_{\mu}^{\nu} - \text{id}) d\mathbb{Q}(\nu)(x) \quad x \in \mathbb{R}^d. \quad (21)$$

Moreover, we denote by $\{\mathcal{F}_k\}_k$ the filtration of the i.i.d. sample $\nu^k \sim \mathbb{Q}$, namely \mathcal{F}_{-1} is the trivial sigma-algebra and \mathcal{F}_{k+1} is the sigma-algebra generated by ν^0, \dots, ν^k and therefore μ_k in (10) is \mathcal{F}_k -measurable.

The next Proposition is needed to prove Theorem 5.

Proposition 6. *For the sequence $(\mu_k)_k$ defined in (10), we have*

$$\mathbb{E}(L(\mu_{k+1}) - L(\mu_k) | \mathcal{F}_k) \leq \gamma_k^2 L(\mu_k) - \gamma_k \|H(\mu_k)\|_{\mathbb{L}^2(\mu_k)}^2. \quad (22)$$

Proof. The arguments are similar to the ones used for the population Wasserstein barycenter iterative scheme in the proof of Proposition 4.6 in [34]. Let us set them for the present problem. Let $\nu \in \text{supp}(\mathbb{Q})$, then $((1 - \gamma_k)\text{id} + \gamma_k S_{\mu_k}^{\nu^k}, S_{\mu_k}^{\nu^k}) \# \mu_k$ belongs to $\Pi(\mu_{k+1}, S_{\mu_k}^{\nu^k} \# \mu_k)$. Therefore we have

$$\begin{aligned} V(\mu_{k+1} | \nu) &\leq W_2^2(\mu_{k+1}, S_{\mu_k}^{\nu^k} \# \mu_k) \quad \text{since } S_{\mu_k}^{\nu^k} \# \mu_k \leq_c \nu \\ &\leq \int \| (1 - \gamma_k)x + \gamma_k S_{\mu_k}^{\nu^k}(x) - S_{\mu_k}^{\nu}(x) \|^2 d\mu_k(x) \\ &= \int \|x - S_{\mu_k}^{\nu}(x)\|^2 d\mu_k(x) - 2\gamma_k \int \langle x - S_{\mu_k}^{\nu}(x), x - S_{\mu_k}^{\nu^k}(x) \rangle d\mu_k(x) \\ &\quad + \gamma_k^2 \int \|x - S_{\mu_k}^{\nu^k}(x)\|^2 d\mu_k(x) \\ &= V(\mu_k | \nu) + \gamma_k^2 V(\mu_k | \nu^k) - 2\gamma_k \int \langle x - S_{\mu_k}^{\nu}(x), x - S_{\mu_k}^{\nu^k}(x) \rangle d\mu_k(x). \end{aligned}$$

Integrating with respect to ν , and divided by 2 we get

$$L(\mu_{k+1}) \leq L(\mu_k) + \frac{\gamma_k^2}{2} V(\mu_k | \nu^k) - \gamma_k \int \langle H(\mu_k)(x), x - S_{\mu_k}^{\nu^k}(x) \rangle d\mu_k(x).$$

We can then take the conditional expectation with respect to the filtration \mathcal{F}_k , knowing that μ_k is \mathcal{F}_k -measurable and that ν^k is independently sampled from \mathbb{Q} , we have

$$\begin{aligned} \mathbb{E}(L(\mu_{k+1}) | \mathcal{F}_k) &\leq L(\mu_k) + \frac{\gamma_k^2}{2} \int V(\mu_k | \nu) d\mathbb{Q}(\nu) - \gamma_k \int \langle H(\mu_k)(x), \int x - S_{\mu_k}^{\nu}(x) d\mathbb{Q}(\nu) \rangle d\mu_k(x) \\ &= L(\mu_k) + \gamma_k^2 L(\mu_k) - \gamma_k \int \langle H(\mu_k)(x), H(\mu_k)(x) \rangle d\mu_k(x) \\ &= (1 + \gamma_k^2) L(\mu_k) - \gamma_k \|H(\mu_k)\|_{\mu_k}^2. \end{aligned}$$

