



HAL
open science

Debiased Sinkhorn barycenters

Hicham Janati, Marco Cuturi, Alexandre Gramfort

► **To cite this version:**

Hicham Janati, Marco Cuturi, Alexandre Gramfort. Debiased Sinkhorn barycenters. ICML 2020 - 37th International Conference on Machine Learning, Jul 2020, Vienna / Virtuel, Austria. hal-03063875

HAL Id: hal-03063875

<https://hal.science/hal-03063875>

Submitted on 14 Dec 2020

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.

Debiased Sinkhorn barycenters

Hicham Janati^{1,2} Marco Cuturi^{3,2} Alexandre Gramfort²

Abstract

Entropy regularization in optimal transport (OT) has been the driver of many recent interests for Wasserstein metrics and barycenters in machine learning. It allows to keep the appealing geometrical properties of the unregularized Wasserstein distance while having a significantly lower complexity thanks to Sinkhorn’s algorithm. However, entropy brings some inherent *smoothing bias*, resulting for example in blurred barycenters. This side effect has prompted an increasing temptation in the community to settle for a slower algorithm such as log-domain stabilized Sinkhorn which breaks the parallel structure that can be leveraged on GPUs, or even go back to unregularized OT. Here we show how this bias is tightly linked to the reference measure that defines the entropy regularizer and propose debiased Wasserstein barycenters that preserve the best of both worlds: fast Sinkhorn-like iterations without entropy smoothing. Theoretically, we prove that the entropic OT barycenter of univariate Gaussians is a Gaussian and quantify its variance bias. This result is obtained by extending the differentiability and convexity of entropic OT to sub-Gaussian measures with unbounded supports. Empirically, we illustrate the reduced blurring and the computational advantage on various applications.

1. Introduction

Comparing, interpolating or averaging probability distributions is an ubiquitous problem in machine learning. Optimal transport (OT) offers an efficient way to do exactly that while taking into account the geometry of the space they live in (Peyré & Cuturi, 2018). Let $\mathcal{P}(\mathbb{R}^d)$ denote the set of probability measures on \mathbb{R}^d . Given some divergence $F : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}$ and weights $(w_k)_k$ such that

¹Inria Saclay, France ²CREST-ENSAE, France ³Google Research, Brain team, France. Correspondence to: Hicham Janati <hicham.janati@inria.fr>.

$\sum_{k=1}^K w_k = 1$, the weighted barycenter of a set of probability measures $(\alpha_k)_k$ can be defined as the Fréchet mean:

$$\alpha_F \stackrel{\text{def}}{=} \arg \min_{\alpha \in \mathcal{P}(\mathbb{R}^d)} \sum_{k=1}^K w_k F(\alpha_k, \alpha) . \quad (1)$$

Here α_F can be thought as a weighted average of distributions. While the $(\alpha_k)_k$ may have a fixed support or known finite supports when working in machine learning applications, the support of α_F may or may not be known. When the latter is unknown a priori, *free support methods* are needed to jointly minimize the objective with respect to both the support and the mass of the distribution (Cuturi & Doucet, 2014). Otherwise, *fixed support methods*, which only optimize weights on known supports, are employed (Benamou et al., 2014). While free support methods are more general and memory efficient, fixed support ones are faster in practice. In this paper, we focus on fixed support methods.

Using the Wasserstein distance as a divergence F , Li & Wang (2006) were the first to propose the Fréchet mean (1) for a clustering application in computer vision. This idea was later adopted by Agueh & Carlier (2011) to formally define Optimal Transport (OT) barycenters. However, the Wasserstein distance is defined through a linear programming problem which does not scale to large datasets. To address this computational issue, some form of regularization is mandatory: either regularize the measures themselves using sliced projections for instances or regularize the OT problem using ℓ_2 (Blondel et al., 2018) or entropy (Cuturi, 2013). While ℓ_2 preserves some of the sparsity of the non-regularized optimal transportation plan, entropy regularization leads to an approximation of the Wasserstein distance that can be solved using a fast and parallelizable GPU-friendly algorithm: the celebrated Sinkhorn’s algorithm (Cuturi, 2013). In the rest of this paper, we will focus on entropic OT. Let C be a non-negative cost function on $\mathbb{R}^d \times \mathbb{R}^d$ such that $C(x, y) = 0 \Leftrightarrow x = y$. For instance, a usual choice is $C(x, y) = \|x - y\|^2$. Entropy regularized OT between $\alpha, \beta \in \mathcal{P}(\mathbb{R}^d)$ with the reference measures $m_1, m_2 \in \mathcal{P}(\mathbb{R}^d)$ is defined as:

$$\text{OT}_\varepsilon^{m_1, m_2}(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\substack{\pi \in \mathcal{P}(\mathbb{R}^d \times \mathbb{R}^d) \\ \pi_{\#1} = \alpha, \pi_{\#2} = \beta}} \int_{\mathbb{R}^d \times \mathbb{R}^d} C d\pi + \varepsilon \text{KL}(\pi | m_1 \otimes m_2) , \quad (2)$$

where $\varepsilon > 0$, $\pi_{\#1}, \pi_{\#2}$ denote the left and right marginals of π respectively, $m_1 \otimes m_2$ is the product measure of m_1 and m_2 , and the relative entropy is defined as:

$$\text{KL}(\pi | m_1 \otimes m_2) \stackrel{\text{def}}{=} \int_{\mathbb{R}^d \times \mathbb{R}^d} \log \left(\frac{d\pi}{d(m_1 \otimes m_2)} \right) d\pi . \quad (3)$$

Naturally, in the discrete case, [Benamou et al. \(2014\)](#) proposed to compute OT barycenters of discrete measures using $F = \text{OT}_\varepsilon^{m_1, m_2}$ with $m_1 = m_2 = \mathcal{U}$, the uniform measure over the finite set on which the measures are defined. Doing so, they showed that the barycenter problem is equivalent to Iterative Bregman Projections (IBP) which are similar to Sinkhorn’s scaling operations. However, entropy regularization leads to an undesirable blurring of the barycenter. While using a very small regularization may appear as an obvious solution, it leads to numerical instabilities that can only be mitigated using log-domain stabilization or full log-domain ‘logsumexp’ operations ([Schmitzer, 2016](#)). This however considerably slows down Sinkhorn’s iterations.

To reduce this entropy bias, several divergences F have been proposed. For instance, [Solomon et al. \(2015\)](#) proposed to modify the IBP algorithm by adding a maximum entropy constraint they called *entropy sharpening*. This leads to a non-convex constraint which does not fit within the IBP framework. [Luise et al. \(2018\)](#) proposed to compute the entropy regularized solution π^* and to evaluate the OT loss (2) without the entropy term KL. This indeed leads to sharper barycenters but can only be estimated via gradient descent, thus requiring a full Sinkhorn loop at each iteration and setting a pre-defined learning rate which can be cumbersome in practice. [Amari et al. \(2019\)](#) proposed a modified entropy regularized divergence OT that can still leverage the fast IBP algorithm of [Benamou et al. \(2014\)](#) but requires a final deconvolution step with the kernel $\exp(-\frac{C}{\varepsilon})$, which is only feasible when ε is small. With this same objective of non-blurred solutions, [Ge et al. \(2019\)](#) even called for a return to the original non-regularized Wasserstein barycenter and proposed an accelerated interior point methods algorithm.

Our main contributions Except ([Ge et al., 2019](#)), all the works proposed above employ the uniform measure as reference, i.e they use $\text{OT}_\varepsilon^{\mathcal{U}} \stackrel{\text{def}}{=} \text{OT}_\varepsilon^{m_1, m_2}$ with $m_1 = m_2 = \mathcal{U}$. The purpose of this paper is to highlight a direct link between the already known entropy bias of the OT barycenter and this particular choice of m_1 and m_2 . This link is illustrated by showing how the choice of m_1 and m_2 impact the barycenter of univariate Gaussians in \mathbb{R}^d . Following ([Ramdas et al., 2017](#); [Genevay et al., 2018](#); [Feydy et al., 2018](#); [Luise et al., 2019](#)), we advocate for using the following Sinkhorn divergence which can be defined without specifying m_1 and m_2 for arbitrary measures $\alpha, \beta \in \mathcal{P}(\mathbb{R}^d)$:

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

The choice of the reference measures m_1 and m_2 has led to different formulations of regularized OT. The main contributions of this paper are twofold. (1) theoretical: we quantify the entropy bias of usual reference measures for univariate Gaussians. Precisely, while the Lebesgue measure ($m_1 = m_2 = \mathcal{L}$) induces a blurring bias and the product measure ($m_1 = \alpha, m_2 = \beta$) induces a shrinking bias, S_ε is actually debiased. (2) empirical: we propose a fast iterative algorithm similar to IBP to compute debiased barycenters. Unlike other gradient-based methods, this fixed point algorithm can be efficiently differentiated with respect to the barycentric weights via backpropagation. This allows one to carry out Wasserstein barycentric projections without entropy blurring. This will be illustrated in the experiments.

In the following section we discuss the different choices of m_1 and m_2 and quantify their induced entropy bias upon the barycenters of univariate Gaussians. In Section 3, we show some useful properties of S_ε (differentiability, convexity) when defined on sub-Gaussian measures with unbounded supports in \mathbb{R}^d which are necessary to prove the theorems of section 2. Next, in Section 4 we turn to computational aspects and provide a fast Sinkhorn-like algorithm for debiased barycenters. We conclude with numerical experiments in Section 5.

2. Reference measure and entropy bias

Notation We denote by $\mathbb{1}$ the vector of ones in \mathbb{R}^n . On matrices, log, exp and the division operator are applied element-wise. We use \odot for the element-wise multiplication between matrices or vectors. On vectors and matrices, the same notation denotes the usual scalar products: for $x, y \in \mathbb{R}^n$, $\langle x, y \rangle = \sum_{i=1}^n x_i y_i$; and for matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n, n}$, $\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i,j=1}^n \mathbf{A}_{ij} \mathbf{B}_{ij}$.

Uniform reference and IBP Let $\mathcal{X} = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ and consider two discrete measures $\alpha = \sum_{i=1}^n \alpha_i \delta_{x_i}$ and $\beta = \sum_{i=1}^n \beta_i \delta_{x_i}$. One can identify α and β with their weights α_i and β_i where $\alpha^\top \mathbb{1} = \beta^\top \mathbb{1}$. Let $\mathbf{C} \in \mathbb{R}_+^{n \times n}$ be the matrix such that $\mathbf{C}_{ij} = C(x_i, x_j)$. The definition of $\text{OT}_\varepsilon^{\mathcal{U}}$ in (2) becomes:

$$\text{OT}_\varepsilon^{\mathcal{U}}(\alpha, \beta) = \min_{\substack{\pi \in \mathbb{R}_+^{n \times n} \\ \pi \mathbb{1} = \alpha, \pi^\top \mathbb{1} = \beta}} \langle \mathbf{C}, \pi \rangle + \varepsilon \text{KL}(\pi | \mathbf{U}) , \quad (4)$$

where \mathbf{U} is the uniform measure on \mathcal{X}^2 given by $\frac{\mathbb{1}\mathbb{1}^\top}{n^2}$. Let \mathbf{K} be the element-wise exponentiated kernel $\exp(-\frac{C}{\varepsilon})$. By adopting the definition $\widetilde{\text{KL}}(\mathbf{A}, \mathbf{B}) = \sum_{i,j} \mathbf{A}_{ij} \log \left(\frac{\mathbf{A}_{ij}}{\mathbf{B}_{ij}} \right) + \mathbf{B}_{ij} - \mathbf{A}_{ij}$ for $\mathbf{A}, \mathbf{B} \in \mathbb{R}_+^{n \times n}$, [Benamou et al. \(2014\)](#) noticed that (4) is equivalent to a Kullback-Leibler projection up to

an additive constant:

$$\text{OT}_\varepsilon^{\mathcal{U}}(\alpha, \beta) = \min_{\substack{\pi \in \mathbb{R}_+^{n \times n} \\ \pi \mathbf{1} = \alpha, \pi^\top \mathbf{1} = \beta}} \varepsilon \widetilde{\text{KL}}(\pi | \mathbf{K}) \quad (5)$$

and proposed the Iterative Bregman Projections (IBP) algorithm to solve the equivalent barycenter problem:

$$\min_{\substack{\pi_1, \dots, \pi_K \\ \pi_k \in \mathcal{C}_k \cap \mathcal{C}'}} \sum_{k=1}^K w_k \widetilde{\text{KL}}(\pi^k | \mathbf{K}) \quad (6)$$

where $\mathcal{C}_k = \{\pi \in \mathbb{R}_+^{n \times n} | \pi \mathbf{1} = \alpha_k\}$ and $\mathcal{C}' = \{\pi \in \mathbb{R}_+^{n \times n} | \exists \alpha \in \Delta_n, \pi_k^\top \mathbf{1} = \alpha, \forall k = 1 \dots K\}$. The IBP algorithm amounts to performing iterative minimization on one constraint set at a time. Each step can be solved in closed form, leading to Sinkhorn-like iterations, see supplementary section D for details on IBP.

Lebesgue reference and smoothing bias As discussed in the introduction, the obtained barycenter $\alpha_{\text{OT}_\varepsilon^{\mathcal{U}}}$ suffers from entropy blurring. To quantify this blur, we turn to Lebesgue continuous measures and consider the Lebesgue measure as a reference by setting $m_1 = m_2 = \mathcal{L}$. We argue that by considering normalized histograms, the discrete formulation (5) provides an approximation of $\text{OT}_\varepsilon^{\mathcal{L}}$ when the number of histogram bins tends to $+\infty$. Indeed, since $\text{OT}_\varepsilon^{\mathcal{L}}$ is defined on Lebesgue-continuous measures, one can identify α, β and π with their density functions. Moreover, if the density functions are positive, the same KL factorization (5) is possible for $\text{OT}_\varepsilon^{\mathcal{L}}$. The following theorem shows that the weighted barycenter of univariate Gaussians is Gaussian with an increased variance. Figure 1 illustrates this smoothing bias using discrete histograms with a grid of 500 bins.

Theorem 1 (Blurring bias of $\text{OT}_\varepsilon^{\mathcal{L}}$). *Let $C(x, y) = (x - y)^2$, $\varepsilon > 0$ and $\varepsilon = 2\varepsilon'^2$. Let $(w_k)_k$ be positive weights that sum to 1. Let \mathcal{N} denote the Gaussian distribution. Assume $\alpha_k \sim \mathcal{N}(\mu_k, \sigma_k^2)$ and let $\bar{\mu} = \sum_k w_k \mu_k$,*

then:

(i) $\alpha_{\text{OT}_\varepsilon^{\mathcal{L}}} \sim \mathcal{N}(\bar{\mu}, S^2)$ where S is a positive solution of the equation: $\sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S^2} = -\varepsilon'^2 + 2S^2$.

(ii) In particular, if all σ_k are equal to some $\sigma > 0$,

then $\alpha_{\text{OT}_\varepsilon^{\mathcal{L}}} \sim \mathcal{N}(\bar{\mu}, \sigma^2 + \varepsilon'^2)$.

PROOF. See section C.3

The product measure and shrinking bias Besides the smoothing bias of the uniform measure, $\text{OT}_\varepsilon^{\mathcal{U}}$ cannot be generalized to a general OT definition for any arbitrary distributions that are non-discrete or non-Lebesgue continuous measures. To go beyond this binary classification

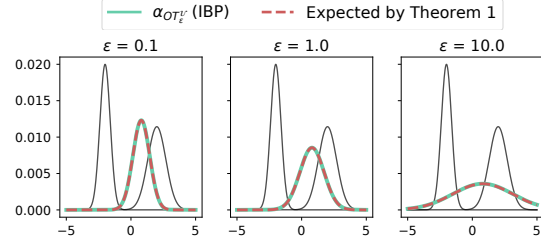


Figure 1. Illustration of theorem 1 with $\mathcal{N}(-2, 0.4)$ and $\mathcal{N}(2, 0.7)$ shown in black, and $(w_1, w_2) = (0.4, 0.6)$. The barycenter $\text{OT}_\varepsilon^{\mathcal{U}}$ matches theoretical expectations and is biased towards blurred distributions.

of probability measures, several authors (Ramdas et al., 2017; Genevay et al., 2018; Feydy et al., 2018) proposed the generic references $m_1 = \alpha, m_2 = \beta$. Indeed, the marginal constraints $\pi_1 = \alpha, \pi_2 = \beta$ imply that the support of π is included in that of $\alpha \otimes \beta$ and the KL term is always well-defined regardless of the nature of α and β . For the sake of convenience, we denote $\text{OT}_\varepsilon^\otimes \stackrel{\text{def}}{=} \text{OT}_\varepsilon^{\alpha, \beta}$. Di Marino & Gerolin (2019) made the following key observation that characterizes the change of reference. For discrete measures α, β :

$$\text{OT}_\varepsilon^{\mathcal{U}}(\alpha, \beta) = \text{OT}_\varepsilon^\otimes(\alpha, \beta) + \varepsilon \text{KL}(\alpha | \mathcal{U}) + \varepsilon \text{KL}(\beta | \mathcal{U}) \quad (7)$$

Similarly, the same identity holds for Lebesgue-continuous measures in $\mathcal{P}(\mathbb{R}^d)$:

$$\text{OT}_\varepsilon^{\mathcal{L}}(\alpha, \beta) = \text{OT}_\varepsilon^\otimes(\alpha, \beta) + \varepsilon \text{KL}(\alpha | \mathcal{L}) + \varepsilon \text{KL}(\beta | \mathcal{L}) \quad (8)$$

The identity (7) unveils another merit of $\text{OT}_\varepsilon^\otimes$ over $\text{OT}_\varepsilon^{\mathcal{U}}$: its corresponding barycenter problem is equivalent to a regularized $\text{OT}_\varepsilon^{\mathcal{U}}$ barycenter with a negative KL penalty. Interestingly, even though ‘ $-\text{KL}$ ’ is concave, $\text{OT}_\varepsilon^\otimes$ remains convex with respect to one of its arguments (Feydy et al., 2018). However, $\text{OT}_\varepsilon^\otimes$ yet suffers from some limitations: (1) $\text{OT}_\varepsilon^\otimes$ cannot be written as a KL projection, thus the fast IBP algorithm is lost; (2) the barycenter $\alpha_{\text{OT}_\varepsilon^\otimes}$ of Gaussians can be a degenerate Gaussian, as demonstrated by Theorem 2 which shows that if ε is large, the barycenter collapses to a Dirac (cf. Figure 3). This phenomenon can however be leveraged as a deconvolution technique: Rigollet & Weed (2018) showed that minimizing $\text{OT}_\varepsilon^\otimes$ is equivalent to maximum-likelihood deconvolution of an additive Gaussian-noise model.

