

Information Geometry for Regularized Optimal Transport and Barycenters of Patterns

Shun-ichi Amari · Ryo Karakida ·
Masafumi Oizumi · Marco Cuturi

Abstract We propose a new divergence on the manifold of probability distributions, building upon the entropic regularization of optimal transportation problems. As shown in (Cuturi, 2013), regularizing the optimal transport problem with an entropic term is known to bring several computational benefits. However, because of that regularization, the resulting quantities do not define a proper distance or divergence between probability distributions. We have recently tried to introduce a family of divergences connecting the Wasserstein distance and the KL divergence from the information geometry point of view (see Amari et al. (2018)). However, that proposal was not able to retain key intuitive aspects of the Wasserstein geometry, such as translation invariances, which play a key role when used in the more general barycenter problem. The divergence we propose in this work is able to retain such properties and admits an intuitive interpretation.

Keywords Wasserstein distance · Kullback-Leibler divergence · Optimal transportation · Barycenter · Shape preservation

1 Introduction

Two major geometrical structures have been introduced on the probability simplex, the manifold of discrete probability distributions. The first one is

Shun-ichi Amari
2-1 Hirosawa, Wako-shi, Saitama, 351-0198, Japan
E-mail: amari@brain.riken.jp

Ryo Karakida
2-3-26 Aomi, Koto-ku, Tokyo, 135-0064, Japan

Masafumi Oizumi
2-8-10 Toranomom, Minato-ku, Tokyo, 105-0001, Japan

Marco Cuturi
CREST, ENSAE, Université Paris-Saclay

based on the principle of parameterization invariance, which requires that the geometry between probability distributions must be invariant under invertible transformations of random variables. That viewpoint is the cornerstone of the theory of information geometry (Amari, 2016), which acts as a foundation for statistical inference. The second direction is grounded on the theory of optimal transport, which exploits prior geometric knowledge on the space in which random variables are valued (Villani, 2003). Computing optimal transport amounts to obtaining a coupling between these two random variables that is optimal in the sense that it has a minimal expected metric cost between the first and second variables. However, computing that solution can be challenging, and is usually carried out by solving a linear program. Cuturi (2013) considered a relaxed formulation of optimal transport, in which the negative-entropy of the coupling is used as a regularizer. We call that approximation of the original optimal transport cost the C function. Entropic regularization provides two major advantages: the regularized optimal transport problem is usually easier and faster to compute than the solution of the linear program, and can be done using Sinkhorn’s algorithm (1964); unlike the original optimal transport geometry, regularized transport distances are differentiable functions of their input, a property which can be exploited in problems arising from pattern classification and clustering (Cuturi and Avis, 2014; Cuturi and Peyré, 2016) as well as more advanced inference tasks that use the C function as an output loss (Frogner et al., 2015; Genevay et al., 2018), a model fitting loss (Rolet et al., 2016) or a way to learn mappings (Courty et al., 2017).

The C function suffers, however, from a few issues. It is neither a distance nor a divergence, notably because comparing a probability measure with itself does not result in a null discrepancy, namely if \mathbf{p} belongs to the simplex, then $C(\mathbf{p}, \mathbf{p}) \neq 0$. More worryingly, the minimizer of $C(\mathbf{p}, \mathbf{q})$ with respect to \mathbf{q} is not reached at $\mathbf{q} = \mathbf{p}$. To solve these issues, we have proposed a first attempt at unifying the information and optimal transport geometrical structures in (Amari et al., 2018). However, the information-geometric divergence introduced in that previous work loses some of the nice properties inherent to the C -function. For example, the C -function can be used to extract a common shape as the barycenter of several patterns (Cuturi and Doucet, 2014), which our former proposal was not able to. Therefore, it is desirable to define a new divergence from C , in the rigorous sense that it is minimized when comparing a measure with itself, and, preferably convex in both arguments, while still retaining the attractive properties of optimal transport.

We propose in this paper such a new divergence between probability distributions \mathbf{p} and \mathbf{q} that is both inspired by optimal transport while incorporating elements of information geometry. Its basic ingredient remains the entropic regularization of optimal transport. We show that the barycenters obtained with that new divergence are more sharply defined than those obtained with the original C -function, still keeping the shape-location decomposition property. We illustrate these new definitions with simple numerical illustrations.

2 C-Function: Entropy-regularized Optimal Transportation Plan

The general transportation problem is concerned with the optimal transportation of commodities, initially distributed according to a distribution \mathbf{p} , so that they end up being distributed as another distribution \mathbf{q} . In its full generality, such a transport can be carried out on a metric manifold, and both \mathbf{p} and \mathbf{q} be continuous measures on that manifold. We consider in this paper the discrete case, where that metric space is of finite size n , namely $X = \{1, 2, \dots, n\}$. We normalize the total amount of commodities such that it sums to 1: \mathbf{p} and \mathbf{q} are therefore probability vectors of size n in the $n - 1$ dimensional simplex,

$$S_{n-1} = \left\{ \mathbf{p} \in \mathbb{R}^n \mid \sum_i p_i = 1, \quad p_i \geq 0 \right\}. \quad (1)$$

We leave out extensions to the continuous case for future work. Let M_{ij} be the cost of transporting a unit of commodity from bin i to bin j , usually defined as a distance (or a suitable power thereof) between points i and j in X . In what follows we only assume that $M_{ij} > 0$ for $i \neq j$ and $M_{ii} = 0$. A transportation plan $\mathbf{P} = (P_{ij}) \in \mathbb{R}_+^{n \times n}$ is a joint (probability) distribution over $X \times X$ which describes at each entry P_{ij} the amount of commodities sent from bin i to bin j . Given a source distribution \mathbf{p} and a target distribution \mathbf{q} , the set of transport plans allowing a transfer from \mathbf{p} to \mathbf{q} is defined as

$$U(\mathbf{p}, \mathbf{q}) = \left\{ \mathbf{P} \in \mathbb{R}_+^{n \times n} : \forall i \leq n, \sum_j P_{ij} = p_i, \quad \forall j \leq n, \sum_i P_{ij} = q_j \right\}. \quad (2)$$

Note that if a matrix is in $U(\mathbf{p}, \mathbf{q})$ then its transpose is in $U(\mathbf{q}, \mathbf{p})$.

