

Course notes on Computational Optimal Transport

Gabriel Peyré
CNRS & DMA
École Normale Supérieure
gabriel.peyre@ens.fr
<https://mathematical-tours.github.io>
www.numerical-tours.com

October 13, 2019

Abstract

These note cours are intended to complement the book [37] with more details on the theory of Optimal Transport. Many parts are extracted from this book, with some additions and re-writing.

Contents

1	Optimal Matching between Point Clouds	2
1.1	Monge Problem between Discrete points	2
1.2	Matching Algorithms	3
2	Monge Problem between Measures	3
2.1	Measures	3
2.2	Push Forward	5
2.3	Monge's Formulation	6
2.4	Existence and Uniqueness of the Monge Map	7
3	Kantorovitch Relaxation	10
3.1	Discrete Relaxation	10
3.2	Relaxation for Arbitrary Measures	13
3.3	Metric Properties	15
4	Sinkhorn	17
4.1	Entropic Regularization for Discrete Measures	17
4.2	General Formulation	19
4.3	Sinkhorn's Algorithm	19
4.4	Convergence	21
5	Dual Problem	22
5.1	Discrete dual	22
5.2	General formulation	23
5.3	c -transforms	24

6	Semi-discrete and W_1	25
6.1	Semi-discrete	25
6.2	W_1	28
6.3	Dual norms (Integral Probability Metrics)	30
6.4	φ -divergences	32
7	Sinkhorn Divergences	34
7.1	Dual of Sinkhorn	34
7.2	Sinkhorn Divergences	35
8	Barycenters	37
8.1	Frechet Mean over the Wasserstein Space	37
8.2	1-D Case	38
8.3	Gaussians Case	38
8.4	Discrete Barycenters	38
8.5	Sinkhorn for barycenters	38
9	Wasserstein Estimation	40
9.1	Wasserstein Loss	40
9.2	Wasserstein Derivatives	41
9.3	Sample Complexity	41
10	Gradient Flows	42
10.1	Optimization over Measures	42
10.2	Particle System and Lagrangian Flows	42
10.3	Wasserstein Gradient Flows	42
10.4	Langevin Flows	42
11	Extensions	43
11.1	Dynamical formulation	43
11.2	Unbalanced OT	43
11.3	Gromov Wasserstein	43
11.4	Quantum OT	44

1 Optimal Matching between Point Clouds

1.1 Monge Problem between Discrete points

Matching problem Given a cost matrix $(\mathbf{C}_{i,j})_{i \in \llbracket n \rrbracket, j \in \llbracket m \rrbracket}$, assuming $n = m$, the optimal assignment problem seeks for a bijection σ in the set $\text{Perm}(n)$ of permutations of n elements solving

$$\min_{\sigma \in \text{Perm}(n)} \frac{1}{n} \sum_{i=1}^n \mathbf{C}_{i, \sigma(i)}. \quad (1)$$

One could naively evaluate the cost function above using all permutations in the set $\text{Perm}(n)$. However, that set has size $n!$, which is gigantic even for small n . In general the optimal σ is non-unique.

1D case If the cost is of the form $\mathbf{C}_{i,j} = h(x_i - y_j)$, where $h : \mathbb{R} \rightarrow \mathbb{R}^+$ is convex (for instance $\mathbf{C}_{i,j} = |x_i - y_j|^p$ for $p \geq 1$), one has that an optimal σ necessarily defines an increasing map $x_i \mapsto x_{\sigma(i)}$, i.e.

$$\forall (i, j), \quad (x_i - y_j)(x_{\sigma(i)} - y_{\sigma(j)}) \geq 0.$$

Indeed, if this property is violated, i.e. there exists (i, j) such that $(x_i - y_j)(x_{\sigma(i)} - y_{\sigma(j)}) < 0$, then one can defines a permutation $\tilde{\sigma}$ by swapping the match, i.e. $\tilde{\sigma}(i) = \sigma(j)$ and $\tilde{\sigma}(j) = \sigma(i)$, with a better cost

$$\sum_i h(x_i - y_{\tilde{\sigma}(i)}) \leq \sum_i h(x_i - y_{\sigma(i)}),$$

because

$$h(x_i - y_{\sigma(j)}) + h(x_j - y_{\sigma(i)}) \leq h(x_i - y_{\sigma(i)}) + h(x_j - y_{\sigma(j)}).$$

So the algorithm to compute an optimal transport (actually all optimal transport) is to sort the points, i.e. find some pair of permutations σ_X, σ_Y such that

$$x_{\sigma_X(1)} \leq x_{\sigma_X(2)} \leq \dots \quad \text{and} \quad y_{\sigma_Y(1)} \leq y_{\sigma_Y(2)} \leq \dots$$

and then an optimal match is mapping $x_{\sigma_X(k)} \mapsto y_{\sigma_Y(k)}$, i.e. an optimal transport is $\sigma = \sigma_Y \circ \sigma_X^{-1}$. The total computational cost is thus $O(n \log(n))$ using for instance quicksort algorithm. Note that if $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is an increasing map, with a change of variable, one can apply this technique to cost of the form $h(|\varphi(x) - \varphi(y)|)$. A typical application is grayscale histogram equalization of the luminance of images.

Note that if h is concave instead of being convex, then the behavior is totally different, and the optimal match actually rather exchange the positions, and in this case there exists an $O(n^2)$ algorithm.

1.2 Matching Algorithms

There exists efficient algorithms to solve the optimal matching problems. The most well known are the hungarian and the auction algorithm, which runs in $O(n^3)$ operations. Their derivation and analysis is however very much simplified by introducing the Kantorovitch relaxation and its associated dual problem. A typical application of these methods is the equalization of the color palette between images, which corresponds to a 3-D optimal transport.

2 Monge Problem between Measures

2.1 Measures

Histograms We will interchangeably the term histogram or probability vector for any element $\mathbf{a} \in \Sigma_n$ that belongs to the probability simplex

$$\Sigma_n \stackrel{\text{def.}}{=} \left\{ \mathbf{a} \in \mathbb{R}_+^n ; \sum_{i=1}^n \mathbf{a}_i = 1 \right\}.$$

Discrete measure, empirical measure A discrete measure with weights \mathbf{a} and locations $x_1, \dots, x_n \in \mathcal{X}$ reads

$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i} \tag{2}$$

where δ_x is the Dirac at position x , intuitively a unit of mass which is infinitely concentrated at location x . Such as measure describes a probability measure if, additionally, $\mathbf{a} \in \Sigma_n$, and more generally a positive measure if each of the “weights” described in vector \mathbf{a} is positive itself. An “empirical” probability distribution is uniform on a point cloud, i.e. $\mathbf{a} = \frac{1}{n} \sum_i \delta_{x_i}$. In practice, it many application is useful to be able to manipulate both the positions x_i (“Lagrangian” discretization) and the weights \mathbf{a}_i (“Eulerian” discretization). Lagrangian modification is usually more powerful (because it leads to adaptive discretization) but it breaks the convexity of most problems.

General measures We consider Borel measures $\alpha \in \mathcal{M}(\mathcal{X})$ on a metric space (\mathcal{X}, d) , i.e. one can compute $\alpha(A)$ for any Borel set A (which can be obtained by applying countable union, countable intersection, and relative complement to open sets). The measure should be finite, i.e. have a finite value on compact set. A Dirac measure δ_x is then define as $\delta_x(A) = 1$ if $x \in A$ and 0 otherwise, and this extend by linearity for discrete measures of the form (2) as

$$\alpha(A) = \sum_{x_i \in A} \mathbf{a}_i$$

We denote $\mathcal{M}_+(\mathcal{X})$ the subset of all positive measures on \mathcal{X} , i.e. $\alpha(A) \geq 0$ (and $\alpha(\mathcal{X}) < +\infty$ for the measure to be finite). The set of probability measures is denoted $\mathcal{M}_+^1(\mathcal{X})$, which means that any $\alpha \in \mathcal{M}_+^1(\mathcal{X})$ is positive, and that $\alpha(\mathcal{X}) = 1$.

Radon measures Using Lebesgue integration, a Borel measure can be used to compute integral of measurable functions (i.e. such that level sets $\{x ; f(x) < t\}$ are Borel sets), and we denote this pairing as

$$\langle f, \alpha \rangle \stackrel{\text{def.}}{=} \int f(x) d\alpha(x).$$

Integration of such a measurable f against a discrete measure α computes a sum

$$\int_{\mathcal{X}} f(x) d\alpha(x) = \sum_{i=1}^n \mathbf{a}_i f(x_i).$$

This can be in particular applied to the subspace of continuous functions which are measurable. Integration against a finite measure on a compact space thus defines a continuous linear form $f \mapsto \int f d\alpha$ on the Banach space of continuous functions $(\mathcal{C}(\mathcal{X}), \|\cdot\|_\infty)$, indeed $|\int f d\alpha| \leq \|f\|_\infty |\alpha(\mathcal{X})|$. On compact spaces, the converse is true, namely that any continuous linear form $\ell : f \mapsto \ell(f)$ on $(\mathcal{C}(\mathcal{X}), \|\cdot\|_\infty)$ is represented as an integral against a measure $\ell(f) = \int f d\alpha$. This is the Riesz-Markov-Kakutani representation theorem, which is often stated that Borel measures can be identified to Radon measures. Radon measures are thus in some sense “less regular” than functions, but more regular than distributions (which are dual to smooth functions). For instance, the derivative of a Dirac is not a measure. This duality pairing $\langle f, \alpha \rangle$ between continuous function and measures will be crucial to develop duality theory for the convex optimization problem we will consider later.

The associated norm, which is the norm of the linear form ℓ , is the so-called total variation norm

$$\|\alpha\|_{TV} = \|\ell\|_{\mathcal{C}(\mathcal{X}) \rightarrow \mathbb{R}} = \sup_{f \in \mathcal{C}(\mathcal{X})} \{ \langle f, \alpha \rangle ; \|f\|_\infty \leq 1 \}.$$

(note that one can remove the $|\cdot|$ in the right hand side, and such a quantity is often called a “dual norm”). One can in fact show that this TV norm is the total mass of the absolute value measure $|\alpha|$. The space $(\mathcal{M}(\mathcal{X}), \|\cdot\|_{TV})$ is a Banach space, which is the dual of $(\mathcal{C}(\mathcal{X}), \|\cdot\|_\infty)$.