\square

Proof of Theorem 5. Having Lemma 4 and Proposition 6, we can proceed in a similar way as in the proof of [34, Theorem 4.7]. Let us first check that the second moments of $(\mu_k)_k$ are a.s. bounded by some constant. Let $\bar{\mu}$ be a weak population barycenter, i.e. $\bar{\mu}$ minimises L defined in (20). We introduce the sequences

$$h_k := L(\mu_k) - L(\bar{\mu}) \quad \text{and} \quad \alpha_k := \prod_{i=1}^{k-1} \frac{1}{1 + \gamma_k^2}.$$

We first notice that $h_k \geq 0$ for all k . From condition (11), the sequence $(\alpha_k)_k$ converges to some $\alpha_\infty > 0$. By Proposition 6, we have

$$\begin{aligned} \mathbb{E}(h_{k+1} - (1 + \gamma_k^2)h_k | \mathcal{F}_k) &\leq \gamma_k^2 L(\bar{\mu}) - \gamma_k \|H(\mu_k)\|_{\mu_k}^2 \leq \gamma_k^2 L(\bar{\mu}) \\ \Rightarrow \mathbb{E}(\alpha_{k+1}h_{k+1} - \alpha_k h_k | \mathcal{F}_k) &\leq \alpha_{k+1} \gamma_k^2 L(\bar{\mu}) \quad \text{by multiplying by } \alpha_{k+1}. \end{aligned} \quad (23)$$

We define

$$\delta_k := \begin{cases} 1 & \text{if } \mathbb{E}(\alpha_{k+1}h_{k+1} - \alpha_k h_k | \mathcal{F}_k) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Therefore

$$\begin{aligned} \sum_{k=1}^{\infty} \mathbb{E}(\delta_k (\alpha_{k+1}h_{k+1} - \alpha_k h_k)) &= \sum_{k=1}^{\infty} \mathbb{E}(\delta_k \mathbb{E}(\alpha_{k+1}h_{k+1} - \alpha_k h_k | \mathcal{F}_k)) \\ &\leq L(\bar{\mu}) \sum_{k=1}^{\infty} \alpha_{k+1} \gamma_k^2 \leq L(\bar{\mu}) \sum_{k=1}^{\infty} \gamma_k^2 < \infty. \end{aligned}$$

Since $h_k \alpha_k \geq 0$, by the quasi-martingale convergence theorem $(h_k \alpha_k)_k$ converges almost surely, but as $(\alpha_k)_k$ converges to α_∞ , then $(h_k)_k$ also converges almost surely to some $h_\infty \geq 0$. Taking expectations in Eq. (23) and summing in k , we get

$$\mathbb{E}(\alpha_{k+1}h_{k+1}) \leq \alpha_0 h_0 + L(\bar{\mu}) \sum_{m=1}^k \alpha_{m+1} \gamma_m^2 \leq C.$$

We then obtain by Fatou's Lemma $\liminf_{k \rightarrow \infty} \mathbb{E}(\alpha_{k+1}h_{k+1}) \geq \mathbb{E}(\liminf_{k \rightarrow \infty} \alpha_{k+1}h_{k+1}) = \mathbb{E}(\alpha_\infty h_\infty)$, and since $\alpha_\infty < \infty$, we have that $\mathbb{E}(h_\infty) < \infty$, so h_∞ is almost surely finite. This implies that $L(\mu_k)$ has a finite a.s. limit, that we call ℓ . Therefore by convexity of W_2 in [41, Theorem 4.8],

$$\frac{1}{2} W_2^2(\mu_k, \int S_{\mu_k}^\nu \# \mu_k d\mathbb{Q}(\nu)) \leq \frac{1}{2} \int W_2^2(\mu_k, S_{\mu_k}^\nu \# \mu_k) d\mathbb{Q}(\nu) = L(\mu_k) \leq \ell + 1$$