The transportation cost of a plan \mathbf{P} is defined as the dot product of \mathbf{P} with the cost matrix $\mathbf{M} = (M_{ij})$,

$$\langle \mathbf{M}, \mathbf{P} \rangle = \sum_{ij} M_{ij} P_{ij}. \quad (3)$$

Its minimum among all feasible plans

$$W(\mathbf{p}, \mathbf{q}) = \min_{\mathbf{P} \in U(\mathbf{p}, \mathbf{q})} \langle \mathbf{M}, \mathbf{P} \rangle \quad (4)$$

is the Wasserstein distance on S_{n-1} parameterized by the ground metric M . Cuturi (2013) studied the regularization of that problem using entropy

$$H(\mathbf{P}) = - \sum_{ij} P_{ij} \log P_{ij}, \quad (5)$$

to consider the problem of minimizing

$$\mathcal{L}(\mathbf{P}) = \langle \mathbf{M}, \mathbf{P} \rangle - \lambda H(\mathbf{P}), \quad (6)$$

where $\lambda > 0$ is a regularization strength. We call its minimum the C -function:

$$C_\lambda(\mathbf{p}, \mathbf{q}) = \min_{\mathbf{P} \in U(\mathbf{p}, \mathbf{q})} \mathcal{L}(\mathbf{P}). \quad (7)$$

The C_λ function is a useful proxy for the Wasserstein distance, with favorable computational properties, and has appeared in several applications as a very useful alternative to information-geometric divergences such as the KL divergence.

The optimal transportation plan is given in the following Theorem (Cuturi and Peyré, 2016; Amari et al., 2018).

Theorem 1 The optimal transportation plan \mathbf{P}_λ^* is given by

$$\mathbf{P}_\lambda^* = [ca_i b_j K_{ij}]_{ij}, \quad (8)$$

$$\mathbf{K} = \left[\exp\left(-\frac{M_{ij}}{\lambda}\right) \right]_{ij}, \quad (9)$$

where c , a normalization constant, and vectors $\mathbf{a}, \mathbf{b} \in S_{n-1}$ are determined from \mathbf{p} and \mathbf{q} such that the sender and receiver conditions (namely marginal conditions) are satisfied.

In our previous paper (Amari et al., 2018) we studied the information geometry of the manifold of optimal transportation plans. We proposed a family of divergences that combine the KL divergence and the Wasserstein distance. However, these divergences are closer in spirit to the KL divergence, and lose therefore some crucial properties of the C -function. We define in this work a new family of divergences directly from the C -function.

3 Divergence Derived from C -Function

The C -function $C_\lambda(\mathbf{p}, \mathbf{q})$ does not satisfy the requirements for a distance or divergence, which would be that for such a function Δ we have that for any $\mathbf{p}, \mathbf{q} \in S_{n-1}$

$$\Delta(\mathbf{p}, \mathbf{q}) \geq \Delta(\mathbf{p}, \mathbf{p}) = 0. \quad (10)$$

In order to find the minimizer \mathbf{q}^* of $C_\lambda(\mathbf{p}, \mathbf{q})$ for a given \mathbf{p} , we use the exponentiated version \mathbf{K}_λ of the cost \mathbf{M} depending on λ given in (9).

We further define its conditional version,

$$\tilde{\mathbf{K}}_\lambda = \left[\frac{K_{ji, \lambda}}{K_{j\cdot}} \right]_{ij}, \quad (11)$$

$$K_{j\cdot} = \sum_i K_{ji}. \quad (12)$$

$\tilde{\mathbf{K}}_\lambda$ is a linear operator from S_{n-1} into S_{n-1} , and we will use for convenience the notation

$$\tilde{\mathbf{q}} = \tilde{\mathbf{K}}_\lambda \mathbf{q}. \quad (13)$$

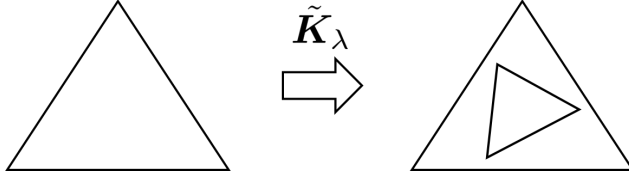


Fig. 1 Shrinkage operator \tilde{K}_λ .

\tilde{K}_λ is a monotonic shrinkage operator mapping S_{n-1} in its interior (see Fig. 1), and $\tilde{K}_\lambda S_{n-1} \subset \tilde{K}_{\lambda'} S_{n-1}$ for $\lambda > \lambda'$. When $\lambda = 0$, \tilde{K}_λ is the identity mapping

$$\tilde{K}_0 \mathbf{q} = \mathbf{q}. \quad (14)$$

As λ tends to infinity, $\tilde{K}_\lambda S_{n-1}$ converges to a single point, the center of S_{n-1} , $\mathbf{1}/n$, where $\mathbf{1} = (1, 1, \dots, 1)^T$, and

$$\tilde{K}_\infty = \frac{1}{n} [\mathbf{1} \cdots \mathbf{1}]. \quad (15)$$

Hence, for any \mathbf{q} , $\tilde{K}_\infty \mathbf{q}$ is the uniform distribution $\mathbf{1}/n$.

Theorem 2 The minimizer of $C_\lambda(\mathbf{p}, \mathbf{q})$ with respect to \mathbf{q} , given \mathbf{p} , is

$$\mathbf{q}^* = \tilde{K}_\lambda \mathbf{p}. \quad (16)$$

Proof By differentiation, we have the equation

$$\partial_{\mathbf{q}} C_\lambda(\mathbf{p}, \mathbf{q}^*) = 0 \quad (17)$$

to determine the minimizer \mathbf{q}^* . This gives the condition

$$b_i = 1, \quad (18)$$

for \mathbf{q}^* (see the duality theorem in (Amari et al., 2018)). Hence, the optimal transportation plan from \mathbf{p} to \mathbf{q}^* is given by

$$P_{ij}^* = ca_i K_{ij}, \quad (19)$$

where suffix λ is omitted to alleviate notations. From

$$\sum_j ca_i K_{ij} = p_i, \quad (20)$$

we have

$$ca_i = \frac{p_i}{K_i}. \quad (21)$$

So

$$\mathbf{q}^* = \tilde{K}_\lambda \mathbf{p}, \quad (22)$$

proving the theorem.

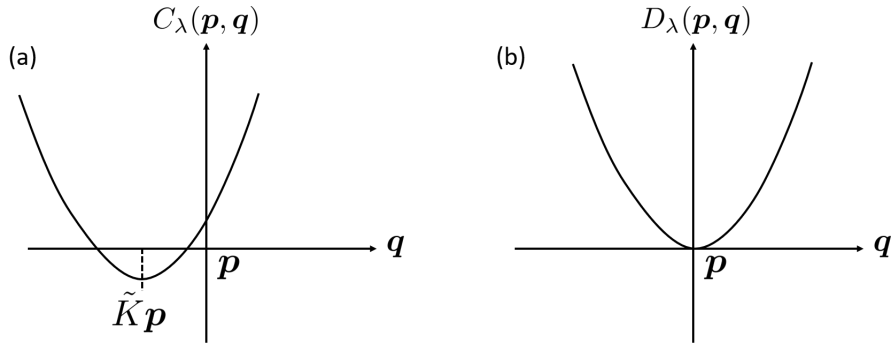


Fig. 2 Comparison of C_λ and D_λ . (a) C_λ as the function of \mathbf{q} . C_λ is minimized when $\mathbf{q} = \tilde{\mathbf{K}}\mathbf{p}$. (c) D_λ as the function of \mathbf{q} . D_λ is minimized when $\mathbf{q} = \mathbf{p}$.

