

Information Geometry of Wasserstein Divergence

Ryo Karakida¹(✉) and Shun-ichi Amari²

¹ National Institute of Advanced Industrial Science and Technology,
Koto-ku, Tokyo 135-0064, Japan

karakida.ryo@aist.go.jp

² RIKEN Brain Science Institute, Wako-shi, Saitama 351-0198, Japan
amari@brain.riken.jp

Abstract. There are two geometrical structures in a manifold of probability distributions. One is invariant, based on the Fisher information, and the other is based on the Wasserstein distance of optimal transportation. We propose a unified framework which connects the Wasserstein distance and the Kullback-Leibler (KL) divergence to give a new information-geometrical theory. We consider the discrete case consisting of n elements and study the geometry of the probability simplex S_{n-1} , the set of all probability distributions over n atoms. The Wasserstein distance is introduced in S_{n-1} by the optimal transportation of commodities from distribution $\mathbf{p} \in S_{n-1}$ to $\mathbf{q} \in S_{n-1}$. We relax the optimal transportation by using entropy, introduced by Cuturi (2013) and show that the entropy-relaxed transportation plan naturally defines the exponential family and the dually flat structure of information geometry. Although the optimal cost does not define a distance function, we introduce a novel divergence function in S_{n-1} , which connects the relaxed Wasserstein distance to the KL-divergence by one parameter.

1 Introduction

Information geometry studies invariant properties of a manifold of probability distributions, which are useful for various applications in statistics, machine learning, signal processing, optimization and others. Two geometrical structures have been introduced from two different backgrounds. One is constructed based on the invariance principle: The geometry is invariant under reversible transformations of random variables. We then have the Fisher information matrix as the unique invariant Riemannian metric (Rao 1945; Chentsov 1982; Amari 2016). Moreover, two dually coupled affine connections are given as invariant connections. These structures are useful for various applications. Another geometrical structure is introduced through the transportation problem. A distribution of commodities in a manifold is transported to another distribution. The transportation with the minimal cost defines a distance between the two distributions, called the Monge-Kantorovich-Wasserstein distance or earth-mover distance. This gives a tool to study the geometry of distributions taking the metric of the supporting manifold into account.

Let $X = \{1, \dots, n\}$ be the support of a probability measure \mathbf{p} . The invariant geometry gives a structure which is invariant under permutations of elements of X . It leads to an efficient estimator in statistical estimation. On the other hand, when we consider a picture over n^2 pixels $X = \{(ij); i, j = 1, \dots, n\}$, neighboring pixels are close. A permutation of X destroys such a neighboring structure, so the invariance should not be required. The Wasserstein distance is responsible for such a structure. Therefore, it is useful for problems having neighboring structure in support X .

An interesting question arises how these two geometrical structures are related. They are useful structures in their own right, but it is intriguing to find a unified framework to include the two. For this purpose in mind, the present paper treats the discrete case over n elements, such that a probability distribution is given by a probability vector $\mathbf{p} = (p_1, \dots, p_n)$ in the probability simplex S_{n-1} , letting a general case of continuous distributions over a manifold to be studied in future.

Cuturi (2013) modified the transportation problem such that the cost is minimized under the entropy constraint. This is called the entropy-relaxed optimal transportation problem. In many applications, his group showed the quasi-distance defined by the entropy-constrained optimal solution gives superior properties to the information-geometric distance such as the KL divergence or the Hellinger distance. As an application, consider a set of normalized histograms over X . A clustering problem categorizes them in some classes such that a class consists of similar histograms. Since a histogram is regarded as an empirical probability distribution, the problem is formulated within the probability simplex S_{n-1} in the discrete case and the distances among supporting pixels play a fundamental role.