for k large enough. Since $\mathbb{Q} \in \mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$, we have that $\int S_{\mu_k}^\nu \# \mu_k d\mathbb{Q}(\nu) \in \mathcal{P}_2(\mathbb{R}^d)$, and the second moments of $(\mu_k)_k$ are a.s. bounded by some constant M . By Markov's inequality, and since closed balls in \mathbb{R}^d are compact, the sequence $(\mu_k)_k$ is a.s. tight. Also, for $q < 2$, we have by Hölder and Chebyshev inequalities

$$\int_{\|x\| > R} \|x\|^2 d\mu_k(x) \leq \frac{1}{R^{1-q/2}} \int \|x\|^2 d\mu_k(x) \leq \frac{M}{R^{1-q/2}},$$

so $(\mu_k)_k$ is a.s. relatively compact in W_q thanks to [40, Theorem 7.12] and

$$\lim_{R \rightarrow \infty} \limsup_{k \rightarrow \infty} \int_{\|x\| > R} \|x\|^2 d\mu_k(x) \leq \lim_{R \rightarrow \infty} \limsup_{k \rightarrow \infty} \frac{M}{R^{1-q/2}} = 0.$$

From (23), we take the expectation and sum over k , then

$$\mathbb{E}(\alpha_{k+1}h_{k+1}) - \alpha_0 h_0 \leq L(\bar{\mu}) \sum_{m=1}^k \alpha_{m+1} \gamma_m^2 - \sum_{m=1}^k \alpha_{m+1} \gamma_m \mathbb{E} \left(\|H(\mu_m)\|_{\mathbb{L}_2(\mu_m)}^2 \right).$$

Taking limit inferior on k , we have by Fatou on the l.h.s. and monotone convergence on the r.h.s.

$$-\infty < \mathbb{E}(\alpha_\infty h_\infty) - \alpha_0 h_0 \leq C - \mathbb{E} \left(\sum_{m=1}^{\infty} \alpha_{m+1} \gamma_m \|H(\mu_m)\|_{\mathbb{L}_2(\mu_m)}^2 \right).$$

In particular, we have

$$\sum_{k=1}^{\infty} \gamma_k \|H(\mu_k)\|_{\mathbb{L}^2(\mu_k)}^2 < +\infty \quad \text{a.s.}$$

Since we assume that $\sum_{k=1}^{\infty} \gamma_k = \infty$, we can follow the arguments in the proof of [34, Theorem 4.7] to conclude the proof, taking also advantage of Proposition 7 hereinafter, which in particular establishes the continuity of the function $\mu \mapsto \|H(\mu)\|_{\mathbb{L}^2(\mu)}^2$.

□

Proposition 7. *The function $\mu \in \mathcal{P}_2(\mathbb{R}^d) \mapsto \|H(\mu)\|_{\mathbb{L}^2(\mu)}^2$ is continuous w.r.t W_2 . Moreover, if $(\rho_m)_m$ in $\mathcal{P}_2(\mathbb{R}^d)$ are uniformly bounded w.r.t W_2 and converges to $\rho \in \mathcal{P}(\mathbb{R}^d)$ w.r.t W_q with $q \in [1, 2)$, we have $\rho \in \mathcal{P}_2(\mathbb{R}^d)$ and $\liminf_{m \rightarrow \infty} \|H(\rho_m)\|_{\mathbb{L}^2(\rho_m)}^2 \geq \|H(\rho)\|_{\mathbb{L}^2(\rho)}^2$.*