We define a new family of divergences $D_\lambda(\mathbf{p}, \mathbf{q})$ that also depend on λ .

Definition 1

$$D_\lambda[\mathbf{p} : \mathbf{q}] = (1 + \lambda) \left(C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{q}) - C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) \right). \quad (23)$$

Figure 2 compares C_λ and D_λ in S_{n-1} . The following theorem is obtained of which proof is given in Appendix I.

Theorem 3 $D_\lambda[\mathbf{p} : \mathbf{q}]$ is a convex function with respect to \mathbf{p} and \mathbf{q} , satisfying the constraints (10). It converges to the Wasserstein distance as $\lambda \rightarrow 0$.

4 Behavior of $\tilde{\mathbf{K}}_\lambda$

The divergence D_λ is defined through $\tilde{\mathbf{K}}_\lambda$. We study properties of $\tilde{\mathbf{K}}_\lambda$, including two limiting cases of $\lambda \rightarrow 0, \infty$.

We first consider the case for small λ to see how $\tilde{\mathbf{K}}_\lambda$ behaves, assuming that X has a graphical structure. We assume that $M_{ii} = 0$, $M_{ij} = 1$ when i is a nearest neighbor of j , and $M_{ij} > 1$, otherwise. By putting

$$\epsilon = \exp\left(-\frac{1}{\lambda}\right), \quad \lambda = -\frac{1}{\log \epsilon}, \quad (24)$$

we have

$$K_{ij} = \exp\left(-\frac{M_{ij}}{\lambda}\right) = \epsilon^{M_{ij}}, \quad (25)$$

$$\tilde{K}_{i|j} = \frac{\epsilon^{M_{ji}}}{\sum_i \epsilon^{M_{ji}}}. \quad (26)$$

Then, by neglecting higher-order terms of ϵ , we have

$$\tilde{K}_{i|j} = \begin{cases} 1 - |N(i)|\epsilon, & i = j, \\ \epsilon, & i \in N(j), \\ 0, & \text{otherwise,} \end{cases} \quad (27)$$

where $N(j)$ is the set of nearest neighbors of j . When X is a plane consisting of $n \times n$ pixels, a neighbor consists of four pixels, $|N(j)| = 4$, except for boundary pixels. We see that $\tilde{K}_{i|j}$ is approximated by the discrete Laplacian operator Δ ,

$$\tilde{\mathbf{K}} = (1 - \epsilon)\mathbf{I} + \epsilon\Delta. \quad (28)$$

This shows that $\tilde{\mathbf{K}}$ is a diffusion operator, flattening pattern \mathbf{q} , that is, shifting \mathbf{q} toward the uniform distribution $\mathbf{1}/n$.

The inverse of $\tilde{\mathbf{K}}$ is

$$\tilde{\mathbf{K}}^{-1} = (1 + \epsilon)\mathbf{I} - \epsilon\Delta. \quad (29)$$

This is the inverse diffusion operator, which sharpens \mathbf{q} by emphasizing larger components.

In order to make clear the character of diffusion without assuming λ is small, we consider a continuous pattern $\mathbf{p} = p(\boldsymbol{\xi})$, $\boldsymbol{\xi} \in \mathbf{R}^n$ and the metric

$$M(\boldsymbol{\xi}, \boldsymbol{\xi}') = |\boldsymbol{\xi} - \boldsymbol{\xi}'|^2. \quad (30)$$

Then,

$$K_\lambda(\boldsymbol{\xi}, \boldsymbol{\xi}') = \exp\left(-\frac{|\boldsymbol{\xi} - \boldsymbol{\xi}'|^2}{\lambda}\right), \quad (31)$$

and we easily have

$$\tilde{K}_\lambda(\boldsymbol{\xi}|\boldsymbol{\xi}') = \frac{1}{(\sqrt{\pi\lambda})^n} \exp\left(-\frac{|\boldsymbol{\xi} - \boldsymbol{\xi}'|^2}{\lambda}\right). \quad (32)$$

This is a diffusion kernel. When \mathbf{p} is a Gaussian distribution

$$p(\boldsymbol{\xi}) = \exp\left(-\frac{1}{2\sigma^2}|\boldsymbol{\xi}|^2\right), \quad (33)$$

we have

$$\tilde{K}_\lambda \mathbf{p} = \exp\left(-\frac{1}{2\tau^2}|\boldsymbol{\xi}|^2\right), \quad (34)$$

with

$$\tau^2 = \sigma^2 + \frac{\lambda}{2}. \quad (35)$$

Hence, $\tilde{K}_\lambda \mathbf{p}$ is Gaussian, where the variance τ^2 is increased by $\lambda/2$, blurring the original \mathbf{p} .

We lastly consider the case when λ is large enough, studying the limiting behavior as $\lambda \rightarrow \infty$.

When λ is large, we expand \mathbf{K}_λ as

$$K_{ij,\lambda} = \exp\left(-\frac{M_{ij}}{\lambda}\right) = 1 - \frac{M_{ij}}{\lambda}, \quad (36)$$

$$K_{j,\lambda} = n\left(1 - \frac{\bar{m}_{j\cdot}}{\lambda}\right), \quad \bar{m}_{j\cdot} = \frac{1}{n} \sum_i M_{ji}, \quad (37)$$

obtaining

$$\tilde{K}_{ij,\lambda} = \frac{K_{ji,\lambda}}{K_{j,\lambda}} = \frac{1}{n} \left(1 - \frac{\tilde{m}_{ij}}{\lambda}\right), \quad \tilde{m}_{ij} = M_{ji} - \bar{m}_{j\cdot}. \quad (38)$$

Hence,

$$\tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p})_i = \frac{1}{n\lambda} \sum_j \tilde{m}_{ij} (q_j - p_j), \quad (39)$$

showing that this is of order $1/\lambda$. Let $\tilde{\mathbf{M}}$ is the moment matrix defined by

$$\tilde{M}_{jk} = \frac{1}{n} \sum_i \tilde{m}_{ij} \tilde{m}_{ik}. \quad (40)$$

Then, we have the following theorem, of which proof is given in Appendix II.

Theorem 4 When λ is large enough,

$$\lim_{\lambda \rightarrow \infty} D_\lambda[\mathbf{p} : \mathbf{q}] = \frac{1}{2} (\mathbf{q} - \mathbf{p})^T \tilde{\mathbf{M}} (\mathbf{q} - \mathbf{p}), \quad (41)$$

which is a squared energy distance defined by the moment matrix $\tilde{\mathbf{M}}$.