Theorem 2 (Shrinking bias of $\text{OT}_\varepsilon^\otimes$). *Let $C(x, y) = (x - y)^2$, $\varepsilon > 0$ and $\varepsilon = 2\varepsilon'^2$. Let $(w_k)_k$ be positive weights that sum to 1. Let \mathcal{N} denote the Gaussian distribution. Assume that $\alpha_k \sim \mathcal{N}(\mu_k, \sigma_k^2)$ and let $\bar{\mu} = \sum_k w_k \mu_k, \bar{\sigma}^2 = \sum_{k=1} w_k \sigma_k^2$,*

(i) if $\varepsilon'^2 < \bar{\sigma}^2$ then $\alpha_{\text{OT}_\varepsilon^\otimes} \sim \mathcal{N}(\bar{\mu}, S^2)$ where S is a positive solution of the equation: $\sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S^2} = \varepsilon'^2 +$

$2S^2$. In particular, if all σ_k are equal to some $\sigma > 0$, then $\alpha_{\text{OT}_\varepsilon^\otimes} \sim \mathcal{N}(\bar{\mu}, \sigma^2 - \varepsilon'^2)$.

(ii) if $\varepsilon'^2 \geq \bar{\sigma}^2$ then $\alpha_{\text{OT}_\varepsilon^\otimes}$ is a Dirac located at $\bar{\mu}$.

PROOF. See section 3.

Debiased barycenters Interestingly, these limitations and significant differences between $\text{OT}_\varepsilon^{\mathcal{U}}$, $\text{OT}_\varepsilon^{\mathcal{L}}$ and $\text{OT}_\varepsilon^\otimes$ disappear when considering the following Sinkhorn divergences:

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

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

Using (7) and (8) it holds:

$$S_\varepsilon(\alpha, \beta) = S_\varepsilon^m(\alpha, \beta), \quad (9)$$

where m is either \mathcal{U} or \mathcal{L} depending on the nature of α and β . Therefore, S_ε is defined on arbitrary probability measures which can be mixtures of continuous measures and Dirac masses. Moreover, Feydy et al. (2018) showed that when the support of the measures is compact and with the additional assumption that C is negative semi-definite, S_ε is differentiable and convex with respect to one of its arguments. In the following section, we generalize the aforementioned statements for measures with unbounded supports in \mathbb{R}^d . The negativity assumption on C holds for instance if $C(x, y) = \|x - y\|^d$ with $0 < d \leq 2$ (Berg et al., 1984, Chapter 3, Cor 3.3) and is the only (cheap) price to pay for a debiased OT divergence. These convexity and differentiability results are essential to prove the debiasing of S_ε stated in Theorem 3 and illustrated in Figure 2.

Theorem 3 (Debiasing of S_ε). *Let $C(x, y) = (x - y)^2$ and $0 < \varepsilon < +\infty$ and $\varepsilon = 2\varepsilon'^2$. Let $(w_k)_k$ be positive weights that sum to 1. Let \mathcal{N} denote the Gaussian distribution. Assume that $\alpha_k \sim \mathcal{N}(\mu_k, \sigma_k^2)$ and let $\bar{\mu} = \sum_k w_k \mu_k$ then:*

(i) $\alpha_{S_\varepsilon} \sim \mathcal{N}(\bar{\mu}, S^2)$ where S is a positive solution S^* of the equation:

$\sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S^2} = \sqrt{\varepsilon'^4 + 4S^4}$. Moreover, given a sorted sequence $\sigma_{(1)} \leq \dots \leq \sigma_{(K)}$, it holds $S^* \in (\sigma_{(0)}, \sigma_{(K)})$.

(ii) In particular, if all σ_k are equal to some $\sigma > 0$, then $\alpha_{S_\varepsilon} \sim \mathcal{N}(\bar{\mu}, \sigma^2)$.

PROOF. See section 3.

Figure 3 shows a comparison of the three barycenters discussed in this section. We intentionally chose Gaussians with equal variances to emphasize two observations: (1) the debiasing of S_ε : the barycenter α_{S_ε} has the same variance

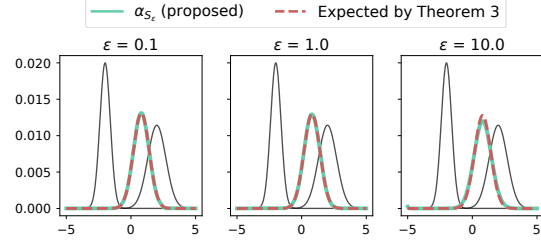


Figure 2. Illustration of theorem 3. Unlike with the uniform measure (Figure 1), the debiased barycenter remains unscathed when increasing ε .

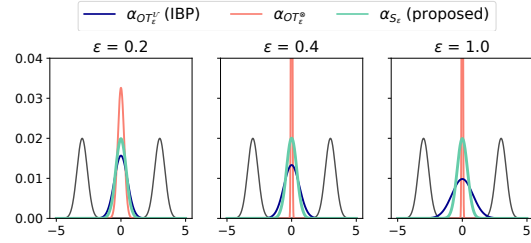


Figure 3. Illustration of the three theorems with $\mathcal{N}(-3, 0.4)$ and $\mathcal{N}(3, 0.4)$ shown in black using uniform weights. Entropy regularization causes a smoothing bias (blue) and a shrinking bias (red). Debiasing with S_ε (cyan) is perfect and independent of ε .

of the input measures for all ε ; (2) the shrinking bias of $\text{OT}_\varepsilon^\otimes$ is significant even for small values of ε .

Besides debiasing, the barycenter α_{S_ε} also comes with a computational advantage. Using the identity (9), we bypass the technical difficulties of the product measure in S_ε and derive an algorithm similar to IBP to compute α_{S_ε} which will be the subject of section 4.

3. S_ε is convex and differentiable on sub-Gaussian measures with unbounded supports

Notation The set of continuous function on \mathbb{R}^d is denoted by $\mathcal{C}(\mathbb{R}^d)$. The set of probability measures with a second order moment is denoted by $\mathcal{P}_2(\mathbb{R}^d)$. For $\alpha \in \mathcal{P}(\mathbb{R}^d)$, $\mathcal{L}_p(\mathbb{R}^d, \alpha)$ denotes the set of continuous functions $\mathbb{R}^d \rightarrow \mathbb{R}$ such that $\int |f|^p d\alpha < +\infty$. Let $f \in \mathcal{L}_1(\mathbb{R}^d, \alpha)$, $g \in \mathcal{L}_1(\mathbb{R}^d, \beta)$ and denote $\langle \alpha, f \rangle = \int_{\mathbb{R}^d} f d\alpha$. The tensor operators \otimes and \oplus denote respectively the mappings $f \otimes g : (x, y) \in \mathbb{R}^d \times \mathbb{R}^d \mapsto f(x).g(y)$ and $f \oplus g : (x, y) \in \mathbb{R}^d \times \mathbb{R}^d \mapsto f(x) + g(y)$.

To prove theorems 2 and 3, we characterize the optimality condition of the barycenter problem. First, we show that $\text{OT}_\varepsilon^\otimes$ and S_ε are convex (w.r.t. one variable) and differentiable. Our differentiability proof is inspired from that of Feydy et al. (2018) where the compactness assumption of the whole \mathcal{X} is replaced with a sub-Gaussian tails assumption.

tion on the measures that allows one to apply Lebesgue's dominated convergence theorem on \mathbb{R}^d . The convexity proof is however novel and is solely based on the dual problem of $\text{OT}_\varepsilon^\otimes$. Proving theorem 1 requires studying $\text{OT}_\varepsilon^\mathcal{L}$ which involves a slightly different dual problem. Since the differences are purely technical, we defer the proof of theorem 1 in the appendix and focus in this section on the product measure $\text{OT}_\varepsilon^\otimes$ and S_ε for the sake of clarity.

Dual problem In this section, we set $C(x, y) = \|x - y\|^2$ with its associated Gaussian kernel $K(x, y) = e^{-\frac{\|x-y\|^2}{\varepsilon}}$. Let $\alpha, \beta \in \mathcal{P}(\mathbb{R}^d)$. We define the linear operators on \mathcal{K} and \mathcal{K}^\top such that $\mathcal{K}(\mu) = \int_{\mathbb{R}^d} K(x, y) d\mu(y)$ and $\mathcal{K}^\top(\mu) = \int_{\mathbb{R}^d} K^\top(x, y) d\mu(x)$ for any non-negative measure $\mu \in \mathcal{M}_+(\mathbb{R}^d)$. Problem (2) has a dual formulation given by:

$$\begin{aligned} \text{OT}_\varepsilon^\otimes(\alpha, \beta) = & \sup_{\substack{f \in \mathcal{L}_1(\mathbb{R}^d, \alpha) \\ g \in \mathcal{L}_1(\mathbb{R}^d, \beta)}} \int_{\mathbb{R}^d} f d\alpha + \int_{\mathbb{R}^d} g d\beta \\ & - \varepsilon \int_{\mathbb{R}^d \times \mathbb{R}^d} \exp\left(\frac{f \oplus g - C}{\varepsilon}\right) d\alpha d\beta + \varepsilon . \end{aligned} \quad (10)$$

if α and β have finite second moments, (10) is well defined and a couple of dual potentials (f, g) are optimal if and only if they are solutions of Sinkhorn's equations (Mena & Weed, 2019):

$$\begin{aligned} e^{\frac{f}{\varepsilon}} \cdot \mathcal{K}(e^{\frac{g}{\varepsilon}} \cdot \beta) &= 1, \quad \alpha - a.e. , \\ e^{\frac{g}{\varepsilon}} \cdot \mathcal{K}^\top(e^{\frac{f}{\varepsilon}} \cdot \alpha) &= 1, \quad \beta - a.e. . \end{aligned} \quad (11)$$

and the optimal transport plan π is given by: $\pi = \exp\left(\frac{f \oplus g - C}{\varepsilon}\right) \cdot (\alpha \otimes \beta)$

Thus, at optimality the integral over $\mathbb{R}^d \times \mathbb{R}^d$ sums to 1 and:

$$\text{OT}_\varepsilon^\otimes(\alpha, \beta) = \int_{\mathbb{R}^d} f d\alpha + \int_{\mathbb{R}^d} g d\beta \quad (12)$$

Symmetric terms $\text{OT}_\varepsilon^\otimes(\alpha, \alpha)$ When $\alpha = \beta$, the symmetry of the problem leads to the existence of a symmetric pair of potentials (h, h) . Indeed, if (f, g) is optimal (g, f) is also optimal. Moreover, since C is symmetric, the optimal transport plan π is also symmetric which leads to $f = g$. Thus the following proposition holds.

Proposition 1. *Let $\alpha \in \mathcal{P}_2(\mathbb{R}^d)$, it holds:*

$$\begin{aligned} \text{OT}_\varepsilon^\otimes(\alpha, \alpha) = & \sup_{h \in \mathcal{L}_1(\mathbb{R}^d, \alpha)} 2 \int_{\mathbb{R}^d} h d\alpha \\ & - \varepsilon \int_{\mathbb{R}^d \times \mathbb{R}^d} \exp\left(\frac{h \oplus h - C}{\varepsilon}\right) d^2\alpha + \varepsilon , \end{aligned} \quad (13)$$

Moreover, the supremum is attained at the unique (by strong concavity; since C is definite negative) autocorrelation potential $h \in \mathcal{L}_1(\mathbb{R}^d, \alpha)$ if and only if h is a solution of $e^{\frac{fh}{\varepsilon}} \cdot \mathcal{K}(e^{\frac{h}{\varepsilon}} \cdot \alpha) = 1$, $\alpha - a.e.$, and at optimality it holds: $\frac{1}{2} \text{OT}_\varepsilon^\otimes(\alpha, \alpha) = \int_{\mathbb{R}^d} h d\alpha$.

Restriction on sub-Gaussians To derive theorems 2 and 3, we show that both $\text{OT}_\varepsilon^\otimes$ and S_ε are convex and differentiable and provide a solution of the first order optimality condition. Notice that the convexity of $\text{OT}_\varepsilon^\otimes$ with respect to α and with respect to β follows immediately from (10) since it corresponds to a supremum of linear functionals. Feydy et al. (2018) showed the differentiability of $\text{OT}_\varepsilon^\otimes$ and the convexity of S_ε on measures with compact supports. On \mathbb{R}^d , more assumptions on α and β are required. Throughout this section we restrict $\text{OT}_\varepsilon^\otimes$ and S_ε to the convex set of sub-Gaussian probability measures:

Assumption 1. *We set $C(x, y) = \|x - y\|^2$ and restrict $\text{OT}_\varepsilon^\otimes$ and S_ε to the set of sub-Gaussian probability measures $\mathcal{G}(\mathbb{R}^d) \stackrel{\text{def}}{=} \{\mu | \exists q > 0, \mathbb{E}_\mu(e^{\frac{\|x\|^2}{2dq^2}}) \leq 2\}$.*

Mena & Weed (2019) showed that if $\alpha, \beta \in \mathcal{G}(\mathbb{R}^d)$, there exists a pair of potentials (f, g) verifying the fixed point equations (11) on the whole space \mathbb{R}^d that are bounded by quadratic functions. This result is key to show the differentiability of $\text{OT}_\varepsilon^\otimes$ on $\mathcal{G}(\mathbb{R}^d)$.

Proposition 2 (Mena & Weed (2019), Prop. 6). *Let $\alpha, \beta \in \mathcal{G}(\mathbb{R}^d)$. There exists a pair of smooth functions (f, g) such that (11) holds on \mathbb{R}^d and $\forall x, y \in \mathbb{R}^d$:*

$$\begin{aligned} -dq^2(1 + \frac{1}{2}(\|x\| + \sqrt{2dq})^2) &\leq \frac{f(x)}{\varepsilon} \leq \frac{1}{2}(\|x\| + \sqrt{2dq})^2 \\ -dq^2(1 + \frac{1}{2}(\|y\| + \sqrt{2dq})^2) &\leq \frac{g(y)}{\varepsilon} \leq \frac{1}{2}(\|y\| + \sqrt{2dq})^2 \end{aligned} \quad (14)$$

Differentiability In the rest of this section, (f, g) denotes a pair of potentials defined by Proposition 2. We say that a function $F : \mathcal{G}(\mathbb{R}^d) \rightarrow \mathbb{R}$ is differentiable at α if there exists $\nabla F(\alpha) \in \mathcal{C}(\mathbb{R}^d)$ such that for any displacement $t\delta\alpha$ with $t > 0$ and $\delta\alpha = \alpha_1 - \alpha_2$ with $\alpha_1, \alpha_2 \in \mathcal{G}(\mathbb{R}^d)$, and:

$$F(\alpha + t\delta\alpha) = F(\alpha) + t\langle \delta\alpha, \nabla F(\alpha) \rangle + o(t) , \quad (15)$$

where $\langle \delta\alpha, \nabla F(\alpha) \rangle = \int_{\mathbb{R}^d} \nabla F(\alpha) d\delta\alpha$.

Proposition 3. *Let $\alpha, \beta \in \mathcal{G}(\mathbb{R}^d)$, and (f, g) their associated pair of dual potentials given by proposition 2. $\text{OT}_\varepsilon^\otimes(\alpha, \cdot)$ is differentiable on sub-Gaussian measures with unbounded supports and its gradient is given by:*

$$\nabla_\beta \text{OT}_\varepsilon^\otimes(\alpha, \beta) = g . \quad (16)$$

SKETCH OF PROOF. The proof is inspired from Feydy et al. (2018) in the case of measures with compact supports. The

difference arises when taking the limit of integrals of the potentials. Thanks to assumption 1, proposition 2 provides an upper bound that allows to conclude by dominated convergence. The full proof is provided in the appendix.

The differentiability of S_ε follows immediately:

Corollary 1. *Let $\alpha, \beta \in \mathcal{G}(\mathbb{R}^d)$, and (f, g) their associated pair of dual potentials given by proposition 2 and h_β the autocorrelation potential associated with β . $S_\varepsilon^\otimes(\alpha, \cdot)$ is differentiable on sub-Gaussian measures with unbounded supports and its gradient is given by:*

$$\nabla_\beta S_\varepsilon^\otimes(\alpha, \beta) = g - h_\beta . \quad (17)$$

Remark 1. It is important to keep in mind that the notion of differentiability (and gradient) of the functions $\text{OT}_\varepsilon^\otimes$ and S_ε differ from the usual Fréchet differentiability. Indeed, the space of probability measures $\mathcal{P}(\mathbb{R}^d)$ has an empty interior in the space of signed Radon measures $\mathcal{M}(\mathbb{R}^d)$. The definition adopted here defines derivatives along feasible directions in $\mathcal{P}(\mathbb{R}^d)$. This is however sufficient to characterize the convexity of S_ε and its stationary points (see appendix A for details).

Convexity Now we turn to showing that S_ε is convex with respect to either one of its arguments separately. To do so, we prove the first order characterization of convexity of a differentiable function $F : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ given by:

$$F(\alpha) \geq F(\alpha') + \langle \alpha - \alpha', \nabla F(\alpha') \rangle , \quad (18)$$

As shown by the proof of the following Lemma, the positivity of K plays a key role in proving the convexity of S_ε .