Recall that the absolute value of a measure is defined as

$$|\alpha|(A) = \sup_{A = \cup_i B_i} \sum_i |\alpha(B_i)|$$

so that for instance if $\alpha = \sum_i \mathbf{a}_i \delta_{x_i}$, $|\alpha| = \sum_i |\mathbf{a}_i| \delta_{x_i}$ and if $d\alpha(x) = \rho dx$ for a positif reference measure dx , then $d|\alpha|(x) = |\rho(x)| dx$.

Relative densities A measure α which is a weighting of another reference one dx is said to have a density, which is denoted $d\alpha(x) = \rho_\alpha(x) dx$ (on \mathbb{R}^d dx is often the Lebesgue measure), often also denoted $\rho_\alpha = \frac{d\alpha}{dx}$, which means that

$$\forall h \in \mathcal{C}(\mathbb{R}^d), \quad \int_{\mathbb{R}^d} h(x) d\alpha(x) = \int_{\mathbb{R}^d} h(x) \rho_\alpha(x) dx.$$

Probabilistic interpretation Radon probability measures can also be viewed as representing the distributions of random variables. A random variable X on \mathcal{X} is actually a map $X : \Omega \rightarrow \mathcal{X}$ from some abstract (often un-specified) probabized space (Ω, \mathbb{P}) , and its distribution is the Radon measure $\alpha \in \mathcal{M}_+^1(\mathcal{X})$ such that $\mathbb{P}(X \in A) = \alpha(A) = \int_A d\alpha(x)$.

2.2 Push Forward

For some continuous map $T : \mathcal{X} \rightarrow \mathcal{Y}$, we define the pushforward operator $T_\# : \mathcal{M}(\mathcal{X}) \rightarrow \mathcal{M}(\mathcal{Y})$. For a Dirac mass, one has $T_\# \delta_x = \delta_{T(x)}$, and this formula is extended to arbitrary measure by linearity. In some sense, moving from T to $T_\#$ is a way to linearize any map at the prize of moving from a (possibly) finite dimensional space \mathcal{X} to the infinite dimensional space $\mathcal{M}(\mathcal{X})$, and this idea is central to many convex relaxation method, most notably Lasserre's relaxation. For discrete measures (2), the pushforward operation consists simply in moving the positions of all the points in the support of the measure

$$T_\# \alpha \stackrel{\text{def.}}{=} \sum_i \mathbf{a}_i \delta_{T(x_i)}.$$

For more general measures, for instance for those with a density, the notion of push-forward plays a fundamental to describe spatial modifications of probability measures. The formal definition reads as follow.

Definition 1 (Push-forward). *For $T : \mathcal{X} \rightarrow \mathcal{Y}$, the push forward measure $\beta = T_\# \alpha \in \mathcal{M}(\mathcal{Y})$ of some $\alpha \in \mathcal{M}(\mathcal{X})$ satisfies*

$$\forall h \in \mathcal{C}(\mathcal{Y}), \quad \int_{\mathcal{Y}} h(y) d\beta(y) = \int_{\mathcal{X}} h(T(x)) d\alpha(x). \quad (3)$$

Equivalently, for any measurable set $B \subset \mathcal{Y}$, one has

$$\beta(B) = \alpha(\{x \in \mathcal{X} ; T(x) \in B\}). \quad (4)$$

Note that $T_\#$ preserves positivity and total mass, so that if $\alpha \in \mathcal{M}_+^1(\mathcal{X})$ then $T_\# \alpha \in \mathcal{M}_+^1(\mathcal{Y})$.

Remark 1 (Push-forward for densities). Explicitly doing the change of variable $x = T(y)$, so that $dx = |\det(T'(y))| dy$ in formula (3) for measures with densities $(\rho_\alpha, \rho_\beta)$ on \mathbb{R}^d (assuming T is smooth and a bijection), one has for all $h \in \mathcal{C}(\mathcal{Y})$

$$\begin{aligned} \int_{\mathcal{Y}} h(y) \rho_\beta(y) dy &= \int_{\mathcal{Y}} h(y) d\beta(y) = \int_{\mathcal{X}} h(T(x)) d\alpha(x) = \int_{\mathcal{X}} h(T(x)) \rho_\alpha(x) dx \\ &= \int_{\mathcal{Y}} h(y) \rho_\alpha(T^{-1}y) \frac{dy}{|\det(T'(T^{-1}y))|}, \end{aligned}$$

which shows that

$$\rho_\beta(y) = \rho_\alpha(T^{-1}y) \frac{1}{|\det(T'(T^{-1}y))|}.$$

Since T is a diffeomorphism, one obtains equivalently

$$\rho_\alpha(x) = |\det(T'(x))| \rho_\beta(T(x)) \quad (5)$$

where $T'(x) \in \mathbb{R}^{d \times d}$ is the Jacobian matrix of T (the matrix formed by taking the gradient of each coordinate of T). This implies, denoting $y = T(x)$

$$|\det(T'(x))| = \frac{\rho_\alpha(x)}{\rho_\beta(y)}.$$

Remark 2 (Probabilistic interpretation). A random variable X , equivalently, is the push-forward of \mathbb{P} by X , $\alpha = X_\# \mathbb{P}$. Applying another push-forward $\beta = T_\# \alpha$ for $T : \mathcal{X} \rightarrow \mathcal{Y}$, following (3), is equivalent to defining another random variable $Y = T(X) : \omega \in \Omega \rightarrow T(X(\omega)) \in \mathcal{Y}$, so that β is the distribution of Y . Drawing a random sample y from Y is thus simply achieved by computing $y = T(x)$ where x is drawn from X .

2.3 Monge's Formulation

Monge problem. Monge problem (1) is extended to the setting of two arbitrary probability measures (α, β) on two spaces $(\mathcal{X}, \mathcal{Y})$ as finding a map $T : \mathcal{X} \rightarrow \mathcal{Y}$ that minimizes

$$\inf_T \left\{ \int_{\mathcal{X}} c(x, T(x)) d\alpha(x) ; T_{\#}\alpha = \beta \right\}. \quad (6)$$

The constraint $T_{\#}\alpha = \beta$ means that T pushes forward the mass of α to β , and makes use of the push-forward operator (3).

For empirical measure with same number $n = m$ of points, one retrieves the optimal matching problem. Indeed, this corresponds to the setting of empirical measures $\alpha = \sum_i \delta_{x_i}$ and $\beta = \sum_i \delta_{y_i}$. In this case, $T_{\#}\alpha = \beta$ necessarily implies that σ is one-to-one, $T : x_i \mapsto x_{\sigma(i)}$, so that

$$\int_{\mathcal{X}} c(x, T(x)) d\alpha(x) = \sum_i c(x_i, x_{\sigma(i)}).$$

In general, an optimal map T solving (6) might fail to exist. In fact, the constraint set $T_{\#}\alpha = \beta$, which is the case for instance if $\alpha = \delta_x$ and β is not a single Dirac. Even if the constraint set is not empty the infimum might not be reached, the most celebrated example being the case of α being distributed uniformly on a single segment and β being distributed on two segments on the two sides.

Monge distance. In the special case $c(x, y) = d^p(x, y)$ where d is a distance, we denote

$$\tilde{\mathcal{W}}_p^p(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_T \left\{ \mathcal{E}_{\alpha}(T) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} d(x, T(x))^p d\alpha(x) ; T_{\#}\alpha = \beta \right\}. \quad (7)$$

If the constraint set is empty, then we set $\tilde{\mathcal{W}}_p^p(\alpha, \beta) = +\infty$. The following proposition shows that quantity defines a distance.

Proposition 1. $\tilde{\mathcal{W}}$ is a distance.

Proof. If $\tilde{\mathcal{W}}_p^p(\alpha, \beta) = 0$ then necessarily the optimal map is Id on the support of α and $\beta = \alpha$. Let us prove that $\tilde{\mathcal{W}}_p^p(\alpha, \beta) \leq \tilde{\mathcal{W}}_p^p(\alpha, \gamma) + \tilde{\mathcal{W}}_p^p(\gamma, \beta)$. If $\tilde{\mathcal{W}}_p^p(\alpha, \beta) = +\infty$, then either $\tilde{\mathcal{W}}_p^p(\alpha, \gamma) = +\infty$ or $\tilde{\mathcal{W}}_p^p(\gamma, \beta) = +\infty$, because otherwise we consider two maps (S, T) such that $S_{\#}\alpha = \gamma$ and $T_{\#}\gamma = \beta$ and then $(T \circ S)_{\#}\alpha = \beta$ so that $\tilde{\mathcal{W}}_p^p(\alpha, \beta) \leq \mathcal{E}_{\alpha}(S \circ T) < +\infty$. So necessarily $\tilde{\mathcal{W}}_p^p(\alpha, \beta) < +\infty$ and we can restrict our attention to the cases where $\tilde{\mathcal{W}}_p^p(\alpha, \gamma) < +\infty$ and $\tilde{\mathcal{W}}_p^p(\gamma, \beta) < +\infty$ because otherwise the inequality is trivial. For any $\varepsilon > 0$, we consider ε -minimizer $S_{\#}\alpha = \gamma$ and $T_{\#}\gamma = \beta$ such that

$$E_{\alpha}(S)^{\frac{1}{p}} \leq \tilde{\mathcal{W}}_p(\alpha, \gamma) + \varepsilon \quad \text{and} \quad E_{\gamma}(T)^{\frac{1}{p}} \leq \tilde{\mathcal{W}}_p(\gamma, \beta) + \varepsilon.$$

Now we have that $(T \circ S)_{\#}\alpha = \beta$, so that one has, using sub-optimality of this map and the triangular inequality

$$\mathcal{W}_p(\alpha, \gamma) \leq \int d(x, T(S(x)))^p d\alpha(x)^{\frac{1}{p}} \leq \int (d(x, S(x)) + d(S(x), T(S(x))))^p d\alpha(x)^{\frac{1}{p}}.$$

The using Minkowski inequality

$$\mathcal{W}_p(\alpha, \gamma) \leq \int d(x, S(x))^p d\alpha(x)^{\frac{1}{p}} + \int d(S(x), T(S(x)))^p d\alpha(x)^{\frac{1}{p}} \leq \mathcal{W}_p(\alpha, \beta) + \mathcal{W}_p(\beta, \gamma) + 2\varepsilon.$$

Letting $\varepsilon \rightarrow 0$ gives the result. \square

2.4 Existence and Uniqueness of the Monge Map

Brenier's theorem. The following celebrated theorem of [12] ensures that in \mathbb{R}^d for $p = 2$, if at least one of the two inputs measures has a density, then Kantorovitch and Monge problems are equivalent.