5 Right Barycenter of Patterns

We consider the barycenter of image patterns represented as probability measures $\mathbf{p} = p(\boldsymbol{\xi})$ on the plane $\boldsymbol{\xi} = (x, y) \in \mathbb{R}^2$ using divergence D_λ . The plane is discretized into a grid of $n \times m$ pixels, and therefore \mathbf{p} is a probability vector of size nm .

Let us consider a family S of patterns, $S = (\mathbf{p}_i)_{i=1,\dots,N}$. A right D -barycenter $\mathbf{q}_D^*(S)$ of these patterns is the minimizer of

$$F_D^r(S, \mathbf{q}) = \sum_i D_\lambda[\mathbf{p}_i : \mathbf{q}]. \quad (42)$$

Cuturi and Doucet (2014) used $C_\lambda(\mathbf{p}, \mathbf{q})$ to define the following barycenter, as a minimizer of

$$F_C(S, \mathbf{q}) = \sum C_\lambda(\mathbf{p}_i, \mathbf{q}). \quad (43)$$

We call such a minimizer the C -barycenter $\mathbf{q}_C^*(S)$. Cuturi and Doucet showed that the C -barycenter can extract for some pattern families S a common shape. In particular, Cuturi and Doucet used a family S of deformed double rings at

various sizes and positions, whose barycenter $\mathbf{q}_C^*(S)$ turns out to be a standard double ring pattern. Other information-theoretic divergences such as the KL-divergence or Heillinger divergence are not able to recover such a common shape. We will show that the right D -barycenter $\mathbf{q}_D^*(S)$ exhibits the same property, but with a sharper solution.

The right D -barycenter minimizes

$$\frac{1}{1+\lambda} \sum D_\lambda(\mathbf{p}_i : \mathbf{q}) = \sum C_\lambda(\mathbf{p}_i, \tilde{\mathbf{K}}_\lambda \mathbf{q}) - \sum C_\lambda(\mathbf{p}_i, \tilde{\mathbf{K}}_\lambda \mathbf{p}_i). \quad (44)$$

The second term of the right-hand side of (44) does not depend on \mathbf{q} , so that it may be deleted for minimization. Let us put

$$\tilde{\mathbf{q}} = \tilde{\mathbf{K}}_\lambda \mathbf{q}. \quad (45)$$

Then, the D -barycenter is derived from the C -barycenter by

$$\mathbf{q}_C^* = \tilde{\mathbf{K}}_\lambda \mathbf{q}_D^*, \quad (46)$$

provided (46) is solvable. In this case,

$$\mathbf{q}_D^* = \tilde{\mathbf{K}}_\lambda^{-1} \mathbf{q}_C^* \quad (47)$$

is a sharpened version of \mathbf{q}_C^* . However, (46) might not be always solvable.

The image $\tilde{\mathbf{K}}_\lambda S_{n-1}$ is a simplex sitting inside S_{n-1} . Since the C -barycenter \mathbf{q}_C^* is not necessarily inside $\tilde{\mathbf{K}}_\lambda S_{n-1}$, we need to solve the D -barycenter problem (44) under the constraint that \mathbf{q} is constrained inside $\tilde{\mathbf{K}}_\lambda S_{n-1}$. When the C -barycenter \mathbf{q}_C^* is inside $\tilde{\mathbf{K}}_\lambda S_{n-1}$, the D -barycenter \mathbf{q}_D^* is simply given by (46), which is more localized or sharper than \mathbf{q}_C^* . When \mathbf{q}_C^* is not inside $\tilde{\mathbf{K}}_\lambda S_{n-1}$, the solution of (44) is on the boundary of the simplex $\tilde{\mathbf{K}}_\lambda S_{n-1}$, which implies that some components of \mathbf{q}_D^* are forced to be equal to 0.

Theorem 5 The right D -barycenter $\mathbf{q}_D^*(S)$ of S is a sharper (more localized) version of the C -barycenter.

Agueh and Carlier (2011) showed that when the ground metric is the quadratic Euclidean distance, the W -barycenter has the property that its shape is determined from the shapes of each element \mathbf{p}_i in S , but does not depend on their location, namely that it is translation invariant. The C -barycenter also inherits this property, and we show that so does the right D -barycenter.

Let us consider a pattern $p(\boldsymbol{\xi}) = p(x, y)$ on the (x, y) -plane, where $\boldsymbol{\xi} = (x, y)$. The center $\boldsymbol{\xi}_p$ of $p(x, y)$ is defined by

$$\boldsymbol{\xi}_p = \int \boldsymbol{\xi} p(\boldsymbol{\xi}) d\boldsymbol{\xi}. \quad (48)$$

Given a $\bar{\boldsymbol{\xi}}$, we define a shift operator $T_{\bar{\boldsymbol{\xi}}}$, which shifts $p(x, y)$ by $\bar{\boldsymbol{\xi}} = (\bar{x}, \bar{y})$,

$$T_{\bar{\boldsymbol{\xi}}} p(\boldsymbol{\xi}) = p(\boldsymbol{\xi} - \bar{\boldsymbol{\xi}}). \quad (49)$$

Let $P(\boldsymbol{\xi}, \boldsymbol{\xi}')$ be a transportation plan from $p(\boldsymbol{\xi})$ to $q(\boldsymbol{\xi})$. When $q(\boldsymbol{\xi})$ is shifted as $T_{\bar{\boldsymbol{\xi}}}q(\boldsymbol{\xi})$, we naturally define the transportation plan $T_{\bar{\boldsymbol{\xi}}}P(\boldsymbol{\xi}, \boldsymbol{\xi}')$,

$$T_{\bar{\boldsymbol{\xi}}}P(\boldsymbol{\xi}, \boldsymbol{\xi}') = \bar{P}(\boldsymbol{\xi}, \boldsymbol{\xi}') = P(\boldsymbol{\xi}, \boldsymbol{\xi}' - \bar{\boldsymbol{\xi}}), \quad (50)$$

which transports $p(\boldsymbol{\xi})$ to $T_{\bar{\boldsymbol{\xi}}}q(\boldsymbol{\xi})$.