Lemma 1. *Let $\alpha, \alpha' \in \mathcal{G}(\mathbb{R}^d)$ and let $h_\alpha, h_{\alpha'}$ denote their respective autocorrelation potentials given by proposition 1. Then if $K(x, y) = e^{-\frac{\|x-y\|^2}{\varepsilon}}$:*

$$\int e^{\frac{h_\alpha(x)}{\varepsilon}} K(x, y) e^{\frac{h_{\alpha'}(y)}{\varepsilon}} d\alpha(x) d\alpha'(y) \leq 1 \quad (19)$$

Proposition 4. *Under assumption (1), S_ε is convex on sub-Gaussian measures with respect to either of its arguments.*

PROOF. Let $\beta \in \mathcal{G}(\mathbb{R}^d)$. Let $\alpha, \alpha' \in \mathcal{G}(\mathbb{R}^d)$. Let (f, g) and (f', g') denote the pair of potentials associated with $\text{OT}_\varepsilon^\otimes(\alpha, \beta)$ and $\text{OT}_\varepsilon^\otimes(\alpha', \beta)$ respectively and for any $\mu \in \mathcal{G}(\mathbb{R}^d)$, let h_μ denote the autocorrelation potential associated with $\text{OT}_\varepsilon^\otimes(\mu, \mu)$. The first order inequality (18) applied to $F = S_\varepsilon(\cdot, \beta)$ is equivalent to:

$$\begin{aligned} (18) &\Leftrightarrow \langle \alpha, f - h_\alpha \rangle + \langle \beta, g - h_\beta \rangle \geq \\ &\langle \alpha', f' - h_{\alpha'} \rangle + \langle \beta, g' - h_\beta \rangle + \langle \alpha - \alpha', f' - h_{\alpha'} \rangle \\ &\Leftrightarrow \langle \alpha, f - h_\alpha \rangle + \langle \beta, g \rangle \geq \langle \beta, g' \rangle + \langle \alpha, f' - h_{\alpha'} \rangle \quad (20) \\ &\Leftrightarrow \langle \alpha, f \rangle + \langle \beta, g \rangle \geq \langle \beta, g' \rangle + \langle \alpha, f' - h_{\alpha'} + h_\alpha \rangle \\ &\Leftrightarrow \text{OT}_\varepsilon^\otimes(\alpha, \beta) \geq \langle \alpha, f' - h_{\alpha'} + h_\alpha \rangle + \langle \beta, g' \rangle \end{aligned}$$

To show the last inequality we use the definition of the dual problem (10) and evaluate the dual function at the suboptimal potentials $(f' - h_{\alpha'} + h_\alpha, g')$. Doing so leads to:

$$\begin{aligned} \text{OT}_\varepsilon^\otimes(\alpha, \beta) &\geq \langle \alpha, f' - h_{\alpha'} + h_\alpha \rangle + \langle \beta, g' \rangle + \varepsilon \\ &- \varepsilon \int_{\mathbb{R}^d \times \mathbb{R}^d} \exp\left(\frac{(f' - h_{\alpha'} + h_\alpha) \oplus g' - C}{\varepsilon}\right) d\alpha d\beta . \end{aligned}$$

To conclude, all we need to show is that,

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \exp\left(\frac{(f' - h_{\alpha'} + h_\alpha) \oplus g' - C}{\varepsilon}\right) d\alpha d\beta \leq 1 \quad (21)$$

By the Fubini-Tonelli theorem, the order of integration is irrelevant. First integrating with respect to β , we use the optimality conditions (11) on the pair (f', g') then on $h_{\alpha'}$:

$$\begin{aligned} B &= \int_{\mathbb{R}^d \times \mathbb{R}^d} \exp\left(\frac{(f' - h_{\alpha'} + h_\alpha) \oplus g' - C}{\varepsilon}\right) d\alpha d\beta \\ &= \int_{\mathbb{R}^d} \exp\left(\frac{h_\alpha - h_{\alpha'}}{\varepsilon}\right) d\alpha \\ &= \int_{\mathbb{R}^d \times \mathbb{R}^d} \exp\left(\frac{h_\alpha \oplus h_{\alpha'} - C}{\varepsilon}\right) d\alpha d\alpha' \end{aligned}$$

Thus, Lemma 1 applies and we have $B \leq 1$. \square

Barycenter of sub-Gaussian distributions. We have shown that $\text{OT}_\varepsilon^\otimes$ and S_ε are convex and differentiable, thus the weighted barycenters $\alpha_{\text{OT}_\varepsilon^\otimes}$ and α_{S_ε} can be characterized by the first order optimality condition as follows. Let (f_k, g_k) denote the potentials associated with $\text{OT}_\varepsilon^\otimes(\alpha_k, \alpha)$ and h_α the autocorrelation potential associated with $\text{OT}_\varepsilon^\otimes(\alpha, \alpha)$. Using the first order characterization of convexity (18), α^* is a global minimizer of the barycenter loss of $\text{OT}_\varepsilon^\otimes$ if and only if for any direction $\beta \in \mathcal{G}(\mathbb{R}^d)$, $\langle \sum_{k=1}^K w_k \nabla_{\alpha^*} \text{OT}_\varepsilon^\otimes(\alpha_k, \alpha^*), \beta - \alpha^* \rangle \geq 0$. This is equivalent to $\sum_{k=1}^K w_k \langle g_k, \beta - \alpha^* \rangle \geq 0$. Similarly, for α_{S_ε} we get the optimality condition $\sum_{k=1}^K w_k \langle g_k - h_{\alpha^*}, \beta - \alpha^* \rangle \geq 0$. We are now ready to summarize the different steps of the proofs of the theorems. For S_ε , we provide solutions of the optimality conditions by considering quadratic potentials and Gaussian barycenters α_{S_ε} . We proceed by identification of the coefficients of the polynomials and the parameters of the barycenters and show that the obtained solutions verify the optimality condition. For $\text{OT}_\varepsilon^\otimes$, we proceed similarly for $2\varepsilon'^2 < \sigma^2$. For $2\varepsilon'^2 \geq 2\sigma^2$, we show directly that for the Dirac measure $\alpha^* = \delta_{\bar{\mu}}$, there exists a set of potentials that verify the optimality condition alongside Sinkhorn's equations. The detailed derivations are provided in the supplementary materials.

4. Fast Sinkhorn-like algorithm

Discrete measures on a finite space The purpose of this section is to derive a fast Sinkhorn-like algorithm to compute α_{S_ε} on a fixed support. Let $\mathcal{X} = \{x_1, \dots, x_n\}$ be a finite grid of size n . With images for instance, each x_i would correspond to a pixel. We identify a probability measure $\alpha = \sum_{i=1}^n \alpha_i \delta_{x_i} \in \mathcal{P}(\mathcal{X})$ with its weights vector $(\alpha_i) \in \mathbb{R}_{++}^n$ such that $\sum_{i=1}^n \alpha_i = 1$. In the rest of this paper, OT_ε and S_ε can be seen as functions operating on the interior of the probability simplex of \mathbb{R}^n denoted by $\Delta_n = \{x \in \mathbb{R}_{++}^n \mid \sum_{i=1}^n x_i = 1\}$. We assume that the cost matrix $\mathbf{C} \in \mathbb{R}_{++}^{n \times n}$ is symmetric negative semi-definite (or equivalently, its associated kernel $\mathbf{K} = e^{-\frac{\mathbf{C}}{\varepsilon}}$ is positive semi-definite). This assumption holds for instance if $\mathbf{C}_{ij} = \|x_i - x_j\|^p$ with $p \in]0, 2]$ (see (Berg et al., 1984, 3, Thm 2.2, Cor 3.3) for both claims)

Debiased barycenters To obtain a fast iterative algorithm for the debiased barycenters α_{S_ε} , we are going to leverage the IBP algorithm through the uniform measure on \mathcal{X} as follows. First, the identity (9) ensures that S_ε is independent of the reference measures. Thus, one can write:

$$S_\varepsilon(\alpha, \beta) = \text{OT}_\varepsilon^{\mathcal{U}}(\alpha, \beta) - \frac{\text{OT}_\varepsilon^{\mathcal{U}}(\alpha, \alpha) + \text{OT}_\varepsilon^{\mathcal{U}}(\beta, \beta)}{2}.$$

Using (5), one can write $\text{OT}_\varepsilon^{\mathcal{U}}(\alpha, \beta)$ as a KL projection. The remaining autocorrelation terms can be replaced by their dual problems to obtain the following proposition. A detailed derivation is provided in appendix E.

Proposition 5. *Let $\alpha_1, \dots, \alpha_K \in \Delta_n$ and $\mathbf{K} = e^{-\frac{\mathbf{C}}{\varepsilon}}$. Let π denote a sequence π_1, \dots, π_K of transport plans in $\mathbb{R}_+^{n \times n}$ and the constraint sets $\mathcal{H}_1 = \{\pi \mid \forall k, \pi_k \mathbf{1} = \alpha_k\}$, and $\mathcal{H}_2 = \{\pi \mid \forall k \forall k', \pi_k^\top \mathbf{1} = \pi_{k'} \mathbf{1}\}$. The barycenter problem $\min_{\alpha \in \Delta_n} \sum_{k=1}^K w_k S_\varepsilon(\alpha_k, \alpha)$ is equivalent to:*

$$\min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2 \\ d \in \mathbb{R}_+^n}} \left[\varepsilon \sum_{k=1}^K w_k \widetilde{\text{KL}}(\pi_k \mid \mathbf{K} \text{diag}(d)) + \frac{\varepsilon}{2} \langle d - \mathbf{1}, \mathbf{K}(d - \mathbf{1}) \rangle \right]. \quad (22)$$

where $\widetilde{\text{KL}}(\mathbf{A}, \mathbf{B}) = \sum_{i,j} \mathbf{A}_{ij} \log \left(\frac{\mathbf{A}_{ij}}{\mathbf{B}_{ij}} \right) + \mathbf{B}_{ij} - \mathbf{A}_{ij}$.

Since $\widetilde{\text{KL}}$ is jointly convex and \mathbf{K} is assumed positive-definite, the objective (22) is convex. Minimizing (22) with respect to π leads to the barycenter problem $\alpha_{\text{OT}_\varepsilon^{\mathcal{U}}}$ (6) with the modified kernel $\mathbf{K} \text{diag}(d)$. This problem can be solved via the fast IBP algorithm. Minimizing with respect to d leads to the Sinkhorn fixed point equation $d = \frac{\sum w_k \pi_k^\top \mathbf{1}}{\mathbf{K}d}$ for which there exists a converging sequence $d_{n+1} \leftarrow \sqrt{\frac{d_n \odot \sum w_k \pi_k^\top \mathbf{1}}{\mathbf{K}d_n}} (\star)$ (Knight et al., 2014). Given

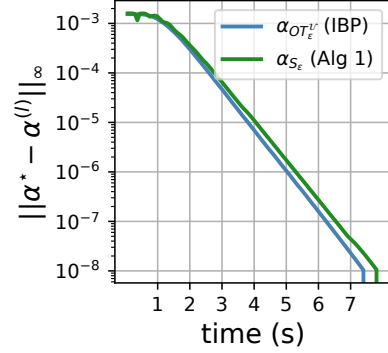


Figure 4. Convergence to the true barycenters of univariate Gaussians $\mathcal{N}(-0.5, 0.1)$ and $\mathcal{N}(0.5, 0.1)$. Algorithm 1 is as fast as IBP with a linear convergence rate.

Algorithm 1 Debiased Sinkhorn Barycenter

Input: $\alpha_1, \dots, \alpha_K, \mathbf{K} = e^{-\frac{\mathbf{C}}{\varepsilon}}$
Output: α_{S_ε}
 Initialize all scalings $(b_k), d$ to $\mathbf{1}$,
repeat
 for $k = 1$ **to** K **do**
 $a_k \leftarrow \left(\frac{\alpha_k}{\mathbf{K}b_k} \right)$
 end for
 $\alpha \leftarrow d \odot \prod_{k=1}^K (\mathbf{K}^\top a_k)^{w_k}$
 for $k = 1$ **to** K **do**
 $b_k \leftarrow \left(\frac{\alpha}{\mathbf{K}^\top a_k} \right)$
 end for
 $d \leftarrow \sqrt{d \odot \left(\frac{\alpha}{\mathbf{K}d} \right)}$
until convergence

that (22) is smooth and convex, alternate minimization – which amounts to perform IBP and (\star) iterations – converges towards its minimum. However, we notice that in practice, either taking one iteration or fully optimizing the subproblems produces the same minimizer. We thus propose to combine one IBP iteration with the update (\star) , which leads to Algorithm 1 (see the appendix for further details on the IBP algorithm). Using the theoretical barycenters of Gaussians given by theorems 1 and 3, we can monitor the convergence to the ground truth (Figure 4). Theoretically, both IBP and algorithm 1 have a $\mathcal{O}(Kn^2)$ complexity per iteration. A convergence proof of IBP can be obtained using alternating Bregman projections (See (Benamou et al., 2014) and the references therein). For Algorithm 1 however, similar techniques were not successful. Proving its convergence will be pursued in future work.

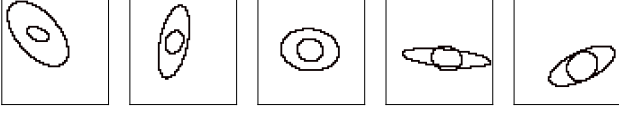


Figure 5. 5 examples of random nested ellipses of size (60×60) used to compute the barycenters of Figure 6.

5. Applications

Now we turn to showing the practical benefits of debiased barycenters in terms of accuracy, speed and performance.

Benchmarks In addition to $\alpha_{OT_\varepsilon^U}$, $\alpha_{OT_\varepsilon^\otimes}$, we evaluate the performance of the following barycenters:

- α_{A_ε} : Sharp barycenters introduced by Luise et al. (2018), where A_ε is defined as: $A_\varepsilon(\alpha, \beta) = \langle \mathbf{C}, \pi_\varepsilon^*(\alpha, \beta) \rangle$. Here $\pi_\varepsilon^*(\alpha, \beta)$ is the primal minimizer of the regularized problem $OT_\varepsilon^U(\alpha, \beta)$, computed via accelerated gradient descent.
- $\alpha_{S_\varepsilon^F}$: Free support barycenters introduced by Luise et al. (2019) that uses the same debiased divergence S_ε , and deals with the free support problem by adding / removing a Dirac particle with Frank-Wolf’s algorithm.
- α_W : The original non-regularized Wasserstein problem solved with interior point methods - using the accelerated MAAIPM algorithm of Ge et al. (2019).

Debiased barycenters of ellipses To demonstrate how debiased barycenters α_{S_ε} reduce smoothing and are computationally competitive with $\alpha_{OT_\varepsilon^U}$, we compare the barycenters of 10 randomly generated nested ellipses displayed in Figure 5. We set the cost matrix \mathbf{C} to the squared Euclidean distance on the unit square and set $\varepsilon = 0.002$. We use the same termination criterion for all methods based on a maximum relative change of the barycenters set to 10^{-5} .

For α_{S_ε} , $\alpha_{OT_\varepsilon^U}$, $\alpha_{OT_\varepsilon^\otimes}$, α_{A_ε} , we use the convolution trick introduced by Solomon et al. (2015) which amounts to computing the kernel operation $\mathbf{K}a$ on a vectorized image a by applying a Gaussian convolution on the rows and the columns of a , thereby reducing the complexity of one Debiased / IBP iteration from $O(n^2)$ to $O(n^{\frac{3}{2}})$.

Figure 5 shows that even though α_{A_ε} and α_W are not blurred compared to $\alpha_{OT_\varepsilon^U}$, they cannot compete computationally with Sinkhorn-like algorithms. The debiased barycenter is sharp and runs in about the same time as $\alpha_{OT_\varepsilon^U}$. Besides, the shrinking bias of OT_ε^\otimes unfolded by theorem 2 is illustrated in the degeneracy of the ellipse $\alpha_{OT_\varepsilon^\otimes}$.

Barycenters of 3D shapes To visually illustrate the impact of the reduced smoothing bias of S_ε , we computed a barycentric interpolation of shapes discretized in a 3D grid of $200 \times 200 \times 200$ voxels. The different inter-

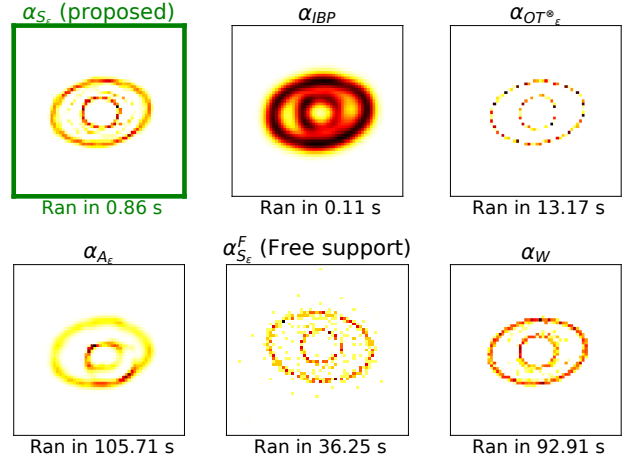


Figure 6. Barycenters of the 10 nested ellipses shown in Figure 5. Results illustrate the reduced blurring of the proposed approach and running times presented below each image demonstrate the computational efficiency. All 6 barycenters were computed on a laptop with an Intel Core i5 3.1 GHz Processor.



Figure 7. Interpolation of two 3D shapes on a $(200)^3$ uniform grid with IBP illustrating a clear blurring bias of OT_ε^U .



Figure 8. Interpolation of two 3D shapes on a $(200)^3$ uniform grid with the proposed Debiased Sinkhorn (Alg 1). The interpolation is sharper and completes in about the same time as figure 7 (5 seconds on a GPU).

polarizations correspond to weights $(w, 1 - w)$ where $w \in [0, 0.25, 0.5, 0.75, 1]$. We set the cost matrix \mathbf{C} to the squared Euclidean distance on the unit cube and set $\varepsilon = 0.01$. Results presented in Figures 7 and 8 using OT_ε^U and S_ε qualitatively demonstrate that S_ε leads to sharper edges, while in both cases it takes a few seconds to compute on a GPU. Again, the kernel operation $\mathbf{K}a$ on a vectorized 3D grid a can be computed via a sequence of 3 Gaussian convolutions on each axis (x, y, z) which reduces the complexity of one Debiased / IBP iteration from $O(n^2)$ to $O(n^{\frac{4}{3}})$.