Theorem 1 (Brenier). *In the case $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ and $c(x, y) = \|x - y\|^2$, if α has a density with respect to the Lebesgue measure, then there exists a unique optimal Monge map T . This map is characterized by being the unique gradient of a convex function $T = \nabla\varphi$ such that $(\nabla\varphi)_\# \alpha = \beta$.*

Its proof requires to study the relaxed Kantorovitch problems and its dual, so we defer it to later (Section 5.3).

Brenier's theorem, stating that an optimal transport map must be the gradient of a convex function, should be examined under the light that a convex function is a natural generalization of the notion of increasing functions in dimension more than one. For instance, the gradient of a convex function is a monotone gradient field in the sense

$$\forall (x, x') \in \mathbb{R}^d \times \mathbb{R}^d, \quad \langle \nabla\varphi(x) - \nabla\varphi(x'), x - x' \rangle \geq 0.$$

Note however that in dimension larger than 1, not all monotone fields are gradient of convex function. For instance, a rotation is monotone but can never be an optimal transport because a gradient field Ax defined by a linear map A is necessarily obtained by a symmetric matrix A . Indeed, such a linear field must be associated to a quadratic form $\varphi(x) = \langle Bx, x \rangle / 2$ and hence $A = \nabla\varphi = (B + B^\top) / 2$. Optimal transport can thus plays an important role to define quantile functions in arbitrary dimensions, which in turn is useful for applications to quantile regression problems [15].

Note also that this theorem can be extended in many directions. The condition that α has a density can be weakened to the condition that it does not give mass to “small sets” having Hausdorff dimension smaller than $d - 1$ (e.g. hypersurfaces). One can also consider costs of the form $c(x, y) = h(x - y)$ where h is a strictly convex smooth function, for for instance $c(x, y) = \|x - y\|^p$ with $1 < p < +\infty$.

Note that Brenier's theorem provides existence and uniqueness, but in general, the map T can be very irregular. Indeed, φ is in general non-smooth, but it is in fact convex and Lipschitz, so that $\nabla\varphi$ is actually well defined α -almost everywhere. Ensuring T to be smooth actually requires the target β to be regular, and more precisely its support must be convex.

If α does not have a density, then T might fail to exists and it should be replaced by a set-valued function included in $\partial\varphi$ which is now the sub-differential of a convex function, which might have singularity on a non-zero measure set. This means that T can “split” the mass by mapping to several locations $T(x) \subset \partial\varphi$. Actually, the condition that $T(x) \subset \partial\varphi(x)$ and $T_\# \alpha = \beta$ implies that the multi-map T defines a solution of Kantorovitch problem that will be studied later.

Monge-Ampère equation. For measures with densities, using (5), one obtains that φ is the unique (up to the addition of a constant) convex function which solves the following Monge-Ampère-type equation

$$\det(\partial^2\varphi(x))\rho_\beta(\nabla\varphi(x)) = \rho_\alpha(x) \tag{8}$$

where $\partial^2\varphi(x) \in \mathbb{R}^{d \times d}$ is the hessian of φ . The convexity constraint forces $\det(\partial^2\varphi(x)) \geq 0$ and is necessary for this equation to have a solution and be well-posed. The Monge-Ampère operator $\det(\partial^2\varphi(x))$ can be understood as a non-linear degenerate Laplacian. In the limit of small displacements, one can consider $\varphi(x) = \|x\|^2 / 2 + \varepsilon\psi$ so that $\nabla\varphi = \text{Id} + \varepsilon\nabla\psi$, one indeed recovers the Laplacian Δ as a linearization since for smooth maps

$$\det(\partial^2\varphi(x)) = 1 + \varepsilon\Delta\psi(x) + o(\varepsilon),$$

where we used the fact that $\det(\text{Id} + \varepsilon A) = 1 + \varepsilon \text{tr}(A) + o(\varepsilon)$.

OT in 1-D. For a measure α on \mathbb{R} , we introduce the cumulative function

$$\forall x \in \mathbb{R}, \quad \mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha, \quad (9)$$

which is a function $\mathcal{C}_\alpha : \mathbb{R} \rightarrow [0, 1]$. Its pseudo-inverse $\mathcal{C}_\alpha^{-1} : [0, 1] \rightarrow \mathbb{R} \cup \{-\infty\}$

$$\forall r \in [0, 1], \quad \mathcal{C}_\alpha^{-1}(r) = \min_x \{x \in \mathbb{R} \cup \{-\infty\} ; \mathcal{C}_\alpha(x) \geq r\}.$$

That function is also called the quantile function of α . The following proposition shows that these defines push-forward toward the uniform distribution \mathcal{U} on $[0, 1]$.

Proposition 2. *One has $(\mathcal{C}_\alpha)_\#^{-1}\mathcal{U} = \alpha$, where \mathcal{U} is the uniform distribution in $[0, 1]$. If α has a density, then $(\mathcal{C}_\alpha)_\#\alpha = \mathcal{U}$.*

Proof. For simplicity, we assume α has a strictly positive density, so that \mathcal{C}_α is a strictly increasing continuous function. Denoting $\gamma \stackrel{\text{def.}}{=} (\mathcal{C}_\alpha)_\#^{-1}\mathcal{U}$ we aim at proving $\gamma = \alpha$, which is equivalent to $\mathcal{C}_\gamma = \mathcal{C}_\alpha$. One has

$$\mathcal{C}_\gamma(x) = \int_{-\infty}^x d\gamma = \int_{\mathbb{R}} 1_{]-\infty, x]} d((\mathcal{C}_\alpha^{-1})_\#\mathcal{U}) = \int_0^1 1_{]-\infty, x]}(\mathcal{C}_\alpha^{-1}(z)) dz = \int_0^1 1_{[0, \mathcal{C}_\alpha(x)]}(z) dz = \mathcal{C}_\alpha(x)$$

where we use the fact that

$$-\infty \leq \mathcal{C}_\alpha^{-1}(z) \leq x \iff 0 \leq z \leq \mathcal{C}_\alpha(x).$$

□

If α has a density, this shows that the map

$$T = \mathcal{C}_\beta^{-1} \circ \mathcal{C}_\alpha \quad (10)$$

satisfies $T_\#\alpha = \beta$.

For the cost $c(x, y) = |x - y|^2$, since this T is increasing (hence the gradient of a convex function since we are in 1-D), by Brenier's theorem, T is the solution to Monge problem (at least if we impose that α has a density, otherwise it might lead to a solution of Kantorovitch problem by properly defining the pseudo-inverse). This closed form formula is also optimal for any cost of the form $h(|x - y|)$ for increasing h . For discrete measures, one cannot apply directly this reasoning (because α does not have a density), but if the measure are uniform on the same number of Dirac masses, then this approach is actually equivalent to the sorting formula.

Plugging this optimal map into the definition of the “Wasserstein” distance (we will see later that this quantity defines a distance), so that for any $p \geq 1$, one has

$$\mathcal{W}_p(\alpha, \beta)^p = \int_{\mathbb{R}} |x - \mathcal{C}_\beta^{-1}(\mathcal{C}_\alpha(x))| d\alpha(x) = \int_0^1 |\mathcal{C}_\alpha^{-1}(r) - \mathcal{C}_\beta^{-1}(r)|^p dr = \|\mathcal{C}_\alpha^{-1} - \mathcal{C}_\beta^{-1}\|_{L^p([0,1])}^p. \quad (11)$$

This formula is still valid for any measure (one can for instance approximate α by a measure with density). This formula means that through the map $\alpha \mapsto \mathcal{C}_\alpha^{-1}$, the Wasserstein distance is isometric to a linear space equipped with the L^p norm. For $p = 2$, the Wasserstein distance for measures on the real line is thus a Hilbertian metric. This makes the geometry of 1-D optimal transport very simple, but also very different from its geometry in higher dimensions, which is not Hilbertian.

For $p = 1$, one even has the simpler formula. Indeed, the previous formula is nothing more than the area between the two graphs of the copula, which can thus be computed by exchanging the role of the two axis, so that

$$\mathcal{W}_1(\alpha, \beta) = \|\mathcal{C}_\alpha - \mathcal{C}_\beta\|_{L^1(\mathbb{R})} = \int_{\mathbb{R}} |\mathcal{C}_\alpha(x) - \mathcal{C}_\beta(x)| dx = \int_{\mathbb{R}} \left| \int_{-\infty}^x d(\alpha - \beta) \right| dx. \quad (12)$$

which shows that \mathcal{W}_1 is a norm (see §?? for the generalization to arbitrary dimensions).

It is possible to define other type of norm which behave similarly (i.e. metrize the convergence in law), for instance $\|\mathcal{C}_\alpha - \mathcal{C}_\beta\|_{L^p(\mathbb{R})}$ define respectively the Wasserstein, Cramer (i.e. Sobolev) and Kolmogorov-Smirnov norms for $p = 1, 2, \infty$.