We study how the transportation cost changes by a shift. As recalled above, we use the squared Euclidean distance as the ground cost,

$$m(\boldsymbol{\xi}, \boldsymbol{\xi}') = \|\boldsymbol{\xi} - \boldsymbol{\xi}'\|^2. \quad (51)$$

Then, the direct cost of transportation is

$$\langle M, P \rangle = \int m(\boldsymbol{\xi}, \boldsymbol{\xi}') P(\boldsymbol{\xi}, \boldsymbol{\xi}') d\boldsymbol{\xi} d\boldsymbol{\xi}'. \quad (52)$$

The cost of the shifted plan is

$$\langle M, T_{\bar{\boldsymbol{\xi}}}P \rangle = \int m(\boldsymbol{\xi}, \boldsymbol{\xi}') P(\boldsymbol{\xi}, \boldsymbol{\xi}' - \bar{\boldsymbol{\xi}}) d\boldsymbol{\xi} d\boldsymbol{\xi}' \quad (53)$$

$$= \int m(\boldsymbol{\xi}, \boldsymbol{\xi}'' + \bar{\boldsymbol{\xi}}) P(\boldsymbol{\xi}, \boldsymbol{\xi}'') d\boldsymbol{\xi} d\boldsymbol{\xi}'' \quad (54)$$

$$= \int \left\{ \|\boldsymbol{\xi} - \boldsymbol{\xi}''\|^2 + \|\bar{\boldsymbol{\xi}}\|^2 - 2\bar{\boldsymbol{\xi}} \cdot (\boldsymbol{\xi} - \boldsymbol{\xi}'') \right\} P(\boldsymbol{\xi}, \boldsymbol{\xi}'') d\boldsymbol{\xi} d\boldsymbol{\xi}'' \quad (55)$$

$$= \langle M, P \rangle + \|\bar{\boldsymbol{\xi}}\|^2 - 2\bar{\boldsymbol{\xi}} \cdot (\boldsymbol{\xi}_p - \boldsymbol{\xi}_q). \quad (56)$$

Note that a shift does not change the entropy

$$H\{P(\boldsymbol{\xi}, \boldsymbol{\xi}')\} = H\{T_{\bar{\boldsymbol{\xi}}}P(\boldsymbol{\xi}, \boldsymbol{\xi}')\}. \quad (57)$$

When $\boldsymbol{\xi}_p = \boldsymbol{\xi}_q$, two patterns \boldsymbol{p} and \boldsymbol{q} are said to be co-centered. For two co-centered $\boldsymbol{p}, \boldsymbol{q}$, we have

$$\mathcal{L}(P_{\boldsymbol{p}, \boldsymbol{q}}) \geq \mathcal{L}\left(P_{\boldsymbol{p}, \bar{\boldsymbol{K}}_\lambda \boldsymbol{p}}\right), \quad (58)$$

where $P_{\boldsymbol{p}, \boldsymbol{q}}$ is a transportation plan sending \boldsymbol{p} to \boldsymbol{q} . Hence, $\boldsymbol{q}^* = \bar{\boldsymbol{K}}_\lambda \boldsymbol{p}$ minimizes the transportation cost among all co-centered \boldsymbol{q} .

We fix $\boldsymbol{\xi}_q$ and search for the optimal shape $\boldsymbol{q}^*(\boldsymbol{\xi})$ located at $\boldsymbol{\xi}_q$ that minimizes the transportation cost $C_\lambda(P_{\boldsymbol{p}, \boldsymbol{q}})$ from \boldsymbol{p} to \boldsymbol{q} . In order to derive the optimal shape, we shift \boldsymbol{q} by $T_{\bar{\boldsymbol{\xi}}}$,

$$\bar{\boldsymbol{\xi}} = \boldsymbol{\xi}_q - \boldsymbol{\xi}_p, \quad (59)$$

such that \boldsymbol{p} and $\bar{\boldsymbol{q}} = T_{\bar{\boldsymbol{\xi}}}q$ are co-centered. Then, for a plan $P_{\boldsymbol{p}, \bar{\boldsymbol{q}}}$, we have $P_{\boldsymbol{p}, \boldsymbol{q}}$ which is the shifted plan of $P_{\boldsymbol{p}, \bar{\boldsymbol{q}}}$ by $-\bar{\boldsymbol{\xi}}$,

$$P_{\boldsymbol{p}, \boldsymbol{q}} = T_{-\bar{\boldsymbol{\xi}}}P_{\boldsymbol{p}, \bar{\boldsymbol{q}}}, \quad (60)$$

$$\boldsymbol{q} = T_{-\bar{\boldsymbol{\xi}}}\bar{\boldsymbol{q}}. \quad (61)$$

We have

$$C_\lambda(T_{\bar{\xi}}P) = C_\lambda(P) + \|\bar{\xi}\|^2 - 2\bar{\xi} \cdot (\xi_p - \xi_q). \quad (62)$$

We have an important relation that $C_\lambda(P_{\mathbf{p},\mathbf{q}})$ is decomposed into a sum of the shape deformation cost among co-centered patterns and the positional transportation cost, as

$$C_\lambda(P_{\mathbf{p},\mathbf{q}}) = C_\lambda(P_{\mathbf{p},\bar{\mathbf{q}}}) + \|\xi_p - \xi_q\|^2 \quad (63)$$

because of

$$\xi_p = \xi_{\bar{q}}. \quad (64)$$

Lemma 1 Given \mathbf{p} , the C_λ -optimal pattern $q(\xi)$ transporting \mathbf{p} to \mathbf{q} is a shifted and blurred version of $p(\xi)$,

$$q^*(\xi) = T_{\xi_p - \xi_q} \tilde{K}_\lambda p(\xi), \quad (65)$$

not depending on the locations ξ_p and ξ_q .

Let $p_1(\xi), \dots, p_n(\xi)$ be n local patterns. We search for their C - and right D -barycenters $q(\xi)$ that minimize

$$F_C(\mathbf{q}) = \sum_{i=1}^n C_\lambda(P_{\mathbf{p}_i,\mathbf{q}}), \quad F_D(\mathbf{q}) = \sum_{i=1}^n D_\lambda(P_{\mathbf{p}_i,\mathbf{q}}), \quad (66)$$

The center of $p_i(\xi)$ is denoted by $\xi_i = \xi_{p_i}$. Before solving the problem, let us respectively shift $p_i(\xi)$ from ξ_i to ξ_0 , a fixed common location, such that all shifted $\bar{p}_i(\xi)$'s are co-centered. Let $q_C^*(\xi)$ and $q_D^*(\xi)$ be the barycenters of all co-centered $\bar{p}_i(\xi)$, which do not depend on the locations of $p_i(\xi)$ but their shapes.

Theorem 6 (Shape-Location Separation Theorem) The barycenters $q_C^*(\xi)$ and $q_D^*(\xi)$ of $p_1(\xi), \dots, p_n(\xi)$ are located at the barycenter of ξ_1, \dots, ξ_n and their shapes are given by the respective barycenters of $\bar{p}_1(\xi), \dots, \bar{p}_n(\xi)$.