We follow the entropy-relaxed framework of Cuturi (2013), Cuturi and Avis (2014), Cuturi and Peyré (2016), etc. and introduce a Lagrangian function which is a linear combination of the transportation cost and the entropy. Given distribution \mathbf{p} of commodities at the sender and \mathbf{q} at the receiver, the optimal transportation plan is the minimizer of the Lagrangian function. We reveal that it is a convex function of \mathbf{p} and \mathbf{q} so it defines a dually flat geometric structure in $S_{n-1} \times S_{n-1}$. The m -flat coordinates are (\mathbf{p}, \mathbf{q}) and their dual, e -flat coordinates $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ are given from the Lagrangian duality of nonlinear optimization problems. The set of the optimal transportation plans is an exponential family with the canonical parameters $(\boldsymbol{\alpha}, \boldsymbol{\beta})$, where the expectation parameters are (\mathbf{p}, \mathbf{q}) . Furthermore, we introduce a novel divergence between \mathbf{p} and \mathbf{q} in S_{n-1} . It connects the relaxed Wasserstein distance to the KL-divergence by a one parameter family. Our divergence will be expected to be useful for practical applications, because a divergence is a general concept including the square of a distance and more flexible admitting non-symmetry between \mathbf{p} and \mathbf{q} .

2 Entropy-Constrained Transportation Problem

Let us consider n terminals $X = (X_1, \dots, X_n)$ at which amounts p_1, \dots, p_n of commodities are stocked. We transport them within X such that amounts

q_1, \dots, q_n are newly stored at X_1, \dots, X_n . We normalize the total amount to be equal to 1, so $\mathbf{p} = (p_1, \dots, p_n)$ and $\mathbf{q} = (q_1, \dots, q_n)$ are regarded as probability distributions in the probability simplex S_{n-1} ,

$$\sum p_i = 1, \quad \sum q_i = 1, \quad p_i > 0, \quad q_i > 0. \tag{1}$$

We consider a transportation plan $\mathbf{P} = (P_{ij})$, where P_{ij} is the amount of commodity transported from X_i to X_j . A plan \mathbf{P} is regarded as a joint probability distribution of commodities flowing from X_i to X_j , satisfying the sender and receiver constraints,

$$\sum_j P_{ij} = p_i, \quad \sum_i P_{ij} = q_j. \tag{2}$$

The set of \mathbf{P} 's satisfying (2) is denoted by $U(\mathbf{p}, \mathbf{q})$.

Let $\mathbf{c} = (c_{ij})$ be the cost matrix, where c_{ij} denotes the cost of transporting one unit of commodities from X_i to X_j .

The transportation cost is defined by

$$c(\mathbf{P}) = \langle \mathbf{c}, \mathbf{P} \rangle = \sum c_{ij} P_{ij}. \tag{3}$$

The Wasserstein distance is defined by the minimal cost of transporting distribution \mathbf{p} at the senders to \mathbf{q} at the receivers,

$$c(\mathbf{p}, \mathbf{q}) = \min_{\mathbf{P} \in U(\mathbf{p}, \mathbf{q})} \langle \mathbf{c}, \mathbf{P} \rangle, \tag{4}$$

where min is taken over all \mathbf{P} satisfying constraints (2). See e.g., Villani (2013).

Given \mathbf{p} and \mathbf{q} , let us consider a special transportation plan \mathbf{P}_D defined by the direct product of \mathbf{p} and \mathbf{q} ,

$$\mathbf{P}_D = \mathbf{p} \otimes \mathbf{q} = (p_i q_j). \tag{5}$$

This plan transports commodities from each sender to the receivers according to the receiver distribution \mathbf{q} , irrespective of \mathbf{c} . The entropy of \mathbf{P}_D ,

$$H(\mathbf{P}_D) = - \sum P_{Dij} \log P_{Dij} = H(\mathbf{p}) + H(\mathbf{q}), \tag{6}$$

is the minimum among all \mathbf{P} 's belonging to $U(\mathbf{p}, \mathbf{q})$, because of $H(\mathbf{P}) \leq H(\mathbf{p}) + H(\mathbf{q})$, where $H(\mathbf{P})$, $H(\mathbf{p})$ and $H(\mathbf{q})$ are the respective entropies and the equality holds for $\mathbf{P} = \mathbf{P}_D$.

We consider a constrained problem of searching for \mathbf{P} that minimizes $\langle \mathbf{c}, \mathbf{P} \rangle$ within a KL-divergence ball centered at \mathbf{P}_D ,

$$KL[\mathbf{P} : \mathbf{P}_D] \leq d \tag{7}$$

for constant d . As d increases, the entropy of \mathbf{P} increases within the ball. This is equivalent to the entropy constrained problem that minimizes a linear combination of the transportation cost $\langle \mathbf{c}, \mathbf{P} \rangle$ and entropy $H(\mathbf{P})$,

$$F_\lambda(\mathbf{P}) = \frac{1}{1 + \lambda} \langle \mathbf{c}, \mathbf{P} \rangle - \frac{\lambda}{1 + \lambda} H(\mathbf{P}) \tag{8}$$

for constant λ (Cuturi 2013). Here, λ is a Lagrangian multiplier and λ becomes smaller as d becomes larger.