OT on 1-D Gaussians We first consider the case where $\alpha = \mathcal{N}(m_\alpha, s_\alpha^2)$ and $\beta = \mathcal{N}(m_\beta, s_\beta^2)$ are two Gaussians in \mathbb{R} . Then one verifies that

$$T(x) = \frac{s_\beta}{s_\alpha}(x - m_\alpha) + m_\beta$$

satisfies $T_\# \alpha = \beta$, furthermore it is the the derivative of the convex function

$$\varphi(x) = \frac{s_\beta}{2s_\alpha}(x - m_\alpha)^2 + m_\beta x,$$

so that according to Brenier's theorem, for the cost $c(x - y) = (x - y)^2$, T is the unique optimal transport, and the associated Monge distance is, after some computation

$$\tilde{W}_2^2(\alpha, \beta) = \int_{\mathbb{R}} \left(\frac{s_\beta}{s_\alpha}(x - m_\alpha) + m_\beta - x \right)^2 d\alpha(x) = (m_\alpha - m_\beta)^2 + (s_\alpha - s_\beta)^2.$$

This formula still holds for Dirac masses, i.e. if $s_\alpha = 0$ or $s_\beta = 0$. The OT geometry of Gaussians is thus the Euclidean distance on the half plane $(m, s) \in \mathbb{R} \times \mathbb{R}_+$. This should be contrasted with the geometry of KL, where singular Gaussians (for which $s = 0$) are infinitely distant.

OT on Gaussians If $\alpha = \mathcal{N}(\mathbf{m}_\alpha, \Sigma_\alpha)$ and $\beta = \mathcal{N}(\mathbf{m}_\beta, \Sigma_\beta)$ are two Gaussians in \mathbb{R}^d , we now look for an affine map

$$T : x \mapsto \mathbf{m}_\beta + A(x - \mathbf{m}_\alpha). \quad (13)$$

This map is the gradient of the convex function $\varphi(x) = \langle \mathbf{m}_\beta, x \rangle + \langle A(x - \mathbf{m}_\alpha), x - \mathbf{m}_\alpha \rangle / 2$ if and only if A is a symmetric positive matrix.

Proposition 3. *One has $T_\# \alpha = \beta$ if and only if*

$$A \Sigma_\alpha A = \Sigma_\beta. \quad (14)$$

Proof. Indeed, one simply has to notice that the change of variables formula (5) is satisfied since

$$\begin{aligned} \rho_\beta(T(x)) &= \det(2\pi \Sigma_\beta)^{-\frac{1}{2}} \exp(-\langle T(x) - \mathbf{m}_\beta, \Sigma_\beta^{-1}(T(x) - \mathbf{m}_\beta) \rangle) \\ &= \det(2\pi \Sigma_\beta)^{-\frac{1}{2}} \exp(-\langle x - \mathbf{m}_\alpha, A^T \Sigma_\beta^{-1} A(x - \mathbf{m}_\alpha) \rangle) \\ &= \det(2\pi \Sigma_\beta)^{-\frac{1}{2}} \exp(-\langle x - \mathbf{m}_\alpha, \Sigma_\alpha^{-1}(x - \mathbf{m}_\alpha) \rangle), \end{aligned}$$

and since T is a linear map we have that

$$|\det T'(x)| = \det A = \left(\frac{\det \Sigma_\beta}{\det \Sigma_\alpha} \right)^{\frac{1}{2}}$$

and we therefore recover $\rho_\alpha = |\det T'| \rho_\beta$ meaning $T_\# \alpha = \beta$. \square

Equation (14) is a quadratic equation on A . Using the square root of positive matrices, which is uniquely defined, one has

$$\Sigma_\alpha^{\frac{1}{2}} \Sigma_\beta \Sigma_\alpha^{\frac{1}{2}} = \Sigma_\alpha^{\frac{1}{2}} A \Sigma_\alpha A \Sigma_\alpha^{\frac{1}{2}} = (\Sigma_\alpha^{\frac{1}{2}} A \Sigma_\alpha^{\frac{1}{2}})^2,$$

so that this equation has a unique solution, given by

$$A = \Sigma_\alpha^{-\frac{1}{2}} \left(\Sigma_\alpha^{\frac{1}{2}} \Sigma_\beta \Sigma_\alpha^{\frac{1}{2}} \right)^{\frac{1}{2}} \Sigma_\alpha^{-\frac{1}{2}} = A^T.$$

Using Brenier's theorem [12], we conclude that T is optimal.

With additional calculations involving first and second order moments of ρ_α , we obtain that the transport cost of that map is

$$\tilde{\mathcal{W}}_2^2(\alpha, \beta) = \|\mathbf{m}_\alpha - \mathbf{m}_\beta\|^2 + \mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2 \quad (15)$$

where \mathcal{B} is the so-called Bures' metric [13] between positive definite matrices (see also [?, 24]),

$$\mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2 \stackrel{\text{def.}}{=} \text{tr} \left(\Sigma_\alpha + \Sigma_\beta - 2(\Sigma_\alpha^{1/2} \Sigma_\beta \Sigma_\alpha^{1/2})^{1/2} \right), \quad (16)$$

where $\Sigma^{1/2}$ is the matrix square root. One can show that \mathcal{B} is a distance on covariance matrices, and that \mathcal{B}^2 is convex with respect to both its arguments. In the case where $\Sigma_\alpha = \text{diag}(r_i)_i$ and $\Sigma_\beta = \text{diag}(s_i)_i$ are diagonals, the Bures metric is the Hellinger distance

$$\mathcal{B}(\Sigma_\alpha, \Sigma_\beta) = \|\sqrt{r} - \sqrt{s}\|_2.$$

3 Kantorovitch Relaxation

3.1 Discrete Relaxation

Monge discrete matching problem is problematic because it cannot be applied when $n \neq m$. One needs to take into account masses $(\mathbf{a}_i, \mathbf{b}_j)$ to handle this more general situation. Monge continuous formulation (6) using push-forward is also problematic because it can be the case that there is no transport map T such that $T_\# \alpha = \beta$, for instance when α is made of a single Dirac to be mapped to several Dirac. Associated to this, it is not symmetric with respect to exchange of α and β (one can map two Diracs to a single one, but not the other way). Also, these are non-convex optimization problem which are not simple to solve numerically.

The key idea of [32] is to relax the deterministic nature of transportation, namely the fact that a source point x_i can only be assigned to another, or transported to one and one location $T(x_i)$ only. Kantorovich proposes instead that the mass at any point x_i be potentially dispatched across several locations. Kantorovich moves away from the idea that mass transportation should be “deterministic” to consider instead a “probabilistic” (or “fuzzy”) transportation, which allows what is commonly known now as “mass splitting” from a source towards several targets. This flexibility is encoded using, in place of a permutation σ or a map T , a coupling matrix $\mathbf{P} \in \mathbb{R}_+^{n \times m}$, where $\mathbf{P}_{i,j}$ describes the amount of mass flowing from bin i (or point x_i) towards bin j (or point x_j), x_i towards y_j in the formalism of discrete measures $\alpha = \sum_i \mathbf{a}_i \delta_{x_i}$, $\beta = \sum_j \mathbf{b}_j \delta_{y_j}$. Admissible couplings are only constrained to satisfy the conservation of mass

$$\mathbf{U}(\mathbf{a}, \mathbf{b}) \stackrel{\text{def.}}{=} \left\{ \mathbf{P} \in \mathbb{R}_+^{n \times m} ; \mathbf{P} \mathbf{1}_m = \mathbf{a} \quad \text{and} \quad \mathbf{P}^T \mathbf{1}_n = \mathbf{b} \right\}, \quad (17)$$

where we used the following matrix-vector notation

$$\mathbf{P} \mathbf{1}_m = \left(\sum_j \mathbf{P}_{i,j} \right)_i \in \mathbb{R}^n \quad \text{and} \quad \mathbf{P}^T \mathbf{1}_n = \left(\sum_i \mathbf{P}_{i,j} \right)_j \in \mathbb{R}^m.$$

The set of matrices $\mathbf{U}(\mathbf{a}, \mathbf{b})$ is bounded, defined by $n + m$ equality constraints, and therefore a convex polytope (the convex hull of a finite set of matrices).

Additionally, whereas the Monge formulation is intrinsically asymmetric, Kantorovich's relaxed formulation is always symmetric, in the sense that a coupling \mathbf{P} is in $\mathbf{U}(\mathbf{a}, \mathbf{b})$ if and only if \mathbf{P}^T is in $\mathbf{U}(\mathbf{b}, \mathbf{a})$.

Kantorovich's optimal transport problem now reads

$$\mathbf{L}_C(\mathbf{a}, \mathbf{b}) \stackrel{\text{def.}}{=} \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \langle \mathbf{C}, \mathbf{P} \rangle \stackrel{\text{def.}}{=} \sum_{i,j} \mathbf{C}_{i,j} \mathbf{P}_{i,j}. \quad (18)$$

This is a linear program, and as is usually the case with such programs, its solutions are not necessarily unique.

Linear programming algorithms The reference algorithms to solve (??) are network simplexes. There exists instances of this method which scale like $O(n^3 \log n)$. Alternative include interior points, which are usually inferior on this particular type of linear program.