Optimal transport barycentric embeddings One of the many machine learning applications of OT barycenters is to compute low-dimensional barycentric embeddings. Introduced by Bonneel et al. (2016), OT barycentric co-

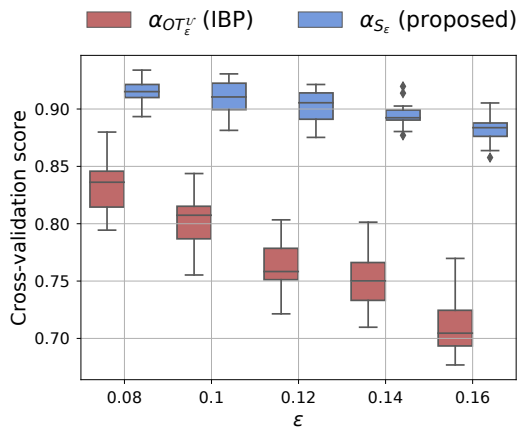


Figure 9. Cross-validation accuracy with 95% confidence intervals obtained on 500 MNIST images using barycentric embedding with S_ϵ or OT_ϵ^U . Debiasing of S_ϵ improves performance. S_ϵ is less sensitive to ϵ .

ordinates are defined as follows. Given a dictionary \mathcal{A} of distributions $\alpha_1, \dots, \alpha_K$ and $w \in \Delta_K$, let $\alpha_F(w) = \arg \min_{\alpha} \sum_{k=1}^K w_k F(\alpha_k, \alpha)$ for some OT divergence F . The OT coordinates \hat{w} of a distribution β are defined as the weights of the barycenter $\alpha_F(w)$ best approximating β for a given divergence. Using a quadratic divergence, it reads: $\hat{w} = \arg \min_{w \in \Delta_K} \|\alpha_F(w) - \beta\|^2$. To leverage the differentiability of the IBP iterations, Bonneel et al. (2016) used the divergence OT_ϵ^U and proposed to substitute the minimizer $\alpha_F(w)$ with the l -th IBP iterate $\alpha_F^{(l)}(w)$. Differentiating the barycenter nets $\alpha_F^{(l)}(w)$ with respect to w can be done via automatic differentiation, while the full minimization can be done using accelerated gradient descent using a soft-max reparametrization. Here we use the ADAM optimizer of the pyTorch library (Paszke et al., 2017). To evaluate the benefits of debiasing, we take 500 samples of the MNIST dataset (LeCun & Cortes, 2010) with 100 instances of each digit (0-1-2-3-4). We select 10% of the dataset (a subset of 50 images; ergo $K=50$) at random as our learning dictionary \mathcal{A} and compute the barycentric coordinates of the remaining 90% subset denoted as \mathcal{D} . Thus, for each image among the 450 samples of \mathcal{D} , we compute the closest (in squared ℓ_2) weighted barycenter of the elements of \mathcal{A} by optimizing over the weights. Thus, each image is represented by a vector of weights $w \in \Delta_K$. Our new embedded dataset is now a table of shape (450×50) . We train a random forest classifier using the Scikit-learn library (Pedregosa et al., 2011) on this learned embedding) and compute a 10-fold cross-validation. Figure 9 displays the accuracy scores for $F = OT_\epsilon^U$ and $F = S_\epsilon$ for 20 different randomized selections of the dictionary \mathcal{A} . The debiased S_ϵ improves accuracy and is less sensitive to the setting of ϵ .

Conclusion

Entropy regularized OT was previously known to induce a bias that can be mitigated using Sinkhorn divergences. Using OT barycenters of Gaussian distributions, we have shown that this entropy bias can be a blur or a shrink depending on the reference measure defining the relative entropy function. We have also extended the convexity and differentiability properties of OT and the Sinkhorn divergence to measures with non-compact supports.

Acknowledgments

MC and HJ acknowledge the support of a chaire d'excellence de l'IDEX Paris Saclay. AG and HJ were supported by the European Research Council Starting Grant SLAB ERC-YStG-676943. We thank Thibault Séjourné and François-Xavier Vialard for fruitful discussions, in particular for pointing out the identity (8). We thank Zikai Ziong for sharing the matlab code and adapting it to our ellipses experiment.

References

- Agueh, M. and Carlier, G. Barycenters in the Wasserstein space. *SIAM*, 43(2):904–924, 2011.
- Amari, S.-i., Karakida, R., Oizumi, M., and Cuturi, M. Information geometry for regularized optimal transport and barycenters of patterns. *Neural computation*, 31(5): 827–848, 2019.
- Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. Iterative bregman projections for regularized transportation problems. *SIAM J. Scientific Computing*, 37, 2014.
- Berg, C., Christensen, J. P. R., and Ressel, P. *Harmonic Analysis on Semigroups*. Springer, Berlin, 1984.
- Blondel, M., Seguy, V., and Rolet, A. Smooth and sparse optimal transport. *international conference on artificial intelligence and statistics*, 2018.
- Bonneel, N., Peyré, G., and Cuturi, M. Wasserstein barycentric coordinates: Histogram regression using optimal transport. *ACM Trans. Graph.*, 35(4), July 2016. ISSN 0730-0301. doi: 10.1145/2897824.2925918. URL <https://doi.org/10.1145/2897824.2925918>.
- Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. Scaling Algorithms for Unbalanced Transport Problems. *arXiv:1607.05816 [math.OA]*, 2017.
- Cuturi, M. Sinkhorn Distances: Lightspeed Computation of Optimal Transport. In *Neural Information Processing Systems*, 2013.

- Cuturi, M. and Doucet, A. Fast computation of wasserstein barycenters. In Xing, E. P. and Jebara, T. (eds.), *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 685–693, Beijing, China, 22–24 Jun 2014. PMLR. URL <http://proceedings.mlr.press/v32/cuturi14.html>.
- Di Marino, S. and Gerolin, A. An Optimal Transport approach for the Schrodinger bridge problem and convergence of Sinkhorn algorithm, 2019.
- Feydy, J., Séjourné, T., Vialard, F.-X., Amari, S.-i., Trounev, A., and Peyré, G. Interpolating between optimal transport and mmd using sinkhorn divergences. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, 10 2018.
- Ge, D., Wang, H., Xiong, Z., and Ye, Y. Interior-point methods strike back: Solving the wasserstein barycenter problem. In *NeurIPS 2019*, 2019.
- Genevay, A., Peyré, G., and Cuturi, M. Learning generative models with sinkhorn divergences. In Storkey, A. and Perez-Cruz, F. (eds.), *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 1608–1617. PMLR, 09–11 Apr 2018.
- Ivan Gentil, Christian Léonard, L. R. About the analogy between optimal transport and minimal entropy. *Annales de la Facult des Sciences de Toulouse, Mathématiques.*, 2017.
- Knight, P. A., Ruiz, D., and Uar, B. A symmetry preserving algorithm for matrix scaling. *SIAM Journal on Matrix Analysis and Applications*, 35, 07 2014. doi: 10.1137/110825753.
- LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Li, J. and Wang, J. Z. Real-time computerized annotation of pictures. In *Proceedings of the 14th ACM International Conference on Multimedia*, MM 06, pp. 911920, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595934472. doi: 10.1145/1180639.1180841. URL <https://doi.org/10.1145/1180639.1180841>.
- Luise, G., Rudi, A., Pontil, M., and Ciliberto, C. Differential properties of sinkhorn approximation for learning with wasserstein distance. In *Advances in Neural Information Processing Systems*, pp. 5859–5870, 2018.
- Luise, G., Salzo, S., Pontil, M., and Ciliberto, C. Sinkhorn barycenters with free support via frank-wolfe algorithm. In *Advances in Neural Information Processing Systems*, 2019.
- Mena, G. and Weed, J. Statistical bounds for entropic optimal transport: sample complexity and the central limit theorem. 2019.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in pytorch. 2017.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Peyré, G. and Cuturi, M. Computational Optimal Transport. *arXiv e-prints*, March 2018.
- Ramdas, A., Trillos, N., and Cuturi, M. On wasserstein two-sample testing and related families of nonparametric tests. *Entropy*, 19(2):47, 2017.
- Rigollet, P. and Weed, J. Entropic optimal transport is maximum-likelihood deconvolution. *Comptes Rendus Mathématique*, 356(11):1228 – 1235, 2018. ISSN 1631-073X. doi: <https://doi.org/10.1016/j.crma.2018.10.010>. URL <http://www.sciencedirect.com/science/article/pii/S1631073X18302802>.
- Schmitzer, B. Stabilized sparse scaling algorithms for entropy regularized transport problems. *SIAM Journal on Scientific Computing*, 41:A1443–A1481, 2016.
- Solomon, J., de Goes, F., Peyré, G., Cuturi, M., Butscher, A., Nguyen, A., Du, T., and Guibas, L. Convolutional Wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Trans. Graph.*, 34(4):66:1–66:11, July 2015. ISSN 0730-0301.
- Sullivan, C. and Kaszynski, A. Pyvista: 3d plotting and mesh analysis through a streamlined interface for the visualization toolkit (vtk). *Journal of Open Source Software*, 4(37):1450, 2019. doi: 10.21105/joss.01450. URL <https://doi.org/10.21105/joss.01450>.

A. Convexity and Optimality condition

In this section we show how the notion of differentiability along feasible directions in $\mathcal{P}(\mathbb{R}^d)$ is enough to characterize convexity and first order optimality conditions. Consider an arbitrary function F on the space of probability measures.

Definition 1. F is said to be differentiable at $\alpha \in \mathcal{P}(\mathbb{R}^d)$, if and only if there exists $\nabla F(\alpha) \in \mathcal{C}(\mathbb{R}^d)$ such that for any displacement $\delta\alpha = \alpha_1 - \alpha_2$ with $\alpha_1, \alpha_2 \in \mathcal{P}(\mathbb{R}^d)$:

$$F(\alpha + t\delta\alpha) = F(\alpha) + t\langle \delta\alpha, \nabla F(\alpha) \rangle + o(t) , \quad (23)$$

where $\langle \eta, \nabla F(\alpha) \rangle = \int_{\mathbb{R}^d} \nabla F(\alpha) d\eta$.

Proposition 6 (convexity). Assume F is differentiable on $\mathcal{P}(\mathbb{R}^d)$. F is convex if and only if for all $\alpha, \alpha' \in \mathcal{P}(\mathbb{R}^d)$:

$$F(\alpha) \geq F(\alpha') + \langle \alpha - \alpha', \nabla F(\alpha') \rangle , \quad (24)$$

PROOF. (\Rightarrow). Assume (24) holds. Let $\lambda \in [0, 1]$ and $\alpha_\lambda = \lambda\alpha + (1 - \lambda)\alpha'$ with arbitrary probability measures α, α' . Applying (24) twice with $\alpha' = \alpha_\lambda$ leads to:

$$\begin{aligned} F(\alpha) &\geq F(\alpha_\lambda) + \langle \alpha - \alpha_\lambda, \nabla F(\alpha_\lambda) \rangle \\ F(\alpha') &\geq F(\alpha_\lambda) + \langle \alpha' - \alpha_\lambda, \nabla F(\alpha_\lambda) \rangle \end{aligned}$$

Multiplying the first equation by λ and the second one by $1 - \lambda$ before summing leads to:

$$\lambda F(\alpha) + (1 - \lambda)F(\alpha') \geq F(\alpha_\lambda).$$

Thus F is convex.

(\Leftarrow). Assume F is convex. Let $\lambda \in (0, 1)$. Convexity implies that:

$$\begin{aligned} F(\lambda\alpha + (1 - \lambda)\alpha') &\leq \lambda F(\alpha) + (1 - \lambda)F(\alpha') \\ \Rightarrow F(\alpha' + \lambda(\alpha - \alpha')) &\leq \lambda F(\alpha) + (1 - \lambda)F(\alpha') \\ \Rightarrow F(\alpha') + \lambda\langle \alpha - \alpha', \nabla F(\alpha') \rangle + o(\lambda) &\leq \lambda F(\alpha) + (1 - \lambda)F(\alpha') \\ \Rightarrow \lambda\langle \alpha - \alpha', \nabla F(\alpha') \rangle + o(\lambda) &\leq \lambda F(\alpha) - \lambda F(\alpha') \\ \Rightarrow \langle \alpha - \alpha', \nabla F(\alpha') \rangle + \frac{o(\lambda)}{\lambda} &\leq F(\alpha) - F(\alpha') \end{aligned}$$

Letting $\lambda \rightarrow 0$ leads to (24). □

Proposition 7 (Optimality condition). Assume F is differentiable and convex on $\mathcal{P}(\mathbb{R}^d)$ then α^* minimizes F if and only if $\langle \nabla F(\alpha^*), \alpha - \alpha^* \rangle \geq 0$.

PROOF. (\Rightarrow) Assume α^* is a minimizer of F . Let $t > 0$. Since $\mathcal{P}(\mathbb{R}^d)$ is convex, we can write for any $\alpha \in \mathcal{P}(\mathbb{R}^d)$:

$$F(\alpha^*) \leq F(\alpha^* + t(\alpha - \alpha^*))$$

For t small enough, we can use (23) on the right-hand side:

$$F(\alpha^*) \leq F(\alpha^*) + t\langle \alpha - \alpha^*, \nabla F(\alpha^*) \rangle + o(t)$$

Dividing by t and letting $t \rightarrow 0$ leads to $\langle \alpha - \alpha^*, \nabla F(\alpha^*) \rangle \geq 0$ for all α .

(\Leftarrow) Assume $\langle \nabla F(\alpha^*), \alpha - \alpha^* \rangle \geq 0$. Proposition 6 applies and (24) allows to conclude that α^* is a minimizer of F . □

B. Proofs of differentiability and convexity

Proof of Lemma 1

Lemma B.1. Let $\alpha, \alpha' \in \mathcal{G}(\mathbb{R}^d)$ and let $h_\alpha, h_{\alpha'}$ denote their respective autocorrelation potentials. Then:

$$\int e^{\frac{h_\alpha(x)}{\varepsilon}} K(x, y) e^{\frac{h_{\alpha'}(y)}{\varepsilon}} d\alpha(x) d\alpha'(y) \leq 1 \quad (25)$$

PROOF. The left side of (25) can be equivalently written using Fubini-Tonelli:

$$\begin{aligned}
 A &= \int e^{\frac{h_{\alpha}(x)}{\varepsilon}} K(x, y) e^{\frac{h_{\alpha'}(y)}{\varepsilon}} d\alpha(x) d\alpha'(y) \\
 &= \langle e^{\frac{h_{\alpha}}{\varepsilon}} \cdot \alpha, \mathcal{K}(e^{\frac{h_{\alpha'}}{\varepsilon}} \alpha') \rangle \\
 &= \langle e^{\frac{h_{\alpha'}}{\varepsilon}} \cdot \alpha', \mathcal{K}^{\top}(e^{\frac{h_{\alpha}}{\varepsilon}} \alpha) \rangle \\
 &= \langle e^{\frac{h_{\alpha'}}{\varepsilon}} \cdot \alpha', \mathcal{K}(e^{\frac{h_{\alpha}}{\varepsilon}} \alpha) \rangle,
 \end{aligned}$$

where the last equality follows from the symmetry of K . Thus we have:

$$A = \frac{1}{2} \langle e^{\frac{h_{\alpha'}}{\varepsilon}} \cdot \alpha', \mathcal{K}(e^{\frac{h_{\alpha}}{\varepsilon}} \alpha) \rangle + \frac{1}{2} \langle e^{\frac{h_{\alpha}}{\varepsilon}} \cdot \alpha, \mathcal{K}(e^{\frac{h_{\alpha'}}{\varepsilon}} \alpha') \rangle \quad (26)$$

Since the optimal transport plans (primal solutions) associated with $\text{OT}_{\varepsilon}^{\otimes}(\alpha, \alpha)$ and $\text{OT}_{\varepsilon}^{\otimes}(\alpha', \alpha')$ integrate to 1, the right side of (25) can be written:

$$1 = \frac{1}{2} \langle e^{\frac{h_{\alpha}}{\varepsilon}} \cdot \alpha, \mathcal{K}(e^{\frac{h_{\alpha}}{\varepsilon}} \alpha) \rangle + \frac{1}{2} \langle e^{\frac{h_{\alpha'}}{\varepsilon}} \cdot \alpha', \mathcal{K}(e^{\frac{h_{\alpha'}}{\varepsilon}} \alpha') \rangle \quad (27)$$

Combining (26) with (27), it holds:

$$1 - A = \frac{1}{2} \langle r, \mathcal{K}(r) \rangle$$

where $r = e^{\frac{h_{\alpha}}{\varepsilon}} \cdot \alpha - e^{\frac{h_{\alpha'}}{\varepsilon}} \cdot \alpha'$. Since K is semi-definite positive, $1 - A \geq 0$. \square

Differentiability of OT_{ε}

Proposition B.1. *Under assumption (1), OT_{ε} is differentiable and its gradient is given by:*

$$\nabla \text{OT}_{\varepsilon}^{\otimes}(\alpha, \beta) = (f, g) \quad (28)$$

Where f and g satisfy the Sinkhorn fixed point system (11) on \mathbb{R}^d .