Proof From (63), we have

$$F_C(\mathbf{q}) = \sum C_\lambda(P_{\mathbf{p}_i,\bar{\mathbf{q}}_i}) + \sum \|\xi_{p_i} - \xi_q\|^2, \quad (67)$$

$$F_D(\mathbf{q}) = \sum C_\lambda(P_{\mathbf{p}_i,\tilde{K}_\lambda \bar{\mathbf{q}}}) + \sum \|\xi_{p_i} - \xi_q\|^2, \quad (68)$$

where $\bar{\mathbf{q}}_i$ is the shifted version of \mathbf{q} to the center of \mathbf{p}_i . Here, the shape and location of \mathbf{q} is separated. Minimizing the first term, we have \mathbf{q}^* which is the respective barycenters of the shapes of co-centered $\mathbf{p}_1, \dots, \mathbf{p}_n$. The second term gives the barycenter of locations ξ_1, \dots, ξ_n .

Corollary 1 When \mathbf{p}_i are shifts of an identical \mathbf{p} , their right D -barycenter has the same shape as the original \mathbf{p} , whereas the C -barycenter is a blurred version of \mathbf{p} ,

$$q_C^* = \mathbf{K}_\lambda \mathbf{p}. \quad (69)$$

We show a simple example where \mathbf{p}_i are shifted \mathbf{p} , a cat shape (Fig. 3a). Its C -barycenter has a blurred shape $\mathbf{K}_\lambda \mathbf{p}$ (Fig. 3b) tending to the uniform distribution as $\lambda \rightarrow \infty$. However, the shape of the right D -barycenter is exactly the same as \mathbf{p} (Fig. 3c).

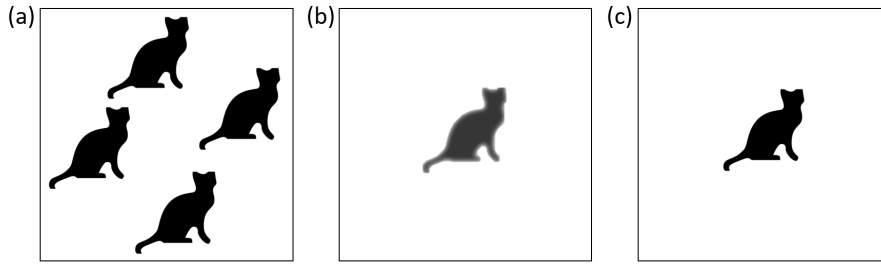


Fig. 3 (a) Cat images. (b) The C -barycenter of panel (a). The C -barycenter has a blurred shape $\mathbf{K}_\lambda \mathbf{p}$ tending to the uniform distribution as $\lambda \rightarrow \infty$. (c) The D -barycenter of panel (a). The shape of the right D -barycenter is exactly the same as the original shape \mathbf{p} .

6 Left Barycenter of Patterns

Since we use the asymmetric divergence D_λ to define a barycenter, we may consider another barycenter by putting the unknown barycenter in the left argument of D_λ .

We consider again a family S of patterns, $S = (\mathbf{q}_i)_{i=1, \dots, N}$. The barycenter \mathbf{p} of these patterns based on divergence D_λ is defined by the minimizer of

$$F_D^l(S, \mathbf{p}) = \sum_i D_\lambda[\mathbf{p} : \mathbf{q}_i], \quad (70)$$

and is called the left D -barycenter. We propose to solve that problem using the accelerated gradient descent approach outlined in (Cuturi and Doucet, 2014), with two differences: all examples \mathbf{q}_i must be smoothed beforehand following an application of $\tilde{\mathbf{K}}_\lambda$, and the gradient incorporates now not only terms resulting from $C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{q}_i)$, but also from $-C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p})$ which, as explained in Equation (73), is simply minus the Kullback-Leibler divergence between \mathbf{p} and the vector $\tilde{\mathbf{K}}_\lambda \mathbf{1}$, that is the entropy of \mathbf{p} plus the dot product between \mathbf{p} and $\log(\tilde{\mathbf{K}}_\lambda \mathbf{1})$. As a result the gradient of $-C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p})$ is equal to, up to a constant term, $-\log(\mathbf{p}) + \log(\tilde{\mathbf{K}}_\lambda \mathbf{1})$, which tends to sharpen further any iterate compared to the simple minimization of the C_λ barycenter. Note that this approach, namely adding the entropy to the regularized Wasserstein barycenter problem, was used in a heuristic way by Solomon et al. (2015) who called it *entropic sharpening*, without proving that the resulting problem was convex. Our work shows that, up to a given strength, the entropic sharpening of regularized Wasserstein barycenters remains a convex problem.

It might be easier to calculate numerically, but it does not have the shape-location separation property. An example of the left D -barycenter of four patterns is shown in Fig. 4.

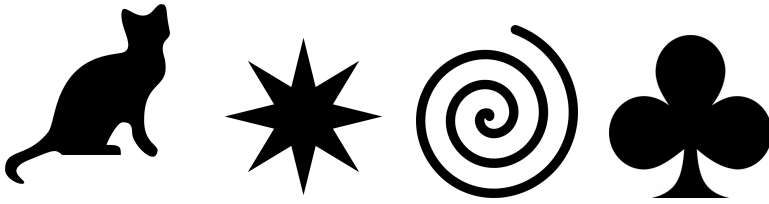


Fig. 4 Four shapes considered in our experiments to compute barycenters according to C_λ or D_λ

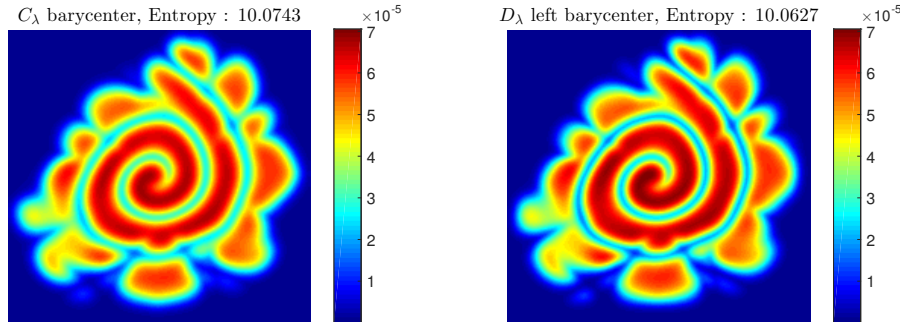


Fig. 5 We consider in this example the squared-Euclidean distance on grids as the ground metric cost. We study the iso-barycenter of the four shapes presented in Figure 4 using two different discrepancy functions, C_λ and D_λ . (*left*) barycenter obtained using Algorithm 2 of Solomon et al. (2015), itself a specialized version of Benamou et al.’s algorithm (2014). We used a 10^{-10} tolerance for the l_1 norm between two successive iterates as a stopping criterion). (*right*) D_λ left-barycenter obtained with our accelerated gradient approach. We use a `jet` colormap to highlight differences in the support of the barycenters. As expected, the entropy of the C_λ barycenter is higher than that of the D_λ left barycenter, since the latter optimization incorporates a penalization term, $-C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p})$, which is equivalent to penalizing the entropy of the solution \mathbf{p} . This difference in entropy results in subtle yet visible differences between the two solutions, with sharper edges for the D_λ left barycenter.