3 Solution of Entropy-Constrained Problem

Since \mathbf{P} satisfies constraints (2), by using Lagrange multipliers α_i, β_j , minimization of (8) is formulated in the Lagrangian form,

$$L_\lambda(\mathbf{P}) = \frac{1}{1+\lambda} \langle \mathbf{c}, \mathbf{P} \rangle - \frac{\lambda}{1+\lambda} H(\mathbf{P}) - \sum_{i,j} (\alpha_i + \beta_j) P_{ij}. \tag{9}$$

Let us fix λ , considering it as a parameter controlling the magnitude of the entropy or the size of the KL-ball. By differentiating (9) with respect to P_{ij} , we have the following solution,

$$P_{ij} = \exp \left\{ -\frac{c_{ij}}{\lambda} + \frac{1+\lambda}{\lambda} (\alpha_i + \beta_j) - 1 \right\}. \tag{10}$$

Let us put

$$a_i = \exp \left(\frac{1+\lambda}{\lambda} \alpha_i \right) \quad b_j = \exp \left(\frac{1+\lambda}{\lambda} \beta_j \right), \quad K_{ij} = \exp \left\{ -\frac{c_{ij}}{\lambda} \right\}, \tag{11}$$

and the optimal solution is written as

$$P_{\lambda ij}^* \propto a_i b_j K_{ij}, \tag{12}$$

where a_i and b_j correspond to the Lagrange multipliers α_i and β_j to be determined from the constraints (2). Note that $2n$ constraints (2) are not independent. Because of $\sum p_i = 1$, we can obtain a_n by $a_n = 1 - \sum_{i \neq n} a_i$. Further, we note that $\mu \mathbf{a}$ and \mathbf{b}/μ give the same answer for any $\mu > 0$, where $\mathbf{a} = (a_i)$ and $\mathbf{b} = (b_j)$. Therefore, the degrees of freedom of \mathbf{a} and \mathbf{b} are $2(n-1)$, which are to be determined from \mathbf{p} and \mathbf{q} of which degrees of freedom are also $2(n-1)$. Therefore, we may choose \mathbf{a} and \mathbf{b} such that they satisfy

$$\sum a_i = 1, \quad \sum b_j = 1. \tag{13}$$

Then, $\mathbf{a}, \mathbf{b} \in S_{n-1}$ and we have the following theorem.

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

$$P_{\lambda ij}^* = c a_i b_j K_{ij}, \tag{14}$$

$$c = \frac{1}{\sum a_i b_j K_{ij}}, \tag{15}$$

where two vectors \mathbf{a} and \mathbf{b} are determined from \mathbf{p} and \mathbf{q} .

Cuturi (2013) obtained the above $P_{\lambda ij}^*$ and applied the Sinkhorn-Knopp algorithm to iteratively compute \mathbf{a} and \mathbf{b} .

The following lemma is useful for later calculations.

Lemma 1. The optimal value

$$\varphi_\lambda(\mathbf{p}, \mathbf{q}) = \min F_\lambda(\mathbf{P}) \quad (16)$$

is given by

$$\varphi_\lambda(\mathbf{p}, \mathbf{q}) = \frac{\lambda}{1+\lambda} \left(\sum p_i \log a_i + \sum q_j \log b_j + \log c \right). \quad (17)$$

Proof. We first calculate $H(\mathbf{P}_\lambda^*)$. Substituting (15) in $H(\mathbf{P}_\lambda^*)$, we have

$$H(\mathbf{P}_\lambda^*) = - \sum_{ij} P_{\lambda ij}^* \left(-\frac{c_{ij}}{\lambda} + \log c a_i b_j \right) \quad (18)$$

$$= \frac{1}{\lambda} \langle \mathbf{c}, \mathbf{P}_\lambda^* \rangle - \sum p_i \log a_i - \sum q_j \log b_j - \log c. \quad (19)$$

Hence, (17) follows.