PROOF. Consider $\alpha, \beta, \alpha_1, \alpha_2, \beta_1, \beta_2 \in \mathcal{G}(\mathbb{R}^d)$ and denote the displacements $\delta\alpha = \alpha_1 - \alpha_2$ and $\delta\beta = \beta_1 - \beta_2$. Let Δ_t denote the ratio of (15):

$$\Delta_t = \frac{\text{OT}_{\varepsilon}^{\otimes}(\alpha_t, \beta_t) - \text{OT}_{\varepsilon}^{\otimes}(\alpha, \beta)}{t}, \quad (29)$$

where $\alpha_t = \alpha + t\delta\alpha$ and $\beta_t = \beta + t\delta\beta$. Similarly to the proof of Proposition 2 of Feydy et al. (2018), we derive a lower and upper bound of Δ_t using suboptimal potentials. On one hand, the pair (f, g) is suboptimal for the dual problem defining $\text{OT}_{\varepsilon}^{\otimes}(\alpha_t, \beta_t)$. Therefore:

$$\begin{aligned}
 \text{OT}_{\varepsilon}^{\otimes}(\alpha_t, \beta_t) &\geq \langle \alpha_t, f \rangle + \langle \beta_t, g \rangle \\
 &\quad - \varepsilon \langle \alpha_t \otimes \beta_t, \exp\left(\frac{f \oplus g - C}{\varepsilon}\right) \rangle + \varepsilon
 \end{aligned}$$

Therefore, (10) and (11) lead to the lower bound:

$$\Delta_t \geq \langle \delta\alpha, f - \varepsilon \rangle + \langle \delta\beta, g - \varepsilon \rangle + o(1)$$

And similarly we get the upper bound:

$$\Delta_t \leq \langle \delta\alpha, f_t - \varepsilon \rangle + \langle \delta\beta, g_t - \varepsilon \rangle + o(1)$$

As $t \rightarrow 0$, $(\alpha_t, \beta) \rightarrow (\alpha, \beta)$. On one hand, Proposition 4 of Mena & Weed (2019) leads to the pointwise convergence of the sequence of potentials (f_t, g_t) towards (f, g) . On the other hand, Proposition 2 implies that there exists $M > 0$ such that $|f_t(x)| \leq M\|x\|^2$ for all $x \in \mathbb{R}^d$. Given that any $\mu \in \mathcal{G}(\mathbb{R}^d)$ has a second order moment, by Lebesgue's dominated convergence we have $\langle \mu, f_t \rangle \rightarrow \langle \mu, f \rangle$. Similarly, $\langle \mu, g_t \rangle \rightarrow \langle \mu, g \rangle$. Finally, since $\langle \delta\alpha, \varepsilon \rangle = \langle \delta\beta, \varepsilon \rangle = 0$, we get as $t \rightarrow 0$, $\Delta_t \rightarrow \langle \delta\alpha, f \rangle + \langle \delta\beta, g \rangle$. Since f and g are smooth (Prop 2) and square-integrable with respect to any $\mu \in \mathcal{G}(\mathbb{R}^d)$, (28) holds for $\nabla \text{OT}_{\varepsilon}^{\otimes}(\alpha, \beta) = (f, g)$. \square

B.1. Differentiability and convexity of $\text{OT}_\varepsilon^\mathcal{L}$

To prove theorem 1 we first need to establish the differentiability and convexity of $\text{OT}_\varepsilon^\mathcal{L}$ on the set of sub-Gaussian measures $\mathcal{G}_\sigma(\mathbb{R}^d)$ which are absolutely continuous with respect to the Lebesgue measure.

Dual problem Let α, β continuous sub-Gaussian measures. Identifying α, β and π with their Lebesgue densities, The OT problem (2) has a dual problem given by:

$$\text{OT}_\varepsilon^\mathcal{L}(\alpha, \beta) = \sup_{f \in \mathcal{L}_1(\alpha), g \in \mathcal{L}_1(\beta)} \langle f, \alpha \rangle + \langle g, \beta \rangle - \varepsilon \int \int \exp\left(\frac{f(x) + g(y) - C(x, y)}{\varepsilon}\right) dx dy + \varepsilon, \quad (30)$$

Notice that the convexity of $\text{OT}_\varepsilon^\mathcal{L}$ follows immediately from (30) since it is a supremum of linear functions in α and β . The optimality conditions are equivalent to the marginal constraints of the primal problem (2). However, they are slightly different than those of $\text{OT}_\varepsilon^\otimes$. Cancelling the gradient of the dual problem leads to the following system (Ivan Gentil, 2017):

$$\begin{aligned} e^{\frac{f}{\varepsilon}} \mathcal{K}(e^{\frac{g}{\varepsilon}}) &= \alpha, \\ e^{\frac{g}{\varepsilon}} \mathcal{K}^\top(e^{\frac{f}{\varepsilon}}) &= \beta, \end{aligned} \quad (31)$$

which in integral form can be written:

$$\begin{aligned} e^{\frac{f(x)}{\varepsilon}} \int e^{\frac{-C(x, y) + g(y)}{\varepsilon}} dy &= \alpha(x) \quad \forall x, \\ e^{\frac{g(y)}{\varepsilon}} \int e^{\frac{-C(y, x) + f(x)}{\varepsilon}} dx &= \beta(y) \quad \forall y, \end{aligned} \quad (32)$$

and the optimal transport plan's density π is given by: $\pi(x, y) = \exp\left(\frac{f(x) + g(y) - C(x, y)}{\varepsilon}\right)$

Thus, at optimality the integral over $\mathbb{R}^d \times \mathbb{R}^d$ sums to 1 and:

$$\text{OT}_\varepsilon^\mathcal{L}(\alpha, \beta) = \langle f, \alpha \rangle + \langle g, \beta \rangle \quad (33)$$

Convexity and Differentiability By using the existence of Lebesgue continuity, one can rewrite the KL in the primal problem such that it holds (Di Marino & Gerolin, 2019):

$$\text{OT}_\varepsilon^\mathcal{L}(\alpha, \beta) = \text{OT}_\varepsilon^\otimes(\alpha, \beta) + \varepsilon \text{KL}(\alpha|\mathcal{L}) + \varepsilon \text{KL}(\beta|\mathcal{L}) \quad (34)$$

We already showed that $\text{OT}_\varepsilon^\otimes$ is convex (w.r.t. to one argument); KL is also convex (even jointly convex). Since the set of Lebesgue-continuous and sub-Gaussian measures is convex, $\text{OT}_\varepsilon^\mathcal{L}$ is also convex with respect to one argument.

Identifying α with its density, we have $E(\alpha) \stackrel{\text{def}}{=} \text{KL}(\alpha, \mathcal{L}) = \int \alpha(x)(\log(\alpha(x)) - 1) dx$. If $\alpha > 0$, then for any feasible displacement $h = h_1 - h_2$ with density functions h_1, h_2 . The functional derivative of E in the direction h is given by: $\left[\frac{dE(\alpha+th)}{dt}\right]_{t=0} = \langle h, \log(\alpha) \rangle$. Thus, in the sense of the directional differentiation (15): $\nabla_\alpha \text{KL}(\alpha, \mathcal{L}) = \log(\alpha)$.

Let (f, g) be a pair of optimal potentials for $\text{OT}_\varepsilon^\otimes(\alpha, \beta)$. Following (B.1) and the differentiability of KL, $\text{OT}_\varepsilon^\mathcal{L}$ is differentiable on the set of sub-Gaussian measures with positive density functions and its gradient is given by: $\nabla_1 \text{OT}_\varepsilon^\mathcal{L}(\alpha, \beta) = f - \varepsilon \log(\alpha)$. By a simple calculation, it is easy to show that $(f - \varepsilon \log(\alpha), g - \varepsilon \log(\beta))$ are actually solutions of the Sinkhorn equations (33). Similarly, given a solution (f_1, g_1) of (30), $(f_1 + \varepsilon \log(\alpha), g_1 + \varepsilon \log(\beta))$ are optimal potentials of $\text{OT}_\varepsilon^\otimes$. Therefore, the following proposition holds:

Proposition B.2. *Let $\alpha, \beta \in \mathcal{G}_\sigma(\mathbb{R}^d)$. If α and β are Lebesgue-continuous with positive density functions, then $\text{OT}_\varepsilon^\mathcal{L}$ is differentiable and it holds:*

$$\nabla \text{OT}_\varepsilon^\mathcal{L}(\alpha, \beta) = (f, g), \quad (35)$$

where (f, g) is a pair of dual potentials verifying the fixed point equations (33).

C. Proofs of the theorems

We first start by showing that the equations verified by the variance of the barycenter have a unique positive solution.

C.1. Fixed point equations Lemmas

For the 3 following lemmas, since the variance S can only be positive, we re-parametrized the equations by replacing S^2 with S for the sake of simplicity.

Lemma C.1. *Under the assumptions of Theorem 1, the equation in S :*

$$\sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S} = -\varepsilon'^2 + 2S \quad (36)$$

has a positive solution.

PROOF. Let $f : S \in \mathbb{R}_+ \rightarrow \sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S} + \varepsilon'^2 - 2S$. Since f is continuous and $f(0) = 2\varepsilon'^2 > 0$ and $\lim_{S \rightarrow -\infty} f(S) = -\infty$, there exists $S > 0$ such that $f(S) = 0$. \square

Lemma C.2. *Under the assumptions of Theorem 2, the equation in S :*

$$\sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S} = \varepsilon'^2 + 2S \quad (37)$$

has a positive solution if and only if $\varepsilon'^2 < \bar{\sigma} = \sum_{k=1} w_k \sigma_k^2$.

PROOF. Let $f : S \in \mathbb{R}_+ \rightarrow \sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S} - \varepsilon'^2 - 2S$. Sufficient condition. since $f(0) = 0$ and $\lim_{S \rightarrow +\infty} f(S) = -\infty$, a positive solution exists if $f'(0) > 0$.

$$\begin{aligned} f'(0) > 0 &\Leftrightarrow \sum_{k=1} \frac{w_k \sigma_k^2}{\varepsilon'^2} - 1 > 0 \\ &\Leftrightarrow \bar{\sigma} > \varepsilon'^2 \end{aligned}$$

Necessary condition. Conversely, we have by Jensen's inequality:

$$\begin{aligned} f(S) &\leq \sqrt{\varepsilon'^4 + 4\bar{\sigma}S} - \varepsilon'^2 - 2S \\ &= 4S \frac{\bar{\sigma} - S - \varepsilon'^2}{\sqrt{\varepsilon'^4 + 4\bar{\sigma}S} + \varepsilon'^2 + 2S} \end{aligned}$$

Therefore, if $\bar{\sigma} \leq \varepsilon'^2$ then $f(S) \leq -\frac{4S^2}{\sqrt{\varepsilon'^4 + 4\bar{\sigma}S} + \varepsilon'^2 + 2S} < 0$ for any $S > 0$. \square

Lemma C.3. *Under the assumptions of Theorem 3, the equation in S :*

$$\sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S} = \sqrt{\varepsilon'^4 + 4S^4} \quad (38)$$

has a positive solution S^* and it holds $S^* \in (\sigma_{(0)}, \sigma_{(K)})$.

PROOF. Let $f : S \in \mathbb{R}_+ \rightarrow \sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S} - \sqrt{\varepsilon'^4 + 4S^4}$. It holds:

$$\begin{aligned} f(S) &\geq \sqrt{\varepsilon'^4 + 4\sigma_{(0)}^2 S} - \sqrt{\varepsilon'^4 + 4S^4} \\ &= \frac{4S(\sigma_{(0)}^2 - S)}{\sqrt{\varepsilon'^4 + 4\sigma_{(0)}^2 S} + \sqrt{\varepsilon'^4 + 4S^4}}. \end{aligned}$$

Thus $f(\sigma_{(0)}^2) \geq 0$. Similarly $f(\sigma_{(K)}) \leq 0$. Thus there exists $S^* \in (\sigma_{(0)}, \sigma_{(K)})$ such that $f(S^*) = 0$. \square

C.2. Proofs of theorems 2 and 3

We turn now to proving theorems 3 and 2. We have shown that both OT_ε and S_ε are convex and differentiable on the convex set of sub-Gaussian measure on \mathbb{R}^d . Thus, the proposition holds:

Proposition 8. *Let $\alpha_1, \dots, \alpha_K \in \mathcal{G}(\mathbb{R}^d)$. Let (w_1, \dots, w_K) be non-negative weights summing to 1. Then:*

$\alpha_{\text{OT}_\varepsilon^\otimes} = \arg \min_\alpha \sum_{k=1}^K w_k \text{OT}_\varepsilon^\otimes(\alpha_k, \alpha)$ if and only if there exists a set of potentials $f_1, \dots, f_K, g_1, \dots, g_K$ such that for any direction $\beta \in \mathcal{G}(\mathbb{R}^d)$ the following equations hold everywhere in \mathbb{R}^d :

$$\begin{cases} e^{\frac{f_k}{\varepsilon}} \mathcal{K}(e^{\frac{g_k}{\varepsilon}} \cdot \alpha_{\text{OT}_\varepsilon^\otimes}) = 1, & e^{\frac{g_k}{\varepsilon}} \mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}} \cdot \alpha_k) = 1, \\ \langle \sum_{k=1}^K w_k g_k, \beta - \alpha_{\text{OT}_\varepsilon^\otimes} \rangle \geq 0 \end{cases} \quad (39)$$

$\alpha_{S_\varepsilon} = \arg \min_\alpha \sum_{k=1}^K w_k S_\varepsilon(\alpha_k, \alpha)$ if and only if there exists a set of potentials $f_1, \dots, f_K, g_1, \dots, g_K, h$ such that for any direction $\beta \in \mathcal{G}(\mathbb{R}^d)$ the following equations hold everywhere in \mathbb{R}^d :

$$\begin{cases} e^{\frac{f_k}{\varepsilon}} \mathcal{K}(e^{\frac{g_k}{\varepsilon}} \cdot \alpha_{S_\varepsilon}) = 1, & e^{\frac{g_k}{\varepsilon}} \mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}} \cdot \alpha_k) = 1, \\ e^{\frac{h}{\varepsilon}} \mathcal{K}(e^{\frac{h}{\varepsilon}} \cdot \alpha_{S_\varepsilon}) = 1, \\ \langle \sum_{k=1}^K w_k g_k - h, \beta - \alpha_{S_\varepsilon} \rangle \geq 0 \end{cases} \quad (40)$$

We solve the systems of equations (39) and (40) by restricting the potentials to quadratic functions. Since the objectives are convex, showing the existence of a solution is sufficient for optimality. We start with the Debiased barycenter (theorem 3).

Theorem C.1 (Debiasing of S_ε). *Let $C(x, y) = (x - y)^2$ and $0 < \varepsilon < +\infty$ and $\varepsilon = 2\varepsilon'^2$. Let (w_k) be positive weights that sum to 1. Let \mathcal{N} denote the Gaussian distribution. Assume that $\alpha_k \sim \mathcal{N}(\mu_k, \sigma_k^2)$ and let $\bar{\mu} = \sum_k w_k \mu_k$, $\bar{\sigma} = \sum_{k=1} w_k \sigma_k^2$ then:*

Then $\alpha_{S_\varepsilon} \sim \mathcal{N}(\bar{\mu}, S^2)$ where S is the unique non-zero solution S^ of the fixed point equation:*

$\sum_{k=1} w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S^2} = \sqrt{\varepsilon'^4 + 4S^4}$. *Moreover, given a sorted sequence $\sigma_{(1)} \leq \dots \leq \sigma_{(K)}$, it holds $S^* \in (\sigma_{(0)}, \sigma_{(K)})$.*

In particular, if all σ_k are equal to some $\sigma > 0$, then $\alpha_{S_\varepsilon} \sim \mathcal{N}(\bar{\mu}, \sigma^2)$.

PROOF. Convexity makes the system equations (39) and (40) sufficient for optimality. Thus, we only need to find a particular solution. We are going to show that there exist a set of quadratic polynomial potentials and Gaussian probability measures satisfying each system. First, let's start with the S_ε barycenter α_{S_ε} .

Consider polynomial potentials of the form $f_k(x) = F_{2,k}x^2 + F_{1,k}x + F_{0,k}$ and $g_k(x) = G_{2,k}x^2 + G_{1,k}x + G_{0,k}$ and $h(x) = H_2x^2 + H_1x + H_0$ for some unknown coefficients $F_{2,k}, F_{1,k}, F_{0,k}, G_{2,k}, G_{1,k}, G_{0,k}, H_2, H_1, H_0 \in \mathbb{R}$, and assume that $\frac{d\alpha_{\text{OT}_\varepsilon}}{d\lambda} = \mathcal{N}(m, S)$. First, we will write the first and second order coefficients as functions of m and S then use the optimality condition to find m and S .