Conclusions

We defined a new divergence between two probability distributions based on the C -function, which is the entropy-regularized cost function (Cuturi, 2013). Although it is useful in many applications, it does not satisfy the criteria of a distance or divergence. We defined a new divergence function D_λ derived from C_λ , which works better than the original C_λ for some problems, in particular, the barycenter problem. We proved that the minimizer of $C_\lambda(\mathbf{p}, \mathbf{q})$ is given by $\tilde{\mathbf{K}}_\lambda \mathbf{p}$, where $\tilde{\mathbf{K}}_\lambda$ is a diffusion operator depending on the base metric \mathbf{M} . We studied properties of $\tilde{\mathbf{K}}_\lambda$ showing how it changes as λ increases, elucidating properties of D_λ .

We applied D_λ to obtain the barycenter of a cluster of image patterns. It is proved that the right D -barycenter keeps a good property that the shape and locations of patterns are separated, which is a merit of the C -function based

barycenter. Moreover, the D -barycenter gives even a sharper shape than the C -barycenter.

We cannot touch upon computational aspects of the D -barycenter, because this is a theoretical paper proposing a new divergence and its properties. We also defined the left D -barycenter, because D_λ is an asymmetric divergence. This is computationally easy to calculate. However, it remains as our future study to explore its properties.

Appendix I: Proof of convexity of D_λ

Let us put the Hessian of the cost function as follows:

$$\mathbf{G}_\lambda = \frac{\partial^2 C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{q})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} = \begin{bmatrix} \mathbf{X} & \mathbf{Y} \\ \mathbf{Y}^T & \mathbf{Z} \end{bmatrix}, \quad (71)$$

expressed by block matrices $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ where $\boldsymbol{\eta} = (\mathbf{p}, \mathbf{q})^T$. Because C_λ is strictly convex (Amari et al., 2018), \mathbf{G}_λ is positive definite and the block component \mathbf{Z} is also regular and positive definite.

By using

$$P_{ij}^* = p_i \frac{K_{ij}}{K_i}, \quad (72)$$

we get

$$C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) = \lambda \sum_{i=1}^n p_i \ln \frac{p_i}{K_i}. \quad (73)$$

Therefore, its Hessian becomes

$$\mathbf{G}'_\lambda = \frac{\partial^2 C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} = \lambda \begin{bmatrix} \text{diag}\left(\frac{1}{p_i}\right) + \frac{1}{p_n} \mathbf{1}\mathbf{1}^T & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}, \quad (74)$$

where \mathbf{O} is the zero matrix and $\text{diag}(p_i)$ represents a diagonal matrix whose ij component is given by $p_i \delta_{ij}$. Let us put $\mathbf{R} = \mathbf{G}_\lambda - \mathbf{G}'_\lambda$. The determinant of \mathbf{R} is given by

$$\det(\mathbf{R}) = \det(\mathbf{Z}) \det(\mathbf{R}'), \quad (75)$$

where we put

$$\mathbf{R}' = \mathbf{X} - \mathbf{Y} \mathbf{Z}^{-1} \mathbf{Y}^T - \lambda \left(\text{diag}\left(\frac{1}{p_i}\right) + \frac{1}{p_n} \mathbf{1}\mathbf{1}^T \right). \quad (76)$$

Because \mathbf{Z} is positive definite, the positive definiteness of \mathbf{R} is equivalent to that of \mathbf{R}' . As derived in (Amari et al., 2018),

$$\mathbf{G}_\lambda^{-1} = \frac{1}{1 + \lambda} \begin{bmatrix} p_i \delta_{ij} - p_i p_j & P_{ij} - p_i q_j \\ P_{ji} - q_i p_j & q_i \delta_{ij} - q_i q_j \end{bmatrix}, \quad (77)$$

and we can represent the (\mathbf{p}, \mathbf{p}) -block part of \mathbf{G}_λ^{-1} by using the block components of \mathbf{G}_λ as follows:

$$(\mathbf{X} - \mathbf{Y}\mathbf{Z}^{-1}\mathbf{Y}^T)^{-1} = \frac{1}{1+\lambda} (\text{diag}(p_i) - \mathbf{p}\mathbf{p}^T). \quad (78)$$

By using the Sherman-Morrison formula, we get

$$\mathbf{X} - \mathbf{Y}\mathbf{Z}^{-1}\mathbf{Y}^T = (1+\lambda) \left(\text{diag} \left(\frac{1}{p_i} \right) + \frac{1}{p_n} \mathbf{1}\mathbf{1}^T \right). \quad (79)$$

Finally, we obtain

$$\mathbf{R}' = \text{diag} \left(\frac{1}{p_i} \right) + \frac{1}{p_n} \mathbf{1}\mathbf{1}^T. \quad (80)$$

Because this \mathbf{R}' is positive definite, \mathbf{R} is also positive definite and D_λ is strictly convex.

By using Eq.(14) and $C_0(\mathbf{p}, \mathbf{p}) = 0$, we can confirm that D_λ converges to the Wasserstein distance as $\lambda \rightarrow 0$.

Appendix II: Proof of D_λ approaching an energy function for large λ

We expand $D_\lambda[\mathbf{p}, \mathbf{q}]$ in terms of $\tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p})$ as

$$\frac{1}{(1+\lambda)} D_\lambda[\mathbf{p}, \mathbf{q}] = C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p} + \tilde{\mathbf{K}}(\mathbf{q} - \mathbf{p})) - C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) \quad (81)$$

$$= \partial_q C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) \cdot \tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p}) + \frac{1}{2} \partial_q \partial_q C(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) \cdot \tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p}) \otimes \tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p}) \quad (82)$$

$$= \frac{1}{2} \partial_q \partial_q C(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) \cdot \tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p}) \otimes \tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p}), \quad (83)$$

because of

$$\partial_q C_\lambda(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) = 0, \quad (84)$$

where \otimes is the tensor product. Higher-order terms are neglected.