Permutation Matrices as Couplings We restrict our attention to the special case $n = m$ and $\mathbf{a}_i = \mathbf{b}_i = 1$ (up to a scaling by $1/n$, these are thus probability measures). In this case one can solve Monge optimal matching problem (1), and it is convenient to re-write it using permutation matrices. For a permutation $\sigma \in \text{Perm}(n)$, we write \mathbf{P}_σ for the corresponding permutation matrix,

$$\forall (i, j) \in \llbracket n \rrbracket^2, \quad (\mathbf{P}_\sigma)_{i,j} = \begin{cases} 1 & \text{if } j = \sigma_i, \\ 0 & \text{otherwise.} \end{cases} \quad (19)$$

We denote the set of permutation matrices as

$$\mathcal{P}_n \stackrel{\text{def.}}{=} \{\mathbf{P}_\sigma ; \sigma \in \text{Perm}(n)\},$$

which is a discrete, hence non-convex, set. One has

$$\langle \mathbf{C}, \mathbf{P}_\sigma \rangle = \sum_{i=1}^n \mathbf{C}_{i, \sigma_i}$$

so that (1) is equivalent to the non-convex optimization problem

$$\min_{\mathbf{P} \in \mathcal{P}_n} \langle \mathbf{C}, \mathbf{P} \rangle.$$

In contrast, one has that $\mathbf{U}(\mathbf{a}, \mathbf{b}) = \mathcal{B}_n$ is equal to the convex set of bistochastic matrices

$$\mathcal{B}_n \stackrel{\text{def.}}{=} \left\{ \mathbf{P} \in \mathbb{R}_+^{n \times n} ; \mathbf{P} \mathbf{1}_n = \mathbf{P}^\top \mathbf{1}_n = \mathbf{1}_n \right\}$$

so that Kantorovitch problem reads

$$\min_{\mathbf{P} \in \mathcal{B}_n} \langle \mathbf{C}, \mathbf{P} \rangle.$$

The set of permutation matrices is strictly included in the set of bistochastic matrices, and more precisely

$$\mathcal{P}_n = \mathcal{B}_n \cap \{0, 1\}^{n \times n}.$$

This shows that one has the following obvious relation between the cost of Monge and Kantorovitch problem

$$\min_{\mathbf{P} \in \mathcal{B}_n} \langle \mathbf{C}, \mathbf{P} \rangle \leq \min_{\mathbf{P} \in \mathcal{P}_n} \langle \mathbf{C}, \mathbf{P} \rangle.$$

We will now show that there is in fact an equality between these two costs, so that both problems are in some sense equivalent.

For this, we will make a detour through more general linear optimization problem of the form $\min_{\mathbf{P} \in \mathcal{C}} \langle \mathbf{C}, \mathbf{P} \rangle$ for some compact convex set \mathcal{C} . We first introduce the notion of extremal point, which are intuitively the vertices of \mathcal{C}

$$\text{Extr}(\mathcal{C}) \stackrel{\text{def.}}{=} \left\{ \mathbf{P} ; \forall (Q, R) \in \mathcal{C}^2, \mathbf{P} = \frac{Q + R}{2} \Rightarrow Q = R \right\}.$$

So to show that $\mathbf{P} \notin \text{Extr}(\mathcal{C})$ it suffices to split \mathbf{P} as $\mathbf{P} = \frac{Q + R}{2}$ with $Q \neq R$ and $(Q, R) \in \mathcal{C}^2$. We will assume the following fundamental result.

Proposition 4. *If \mathcal{C} is compact, then $\text{Extr}(\mathcal{C}) \neq \emptyset$.*

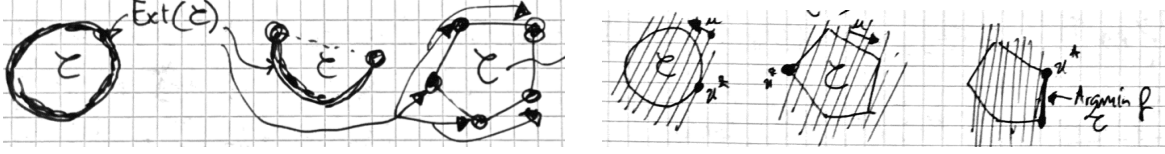


Figure 1: Left: extremal points of a convex set. Right: the solution of a convex program is a convex set.

The fact that \mathcal{C} is compact is crucial, for instance the set $\{(x, y) \in \mathbb{R}_+^2 ; xy \geq 1\}$ has no extremal point.

We can now use this result to show the following fundamental result, namely that there is always a solution to a linear program which is an extremal point. Note that of course the set of solution (which is non-empty because one minimizes a continuous function on a compact) might not be a singleton.

Proposition 5. *If \mathcal{C} is compact, then*

$$\text{Extr}(\mathcal{C}) \cap \left(\underset{\mathbf{P} \in \mathcal{C}}{\text{argmin}} \langle \mathbf{C}, \mathbf{P} \rangle \right) \neq \emptyset.$$

Proof. One consider $\mathcal{S} \stackrel{\text{def.}}{=} \underset{\mathbf{P} \in \mathcal{C}}{\text{argmin}} \langle \mathbf{C}, \mathbf{P} \rangle$. We first note that \mathcal{S} is convex (as always for an argmin) and compact, because \mathcal{C} is compact and the objective function is continuous, so that $\text{Extr}(\mathcal{S}) \neq \emptyset$. We will show that $\text{Extr}(\mathcal{S}) \subset \text{Extr}(\mathcal{C})$. **[ToDo: finish]** \square

The following theorem states that the extremal points of bistochastic matrices are the permutation matrices. It implies as a corollary that the cost of Monge and Kantorovitch are the same, and that they share a common solution.

Theorem 2 (Birkhoff and von Neumann). *One has $\text{Extr}(\mathcal{B}_n) = \mathcal{P}_n$.*

Proof. We first show the simplest inclusion $\mathcal{P}_n \subset \text{Extr}(\mathcal{B}_n)$. Indeed it follows from the fact that $\text{Extr}([0, 1]) = \{0, 1\}$. Take $\mathbf{P} \in \mathcal{P}_n$, if $\mathbf{P} = (Q + R)/2$ with $Q_{i,j}, R_{i,j} \in [0, 1]$, since $\mathbf{P}_{i,j} \in \{0, 1\}$ then necessarily $Q_{i,j} = R_{i,j} \in \{0, 1\}$.

Now we show $\text{Extr}(\mathcal{B}_n) \subset \mathcal{P}_n$ by showing that $\mathcal{P}_n^c \subset \text{Extr}(\mathcal{B}_n)^c$ where the complementary are computed inside the larger set \mathcal{B}_n . So picking $\mathbf{P} \in \mathcal{B}_n \setminus \mathcal{P}_n$, we need to split $\mathbf{P} = (Q + R)/2$ where Q, R are distinct bistochastic matrices. As shown on figure 2, \mathbf{P} defines a partite graph linking two sets of n vertices. This graph is composed of isolated edge when $\mathbf{P}_{i,j} = 1$ and connected edges corresponding to $0 < \mathbf{P}_{i,j} < 1$. If i is such a connected vertex on the left (similarly for j on the right), because $\sum_j \mathbf{P}_{i,j} = 1$, there is necessarily at least two edges (i, j_1) and (i, j_2) emating from it (similarly on the right there are at least two converging edges (i_1, j) and (i_2, j)). This means that by following these connexions, one necessarily can extract a cycle (if not, one could alway extend it by the previous remarks) of the form

$$(i_1, j_1, i_2, j_2, \dots, i_p, j_p), \quad \text{i.e.} \quad i_{p+1} = i_1.$$

We assume this cycle is the shortest one among all this (finite) ensemble of cycle. Along this cycle, the left-right and right-left edges satisfy

$$0 < \mathbf{P}_{i_s, j_s}, \mathbf{P}_{j_s, i_{s+1}} < 1.$$

The $(i_s)_s$ and $(j_s)_s$ are also all distincts because the cycle is the shortest. Lets pick

$$\varepsilon \stackrel{\text{def.}}{=} \min_{0 \leq s \leq p} \{\mathbf{P}_{i_s, j_s}, \mathbf{P}_{j_s, i_{s+1}}, 1 - \mathbf{P}_{i_s, j_s}, 1 - \mathbf{P}_{j_s, i_{s+1}}\}$$

so that $0 < \varepsilon < 1$. As shown on Figure 2, right, we split the graph in two set of edges, left-right and right-left

$$\mathcal{A} \stackrel{\text{def.}}{=} \{(i_s, j_s)\}_{s=1}^p \quad \text{and} \quad \mathcal{B} \stackrel{\text{def.}}{=} \{(j_s, i_{s+1})\}_{s=1}^p.$$

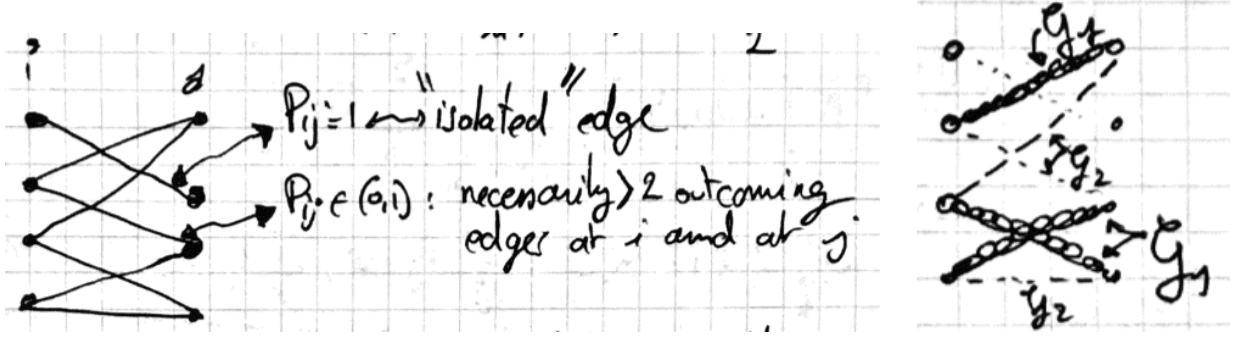


Figure 2: Left: the support of the coupling \mathbf{P} defines a bipartite graph. Right: splitting of this graph in two set of edges.

We define then two matrices as

$$Q_{i,j} \stackrel{\text{def.}}{=} \begin{cases} \mathbf{P}_{i,j} & \text{if } (i,j) \notin \mathcal{A} \cup \mathcal{B}, \\ \mathbf{P}_{i,j} + \varepsilon/2 & \text{if } (i,j) \in \mathcal{A}, \\ \mathbf{P}_{i,j} - \varepsilon/2 & \text{if } (i,j) \in \mathcal{B}, \end{cases} \quad \text{and} \quad R_{i,j} \stackrel{\text{def.}}{=} \begin{cases} \mathbf{P}_{i,j} & \text{if } (i,j) \notin \mathcal{A} \cup \mathcal{B}, \\ \mathbf{P}_{i,j} - \varepsilon/2 & \text{if } (i,j) \in \mathcal{A}, \\ \mathbf{P}_{i,j} + \varepsilon/2 & \text{if } (i,j) \in \mathcal{B}, \end{cases}$$

Because of the choice of ε , one has $0 \leq Q_{i,j}, R_{i,j} \leq 1$. Because each left-right edge in \mathcal{A} is associated to a right-left edge in \mathcal{B} , (and the other way) the sum constraint on the row (and on the column) is maintain, so that $U, V \in \mathcal{B}_n$. Finally, note that $\mathbf{P} = (P + Q)/2$. \square

By putting together Proposition 5 and Theorem 2, one obtains that for the discrete optimal problem with empirical measures, Monge and Kantoritch problems are equivalent.