Sufficient optimality condition Let $\beta \in \mathcal{G}(\mathbb{R}^d)$. And let $M_r(\beta)$ denote the r -th moment of β . For any real sequence y_1, \dots, y_k , let \bar{y} denote its weighted average $\sum_{k=1}^K w_k y_k$. The optimality condition reads:

$$\begin{aligned} & \left\langle \sum_{k=1}^K w_k g_k - h, \beta - \alpha_{S_\varepsilon} \right\rangle \geq 0 \\ & \Leftrightarrow (\bar{G}_2 - H_2)(M_2(\beta) - M_2(\alpha_{S_\varepsilon})) + (\bar{G}_1 - H_1)(M_1(\beta) - M_1(\alpha_{S_\varepsilon})) + (\bar{G}_0 - H_0)(M_0(\beta) - M_0(\alpha_{S_\varepsilon})) \geq 0 \\ & \Leftrightarrow (\bar{G}_2 - H_2)(M_2(\beta) - M_2(\alpha_{S_\varepsilon})) + (\bar{G}_1 - H_1)(M_1(\beta) - M_1(\alpha_{S_\varepsilon})) \geq 0 \end{aligned}$$

Where the last inequality follows from $M_0(\beta) = M_0(\alpha_{S_\varepsilon}) = \int d\alpha_{S_\varepsilon} = 1$. Thus, the 0-order coefficients are irrelevant for optimality. We are going to show that there exist a set of coefficients such that the following sufficient conditions for

optimality hold:

$$\bar{G}_2 \stackrel{\text{def}}{=} \sum_{k=1}^K w_k G_{2,k} = H_2 \quad (41)$$

$$\bar{G}_1 \stackrel{\text{def}}{=} \sum_{k=1}^K w_k G_{1,k} = H_1 \quad (42)$$

Kernel integration Dropping the k exponent for the sake of convenience, let's carefully derive the integral $\mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}} \cdot \alpha_k)$:

$$\begin{aligned} \mathcal{K}^\top(e^{\frac{f}{\varepsilon}} \alpha_k)(x) &= \int K(x, y) e^{\frac{f(y)}{2\varepsilon'^2}} \frac{d\alpha}{d\lambda}(y) dy \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int \exp\left(\frac{-(x-y)^2 + f(y)}{2\varepsilon'^2} - \frac{(y-\mu)^2}{2\sigma^2}\right) dy \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int \exp\left(\underbrace{\left[\frac{F_2-1}{2\varepsilon'^2} - \frac{1}{2\sigma^2}\right]}_A y^2 + \underbrace{\left[\frac{F_1}{2\varepsilon'^2} + \frac{x}{\varepsilon'^2} + \frac{\mu}{\sigma^2}\right]}_{Z(x)} y + \left[\frac{F_0-x^2}{2\varepsilon'^2} - \frac{\mu^2}{2\sigma^2}\right]\right) dy \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{F_0-x^2}{2\varepsilon'^2} - \frac{\mu^2}{2\sigma^2}\right) \int \exp\left(A \left[y^2 + \frac{Z(x)}{A} y\right]\right) dy \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{F_0-x^2}{2\varepsilon'^2} - \frac{\mu^2}{2\sigma^2}\right) \int \exp\left(A \left[y + \frac{Z(x)}{2A}\right]^2 - \frac{Z(x)^2}{4A}\right) dy \\ &= \frac{1}{\sigma} \exp\left(\frac{F_0-x^2}{2\varepsilon'^2} - \frac{\mu^2}{2\sigma^2} - \frac{Z(x)^2}{4A}\right) \underbrace{\frac{1}{\sqrt{2\pi}} \int \exp\left(A \left[y + \frac{Z(x)}{2A}\right]^2\right) dy}_I \end{aligned}$$

For the fourth equality to be sound, we need $A \neq 0$, and for the integral I to be finite, we need $A < 0$ which is equivalent to:

$$F_2 < 1 + \frac{\varepsilon'^2}{\sigma^2}. \quad (43)$$

In that case, $I = \frac{1}{\sqrt{-2A}}$ thus:

$$\begin{aligned} \mathcal{K}^\top(e^{\frac{f}{\varepsilon}} \alpha_k)(x) &= \frac{1}{\sigma\sqrt{-2A}} \exp\left(\frac{F_0-x^2}{2\varepsilon'^2} - \frac{\mu^2}{2\sigma^2} - \frac{Z(x)^2}{4A}\right) \\ &= \frac{1}{\sigma\sqrt{-2A}} \exp\left(\left[-\frac{1}{2\varepsilon'^2} - \frac{1}{4A\varepsilon'^4}\right] x^2 - \left[\frac{F_1}{4A\varepsilon'^4} + \frac{\mu}{2A\sigma^2\varepsilon'^2}\right] x - \frac{\mu^2}{2\sigma^2} + \frac{F_0}{2\varepsilon'^2} - \frac{\left[\frac{F_1}{2\varepsilon'^2} + \frac{\mu}{\sigma^2}\right]^2}{4A}\right) \\ &= \exp\left(\left[-\frac{1}{2\varepsilon'^2} - \frac{1}{4A\varepsilon'^4}\right] x^2 - \left[\frac{F_1}{4A\varepsilon'^4} + \frac{\mu}{2A\sigma^2\varepsilon'^2}\right] x - \frac{\mu^2}{2\sigma^2} + \frac{F_0}{2\varepsilon'^2} - \frac{\left[\frac{F_1}{2\varepsilon'^2} + \frac{\mu}{\sigma^2}\right]^2}{4A} - \log(\sigma\sqrt{-2A})\right) \end{aligned}$$

Sinkhorn equations Using the first Sinkhorn equation $e^{\frac{g_k}{\varepsilon}} \cdot \mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}} \cdot \alpha_k) = 1$ we get by identification, for all k , with :

$$A_k = \frac{F_{2,k} - 1}{2\varepsilon'^2} - \frac{1}{2\sigma_k^2} \quad (44)$$

$$\begin{cases} \frac{G_{2,k-1}}{\varepsilon'^2} - \frac{1}{2A_k\varepsilon'^4} = 0 & [i] \\ \frac{G_{1,k}}{\varepsilon'^2} - \frac{F_{1,k}}{2A_k\varepsilon'^4} - \frac{\mu_k}{A_k\sigma_k^2\varepsilon'^2} = 0 & [ii] \\ \frac{G_{0,k}}{\varepsilon'^2} - \frac{\mu_k^2}{\sigma_k^2} + \frac{F_{0,k}}{\varepsilon'^2} - \frac{\left[\frac{F_{1,k}}{2\varepsilon'^2} + \frac{\mu_k}{\sigma_k^2}\right]^2}{2A_k} - \log(-2\sigma_k^2 A_k) = 0 & [iii] \end{cases}$$

Similarly, since the equations are symmetric, $e^{\frac{f_k}{\varepsilon}} \mathcal{K}(e^{\frac{g_k}{\varepsilon}} \cdot \alpha_{OT_\varepsilon}) = 1$ with $G_{2,k} < 1 + \frac{\varepsilon'^2}{S^2}$ and:

$$B_k = \frac{G_{2,k} - 1}{2\varepsilon'^2} - \frac{1}{2S^2} \quad (45)$$

lead to:

$$\begin{cases} \frac{F_{2,k}-1}{\varepsilon'^2} - \frac{1}{2B_k\varepsilon'^4} = 0 & [j] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{m}{B_kS^2\varepsilon'^2} = 0 & [jj] \\ \frac{F_{0,k}}{\varepsilon'^2} - \frac{m^2}{S^2} + \frac{G_{0,k}}{\varepsilon'^2} - \frac{\left[\frac{G_{1,k}}{2\varepsilon'^2} + \frac{m}{S^2}\right]^2}{2B_k} - \log(-2S^2B_k) = 0 & [jjj] \end{cases}$$

Second order coefficients G_2, F_2 Let's rewrite [i] and [j] separately:

$$\begin{aligned} & \begin{cases} 2B_k + \frac{1}{S^2} - \frac{1}{2A_k\varepsilon'^4} = 0 \\ 2A_k + \frac{1}{\sigma_k^2} - \frac{1}{2B_k\varepsilon'^4} = 0 \end{cases} \\ \Leftrightarrow & \begin{cases} A_kB_k + \frac{A_k}{2S^2} - \frac{1}{4\varepsilon'^4} = 0 \\ A_kB_k + \frac{B_k}{2\sigma_k^2} - \frac{1}{4\varepsilon'^4} = 0 \end{cases} \\ \Leftrightarrow & \begin{cases} \frac{B_k}{\sigma_k^2} = \frac{A_k}{S^2} \\ A_kB_k + \frac{B_k}{2\sigma_k^2} - \frac{1}{4\varepsilon'^4} = 0 \end{cases} \\ \Leftrightarrow & \begin{cases} \frac{B_k}{\sigma_k^2} = \frac{A_k}{S^2} \\ B_k^2 + \frac{B_k}{2S^2} - \frac{\sigma_k^2}{4\varepsilon'^4S^2} = 0 \end{cases} \end{aligned}$$

The roots of the polynomial above are: $-\frac{1}{4S^2} \pm \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}}$. The constraint $B_k < 0$ eliminates the positive solution and it holds:

$$B_k = -\frac{1}{4S^2} - \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}} \quad (46)$$

$$A_k = \frac{S^2}{\sigma_k^2} B_k \quad (47)$$

First order coefficients G_1, F_1 Let's rewrite [ii] and [jj] separately:

$$\begin{aligned} & \begin{cases} \frac{G_{1,k}}{\varepsilon'^2} - \frac{F_{1,k}}{2A_k\varepsilon'^4} - \frac{\mu_k}{A_k\sigma_k^2\varepsilon'^2} = 0 & [ii] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{m}{B_kS^2\varepsilon'^2} = 0 & [jj] \end{cases} \\ \Leftrightarrow & \begin{cases} 2A_kG_{1,k} - \frac{F_{1,k}}{\varepsilon'^2} - \frac{2\mu_k}{\sigma_k^2} = 0 & [ii] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{m}{B_kS^2\varepsilon'^2} = 0 & [jj] \end{cases} \\ \Leftrightarrow & \begin{cases} \left(2A_k - \frac{1}{2B_k\varepsilon'^4}\right)G_{1,k} - \frac{2\mu_k}{\sigma_k^2} - \frac{m}{B_kS^2\varepsilon'^2} = 0 & [ii] + [jj] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{m}{B_kS^2\varepsilon'^2} = 0 & [jj] \end{cases} \end{aligned}$$

The equations above between A_k and B_k lead to $2A_k - \frac{1}{2B_k\varepsilon'^4} = -\frac{1}{\sigma_k^2}$ and the second order polynomial equation in B_k leads to $\frac{\sigma_k^2}{B_kS^2\varepsilon'^2} = 4\varepsilon'^2B_k + 2\varepsilon'^2\frac{1}{S^2}$. Therefore, [ii] + [jj] can be written:

$$\begin{aligned} & \frac{1}{\sigma_k^2}G_{1,k} + \frac{2\mu_k}{\sigma_k^2} + \frac{m}{B_kS^2\varepsilon'^2} = 0 \\ \Rightarrow & G_{1,k} + 2\mu_k + m\left(4\varepsilon'^2B_k + 2\varepsilon'^2\frac{1}{S^2}\right) = 0 \\ \Rightarrow & G_{1,k} + 2\mu_k + 2m(G_{2,k} - 1) = 0 \end{aligned}$$

Using [jj] we recover the first order coefficients $F_{1,k}$ and $G_{1,k}$ as function of m :

$$\begin{aligned} G_{1,k} + 2\mu_k + 2m(G_{2,k} - 1) &= 0 \\ F_{1,k} + 2m + 2\mu_k(F_{2,k} - 1) &= 0 \end{aligned} \quad (48)$$

Sinkhorn auto-correlation equation Similarly, the auto-correlation equation $e^{\frac{h}{\varepsilon}} \cdot \mathcal{K}^\top(e^{\frac{h}{\varepsilon}} \cdot \alpha_k) = 1$ leads to the same system of equations (equal dual potentials), with $H_2 < 1 + \frac{\varepsilon'^2}{S^2}$ and:

$$C = \frac{H_2 - 1}{2\varepsilon'^2} - \frac{1}{2S^2} < 0 \quad (49)$$

$$\begin{cases} \frac{H_2 - 1}{\varepsilon'^2} - \frac{1}{2C\varepsilon'^4} = 0 & [a] \\ \frac{H_1}{\varepsilon'^2} - \frac{H_1}{2C\varepsilon'^4} - \frac{\mu}{CS^2\varepsilon'^2} = 0 & [b] \\ \frac{H_0}{\varepsilon'^2} - \frac{\mu^2}{S^2} + \frac{H_0}{\varepsilon'^2} - \frac{\left[\frac{H_1}{2\varepsilon'^2} + \frac{\mu}{S^2}\right]^2}{2C} - \log(-2S^2C) = 0 & [c] \end{cases}$$

Isolating [a] we get: $2C + \frac{1}{S^2} - \frac{1}{2C\varepsilon'^4} = 0$. Again, the only negative root of [a] is given by:

$$C = -\frac{1}{4S^2} - \sqrt{\frac{1}{16S^4} + \frac{1}{4\varepsilon'^4}} \quad (50)$$

and similarly to (48), we also get the link between H_1 and H_2 :

$$H_1 + 2mH_2 = 0 \quad (51)$$

Optimality condition and identifying σ and μ Using the definition of B_2 (45) and (49) and then with their closed form formulas (46) and (50), the first sufficient optimality condition (41) reads:

$$\begin{aligned} \sum_{k=1}^K w_k G_{2,k} = H_2 &\Rightarrow \sum_{k=1}^K w_k B_{2,k} = \frac{H_2 - 1}{2\varepsilon'^2} - \frac{1}{2S^2} \\ &\Rightarrow \sum_{k=1}^K w_k B_{2,k} = C \\ &\Rightarrow \sum_{k=1}^K w_k \left(\frac{1}{4S^2} + \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}} \right) = \frac{1}{4S^2} + \sqrt{\frac{1}{16S^4} + \frac{1}{4\varepsilon'^4}} \\ &\Rightarrow \sum_{k=1}^K w_k \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}} = \sqrt{\frac{1}{16S^4} + \frac{1}{4\varepsilon'^4}} \\ &\Rightarrow \sum_{k=1}^K w_k \sqrt{4\sigma_k^2 S^2 + \varepsilon'^4} = \sqrt{4S^4 + \varepsilon'^4} \end{aligned}$$

Lemma C.3 guarantees that the fixed point equation above possesses a unique positive solution S .

The second sufficient optimality condition (41) combined with the equations on G_1 , F_1 (48) and H_1 (51) lead to identifying m :

$$\sum_{k=1}^K w_k G_{1,k} = H_1 \Rightarrow m = \sum_{k=1}^K w_k \mu_k$$

Identifying the offset coefficients F_0 , G_0 , H_0 Since now m and S are known and unique, all the first and second order coefficients $F_{2,k}$, $G_{2,k}$, H_2 , $F_{1,k}$, $G_{1,k}$, H_1 are uniquely determined. H_0 follows immediately from [c]. Finding $F_{0,k}$ and $G_{0,k}$ can be done up to an additive constant. Adding [iii] and [jjj] leads to a closed form expression on $F_{0,k} + G_{0,k}$. Since the optimality condition does not depend on H_0 , F_0 and G_0 , one may simply set $F_{0,k}$ to 0, and solve $G_{0,k}$ exactly. \square

Theorem C.2 (Shrinking bias of $\text{OT}_\varepsilon^\otimes$). Let $C(x, y) = (x - y)^2$ and $0 < \varepsilon < +\infty$ and $\varepsilon = 2\varepsilon'^2$. Let (w_k) be positive weights that sum to 1. Let \mathcal{N} denote the Gaussian distribution. Assume that $\alpha_k \sim \mathcal{N}(\mu_k, \sigma_k^2)$ and let $\bar{\mu} = \sum_k w_k \mu_k$, $\bar{\sigma} = \sum_{k=1}^K w_k \sigma_k^2$ then:

if $\varepsilon'^2 < \bar{\sigma}$ then $\alpha_{\text{OT}_\varepsilon^\otimes} \sim \mathcal{N}(\bar{\mu}, S^2)$ where S is the unique non-zero solution of the fixed point equation: $\sum_{k=1}^K w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S^2} = \varepsilon'^2 + 2S^2$. In particular, if all σ_k are equal to some $\sigma > 0$, then $\alpha_{\text{OT}_\varepsilon^\otimes} \sim \mathcal{N}(\bar{\mu}, \sigma^2 - \varepsilon'^2)$.

if $\varepsilon'^2 \geq \bar{\sigma}$ then $\alpha_{\text{OT}_\varepsilon^\otimes}$ is a Dirac distribution located at $\bar{\mu}$.

PROOF. When $\varepsilon'^2 < \bar{\sigma}$, the same proof of theorem C.1 applies. Convexity makes the system equations (39) sufficient for optimality. Thus, we only need to find a *particular* solution. We are going to show that there exist a set of quadratic polynomial potentials and Gaussian probability measures satisfying each system.

Consider polynomial potentials of the form $f_k(x) = F_{2,k}x^2 + F_{1,k}x + F_{0,k}$ and $g_k(x) = G_{2,k}x^2 + G_{1,k}x + G_{0,k}$ for some unknown coefficients $F_{2,k}, F_{1,k}, F_{0,k}, G_{2,k}, G_{1,k}, G_{0,k} \in \mathbb{R}$, and assume that $\frac{d\alpha_{\text{OT}_\varepsilon}}{d\lambda} = \mathcal{N}(m, S)$. First, we will write the first and second order coefficients as functions of m and S then use the optimality condition to find m and S .

Sufficient optimality condition Let $\beta \in \mathcal{P}_2(\mathbb{R}^d)$. And let $M_r(\beta)$ denote the r -th moment of β . For any real sequence y_1, \dots, y_k , let \bar{y} denote its weighted average $\sum_{k=1}^K w_k y_k$. The optimality condition reads:

$$\begin{aligned} \left\langle \sum_{k=1}^K w_k g_k, \beta - \alpha_{S_\varepsilon} \right\rangle &\geq 0 \\ \Leftrightarrow \bar{G}_2(M_2(\beta) - M_2(\alpha_{S_\varepsilon})) + \bar{G}_1(M_1(\beta) - M_1(\alpha_{S_\varepsilon})) + \bar{G}_0(M_0(\beta) - M_0(\alpha_{S_\varepsilon})) &\geq 0 \\ \Leftrightarrow \bar{G}_2(M_2(\beta) - M_2(\alpha_{S_\varepsilon})) + \bar{G}_1(M_1(\beta) - M_1(\alpha_{S_\varepsilon})) &\geq 0 \end{aligned}$$

Where the last inequality follows from $M_0(\beta) = M_0(\alpha_{S_\varepsilon}) = \int d\alpha_{S_\varepsilon} = 1$. Thus, the 0-order coefficients are irrelevant for optimality.

1. Case 1: if $\varepsilon'^2 < \bar{\sigma}$: We are going to show that there exist a set of coefficients such that the following sufficient conditions for optimality hold:

$$\bar{G}_2 \stackrel{\text{def}}{=} \sum_{k=1}^K w_k G_{2,k} = 0 \quad (52)$$

$$\bar{G}_1 \stackrel{\text{def}}{=} \sum_{k=1}^K w_k G_{1,k} = 0 \quad (53)$$

The Sinkhorn system on f_k and g_k is the same as in the proof above. Thus, the same equations still hold. The first differences arise when using the optimality condition (52):

Optimality condition and identifying σ and μ Using the definition of B_2 (45) and its closed form formulas (46), the first sufficient optimality condition (52) reads:

$$\begin{aligned} \sum_{k=1}^K w_k G_{2,k} = 0 &\Rightarrow \sum_{k=1}^K w_k B_{2,k} = -\frac{1}{2\varepsilon'^2} - \frac{1}{2S^2} \\ &\Rightarrow \sum_{k=1}^K w_k \left(\frac{1}{4S^2} + \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}} \right) = \frac{1}{2\varepsilon'^2} + \frac{1}{2S^2} \\ &\Rightarrow \sum_{k=1}^K w_k \sqrt{4\sigma_k^2 S^2 + \varepsilon'^4} = 2S^2 + \varepsilon'^2 \end{aligned}$$

Lemma C.2 guarantees that the fixed point equation above possesses a unique positive solution S when $\varepsilon'^2 < \bar{\sigma} = \sum_{k=1}^K w_k \sigma_k^2$.