Proof. Let us first assume that $(\rho_m)_m, \rho$ in $\mathcal{P}_2(\mathbb{R}^d)$ are such that $W_2(\rho_m, \rho) \rightarrow 0$. We want to prove that

$$\|H(\rho_m)\|_{\mathbb{L}^2(\rho_m)}^2 \rightarrow \|H(\rho)\|_{\mathbb{L}^2(\rho)}^2 \quad (24)$$

when $m \rightarrow \infty$. Consider the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and r.v.'s $(X_m)_m$ and X constructed in Lemma 6.ii), and recall that, for each $\nu \in \mathcal{P}_2(\mathbb{R}^d)$, the r.v.'s $(X_m, S_{\rho_m}^\nu(X_m))$ have law $(\text{id}, S_{\rho_m}^\nu) \# \rho_m$ for each m and converge in $\mathbb{L}^2(\Omega, \mathcal{F}, \mathbb{P})$ to the r.v. $(X, S_\rho^\nu(X))$, which has the law $(\text{id}, S_\rho^\nu) \# \rho$. We next extend this construction in order to suitably randomise ν . More precisely, we enlarge the probability space $(\Omega, \mathcal{F}, \mathbb{P})$ to the product space $(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}}) = (\Omega \times \mathcal{P}_2(\mathbb{R}^d), \mathcal{F} \otimes \mathcal{B}(\mathcal{P}_2(\mathbb{R}^d)), \mathbb{P} \otimes \mathbb{Q})$, that is, we add an independent random variable, called ν , taking values in $\mathcal{P}_2(\mathbb{R}^d)$ and which has distribution \mathbb{Q} .

Thanks to the measurability of the mappings $(x, \nu) \mapsto S_{\rho_m}^\nu$ and $(x, \nu) \mapsto S_\rho^\nu$ proven in Lemma 3, by replacing ν by ν in the previous objects we obtain random vectors $(X_m, S_{\rho_m}^\nu(X_m))$ and $(X, S_\rho^\nu(X))$ defined in $(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}})$ which have, conditionally on $\{\nu = \nu\}$, the laws $(\text{id}, S_{\rho_m}^\nu) \# \rho_m$ and $(\text{id}, S_\rho^\nu) \# \rho$ respectively. Moreover, ν is independent of the r.v. X, X_1, \dots, X_m under $\bar{\mathbb{P}}$.

Now, by conditioning on $\{\nu = \nu\}$, using the convergence result in Lemma 6.ii) and the dominated convergence Theorem, we can easily check that $((X_m, S_{\rho_m}^\nu(X_m)))_m$ converges to $(X, S_\rho^\nu(X))$ in $\bar{\mathbb{P}}$ -probability. Furthermore, one can integrate w.r.t. \mathbb{Q} the bound (18) obtained for fixed ν in the proof of Lemma 5 and, denoting $\bar{\mathbb{E}}$ the expectation with respect to $\bar{\mathbb{P}}$, deduce that

$$\begin{aligned} \sup_m \bar{\mathbb{E}} \left(\|S_{\rho_m}^\nu(X_m)\|^2 \mathbf{1}_{\{\|S_{\rho_m}^\nu(X_m)\|^2 \geq M\}} \right) &\leq \int \int_{\{\|x\|^2 \geq K\}} \|x\|^2 d\nu(x) \mathbb{Q}(d\nu) \\ &\quad + \frac{K}{M} \int \int \|x\|^2 d\nu(x) \mathbb{Q}(d\nu), \end{aligned}$$

for each $M, K \geq 0$, where the r.h.s. is finite since $\mathbb{Q} \in \mathcal{P}_2(\mathcal{P}_2(\mathbb{R}^d))$. It follows that the sequence $((X_m, S_{\rho_m}^\nu(X_m)))_m$ has uniformly integrable second moments, and therefore converges also in $L^2(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}})$ to $(X, S_\rho^\nu(X))$, thanks to the Vitali convergence theorem.