When λ is large, we expand $\tilde{\mathbf{K}}_\lambda$ as

$$K_{ij,\lambda} = \exp \left\{ -\frac{M_{ij}}{\lambda} \right\} = 1 - \frac{M_{ij}}{\lambda}, \quad (85)$$

$$K_{j\cdot,\lambda} = n \left(1 - \frac{\tilde{m}_{\cdot j}}{\lambda} \right), \quad \tilde{m}_{j\cdot} = \frac{1}{n} \sum_i M_{ji}, \quad (86)$$

obtaining

$$\tilde{K}_{i|j,\lambda} = \frac{K_{ji}}{K_{j\cdot}} = \frac{1}{n} \left(1 - \frac{\tilde{m}_{ij}}{\lambda} \right), \quad \tilde{m}_{ij} = M_{ji} - \tilde{m}_{j\cdot}. \quad (87)$$

Hence

$$\tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p})_i = \frac{1}{n\lambda} \sum_j \tilde{m}_{ij} (q_j - p_j). \quad (88)$$

We have

$$\partial_{\mathbf{q}} \partial_{\mathbf{q}} C(\mathbf{p}, \tilde{\mathbf{K}}_\lambda \mathbf{p}) = \partial_{\mathbf{q}} \partial_{\mathbf{q}} C_\lambda(\mathbf{p}, \mathbf{1}) + O\left(\frac{1}{\lambda}\right). \quad (89)$$

So we calculate $\partial_{\mathbf{q}} \partial_{\mathbf{q}} C_\lambda(\mathbf{p}, \tilde{\mathbf{q}})$ when

$$\tilde{q}_i = \left(\tilde{\mathbf{K}}_\lambda \mathbf{p}\right)_i = \frac{1}{n} + \varepsilon_i \quad (90)$$

and expand it up to $O(\varepsilon^2)$.

The optimal transportation plan for $\mathbf{p} \rightarrow \mathbf{q}$ is

$$p_{ij}^* = ca_i b_j \quad (91)$$

when $\lambda \rightarrow \infty$, because $K_{ij,\lambda} = 1$. Hence,

$$P_{ij}^* = p_i q_j. \quad (92)$$

We have already obtained the inverse of $\partial_{\boldsymbol{\eta}} \partial_{\boldsymbol{\eta}} C_\lambda(\mathbf{p}, \mathbf{q})$ in (45) of the previous paper

$$\mathbf{G}_\lambda^{-1} = \frac{1}{1 + \lambda} \begin{bmatrix} p_i \delta_{ij} - p_i p_j & P_{ij}^* - p_i q_j \\ P_{ji}^* - q_i p_j & q_i \delta_{ij} - q_i q_j \end{bmatrix}. \quad (93)$$

Hence, it is block-diagonal ($P_{ij}^* - p_i q_j = 0$), and the (\mathbf{q}, \mathbf{q}) -part of \mathbf{G}_λ ($\lambda \rightarrow \infty$) is

$$\mathbf{G}_{\lambda, \mathbf{q}\mathbf{q}} = (1 + \lambda) [q_i \delta_{ij} - q_i q_j]^{-1} \quad (94)$$

$$= (1 + \lambda) \left(\text{diag} \left(\frac{1}{q_i} \right) + \frac{1}{q_n} \mathbf{1}\mathbf{1}^T \right). \quad (95)$$

In our case of $\mathbf{q} = \mathbf{1}/n$,

$$\mathbf{G} = \partial_{\mathbf{q}} \partial_{\mathbf{q}} C = (1 + \lambda)(n\delta_{ij} + n). \quad (96)$$

We finally calculate

$$\begin{aligned} & \mathbf{G} \cdot \tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p}) \otimes \tilde{\mathbf{K}}_\lambda(\mathbf{q} - \mathbf{p}) \\ &= \frac{(1 + \lambda)n}{n^2 \lambda^2} \sum \tilde{m}_{ij} \tilde{m}_{ik} (q_j - p_j) (q_k - p_k) + \frac{(1 + \lambda)n}{n^2 \lambda^2} \left\{ \sum_{ij} \tilde{m}_{ij} (q_j - p_j) \right\}^2 \\ &= \frac{1 + \lambda}{\lambda^2} (\mathbf{q} - \mathbf{p})^T \tilde{\mathbf{M}} (\mathbf{q} - \mathbf{p}), \end{aligned}$$

Note that $\sum_i \tilde{m}_{ij} = 0$.

Bibliography

- M Agueh and G Carlier. Barycenters in the wasserstein space. SIAM J. on Mathematical Analysis, 43:904–924, 2011.
- Shun-ichi Amari. Information Geometry and Its Applications. Springer, 2016.
- Shun-ichi Amari, Ryo Karakida, and Masafumi Oizumi. Information geometry connecting wasserstein distance and kullback–leibler divergence via the entropy-relaxed transportation problem. Information Geometry, 2018. doi: 10.1007/s41884-018-0002-8.
- J.-D. Benamou, G. Carlier, M. Cuturi, L. Nenna, and G. Peyré. Iterative Bregman projections for regularized transportation problems. arXiv preprint arXiv:1412.5154, 2014.
- N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy. Optimal transport for domain adaptation. IEEE Transactions on Pattern Analysis and Machine Intelligence, 39(9):1853–1865, Sept 2017. ISSN 0162-8828. doi: 10.1109/TPAMI.2016.2615921.
- M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems, pages 2292–2300, 2013.
- M. Cuturi and D. Avis. Ground metric learning. The Journal of Machine Learning Research, 2014.
- M. Cuturi and A. Doucet. Fast computation of Wasserstein barycenters. In Proceedings of the 31st International Conference on Machine Learning (ICML-14), pages 685–693, 2014.
- Marco Cuturi and Gabriel Peyré. A smoothed dual approach for variational wasserstein problems. SIAM Journal on Imaging Sciences, 9(1):320–343, 2016.
- Charlie Frogner, Chiyuan Zhang, Hossein Mobahi, Mauricio Araya, and Tomaso A Poggio. Learning with a wasserstein loss. In Advances in Neural Information Processing Systems, pages 2053–2061, 2015.
- Aude Genevay, Gabriel Peyre, and Marco Cuturi. Learning generative models with sinkhorn divergences. In Amos Storkey and Fernando Perez-Cruz, editors, Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics, volume 84 of Proceedings of Machine Learning Research, pages 1608–1617, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR.
- Antoine Rolet, Marco Cuturi, and Gabriel Peyr. Fast dictionary learning with a smoothed wasserstein loss. In Arthur Gretton and Christian C. Robert, editors, Proceedings of the 19th International Conference on Artificial Intelligence and Statistics, volume 51 of Proceedings of Machine Learning Research, pages 630–638, Cadiz, Spain, 09–11 May 2016. PMLR.
- Richard Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices. The annals of mathematical statistics, 35(2): 876–879, 1964.
- J. Solomon, F. de Goes, G. Peyré, M. Cuturi, A. Butscher, A. Nguyen, T. Du, and L. Guibas. Convolutional Wasserstein distances: Efficient optimal trans-

portation on geometric domains. ACM Transactions on Graphics (Proc. SIGGRAPH 2015), 2015.

Cédric Villani. Topics in Optimal Transportation, volume 58. AMS Graduate Studies in Mathematics, 2003.