The second sufficient optimality condition (52) combined with the equations on G_1, F_1 (48) lead to identifying m :

$$\sum_{k=1}^K w_k G_{1,k} = 0 \Rightarrow m = \sum_{k=1}^K w_k \mu_k$$

Identifying the offset coefficients F_0, G_0 Since now m and S are known and unique, all the first and second order coefficients $F_{2,k}, G_{2,k}, F_{1,k}, G_{1,k}$ are uniquely determined. Finding $F_{0,k}$ and $G_{0,k}$ can be done up to an additive constant. Adding [iii] and [jjj] leads to a closed form expression on $F_{0,k} + G_{0,k}$. Since the optimality condition does not depend on H_0, F_0 and G_0 , one may simply set $F_{0,k}$ to 0, and solve $G_{0,k}$ exactly.

2. Case 2: if $\varepsilon'^2 \geq \bar{\sigma}$: We are going to show that there exist a set of potentials such that the Dirac at $\bar{\mu} = \sum_{k=1}^K w_k \mu_k$ verifies the optimality conditions (39). Let's simplify the optimality condition for a Dirac minimizer $\alpha_{S_\varepsilon} = \delta_{\bar{\mu}}$

$$\begin{aligned} & \left\langle \sum_{k=1}^K w_k g_k, \beta - \alpha_{S_\varepsilon} \right\rangle \geq 0 \\ & \Leftrightarrow \bar{G}_2 M_2(\beta) + \bar{G}_1 M_1(\beta) - \bar{G}_2 \bar{\mu}^2 - \bar{G}_1 \bar{\mu} \geq 0 \end{aligned}$$

However since for any measure $\beta, M_1(\beta)^2 \leq M_2(\beta)$, the following condition is sufficient for optimality:

$$(\forall x \in \mathbb{R}) \quad \bar{G}_2 x^2 + \bar{G}_1 x - \bar{G}_2 \bar{\mu}^2 - \bar{G}_1 \bar{\mu} \geq 0 \quad (54)$$

Sinkhorn equations Using the second Sinkhorn equation with $\alpha_{S_\varepsilon} = \delta_{\bar{\mu}}$ given by $e^{\frac{f_k}{\varepsilon}} \cdot \mathcal{K}(e^{\frac{g_k}{\varepsilon}} \cdot \alpha_{\text{OT}_\varepsilon}) = 1$:

$$\begin{cases} F_{2,k} - 1 = 0 & [j] \\ F_{1,k} + 2m = 0 & [jj] \\ F_{0,k} - \bar{\mu}^2 + G_{2,k} \bar{\mu}^2 + G_{1,k} \bar{\mu} + G_{0,k} = 0 & [jjj] \end{cases}$$

Using the first Sinkhorn equation $e^{\frac{g_k}{\varepsilon}} \cdot \mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}} \cdot \alpha_k) = 1$ we get by identification, for all k , with :

$$A_k = \frac{F_{2,k} - 1}{2\varepsilon'^2} - \frac{1}{2\sigma_k^2} \quad (55)$$

$$\begin{cases} \frac{G_{2,k} - 1}{\varepsilon'^2} - \frac{1}{2A_k \varepsilon'^4} = 0 & [i] \\ \frac{G_{1,k}}{\varepsilon'^2} - \frac{F_{1,k}}{2A_k \varepsilon'^4} - \frac{\mu_k}{A_k \sigma_k^2 \varepsilon'^2} = 0 & [ii] \\ \frac{G_{0,k}}{\varepsilon'^2} - \frac{\mu_k^2}{\sigma_k^2} + \frac{F_{0,k}}{\varepsilon'^2} - \frac{\left[\frac{F_{1,k} + \mu_k}{2\varepsilon'^2 + \sigma_k^2} \right]^2}{2A_k} - \log(-2\sigma_k^2 A_k) = 0 & [iii] \end{cases}$$

Combining both systems:

$$\begin{cases} F_{2,k} = 1 & [j] \\ F_{1,k} = -2\bar{\mu} & [jj] \\ G_{2,k} = 1 - \frac{\sigma_k^2}{\varepsilon'^2} & [i] \\ G_{1,k} = \frac{2\bar{\mu}\sigma_k^2}{\varepsilon'^2} - 2\mu_k & [ii] \\ \frac{G_{0,k}}{\varepsilon'^2} - \frac{\mu_k^2}{\sigma_k^2} + \frac{F_{0,k}}{\varepsilon'^2} + \sigma_k^2 \left[\frac{\bar{\mu}}{\varepsilon'^2} - \frac{\mu_k}{\sigma_k^2} \right]^2 = 0 & [iii] \end{cases} \Rightarrow \begin{cases} \bar{G}_2 = 1 - \frac{\bar{\sigma}}{\varepsilon'^2} \\ \bar{G}_1 = 2\bar{\mu} \left(\frac{\bar{\sigma}}{\varepsilon'^2} - 1 \right) \end{cases}$$

[iii] can be simplified, it leads to: $G_{0,k} + F_{0,k} + \frac{\sigma_k^2 \bar{\mu}^2}{\varepsilon'^2} - 2\bar{\mu} \mu_k = 0$. Similarly, [jjj] can be written: $F_{0,k} + G_{0,k} - \bar{\mu}^2 \frac{\sigma_k^2}{\varepsilon'^2} + \left(\frac{2\bar{\mu}\sigma_k^2}{\varepsilon'^2} - 2\mu_k \right) \bar{\mu} = 0$ which are equivalent.

Using the assumption $\varepsilon'^2 > \bar{\sigma}$, the optimality condition (54) is equivalent to:

$$\begin{aligned}
 & (\forall x \in \mathbb{R}) \quad \bar{G}_2 x^2 + \bar{G}_1 x - \bar{G}_2 \bar{\mu}^2 - \bar{G}_1 \bar{\mu} \geq 0 \\
 \Leftrightarrow & (\forall x \in \mathbb{R}) \quad \left(1 - \frac{\bar{\sigma}}{\varepsilon'^2}\right) x^2 + 2\bar{\mu} \left(\frac{\bar{\sigma}}{\varepsilon'^2} - 1\right) x - \left(1 - \frac{\bar{\sigma}}{\varepsilon'^2}\right) \bar{\mu}^2 - 2\bar{\mu} \left(\frac{\bar{\sigma}}{\varepsilon'^2} - 1\right) \bar{\mu} \geq 0 \\
 \Leftrightarrow & (\forall x \in \mathbb{R}) \quad x^2 - 2\bar{\mu}x + \bar{\mu}^2 \geq 0 \\
 \Leftrightarrow & (\forall x \in \mathbb{R}) \quad (x - \bar{\mu})^2 \geq 0
 \end{aligned}$$

Thus, the optimality condition holds.

As before, $F_{0,k}, G_{0,k}$ can be determined up to a constant using [iii]. \square

C.3. Proof of theorem 1

We showed in the supplementary section B.1 that $\text{OT}_\varepsilon^\mathcal{L}$ is convex and differentiable on the set of sub-Gaussian measures with positive densities with respect to the Lebesgue measure. Thus, the barycenter $\alpha_{\text{OT}_\varepsilon^\mathcal{L}}$ can be characterized by the first order optimality condition:

Proposition C.1. *Let $\alpha_1, \dots, \alpha_K \in \mathcal{G}(\mathbb{R}^d)$ be Lebesgue-continuous measures with positive density functions. Let (w_1, \dots, w_K) be non-negative weights summing to 1. Then:*

$\alpha_{\text{OT}_\varepsilon^\mathcal{L}} = \arg \min_\alpha \sum_{k=1}^K w_k \text{OT}_\varepsilon^\mathcal{L}(\alpha_k, \alpha)$ if and only if there exists a set of potentials $f_1, \dots, f_K, g_1, \dots, g_K$ such that for any feasible (Lebesgue-continuous) direction $\beta \in \mathcal{G}(\mathbb{R}^d)$ the following equations hold everywhere in \mathbb{R}^d , identifying the measures with their density functions:

$$\begin{cases} e^{\frac{f_k}{\varepsilon}} \mathcal{K}(e^{\frac{g_k}{\varepsilon}}) = \alpha_k, & e^{\frac{g_k}{\varepsilon}} \mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}}) = \alpha_{\text{OT}_\varepsilon^\mathcal{L}}, \\ \langle \sum_{k=1}^K w_k g_k, \beta - \alpha_{\text{OT}_\varepsilon^\mathcal{L}} \rangle \geq 0 \end{cases} \quad (56)$$

Theorem C.3 (Blurring bias of $\text{OT}_\varepsilon^\mathcal{L}$). *Let $C(x, y) = (x - y)^2$ and $0 < \varepsilon < +\infty$ and $\varepsilon = 2\varepsilon'^2$. Let (w_k) be positive weights that sum to 1. Let \mathcal{N} denote the Gaussian distribution. Assume $\alpha_k \sim \mathcal{N}(\mu_k, \sigma_k^2)$ and let $\bar{\mu} = \sum_k w_k \mu_k$,*

then $\alpha_{\text{OT}_\varepsilon^\mathcal{L}} \sim \mathcal{N}(\bar{\mu}, S^2)$ where S is the unique solution of the fixed point equation:

$$\sum_{k=1}^K w_k \sqrt{\varepsilon'^4 + 4\sigma_k^2 S^2} = -\varepsilon'^2 + 2S^2.$$

In particular, if all σ_k are equal to some $\sigma > 0$, then

$$\text{then } \alpha_{\text{OT}_\varepsilon^\mathcal{L}} \sim \mathcal{N}(\bar{\mu}, \sigma^2 + \varepsilon'^2).$$

PROOF. The proof of this theorem is technically identical to that of 3. Except that the Sinkhorn-equations are slightly different.

Consider polynomial potentials of the form $f_k(x) = F_{2,k}x^2 + F_{1,k}x + F_{0,k}$ and $g_k(x) = G_{2,k}x^2 + G_{1,k}x + G_{0,k}$ for some unknown coefficients $F_{2,k}, F_{1,k}, F_{0,k}, G_{2,k}, G_{1,k}, G_{0,k} \in \mathbb{R}$, and assume that $\frac{d\alpha_\mathcal{L}}{d\lambda} = \mathcal{N}(m, S)$. First, we will write the first and second order coefficients as functions of m and S then use the optimality condition to find m and S .

Sufficient optimality condition Let a continuous measure $\beta \in \mathcal{G}(\mathbb{R}^d)$ identified with a positive density function. And let $M_r(\beta)$ denote the r -th moment of β . For any real sequence y_1, \dots, y_K , let \bar{y} denote its weighted average $\sum_{k=1}^K w_k y_k$. The optimality condition reads:

$$\begin{aligned}
 & \left\langle \sum_{k=1}^K w_k g_k, \beta - \alpha_\mathcal{L} \right\rangle \geq 0 \\
 \Leftrightarrow & \bar{G}_2(M_2(\beta) - M_2(\alpha_\mathcal{L})) + \bar{G}_1(M_1(\beta) - M_1(\alpha_\mathcal{L})) + \bar{G}_0(M_0(\beta) - M_0(\alpha_\mathcal{L})) \geq 0 \\
 \Leftrightarrow & \bar{G}_2(M_2(\beta) - M_2(\alpha_\mathcal{L})) + \bar{G}_1(M_1(\beta) - M_1(\alpha_\mathcal{L})) \geq 0
 \end{aligned}$$

Where the last inequality follows from $M_0(\beta) = M_0(\alpha_{\mathcal{L}}) = \int d\alpha_{\mathcal{L}} = 1$. Thus, the 0-order coefficients are irrelevant for optimality. We are going to show that there exist a set of coefficients such that the following sufficient conditions for optimality hold:

$$\bar{G}_2 \stackrel{\text{def}}{=} \sum_{k=1}^K w_k G_{2,k} \quad (57)$$

$$\bar{G}_1 \stackrel{\text{def}}{=} \sum_{k=1}^K w_k G_{1,k} \quad (58)$$

Kernel integration Dropping the k exponent for the sake of convenience, let's carefully derive the integral $\mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}})$:

$$\begin{aligned} \mathcal{K}^\top e^{\frac{f}{\varepsilon}} &= \int K(x, y) e^{\frac{f(y)}{2\varepsilon'^2}} dy \\ &= \int \exp\left(\frac{-(x-y)^2 + f(y)}{2\varepsilon'^2}\right) dy \\ &= \int \exp\left(\underbrace{\left[\frac{F_2 - 1}{2\varepsilon'^2}\right]}_A y^2 + \underbrace{\left[\frac{F_1}{2\varepsilon'^2} + \frac{x}{\varepsilon'^2}\right]}_{Z(x)} y + \left[\frac{F_0 - x^2}{2\varepsilon'^2}\right]\right) dy \\ &= \exp\left(\frac{F_0 - x^2}{2\varepsilon'^2}\right) \int \exp\left(A \left[y^2 + \frac{Z(x)}{A} y\right]\right) dy \\ &= \exp\left(\frac{F_0 - x^2}{2\varepsilon'^2}\right) \int \exp\left(A \left[y + \frac{Z(x)}{2A}\right]^2 - \frac{Z(x)^2}{4A}\right) dy \\ &= \exp\left(\frac{F_0 - x^2}{2\varepsilon'^2} - \frac{Z(x)^2}{4A}\right) \underbrace{\int \exp\left(A \left[y + \frac{Z(x)}{2A}\right]^2\right) dy}_I \end{aligned}$$

For the fourth equality to be sound, we need $A \neq 0$, and for the integral I to be finite, we need $A < 0$ which is equivalent to:

$$F_2 < 1 \quad (59)$$

In that case, $I = \sqrt{-\frac{\pi}{A}}$ thus:

$$\begin{aligned} \mathcal{K}^\top(e^{\frac{f}{\varepsilon}})(x) &= \sqrt{-\frac{\pi}{A}} \exp\left(\frac{F_0 - x^2}{2\varepsilon'^2} - \frac{Z(x)^2}{4A}\right) \\ &= \sqrt{-\frac{\pi}{A}} \exp\left(\left[-\frac{1}{2\varepsilon'^2} - \frac{1}{4A\varepsilon'^4}\right] x^2 - \left[\frac{F_1}{4A\varepsilon'^4}\right] x + \frac{F_0}{2\varepsilon'^2} - \frac{\left[\frac{F_1}{2\varepsilon'^2}\right]^2}{4A}\right) \\ &= \exp\left(\left[-\frac{1}{2\varepsilon'^2} - \frac{1}{4A\varepsilon'^4}\right] x^2 - \left[\frac{F_1}{4A\varepsilon'^4}\right] x + \frac{F_0}{2\varepsilon'^2} - \frac{\left[\frac{F_1}{2\varepsilon'^2}\right]^2}{4A} - \log\left(\sqrt{-\frac{A}{\pi}}\right)\right) \end{aligned}$$

Sinkhorn equations Using the first Sinkhorn equation $e^{\frac{g_k}{\varepsilon}} \mathcal{K}^\top(e^{\frac{f_k}{\varepsilon}}) = \alpha_{\mathcal{L}}$ we get by identification, for all k , with :

$$A_k = \frac{F_{2,k} - 1}{2\varepsilon'^2} \quad (60)$$

$$\begin{cases} \frac{G_{2,k-1}}{\varepsilon'^2} - \frac{1}{2A_k \varepsilon'^4} + \frac{1}{S^2} = 0 & [i] \\ \frac{G_{1,k}}{\varepsilon'^2} - \frac{F_{1,k}}{2A_k \varepsilon'^4} - \frac{2m}{S^2} = 0 & [ii] \\ \frac{G_{0,k}}{\varepsilon'^2} + \frac{m^2}{S^2} + \frac{F_{0,k}}{\varepsilon'^2} - \frac{\left[\frac{F_{1,k}}{2\varepsilon'^2}\right]^2}{2A_k} - \log\left(-\frac{A_k}{\pi}\right) + \log(2\pi S^2) = 0 & [iii] \end{cases}$$

Similarly, since the equations are symmetric, $e^{\frac{f_k}{\varepsilon}} \mathcal{K} e^{\frac{g_k}{\varepsilon}} = \alpha_k$) with $G_{2,k} < 1$ and:

$$B_k = \frac{G_{2,k} - 1}{2\varepsilon'^2} \quad (61)$$

lead to:

$$\begin{cases} \frac{F_{2,k}-1}{\varepsilon'^2} - \frac{1}{2B_k\varepsilon'^4} + \frac{1}{\sigma_k^2} = 0 & [i] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{2\mu_k}{\sigma_k^2} = 0 & [ii] \\ \frac{F_{0,k}}{\varepsilon'^2} + \frac{\mu_k^2}{\sigma_k^2} + \frac{G_{0,k}}{\varepsilon'^2} - \frac{[\frac{G_{1,k}}{2\varepsilon'^2}]^2}{2B_k} - \log(-\frac{B_k}{\pi}) + \log(2\pi\sigma_k^2) = 0 & [iii] \end{cases}$$

Second order coefficients G_2, F_2 Let's rewrite [i] and [j] separately:

$$\begin{aligned} & \begin{cases} 2B_k + \frac{1}{S^2} - \frac{1}{2A_k\varepsilon'^4} = 0 \\ 2A_k + \frac{1}{\sigma_k^2} - \frac{1}{2B_k\varepsilon'^4} = 0 \end{cases} \\ \Leftrightarrow & \begin{cases} A_k B_k + \frac{A_k}{2S^2} - \frac{1}{4\varepsilon'^4} = 0 \\ A_k B_k + \frac{B_k}{2\sigma_k^2} - \frac{1}{4\varepsilon'^4} = 0 \end{cases} \\ \Leftrightarrow & \begin{cases} \frac{B_k}{\sigma_k^2} = \frac{A_k}{S^2} \\ A_k B_k + \frac{B_k}{2\sigma_k^2} - \frac{1}{4\varepsilon'^4} = 0 \end{cases} \\ \Leftrightarrow & \begin{cases} \frac{B_k}{\sigma_k^2} = \frac{A_k}{S^2} \\ B_k^2 + \frac{B_k}{2S^2} - \frac{\sigma_k^2}{4\varepsilon'^4 S^2} = 0 \end{cases} \end{aligned}$$