4 Exponential Family of Optimal Transportation Plans

A transportation plan \mathbf{P} is a probability distribution over branches (i, j) connecting terminals i and j . Let x denote branches and $\delta_{ij}(x) = 1$ when x is (i, j) and 0 otherwise. Then \mathbf{P} is a probability distribution of random variable x ,

$$P(x) = \sum_{i,j=1}^n P_{ij} \delta_{ij}(x). \quad (20)$$

By introducing new parameters

$$\theta^{ij} = \log \frac{P_{ij}}{P_{nn}}, \quad \boldsymbol{\theta} = (\theta^{ij}), \quad (21)$$

it is rewritten in a parameterized form as

$$P(x, \boldsymbol{\theta}) = \exp \left\{ \sum_{i,j} \theta^{ij} \delta_{ij}(x) + \log P_{nn} \right\}. \quad (22)$$

This shows that the set of transportation plans is an exponential family, where θ^{ij} are the canonical parameters and $\eta_{ij} = P_{ij}$ the expectation parameters. They form an $(n^2 - 1)$ -dimensional manifold denoted by S_{TP} , because $\theta^{nn} = 0$.

An optimal transportation plan is specified by $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ in (10), $\boldsymbol{\alpha} = (\alpha_i)$, $\boldsymbol{\beta} = (\beta_j)$ which are determined from (\mathbf{p}, \mathbf{q}) . It is written as

$$P(x, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \exp \left[\sum_{i,j} \left\{ \frac{\lambda+1}{\lambda} (\alpha_i + \beta_j) - \frac{c_{ij}}{\lambda} \right\} \delta_{ij}(x) - \frac{\lambda+1}{\lambda} \psi \right], \quad (23)$$

where

$$\psi(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\lambda}{1+\lambda} \log \sum_{i,j} \exp \left\{ \frac{\lambda+1}{\lambda} (\alpha_i + \beta_j) - \frac{c_{ij}}{\lambda} \right\} \quad (24)$$

is the normalization factor. By putting

$$\theta^{ij} = \frac{1+\lambda}{\lambda} (\alpha_i + \beta_j) - \frac{c_{ij}}{\lambda}, \quad (25)$$

we see that the set S_{OTP} of optimal transformation plans is a submanifold of S_{TP} . Because (25) is linear in $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$, S_{OTP} itself is an exponential family, where the canonical parameters are $(1+\lambda)/\lambda$ times $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ and the expectation parameters are $(\mathbf{p}, \mathbf{q}) \in S_{n-1} \times S_{n-1}$, since

$$E \left[\sum_j \delta_{ij}(x) \right] = p_i, \quad (26)$$

$$E \left[\sum_i \delta_{ij}(x) \right] = q_j. \quad (27)$$

Since each of $\mathbf{p}, \mathbf{q} \in S_{n-1}$ has $n-1$ degrees of freedom, S_{OPT} is a $2(n-1)$ -dimensional dually flat manifold. We may put $\alpha_n = \beta_n = 0$ without loss of generality, which correspond to putting $a_n = b_n = 1$ instead of $\sum a_i = \sum b_j = 1$.

We calculate the relaxed cost function $\varphi_\lambda(\mathbf{p}, \mathbf{q})$ corresponding to $\mathbf{P}(\boldsymbol{\alpha}, \boldsymbol{\beta})$. We then have

$$\varphi_\lambda(\mathbf{p}, \mathbf{q}) = \frac{1}{1+\lambda} \langle \mathbf{c}, \mathbf{P} \rangle + \frac{\lambda}{1+\lambda} \sum_{i,j} P_{ij} \left\{ \frac{1+\lambda}{\lambda} (\alpha_i + \beta_j) - \frac{c_{ij}}{\lambda} - \frac{1+\lambda}{\lambda} \psi_\lambda \right\} \quad (28)$$

$$= \mathbf{p} \cdot \boldsymbol{\alpha} + \mathbf{q} \cdot \boldsymbol{\beta} - \psi_\lambda(\boldsymbol{\alpha}, \boldsymbol{\beta}). \quad (29)$$

When we use new notations $\boldsymbol{\eta} = (\mathbf{p}, \mathbf{q})^T$, $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\beta})^T$, we have

$$\psi_\lambda(\boldsymbol{\theta}) + \varphi_\lambda(\boldsymbol{\eta}) = \boldsymbol{\theta} \cdot \boldsymbol{\eta}, \quad (30)$$

which is the Legendre relation between $\boldsymbol{\theta}$ and $\boldsymbol{\eta}$. Thus, we have the following theorem.