We observe now that $\bar{\mathbb{E}}(S_{\rho_m}^\nu(X_m)|X_m) = \int S_{\rho_m}^\nu(X_m) d\mathbb{Q}(\nu)$ and $\bar{\mathbb{E}}(S_\rho^\nu(X)|X) = \int S_\rho^\nu(X) d\mathbb{Q}(\nu)$, $\bar{\mathbb{P}}$ - a.s., Moreover, if \mathcal{F}_∞ denotes the σ -algebra generated by (X_1, X_2, \dots) , one has $\bar{\mathbb{E}}(S_{\rho_m}^\nu(X_m)|\mathcal{F}_\infty) = \bar{\mathbb{E}}(S_{\rho_m}^\nu(X_m)|X_m)$ and $\bar{\mathbb{E}}(S_\rho^\nu(X)|\mathcal{F}_\infty) = \bar{\mathbb{E}}(S_\rho^\nu(X)|X)$. Using the continuity in $L^2(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}})$ of the conditional expectation with respect to \mathcal{F}_∞ , we deduce that

$$X_m - \bar{\mathbb{E}}(S_{\rho_m}^\nu(X_m)|X_m) \rightarrow X - \bar{\mathbb{E}}(S_\rho^\nu(X)|X) \quad (25)$$

in $L^2(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}})$. We conclude that $\bar{\mathbb{E}}\|X_m - \bar{\mathbb{E}}(S_{\rho_m}^\nu(X_m)|X_m)\|^2 \rightarrow \bar{\mathbb{E}}\|X - \bar{\mathbb{E}}(S_\rho^\nu(X)|X)\|^2$ as $m \rightarrow \infty$, which is exactly the required convergence (24).

Let us now assume that $(\rho_m)_m$ in $\mathcal{P}_2(\mathbb{R}^d)$ are uniformly bounded w.r.t W_2 and converge to $\rho \in \mathcal{P}(\mathbb{R}^d)$ w.r.t W_q with $q \in [1, 2)$. The previous arguments can be easily adapted to show that convergence (25) holds in $L^q(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}})$. Moreover, it can be easily checked that for every $M > 0$, the mapping

$Y \mapsto \bar{\mathbb{E}}(\|Y\|^2 \wedge M)$ is continuous in $L^q(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}})$. It follows that

$$\begin{aligned} \liminf_{m \rightarrow \infty} \bar{\mathbb{E}}(\|X_m - \bar{\mathbb{E}}(S_{\rho_m}^\nu(X_m)|X_m)\|^2) &\geq \liminf_{m \rightarrow \infty} \bar{\mathbb{E}}(\|X_m - \bar{\mathbb{E}}(S_{\rho_m}^\nu(X_m)|X_m)\|^2 \wedge M) \\ &= \bar{\mathbb{E}}(\|X - \bar{\mathbb{E}}S_\rho^\nu(X)|X\|^2 \wedge M) \end{aligned}$$

Letting $M \rightarrow \infty$ and using monotone convergence in the last term, the stated property follows. \square

E Numerical results

E.1 Proximal algorithm for the computation of the OWT plan

This section is dedicated to the resolution of the OWT problem. Let $\mu = \sum_{i=1}^r a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^m b_j \delta_{y_j}$, be two discrete measures, the OWT problem boils down to solving

$$\min_{\pi \in \mathbb{R}^{r \times m}} \underbrace{\sum_{i=1}^r a_i \|x_i - \left(\frac{\pi \mathbf{y}}{\mathbf{a}}\right)_i\|^2}_{f(\pi)} + \underbrace{1_{\Pi(\mu, \nu)}(\pi)}_{g(\pi)}, \quad (26)$$

where 1_C is the indicator function of the set C i.e.

$$1_C(\pi) = \begin{cases} \pi & \text{if } \pi \in C \\ \infty & \text{otherwise.} \end{cases}$$

The proximal algorithm to solve Eq. (26) then reads:

$$\pi^{\ell+1} = \text{prox}_{\theta_\ell g}(\pi^\ell - \theta^\ell \nabla f(\pi^\ell)). \quad (27)$$