Corollary 1 (Kantorovich for matching). *If $m = n$ and $\mathbf{a} = \mathbf{b} = \mathbf{1}_n$, then there exists an optimal solution for Problem (??) \mathbf{P}_{σ^*} , which is a permutation matrix associated to an optimal permutation $\sigma^* \in \text{Perm}(n)$ for Problem (1).*

The following proposition shows that these problems result in fact in the same optimum, namely that one can always find a permutation matrix that minimizes Kantorovich's problem (??) between two uniform measures $\mathbf{a} = \mathbf{b} = \mathbf{1}_n/n$, which shows that the Kantorovich relaxation is *tight* when considered on assignment problems.

3.2 Relaxation for Arbitrary Measures

Continuous couplings. The definition of \mathcal{L}_c in (18) is extended to arbitrary measures by considering couplings $\pi \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y})$ which are joint distributions over the product space. The marginal constraint $\mathbf{P}\mathbf{1}_m = \mathbf{a}, \mathbf{P}\mathbf{1}_n = \mathbf{b}$ must be replaced by “integrated” versions, which are written $\pi_1 = \alpha$ and $\pi_2 = \beta$, where $(\pi_1, \pi_2) \in \mathcal{M}(\mathcal{X}) \times \mathcal{M}(\mathcal{Y})$ are the two marginals. They are defined as $\pi_1 \stackrel{\text{def.}}{=} P_{1\#}\pi$ and $\pi_2 \stackrel{\text{def.}}{=} P_{2\#}\pi$ the two marginals of π , which are defined using push-forward by the projectors $P_1(x, y) = x$ and $P_2(x, y) = y$.

A heuristic way to understand the marginal constraint $\pi_1 = \alpha$ and $\pi_2 = \beta$, which mimics the discrete case where one sums along the rows and columns is to write

$$\int_{\mathcal{Y}} d\pi(x, y) = d\alpha(x) \quad \text{and} \quad \int_{\mathcal{X}} d\pi(x, y) = d\beta(y),$$

and the mathematically rigorous way to write this, which corresponds to the change of variables formula, is

$$\forall (f, g) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y}), \quad \int_{\mathcal{X} \times \mathcal{Y}} f(x) d\pi(x, y) = \int_{\mathcal{X}} f d\alpha \quad \text{and} \quad \int_{\mathcal{X} \times \mathcal{Y}} d\pi(x, y) = \int_{\mathcal{Y}} g d\beta.$$

Using (4), these marginal constraints are also equivalent to imposing that $\pi(A \times \mathcal{Y}) = \alpha(A)$ and $\pi(\mathcal{X} \times B) = \beta(B)$ for sets $A \subset \mathcal{X}$ and $B \subset \mathcal{Y}$.

In the general case, the mass conservation constraint (17) should thus be rewritten as a marginal constraint on joint probability distributions

$$\mathcal{U}(\alpha, \beta) \stackrel{\text{def.}}{=} \{ \pi \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y}) ; \pi_1 = \alpha \text{ and } \pi_2 = \beta \}. \quad (20)$$

The discrete case, when $\alpha = \sum_i \mathbf{a}_i \delta_{x_i}$, $\beta = \sum_j \mathbf{a}_j \delta_{x_j}$, the constraint $\pi_1 = \alpha$ and $\pi_2 = \beta$ necessarily imposes that π is discrete, supported on the set $\{(x_i, y_j)\}_{i,j}$, and thus has the form $\pi = \sum_{i,j} \mathbf{P}_{i,j} \delta_{(x_i, y_j)}$. The discrete formulation is thus a special case (and not some sort of approximation) of the continuous formulation.

Continuous Kantorovitch problem. The Kantorovich problem (18) is then generalized as

$$\mathcal{L}_c(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi \in \mathcal{U}(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y). \quad (21)$$

This is an infinite-dimensional linear program over a space of measures.

On compact domain $(\mathcal{X}, \mathcal{Y})$, (21) always has a solution, because using the weak-* topology (so called weak topology of measures), the set of measure is compact, and a linear function with a continuous $c(x, y)$ is weak-* continuous. And the set of constraint is non empty, taking $\alpha \otimes \beta$. On non compact domain, one needs to impose moment condition on α and β .

Probabilistic interpretation. If we denote $X \sim \alpha$ the fact that the law of a random vector X is the probability distribution α , then the marginal constraint appearing in (21) is simply that π is the law of a couple (X, Y) and that its coordinates X and Y have laws α and β . The coupling π encodes the statistical dependency between X and Y . For instance, $\pi = \alpha \otimes \beta$ means that X and Y are independent, and it is unlikely that such a coupling is optimal. Indeed as stated by Brenier's theorem, optimal coupling for a square Euclidean loss on contrary describe totally dependent variable.

With this remark, problem (21) reads equivalently

$$\mathcal{L}_c(\alpha, \beta) = \min_{X \sim \alpha, Y \sim \beta} \mathbb{E}(c(X, Y)). \quad (22)$$

Monge-Kantorovitch equivalence. The proof of Brenier theorem 1 (detailed in Section 5.3) to prove the existence of a Monge map actually studies Kantorovitch relaxation, and proves that this relaxation is tight in the sense that it has the same cost as Monge problem.

Indeed, if α has a density and we denote $T = \nabla \varphi$ the unique optimal transport, then the coupling

$$\pi = (\text{Id}, T)_\# \alpha \quad \text{i.e.} \quad \forall h \in \mathcal{C}(\mathcal{X} \times \mathcal{Y}), \quad \int_{\mathcal{X} \times \mathcal{Y}} h d\pi = \int_{\mathcal{X}} h(x, T(x)) d\alpha(x)$$

is optimal. In term of random vector, denoting (X, Y) a random vector with law π , it means that any such optimal random vector satisfies $Y = T(X)$ where $X \sim \alpha$ (and of course $T(X) \sim \beta$ by the marginal constraint).

This key result is similar to Birkoff-von-Neumann Theorem 1 in the sense that it provides conditions ensuring the equivalence between Monge and Kantorovitch problems (note however that Birkoff-von-Neumann does not implies uniqueness). Note however that the settings are radically difference (one is fully discrete while the other requires the sources to be “continuous”, i.e. to have a density).

3.3 Metric Properties

OT defines a distance. An important feature of OT is that it defines a distance between histograms and probability measures as soon as the cost matrix satisfies certain suitable properties. Indeed, OT can be understood as a canonical way to lift a ground distance between points to a distance between histogram or measures.

Proposition 6. *We suppose $n = m$, and that for some $p \geq 1$, $\mathbf{C} = \mathbf{D}^p = (\mathbf{D}_{i,j}^p)_{i,j} \in \mathbb{R}^{n \times n}$ where $\mathbf{D} \in \mathbb{R}_+^{n \times n}$ is a distance on $\llbracket n \rrbracket$, i.e.*

1. $\mathbf{D} \in \mathbb{R}_+^{n \times n}$ is symmetric;
2. $\mathbf{D}_{i,j} = 0$ if and only if $i = j$;
3. $\forall (i, j, k) \in \llbracket n \rrbracket^3, \mathbf{D}_{i,k} \leq \mathbf{D}_{i,j} + \mathbf{D}_{j,k}$.

Then

$$W_p(\mathbf{a}, \mathbf{b}) \stackrel{\text{def.}}{=} L_{\mathbf{D}^p}(\mathbf{a}, \mathbf{b})^{1/p} \quad (23)$$

(note that W_p depends on \mathbf{D}) defines the p -Wasserstein distance on Σ_n , i.e. W_p is symmetric, positive, $W_p(\mathbf{a}, \mathbf{b}) = 0$ if and only if $\mathbf{a} = \mathbf{b}$, and it satisfies the triangle inequality

$$\forall \mathbf{a}, \mathbf{a}', \mathbf{b} \in \Sigma_n, \quad W_p(\mathbf{a}, \mathbf{b}) \leq W_p(\mathbf{a}, \mathbf{a}') + W_p(\mathbf{a}', \mathbf{b}).$$

Proof. Symmetry and definiteness of the distance are easy to prove: since $\mathbf{C} = \mathbf{D}^p$ has a null diagonal, $W_p(\mathbf{a}, \mathbf{a}) = 0$, with corresponding optimal transport matrix $\mathbf{P}^* = \text{diag}(\mathbf{a})$; by the positivity of all off-diagonal elements of \mathbf{D}^p , $W_p(\mathbf{a}, \mathbf{b}) > 0$ whenever $\mathbf{a} \neq \mathbf{b}$ (because in this case, an admissible coupling necessarily has a non-zero element outside the diagonal); by symmetry of \mathbf{D}^p , $W_p(\mathbf{a}, \mathbf{b}) = 0$ is itself a symmetric function.