The roots of the polynomial above are: $-\frac{1}{4S^2} \pm \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}}$. The constraint $B_k < 0$ eliminates the positive solution and it holds:

$$B_k = -\frac{1}{4S^2} - \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}} \quad (62)$$

$$A_k = \frac{S^2}{\sigma_k^2} B_k \quad (63)$$

First order coefficients G_1, F_1 Let's rewrite [ii] and [jj] separately:

$$\begin{aligned} & \begin{cases} \frac{G_{1,k}}{\varepsilon'^2} - \frac{F_{1,k}}{2A_k\varepsilon'^4} - \frac{2m}{S^2} = 0 & [ii] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{2\mu_k}{\sigma_k^2} = 0 & [jj] \end{cases} \\ \Leftrightarrow & \begin{cases} 2A_k G_{1,k} - \frac{F_{1,k}}{\varepsilon'^2} - \frac{4A_k\varepsilon'^2 m}{S^2} = 0 & [ii] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{2\mu_k}{\sigma_k^2} = 0 & [jj] \end{cases} \\ \Leftrightarrow & \begin{cases} \left(2A_k - \frac{1}{2B_k\varepsilon'^4}\right) G_{1,k} - \frac{4A_k\varepsilon'^2 m}{S^2} - \frac{2\mu_k}{\sigma_k^2} = 0 & [ii] + [jj] \\ \frac{F_{1,k}}{\varepsilon'^2} - \frac{G_{1,k}}{2B_k\varepsilon'^4} - \frac{2\mu_k}{\sigma_k^2} = 0 & [jj] \end{cases} \end{aligned}$$

The equations above between A_k and B_k lead to $2A_k - \frac{1}{2B_k\varepsilon'^4} = -\frac{1}{\sigma_k^2}$ and $A_k = \frac{S^2}{\sigma_k^2} B_k$ leads to:

$$\begin{aligned} & \frac{1}{\sigma_k^2} G_{1,k} + \frac{2\mu_k}{\sigma_k^2} + \frac{4B_k\varepsilon'^2 m}{\sigma_k^2} = 0 \\ \Rightarrow & G_{1,k} + 2\mu_k + 2m(G_{2,k} - 1) = 0 \end{aligned}$$

Using [jj] we recover the first order coefficients $F_{1,k}$ and $G_{1,k}$ as function of m :

$$\begin{aligned} & G_{1,k} + 2\mu_k + 2m(G_{2,k} - 1) = 0 \\ & F_{1,k} + 2m + 2\mu_k(F_{2,k} - 1) = 0 \end{aligned} \quad (64)$$

Optimality condition and identifying σ and μ Using the definition of B_2 (61) and then with their closed form formulas (62), the first sufficient optimality condition (57) reads:

$$\begin{aligned}
 \sum_{k=1}^K w_k G_{2,k} = 0 &\Rightarrow \sum_{k=1}^K w_k B_{2,k} = -\frac{1}{2\varepsilon'^2} \\
 &\Rightarrow \sum_{k=1}^K w_k \left(\frac{1}{4S^2} + \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}} \right) = \frac{1}{2\varepsilon^2} \\
 &\Rightarrow \sum_{k=1}^K w_k \sqrt{\frac{1}{16S^4} + \frac{\sigma_k^2}{4S^2\varepsilon'^4}} = \frac{1}{2\varepsilon^2} - \frac{1}{4S^2} \\
 &\Rightarrow \sum_{k=1}^K w_k \sqrt{4\sigma_k^2 S^2 + \varepsilon'^4} = 2S^2 - \varepsilon'^2
 \end{aligned}$$

Lemma C.1 guarantees that the fixed point equation above possesses a unique positive solution S .

The second sufficient optimality condition (57) combined with the equations on G_1, F_1 (64) lead to identifying m :

$$\sum_{k=1}^K w_k G_{1,k} = 0 \Rightarrow m = \sum_{k=1}^K w_k \mu_k$$

Identifying the offset coefficients F_0, G_0 Since now m and S are known and unique, all the first and second order coefficients $F_{2,k}, G_{2,k}, F_{1,k}, G_{1,k}$ are uniquely determined. Finding $F_{0,k}$ and $G_{0,k}$ can be done up to an additive constant. Adding [iii] and [jjj] leads to a closed form expression on $F_{0,k} + G_{0,k}$. Since the optimality condition does not depend on F_0 and G_0 , one may simply set $F_{0,k}$ to 0, and solve $G_{0,k}$ exactly. \square

D. The IBP algorithm

Computing the OT barycenter with the divergence $\text{OT}_\varepsilon^{\mathcal{U}}$ can be shown to be equivalent to the KL projection problem:

$$\min_{\substack{\pi_1, \dots, \pi_K \\ \pi_k \in \mathcal{C}_k \cap \mathcal{C}'}} \sum_{k=1}^K w_k \widetilde{\text{KL}}(\pi^k | \mathbf{K}) , \quad (\text{IBP})$$

where $\mathcal{C}_k = \{\pi \in \mathbb{R}_+^{n \times n} | \pi \mathbf{1} = \alpha_k\}$ and $\mathcal{C}' = \{\pi \in \mathbb{R}_+^{n \times n} | \exists \alpha \in \Delta_n, \pi_k^\top \mathbf{1} = \alpha, \forall k = 1 \dots K\}$. The IBP algorithm amounts to performing iterative minimization on one constraint set \mathcal{C} at a time. Each step can be solved in closed form, leading to Sinkhorn-like iterations. By combining both iterations, one can write every iterate of the transport plan as $\pi^{(l)} = \text{diag}(\mathbf{a}^{(l)}) \mathbf{K} \text{diag}(\mathbf{b}^{(l)})$ and perform the scaling operations on the variables \mathbf{a}, \mathbf{b} given in algorithm 2.

Algorithm 2 IBP algorithm (Benamou et al., 2014; Chizat et al., 2017)

Input: $\alpha_1, \dots, \alpha_K, \mathbf{K} = e^{-\frac{c}{\varepsilon}}$

Output: $\alpha_{\text{OT}_\varepsilon^{\mathcal{U}}}$

Initialize all scalings (b_k) to $\mathbf{1}$,

repeat

for $k = 1$ **to** K **do**

$$a_k \leftarrow \left(\frac{\alpha_k}{\mathbf{K} b_k} \right)$$

end for

$$\alpha \leftarrow \prod_{k=1}^K (\mathbf{K}^\top a_k)^{w_k}$$

for $k = 1$ **to** K **do**

$$b_k \leftarrow \left(\frac{\alpha}{\mathbf{K}^\top a_k} \right)$$

end for

until convergence

E. Additional Proofs

Proof of proposition 5

Proposition E.1. *Let $\alpha_1, \dots, \alpha_K \in \Delta_n$ and $\mathbf{K} = e^{-\frac{\mathbf{C}}{\varepsilon}}$. Let π denote a sequence π_1, \dots, π_K of transport plans in $\mathbb{R}_+^{n \times n}$ and the constraint sets $\mathcal{H}_1 = \{\pi | \forall k, \pi_k \mathbf{1} = \alpha_k\}$, and $\mathcal{H}_2 = \{\pi | \forall k \forall k', \pi_k^\top \mathbf{1} = \pi_{k'} \mathbf{1}\}$. The barycenter problem $\min_{\alpha \in \Delta_n} \sum_{k=1}^K w_k S_\varepsilon(\alpha_k, \alpha)$ is equivalent to:*

$$\min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2 \\ d \in \mathbb{R}_+^n}} \left[\varepsilon \sum_{k=1}^K w_k \text{KL}(\pi_k | \mathbf{K} \text{diag}(d)) + \frac{\varepsilon}{2} \langle d - \mathbf{1}, \mathbf{K}(d - \mathbf{1}) \rangle \right]. \quad (65)$$

PROOF. The barycenter problem of S_ε only depends on $\text{OT}_\varepsilon^\mathcal{U}(\alpha, \beta) - \frac{1}{2}(\text{OT}_\varepsilon(\alpha, \alpha))$. Let's rewrite this expression using the IBP formulation and duality. the IBP formulation (5) is explicitly given by:

$$\text{OT}_\varepsilon^\mathcal{U}(\alpha, \beta) = \min_{\substack{\pi \in \mathbb{R}_+^{n \times n} \\ \pi \mathbf{1} = \alpha, \pi^\top \mathbf{1} = \beta}} \varepsilon \widetilde{\text{KL}}(\pi | \mathbf{K}) - \varepsilon \sum_{i,j} \mathbf{K}_{ij} \quad (66)$$

And the autocorrelation term can be expressed via its dual problem:

$$\text{OT}_\varepsilon^\mathcal{U}(\alpha, \alpha) = \max_{h \in \mathbb{R}^n} 2 \langle h, \alpha \rangle - \varepsilon \langle e^{\frac{h}{\varepsilon}}, \mathbf{K} e^{\frac{h}{\varepsilon}} \rangle - \varepsilon \sum_{i,j} \mathbf{K}_{ij} \quad (67)$$

$$= \max_{d \in \mathbb{R}_+^n} 2 \langle \varepsilon \log(d), \alpha \rangle - \varepsilon \langle d, \mathbf{K}d \rangle - \varepsilon \sum_{i,j} \mathbf{K}_{ij} \quad (68)$$

$$= - \min_{d \in \mathbb{R}_+^n} -2 \langle \varepsilon \log(d), \alpha \rangle + \varepsilon \langle d, \mathbf{K}d \rangle + \varepsilon \sum_{i,j} \mathbf{K}_{ij} \quad (69)$$

Moreover, on the constraint set $\mathcal{H}_1 \cap \mathcal{H}_2$, it holds $\alpha = \pi_k^\top \mathbf{1}$ for all k . Thus, denoting $\mathcal{H}_2(\alpha) = \{\pi | \forall k \forall k', \pi_k^\top \mathbf{1} = \alpha\}$ the following can be written:

$$\begin{aligned}
 & \arg \min_{\alpha \in \Delta_n} \sum_{k=1}^K w_k S_\varepsilon(\alpha_k, \alpha) \\
 &= \arg \min_{\alpha \in \Delta_n} \min_{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2(\alpha)} \sum_{k=1}^K w_k \varepsilon \widetilde{\text{KL}}(\pi_k | \mathbf{K}) + \min_{d \in \mathbb{R}_+^n} -\langle \varepsilon \log(d), \alpha \rangle + \frac{1}{2} \varepsilon \langle d, \mathbf{K}d \rangle - \frac{1}{2} \varepsilon \sum_{i,j} \mathbf{K}_{ij} \\
 &= \arg \min_{\alpha \in \Delta_n} \min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2(\alpha) \\ d \in \mathbb{R}_+^n}} \sum_{k=1}^K w_k \left(\varepsilon \widetilde{\text{KL}}(\pi_k | \mathbf{K}) - \langle \varepsilon \log(d), \alpha \rangle \right) + \frac{1}{2} \varepsilon \langle d, \mathbf{K}d \rangle - \frac{1}{2} \varepsilon \sum_{i,j} \mathbf{K}_{ij} \\
 &= \arg \min_{\alpha \in \Delta_n} \min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2(\alpha) \\ d \in \mathbb{R}_+^n}} \sum_{k=1}^K w_k \left(\varepsilon \widetilde{\text{KL}}(\pi_k | \mathbf{K}) - \langle \varepsilon \log(d), \pi_k^\top \mathbf{1} \rangle \right) + \frac{1}{2} \varepsilon \langle d, \mathbf{K}d \rangle - \frac{1}{2} \varepsilon \sum_{i,j} \mathbf{K}_{ij} \\
 &= \min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2 \\ d \in \mathbb{R}_+^n}} \sum_{k=1}^K w_k \left(\varepsilon \widetilde{\text{KL}}(\pi_k | \mathbf{K}) - \langle \varepsilon \log(d), \pi_k^\top \mathbf{1} \rangle \right) + \frac{1}{2} \varepsilon \langle d, \mathbf{K}d \rangle - \frac{1}{2} \varepsilon \sum_{i,j} \mathbf{K}_{ij} \\
 &= \min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2 \\ d \in \mathbb{R}_+^n}} \sum_{k=1}^K w_k \left(\varepsilon \widetilde{\text{KL}}(\pi_k | \mathbf{K} \text{diag}(d)) - \varepsilon \langle \mathbf{K}d, \mathbf{1} \rangle + \varepsilon \sum_{ij} \mathbf{K}_{ij} \right) + \frac{1}{2} \varepsilon \langle d, \mathbf{K}d \rangle - \frac{1}{2} \varepsilon \sum_{i,j} \mathbf{K}_{ij} \\
 &= \min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2 \\ d \in \mathbb{R}_+^n}} \sum_{k=1}^K \varepsilon w_k \widetilde{\text{KL}}(\pi_k | \mathbf{K} \text{diag}(d)) - \varepsilon \langle \mathbf{K}d, \mathbf{1} \rangle + \frac{1}{2} \varepsilon \langle d, \mathbf{K}d \rangle + \frac{1}{2} \varepsilon \sum_{ij} \mathbf{K}_{ij} \\
 &= \min_{\substack{\pi \in \mathcal{H}_1 \cap \mathcal{H}_2 \\ d \in \mathbb{R}_+^n}} \sum_{k=1}^K \varepsilon w_k \text{KL}(\pi_k | \mathbf{K} \text{diag}(d)) + \frac{\varepsilon}{2} \langle d - \mathbf{1}, \mathbf{K}(d - \mathbf{1}) \rangle .
 \end{aligned}$$

□

F. Supplementary details on experiments

F.1. Barycenters of nested ellipses

We simulate each ellipse by generating random major and minor radii with a moving a center from the top left quarter corner to the bottom right quarter corner. The box constraints of the random generators of the radii are manually picked so that ellipses are more likely to be nested with an asymmetric surrounding ellipse (see supplementary code). The full list of 10 images used to compute the barycenters is displayed in Figure 10. Each image has 60×60 pixels. The ground OT cost function is the squared Euclidean cost over the unit square $[0, 1]^2$. For entropy regularized distances (All except W), we set ε to the lowest value guaranteeing no numerical instabilities in Sinkhorn's algorithm (this was particularly an issue for *Sharp barycenters* α_{A_ε} of Luise et al. (2018)). Now we detail the algorithm used for each divergence F defining each barycenter α_F of the experiment in Figure 6:

1. $\text{OT}_\varepsilon^{\mathcal{U}}$: OT with the uniform measure; computed using the IBP algorithm (Algorithm 2).
2. S_ε : Proposed debiased divergence; computed using the proposed algorithm (Algorithm 1).
3. $\text{OT}_\varepsilon^{\otimes}$: Computed using iterative IBP in minimization-majorization alternative algorithm. With (7), one can linearize the concave negative KL penalty and solve the resulting problem using IBP iteratively and then update the KL term etc. This leads to a series of nested IBP loops.
4. A_ε : Sharp barycenters introduced by Luise et al. (2018). Solved using accelerated gradient descent. This method required considerable manual effort to tune the learning rate in order to get an acceptable barycenter and was more prone to numerical instabilities.
5. Free support barycenters with S_ε : introduced by Luise et al. (2019), we used the online Python code provided by the authors which amounts to add or remove a dirac particle at each iteration and update their weights using Frank-Wolf's algorithm. The algorithm is stopped when no particles are created / removed.
6. W : non regularized Wasserstein distance. We used the accelerated interior point methods introduced by Ge et al. (2019)

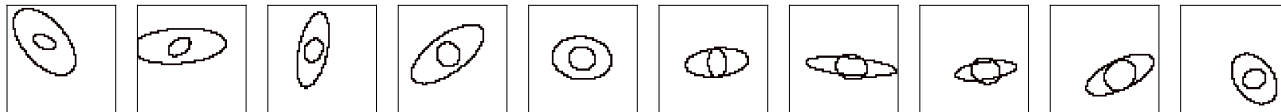


Figure 10. All 10 nested ellipses images used to compute the barycenters of Figure 6.



Figure 11. Input meshes used to compute the barycenters of 3D meshes.

with the online matlab implementation provided by the authors.

F.2. Barycenters of 3D shapes

The original 3D shapes (tore and rabbit) are taken from the `PyVista` (Sullivan & Kaszynski, 2019) Python library. We preprocess the original meshes as follows. Each mesh is smoothed by 100 iterations of a Laplacian operator then the coordinates are centered and rescaled to fit within 95% of the cube $(-1, 1)^3$. We sample 3D histograms of both meshes on a uniform 3D grid of size 200^3 . Both histograms are normalized and regularized by adding a 10^{-10} weight to avoid numerical errors. We set the lowest stable regularization $\varepsilon = 0.01$ for the ground cost defined as the squared Euclidean distance over the $(-1, 1)^3$ cube. We compute weighted barycenters with the IBP algorithm 2 and the proposed debiased Sinkhorn barycenter algorithm 1. For each method, we use the weights $(w, 1 - w)$ for $w \in [0, 0.25, 0.5, 0.75, 1.]$. The original meshes are shown in Figure 11.

F.3. OT barycentric embeddings

We use the Python library `Torchvision` that provides a fetch method to download the MNIST dataset. We first filter the data by keeping the labels (0, 1, 2, 3, 4) then select the first 500 samples. This constitutes the global dataset of the experiment. Then we randomly select $K = 50$ samples that will be considered as our learning dictionary \mathcal{A} . Then for each sample (image) β in the remaining 450 samples, we compute the weights $w \in \Sigma_K$ minimizing $\|\alpha_F(w) - \beta\|^2$ where $\alpha_F(w)$ is the weighted barycenter of the dictionary \mathcal{A} . This leads to an embedding of 450 MNIST samples in a space of dimension K . We then use this embedding to train a Random Forest Classifier with 100 estimators using `scikit-learn`'s default parameters (version 0.21.3).