As $\Pi(\mu, \nu)$ is a closed non-empty convex set, the proximal operator of g reduces to the Euclidean projection onto $\Pi(\mu, \nu)$:

$$\text{proj}_{\Pi(\mu, \nu)}(P) = \arg \min_{\pi \in \mathbb{R}^{r \times m}} \|P - \pi\|^2 = \arg \min_{\pi \in \mathbb{R}^{r \times m}} \langle \pi, -P \rangle + \frac{1}{2} \|\pi\|^2$$

where $\|\cdot\|$ is the Frobenius norm. This projection problem can be solved by Dykstra's algorithm with alternate Bregman projections [19] or by stochastic dual approaches of OT regularised by an \mathbb{L}_2 norm [35]. This method is summarised in Algorithm 3. In particular, we used an accelerated version of Eq. (27) via FISTA [11] (with $\omega_\ell \in [0, 1)$ an extrapolation parameter and θ_ℓ the usual stepsize chosen by a line search) in order to compute the optimal plan π_μ^ν in the weak transport problem.

The optimal barycentric projection is then given by $S_\mu^\nu = \frac{\pi_\mu^\nu \mathbf{y}}{\mathbf{a}}$. We initialised the algorithm with a random matrix whose elements sum to 1. Observe that, from Algorithm 1, the K optimal barycentric projection computations can be parallelised for each step n .

Algorithm 3: Computation of the optimal weak plan

Output: π_μ^ν ;

Input: $\mu = \sum_{i=1}^r a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^m b_j \delta_{y_j}$;

Initialise π_0 random matrix;

while not converge **do**

| $P_{\ell+1} := \pi_\ell + \omega_\ell(\pi_\ell - \pi_{\ell-1})$;

| $\pi_{\ell+1} := \text{proj}_{\Pi(\mu, \nu)}(P_{\ell+1} - \theta_\ell \nabla f(P_{\ell+1}))$;

end

With respect to the efficiency of this algorithm, Figure 7 shows a comparison of different settings for Eq. (27) in order to compute an optimal weak transport plan. For that purpose, we considered two discrete distributions μ and ν each constructed from $r = m = 10, 100$ and 250 samples of two dimensional Gaussian measures. We illustrate the convergence for both the standard and accelerated versions of the proximal algorithm, as well as for the projection into $\Pi(\mu, \nu)$ via Dykstra's algorithm or the stochastic dual approach. As expected, the accelerated version of Eq. (27) converges faster than the classical proximal algorithm, and the projection step is more stable with Dykstra's algorithm. Moreover, the smaller the number of support points, the faster the convergence. We have also noted that the random initialisation does not affect the convergence towards the minimiser of Eq. (26).

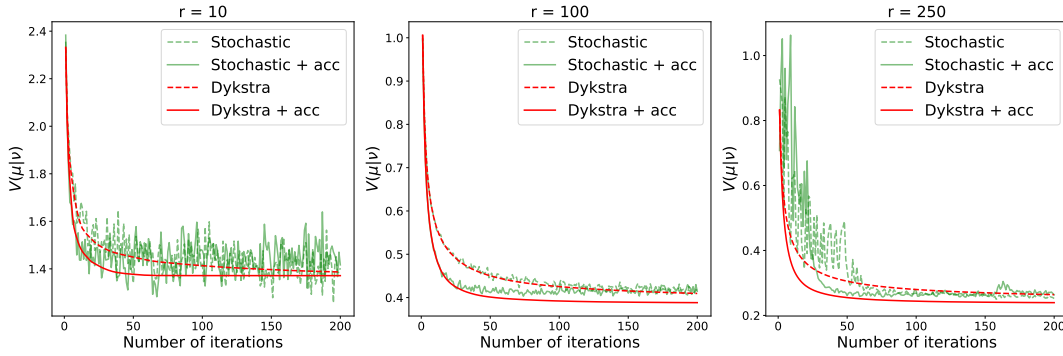


Figure 7: Convergence of the algorithm (3) in several settings for measures μ and ν supported on $r = m$ points.