To prove the triangle inequality of Wasserstein distances for arbitrary measures, [44, Theorem 7.3] uses the gluing lemma, which stresses the existence of couplings with a prescribed structure. In the discrete setting, the explicit construction of this glued coupling is simple. Let $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \Sigma_n$. Let \mathbf{P} and \mathbf{Q} be two optimal solutions of the transport problems between \mathbf{a} and \mathbf{b} , and \mathbf{b} and \mathbf{c} respectively. We define $\bar{\mathbf{b}}_j \stackrel{\text{def.}}{=} \mathbf{b}_j$ if $\mathbf{b}_j > 0$ and set otherwise $\bar{\mathbf{b}}_j = 1$ (or actually any other value). We then define

$$\mathbf{S} \stackrel{\text{def.}}{=} \mathbf{P} \text{diag}(1/\bar{\mathbf{b}}) \mathbf{Q} \in \mathbb{R}_+^{n \times n}.$$

We remark that $\mathbf{S} \in U(\mathbf{a}, \mathbf{c})$ because

$$\mathbf{S} \mathbf{1}_n = \mathbf{P} \text{diag}(1/\bar{\mathbf{b}}) \mathbf{Q} \mathbf{1}_n = \mathbf{P}(\mathbf{b}/\bar{\mathbf{b}}) = \mathbf{P} \mathbf{1}_{\text{Supp}(\mathbf{b})} = \mathbf{a}$$

where we denoted $\mathbf{1}_{\text{Supp}(\mathbf{b})}$ the indicator of the support of \mathbf{b} , and we use the fact that $\mathbf{P} \mathbf{1}_{\text{Supp}(\mathbf{b})} = \mathbf{P} \mathbf{1} = \mathbf{b}$ because necessarily $\mathbf{P}_{i,j} = 0$ for $j \notin \text{Supp}(\mathbf{b})$. Similarly one verifies that $\mathbf{S}^\top \mathbf{1}_n = \mathbf{c}$.

The triangle inequality follows from

$$\begin{aligned}
W_p(\mathbf{a}, \mathbf{c}) &= \left(\min_{\mathbf{P} \in U(\mathbf{a}, \mathbf{c})} \langle \mathbf{P}, \mathbf{D}^p \rangle \right)^{1/p} \leq \langle \mathbf{S}, \mathbf{D}^p \rangle^{1/p} \\
&= \left(\sum_{ik} \mathbf{D}_{ik}^p \sum_j \frac{\mathbf{P}_{ij} \mathbf{Q}_{jk}}{\bar{\mathbf{b}}_j} \right)^{1/p} \leq \left(\sum_{ijk} (\mathbf{D}_{ij} + \mathbf{D}_{jk})^p \frac{\mathbf{P}_{ij} \mathbf{Q}_{jk}}{\bar{\mathbf{b}}_j} \right)^{1/p} \\
&\leq \left(\sum_{ijk} \mathbf{D}_{ij}^p \frac{\mathbf{P}_{ij} \mathbf{Q}_{jk}}{\bar{\mathbf{b}}_j} \right)^{1/p} + \left(\sum_{ijk} \mathbf{D}_{jk}^p \frac{\mathbf{P}_{ij} \mathbf{Q}_{jk}}{\bar{\mathbf{b}}_j} \right)^{1/p} \\
&= \left(\sum_{ij} \mathbf{D}_{ij}^p \mathbf{P}_{ij} \sum_k \frac{\mathbf{Q}_{jk}}{\bar{\mathbf{b}}_j} \right)^{1/p} + \left(\sum_{jk} \mathbf{D}_{jk}^p \mathbf{Q}_{jk} \sum_i \frac{\mathbf{P}_{ij}}{\bar{\mathbf{b}}_j} \right)^{1/p} \\
&= \left(\sum_{ij} \mathbf{D}_{ij}^p \mathbf{P}_{ij} \right)^{1/p} + \left(\sum_{jk} \mathbf{D}_{jk}^p \mathbf{Q}_{jk} \right)^{1/p} \\
&= W_p(\mathbf{a}, \mathbf{b}) + W_p(\mathbf{b}, \mathbf{c}).
\end{aligned}$$

The first inequality is due to the suboptimality of \mathbf{S} , the second is the usual triangle inequality for elements in \mathbf{D} , and the third comes from Minkowski's inequality. \square

Proposition 6 generalizes from histogram to arbitrary measures that need not be discrete.

Proposition 7. *We assume $\mathcal{X} = \mathcal{Y}$, and that for some $p \geq 1$, $c(x, y) = d(x, y)^p$ where d is a distance on \mathcal{X} , i.e.*

- (i) $d(x, y) = d(y, x) \geq 0$;
- (ii) $d(x, y) = 0$ if and only if $x = y$;
- (iii) $\forall (x, y, z) \in \mathcal{X}^3, d(x, z) \leq d(x, y) + d(y, z)$.

Then

$$\mathcal{W}_p(\alpha, \beta) \stackrel{\text{def.}}{=} \mathcal{L}_{d^p}(\alpha, \beta)^{1/p} \quad (24)$$

(note that \mathcal{W}_p depends on d) defines the p -Wasserstein distance on \mathcal{X} , i.e. \mathcal{W}_p is symmetric, positive, $\mathcal{W}_p(\alpha, \beta) = 0$ if and only if $\alpha = \beta$, and it satisfies the triangle inequality

$$\forall (\alpha, \beta, \gamma) \in \mathcal{M}_+^1(\mathcal{X})^3, \quad \mathcal{W}_p(\alpha, \gamma) \leq \mathcal{W}_p(\alpha, \beta) + \mathcal{W}_p(\beta, \gamma).$$

This distance \mathcal{W}_p defined though Kantorovitch problem (24) should be contrasted with the distance $\tilde{\mathcal{W}}$ obtained using Monge's problem (7). **Kantorovitch distance is always finite, while Monge's one might be infinite if the constraint set $\{T; T_{\#}\alpha = \beta\}$ is empty.** In fact, one can show that as soon as this constraint set is non-empty, and even if no optimal T exists, then one has $\mathcal{W}_p = \tilde{\mathcal{W}}_p$, which is a non-trivial result. Kantorovitch distance should thus be seen as a (convex) relaxation of Monge's distance, which behave in a much nicer way, as we will explore next (it is continuous with respect to the convergence in law topology).

Convergence in law topology. Let us first note that on a compact space, all \mathcal{W}_p distance defines the same topology (although they are not equivalent, the notion of converging sequence is the same).

Proposition 8. *On a compact space \mathcal{X} , one has for $p \leq q$*

$$\mathcal{W}_p(\alpha, \beta) \leq \mathcal{W}_q(\alpha, \beta) \leq \text{diam}(\mathcal{X})^{\frac{q-p}{q}} \mathcal{W}_p(\alpha, \beta)^{\frac{q}{p}}$$

Proof. The left inequality follows from Jensen inequality, $\varphi(\int c(x, y) d\pi(x, y)) \leq \int \varphi(c(x, y)) d\pi(x, y)$, applied to any probability distribution π and to the convex function $\varphi(r) = r^{q/p}$ to $c(x, y) = \|x - y\|^p$, so that one gets

$$\left(\int \|x - y\|^p d\pi(x, y) \right)^{\frac{q}{p}} \leq \int \|x - y\|^q d\pi(x, y).$$

The right inequality follows from

$$\|x - y\|^q \leq \text{diam}(\mathcal{X})^{q-p} \|x - y\|^p.$$

□

The Wasserstein distance \mathcal{W}_p has many important properties, the most important one being that it is a weak distance, *i.e.* it allows to compare singular distributions (for instance discrete ones) and to quantify spatial shift between the supports of the distributions. This corresponds to the notion of weak* convergence.

Definition 2 (Weak* topology). $(\alpha_k)_k$ converges weakly* to α in $\mathcal{M}_+^1(\mathcal{X})$ (denoted $\alpha_k \rightharpoonup \alpha$) if and only if for any continuous function $f \in \mathcal{C}(\mathcal{X})$, $\int_{\mathcal{X}} f d\alpha_k \rightarrow \int_{\mathcal{X}} f d\alpha$.

In term of random vectors, if $X_n \sim \alpha_n$ and $X \sim \alpha$ (not necessarily defined on the same probability space), the weak* convergence corresponds to the convergence in law of X_n toward X .

Definition 3 (Strong topology). The simplest distance on Radon measures is the total variation norm, which is the dual norm of the L^∞ norm on $\mathcal{C}(\mathcal{X})$ and whose topology is often called the “strong” topology

$$\|\alpha - \beta\|_{TV} \stackrel{\text{def.}}{=} \sup_{\|f\|_\infty \leq 1} \int f d(\alpha - \beta) = |\alpha - \beta|(\mathcal{X})$$

where $|\alpha - \beta|(\mathcal{X})$ is the mass of the absolute value of the difference measure. When $\alpha - \beta = \rho dx$ has a density, then $\|\alpha - \beta\|_{TV} = \int |\rho(x)| dx = \|\rho\|_{L^1(dx)}$ is the L^1 norm associated to dx . When $\alpha - \beta = \sum_i u_i \delta_{z_i}$ is discrete, then $\|\alpha - \beta\|_{TV} = \sum_i |u_i| = \|u\|_{\ell^1}$ is the discrete ℓ^1 norm.

In the special case of Diracs, having $\int f d\delta_{x_n} = f(x_n) \rightarrow \int f d\delta_x = f(x)$ for any continuous f is equivalent to $x_n \rightarrow x$. One can then contrast the strong topology with the Wasserstein distance, if $x_n \neq x$,

$$\|\delta_{x_n} - \delta_x\|_{TV} = 2 \quad \text{and} \quad W_p(\delta_{x_n}, \delta_x) = d(x_n, x).$$

This shows that for the strong topology, Diracs never converge, while they do converge for the Wasserstein distance. In fact it is a powerful property of the Wasserstein distance, which is regular with respect to the weak* topology, and metrizes it.

Proposition 9. If \mathcal{X} is compact, $\alpha_k \rightharpoonup \alpha$ if and only if $W_p(\alpha_k, \alpha) \rightarrow 0$.

The proof of this proposition requires the use of duality, and is delayed to later, see Proposition 2. On non-compact spaces, one needs also to impose the convergence of the moments up to order p . Note that there exists alternative distances which also metrize weak convergence. The simplest one are Hilbertian kernel norms, which are detailed in Section 6.3.