Theorem 2. The relaxed cost function φ_λ and the free energy (cumulant generating function) ψ_λ of the exponential family are both convex, connected by the Legendre transformation,

$$\boldsymbol{\theta} = \nabla_{\boldsymbol{\eta}} \varphi_\lambda(\boldsymbol{\eta}), \quad \boldsymbol{\eta} = \nabla_{\boldsymbol{\theta}} \psi_\lambda(\boldsymbol{\theta}), \quad (31)$$

$$\boldsymbol{\alpha} = \nabla_{\mathbf{p}} \varphi_\lambda(\mathbf{p}, \mathbf{q}), \quad \boldsymbol{\beta} = \nabla_{\mathbf{q}} \varphi_\lambda(\mathbf{p}, \mathbf{q}), \quad (32)$$

$$\mathbf{p} = \nabla_{\boldsymbol{\alpha}} \psi_\lambda(\boldsymbol{\alpha}, \boldsymbol{\beta}), \quad \mathbf{q} = \nabla_{\boldsymbol{\beta}} \psi_\lambda(\boldsymbol{\alpha}, \boldsymbol{\beta}). \quad (33)$$

The Riemannian metric \mathbf{G}_λ is given to $\mathcal{S}_{n-1} \times \mathcal{S}_{n-1}$ by

$$\mathbf{G}_\lambda = \nabla_\eta \nabla_\eta \varphi_\lambda(\boldsymbol{\eta}) \quad (34)$$

in the $\boldsymbol{\eta}$ -coordinate system (\mathbf{p}, \mathbf{q}) . Its inverse is

$$\mathbf{G}_\lambda^{-1} = \nabla_\theta \nabla_\theta \psi_\lambda(\boldsymbol{\theta}). \quad (35)$$

In addition, we can calculate \mathbf{G}_λ^{-1} explicitly from (24).

Theorem 3. The Fisher information matrix \mathbf{G}_λ^{-1} in the $\boldsymbol{\theta}$ -coordinate system is given by

$$\mathbf{G}_\lambda^{-1} = \frac{1+\lambda}{\lambda} \left\{ \begin{bmatrix} \text{diag}(\mathbf{p}) & \mathbf{P} \\ \mathbf{P}^T & \text{diag}(\mathbf{q}) \end{bmatrix} - \begin{bmatrix} \mathbf{p} \\ \mathbf{q} \end{bmatrix} \begin{bmatrix} \mathbf{p}^T & \mathbf{q}^T \end{bmatrix} \right\}, \quad (36)$$

or in the component form as

$$\mathbf{G}_\lambda^{-1} = \frac{1+\lambda}{\lambda} \begin{bmatrix} p_i \delta_{ij} - p_i p_j & P_{ij} - p_i q_j \\ P_{ij} - p_i q_j & q_i \delta_{ij} - q_i q_j \end{bmatrix}. \quad (37)$$

Remark 1. The \mathbf{p} -part of \mathbf{G}_λ^{-1} is a scalar multiple of the Fisher information of \mathbf{p} in \mathcal{S}_{n-1} in the e -coordinate system. So is the \mathbf{q} -part. They are independent of the cost matrix c_{ij} , but the off-diagonal blocks of \mathbf{G}_λ^{-1} depend on it.

Remark 2. The \mathbf{p} -part of \mathbf{G}_λ is not equal to the Fisher information of \mathbf{p} in the m -coordinate system. It is the \mathbf{p} -part of the inverse of \mathbf{G}_λ^{-1} , depending on \mathbf{q} , too.