E.2 Additional experiments

Gaussian distributions As in Section 5 of [5], we computed a weak barycenter between two 2D centered ellipses $E(\Sigma_i) = \{s \in \mathbb{R}^2 : s^t \Sigma_i^{-1} s = 1\}$ with covariances matrices

$$\Sigma_1 = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{and} \quad \Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix},$$

by considering 300 random observations for each ellipse. We then executed the iterations of Algorithm 1 until the difference of the objective function (i.e., the sum in Eq. (5)) between two successive iterations was smaller than $1e - 5$. This occurred at the 8th iteration, and the resulting weak barycenter was a circle within both ellipses. As we have access to the value of the weak barycenter problem (see Eq. (7)), we also compared the value of the objective function at the 8th iteration (that is $3.62e - 4$) to $\frac{1}{2} \sum_{i=1}^2 \|\mathbb{E}(Y_i)\|^2 - \|\frac{1}{2} \sum_{i=1}^2 \mathbb{E}(Y_i)\|^2$, with a plug-in estimator for $\mathbb{E}(Y_i)$. The approximated objective was equal to $3.21e - 4$, therefore, Algorithm 1 gave a satisfactory optimised weak barycenter.

Ellipse distributions ($r = 100$ & $K = 15$). We considered ellipse distributions with random center in $(-5, 5)$, random semi-major and semi-minor axes in (6, 14). The results are presented in Fig. 8, where the same conclusions as in the Gaussian examples hold.

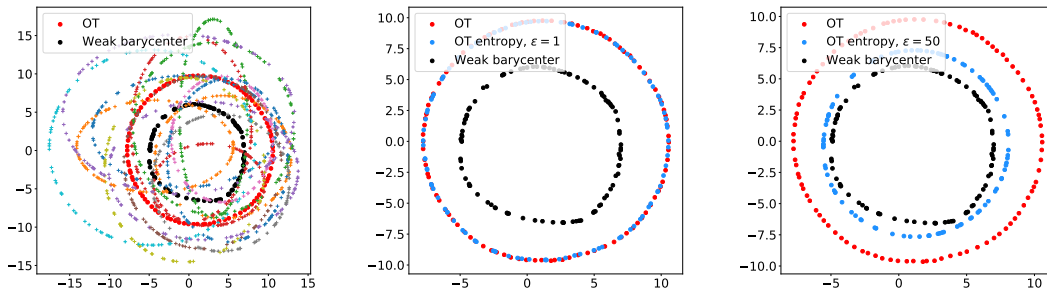


Figure 8: (left) Ellipse distributions and their OWT (black) and OT (red) barycenters computed with Algorithm 2. (center & right) Illustration of the weak (black), OT (red) and OT Sinkhorn (blue) barycenters for different values of $\epsilon = 1, 50$.

Pair-of-ellipses ($r \in (200, 300)$ & $K = 10$). In the same fashion, we considered distributions supported on two ellipses with random centers in $(-5, 5)$, random semi-major and semi-minor axes in (1, 7) and (7, 13) respectively. Fig. 9 shows the distributions (left) as well as the OT and OWT

barycenters (right) computed from random samples of the distributions. Observe that, once again, the weak barycenter better preserved the structure of the distributions when computing Algorithm 2.

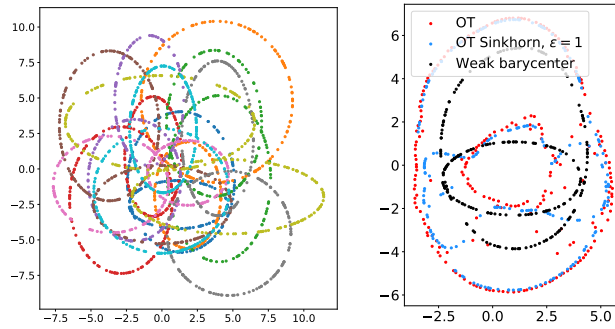


Figure 9: (left) Distributions supported on a pair-of-squares. (right) OWT (black), OT (red) and OT Sinkhorn for $\epsilon = 1$ (blue) barycenters computed with Algorithm 2.