5 λ -Divergence in \mathcal{S}_{n-1}

The relaxed Wasserstein distance $\varphi_\lambda(\mathbf{p} : \mathbf{q})$ does not satisfy a criterion of divergence, i.e. $\varphi_\lambda(\mathbf{p} : \mathbf{p}) \neq 0$, because $\varphi_\lambda(\mathbf{p} : \mathbf{q})$ is minimized at $\mathbf{q} \neq \mathbf{p}$ in general. In contrast, the original Wasserstein distance literally satisfies the criteria of distance and those of divergence. To recover the property of divergence in the relaxed form, we introduce a canonical divergence between two transportation plans (\mathbf{p}, \mathbf{p}) and (\mathbf{p}, \mathbf{q}) , which is composed of the Legendre pair of the convex functions φ_λ and ψ_λ (Amari 2016):

$$D_\lambda[\mathbf{p} : \mathbf{q}] = \psi_\lambda(\boldsymbol{\alpha}, \boldsymbol{\beta}) + \varphi_\lambda(\mathbf{p}, \mathbf{p}) - \boldsymbol{\alpha} \cdot \mathbf{p} - \boldsymbol{\beta} \cdot \mathbf{p}, \quad (38)$$

where $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ corresponds to (\mathbf{p}, \mathbf{q}) . We call this a λ -divergence $D_\lambda[\mathbf{p} : \mathbf{q}]$ in \mathcal{S}_{n-1} from \mathbf{p} to \mathbf{q} . It connects the Wasserstein distance and the KL-divergence in the following way.

This λ -divergence can be transformed into a Bregman-like divergence with the relaxed cost function φ_λ (not a Bregman divergence constructed from a convex function of a single variable \mathbf{q}):

$$D_\lambda[\mathbf{p} : \mathbf{q}] = \varphi_\lambda(\mathbf{p}, \mathbf{p}) - \varphi_\lambda(\mathbf{p}, \mathbf{q}) - \langle \nabla_{\mathbf{q}} \varphi_\lambda(\mathbf{p}, \mathbf{q}), \mathbf{p} - \mathbf{q} \rangle. \quad (39)$$

As easily confirmed by substituting (14) to (38), the λ -divergence is equivalent to the KL-divergence between the two transportation plans, up to a constant factor:

$$D_\lambda[\mathbf{p} : \mathbf{q}] = \frac{\lambda}{1 + \lambda} KL[\mathbf{P}' : \mathbf{P}], \quad (40)$$

where \mathbf{P}' and \mathbf{P} are the optimal plans from \mathbf{p} to \mathbf{p} and \mathbf{p} to \mathbf{q} , respectively. It is easy to see that $D_\lambda[\mathbf{p} : \mathbf{q}]$ satisfies the criteria of divergence. However, it is not dually flat in general.

Let us consider the case of $\lambda \rightarrow \infty$. Then,

$$\mathbf{P}' = (p_i p_j), \quad \mathbf{P} = (p_i q_j), \quad (41)$$

and hence

$$D_\lambda[\mathbf{p} : \mathbf{q}] = KL[\mathbf{p} : \mathbf{q}], \quad (42)$$

converging to the KL-divergence of S_{n-1} .

6 Conclusions

We have opened a new way of studying the geometry of probability distributions. We showed that the entropy-relaxed transportation plan in a probability simplex naturally defines the exponential family and the dually flat structure of information geometry. We also introduced a one-parameter family which connects the relaxed Wasserstein distance to the KL-divergence.

It remains as future problems to extend the information geometry of the relaxed Wasserstein distance into a general case of continuous distributions on a metric manifold. Another direction of research is to study the geometrical properties of the manifold through the new family of λ -divergence and to apply it to various practical applications, where some modifications of D_λ might be useful.

References

- Amari, S.: Information Geometry and Its Applications. Springer, Tokyo (2016)
- Chentsov, N.N.: Statistical Decision Rules and Optimal Inference. Nauka (1972). Translated in English. AMS (1982)
- Cuturi, M.: Sinkhorn distances: light speed computation of optimal transport. In: Advances in Neural Information Processing Systems, pp. 2292–2300 (2013)
- Cuturi, M., Avis, D.: Ground metric learning. J. Mach. Learn. Res. **15**, 533–564 (2014)
- Cuturi, M., Peyré, G.: A smoothed dual formulation for variational Wasserstein problems. SIAM J. Imaging Sci. **9**, 320–343 (2016)
- Rao, C.R.: Information and accuracy attainable in the estimation of statistical parameters. Bull. Calcutta Math. Soc. **37**, 81–91 (1945)
- Villani, C.: Topics in Optimal Transportation. Graduate Studies in Math. AMS (2013)