Applications and implications Applications for having a geometric distance : barycenters, shape registration loss functions, density fitting

4 Sinkhorn

4.1 Entropic Regularization for Discrete Measures

Relative entropy The Kullback-Leibler divergence is defined as

$$\text{KL}(\mathbf{P}|\mathbf{Q}) \stackrel{\text{def.}}{=} \sum_{i,j} \mathbf{P}_{i,j} \log \left(\frac{\mathbf{P}_{i,j}}{\mathbf{Q}_{i,j}} \right) - \mathbf{P}_{i,j} + \mathbf{Q}_{i,j}. \quad (25)$$

with the convention $0 \log(0) = 0$ and $\mathbf{KL}(\mathbf{P}|\mathbf{Q}) = +\infty$ if there exists some (i, j) such that $\mathbf{Q}_{i,j} = 0$ but $\mathbf{P}_{i,j} \neq 0$. The special case $\mathbf{KL}(\mathbf{P}|\mathbf{1})$ corresponds to minus the Shannon-Boltzmann entropy. The function $\mathbf{KL}(\cdot|\mathbf{Q})$ is strongly convex, because its hessian is $\partial^2 \mathbf{KL}(\mathbf{P}|\mathbf{Q}) = \text{diag}(1/\mathbf{P}_{i,j})$ and $\mathbf{P}_{i,j} \leq 1$.

\mathbf{KL} is a particular instance (and actually the unique case) of both a φ -divergence (as defined in Section ??) and a Bregman divergence. This unique property is at the heart of the fact that this regularization leads to elegant algorithms and a tractable mathematical analysis. One thus has $\mathbf{KL}(\mathbf{P}|\mathbf{Q}) \geq 0$ and $\mathbf{KL}(\mathbf{P}|\mathbf{Q}) = 0$ if and only if $\mathbf{P} = \mathbf{Q}$.

Entropic Regularization for Discrete Measures. The idea of the entropic regularization of optimal transport is to use \mathbf{KL} as a regularizing function to obtain approximate solutions to the original transport problem (??):

$$L_{\mathbf{C}}^{\varepsilon}(\mathbf{a}, \mathbf{b}) \stackrel{\text{def.}}{=} \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \mathbf{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}). \quad (26)$$

Here we used as a reference measure for the relative entropy $\mathbf{a} \otimes \mathbf{b} = (\mathbf{a}_i \mathbf{b}_j)_{i,j}$. This choice of normalization, specially in this discrete setting, has no importance for the selection of the optimal \mathbf{P} since it only affects the objective by a constant, indeed for $\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})$, one has

$$\mathbf{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) = \mathbf{KL}(\mathbf{P}|\mathbf{a}' \otimes \mathbf{b}') + \mathbf{KL}(\mathbf{a}' \otimes \mathbf{b}'|\mathbf{a} \otimes \mathbf{b})$$

[ToDo: check this]. This choice of normalization is however important to deal with situation where the support of \mathbf{a} and \mathbf{b} can change, and in particular when later we will deal with possibly continuous distribution. It also affect the values of the cost $L_{\mathbf{C}}^{\varepsilon}(\mathbf{a}, \mathbf{b})$ and this normalization will be instrumental to define a proper Sinkhorn divergence.

Smoothing effect. Since the objective is a ε -strongly convex function, problem 26 has a unique optimal solution. As studied in Section ??, this smoothing, beyond providing uniqueness, actually leads to $L_{\mathbf{C}}^{\varepsilon}(\mathbf{a}, \mathbf{b})$ being a smooth function of \mathbf{a}, \mathbf{b} and \mathbf{C} . The effect of the entropy is to act as a barrier function for the positivity constraint. As we will show next, this forces the solution \mathbf{P} to be strictly positive on the support of $\mathbf{a} \otimes \mathbf{b}$.

One has the following convergence property.

Proposition 10 (Convergence with ε). *The unique solution \mathbf{P}_{ε} of (26) converges to the optimal solution with maximal entropy within the set of all optimal solutions of the Kantorovich problem, namely*

$$\mathbf{P}_{\varepsilon} \xrightarrow{\varepsilon \rightarrow 0} \underset{\mathbf{P}}{\text{argmin}} \{ \mathbf{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) ; \mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b}), \langle \mathbf{P}, \mathbf{C} \rangle = L_{\mathbf{C}}(\mathbf{a}, \mathbf{b}) \} \quad (27)$$

so that in particular

$$L_{\mathbf{C}}^{\varepsilon}(\mathbf{a}, \mathbf{b}) \xrightarrow{\varepsilon \rightarrow 0} L_{\mathbf{C}}(\mathbf{a}, \mathbf{b}).$$

One has

$$\mathbf{P}_{\varepsilon} \xrightarrow{\varepsilon \rightarrow \infty} \mathbf{a} \otimes \mathbf{b}. \quad (28)$$

Proof. Case $\varepsilon \rightarrow 0$. We consider a sequence $(\varepsilon_{\ell})_{\ell}$ such that $\varepsilon_{\ell} \rightarrow 0$ and $\varepsilon_{\ell} > 0$. We denote \mathbf{P}_{ℓ} the solution of (26) for $\varepsilon = \varepsilon_{\ell}$. Since $\mathbf{U}(\mathbf{a}, \mathbf{b})$ is bounded, we can extract a sequence (that we do not relabel for sake of simplicity) such that $\mathbf{P}_{\ell} \rightarrow \mathbf{P}^*$. Since $\mathbf{U}(\mathbf{a}, \mathbf{b})$ is closed, $\mathbf{P}^* \in \mathbf{U}(\mathbf{a}, \mathbf{b})$. We consider any \mathbf{P} such that $\langle \mathbf{C}, \mathbf{P} \rangle = L_{\mathbf{C}}(\mathbf{a}, \mathbf{b})$. By optimality of \mathbf{P} and \mathbf{P}_{ℓ} for their respective optimization problems (for $\varepsilon = 0$ and $\varepsilon = \varepsilon_{\ell}$), one has

$$0 \leq \langle \mathbf{C}, \mathbf{P}_{\ell} \rangle - \langle \mathbf{C}, \mathbf{P} \rangle \leq \varepsilon_{\ell} (\mathbf{KL}(\mathbf{P}_{\ell}|\mathbf{a} \otimes \mathbf{b}) - \mathbf{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b})). \quad (29)$$

Since \mathbf{H} is continuous, taking the limit $\ell \rightarrow +\infty$ in this expression shows that $\langle \mathbf{C}, \mathbf{P}^* \rangle = \langle \mathbf{C}, \mathbf{P} \rangle$ so that \mathbf{P}^* is a feasible point of (27). Furthermore, dividing by ε_{ℓ} in (29) and taking the limit shows that $\mathbf{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) \leq \mathbf{KL}(\mathbf{P}^*|\mathbf{a} \otimes \mathbf{b})$, which shows that \mathbf{P}^* is a solution of (27). Since the solution \mathbf{P}_0^* to this program is unique by strict convexity of $\mathbf{KL}(\cdot|\mathbf{a} \otimes \mathbf{b})$, one has $\mathbf{P}^* = \mathbf{P}_0^*$, and the whole sequence is converging.

Case $\varepsilon \rightarrow +\infty$. Evaluating at $\mathbf{a} \otimes \mathbf{b}$ the energy, one has

$$\langle \mathbf{C}, \mathbf{P}_\varepsilon \rangle + \varepsilon \text{KL}(\mathbf{P}_\varepsilon | \alpha \otimes \beta) \leq \langle \mathbf{C}, \alpha \otimes \beta \rangle + \varepsilon \times 0$$

and since $\langle \mathbf{C}, \mathbf{P}_\varepsilon \rangle \geq 0$, this leads to

$$\text{KL}(\mathbf{P}_\varepsilon | \alpha \otimes \beta) \leq \varepsilon^{-1} \langle \mathbf{C}, \alpha \otimes \beta \rangle \leq \frac{\|\mathbf{C}\|_\infty}{\varepsilon}$$

so that $\text{KL}(\mathbf{P}_\varepsilon | \alpha \otimes \beta) \rightarrow 0$ and thus $\mathbf{P}_\varepsilon \rightarrow \alpha \otimes \beta$ since KL is a valid divergence. \square

4.2 General Formulation

One can consider arbitrary measures by replacing the discrete entropy by the relative entropy with respect to the product measure $d\alpha \otimes d\beta(x, y) \stackrel{\text{def.}}{=} d\alpha(x)d\beta(y)$, and propose a regularized counterpart to (21) using

$$\mathcal{L}_c^\varepsilon(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi \in \mathcal{U}(\alpha, \beta)} \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta) \quad (30)$$

where the relative entropy is a generalization of the discrete Kullback-Leibler divergence (25)

$$\text{KL}(\pi | \xi) \stackrel{\text{def.}}{=} \int_{X \times Y} \log \left(\frac{d\pi}{d\xi}(x, y) \right) d\pi(x, y) + \int_{X \times Y} (d\xi(x, y) - d\pi(x, y)), \quad (31)$$

and by convention $\text{KL}(\pi | \xi) = +\infty$ if π does not have a density $\frac{d\pi}{d\xi}$ with respect to ξ . It is important to realize that the reference measure $\alpha \otimes \beta$ chosen in (30) to define the entropic regularizing term $\text{KL}(\cdot | \alpha \otimes \beta)$ plays no specific role, only its support matters. This problem is often referred to as the “static Schrödinger problem”, since π is intended to model the most likely coupling between particles of gas which can be only observed at two different times (it is the so-called lazy gas model). The parameter ε controls the temperature of the gas, and particles do not move in deterministic straight line as in optimal transport for the Euclidean cost, but rather according to a stochastic Brownian bridge.

Remark 3 (Probabilistic interpretation). If $(X, Y) \sim \pi$ have marginals $X \sim \alpha$ and $Y \sim \beta$, then $\text{KL}(\pi | \alpha \otimes \beta) = \mathcal{I}(X, Y)$ is the mutual information of the couple, which is 0 if and only if X and Y are independent. The entropic problem (30) is thus equivalent to

$$\min_{(X, Y), X \sim \alpha, Y \sim \beta} \mathbb{E}(c(X, Y)) + \varepsilon \mathcal{I}(X, Y).$$

Using a large ε thus enforces the optimal coupling to describe independent variables, while, according to Brenier’s theorem, small ε rather imposes a deterministic dependency between the couple according to a Monge map.

4.3 Sinkhorn’s Algorithm

The following proposition shows that the solution of (26) has a specific form, which can be parameterized using $n + m$ variables. That parameterization is therefore essentially dual, in the sense that a coupling \mathbf{P} in $\mathcal{U}(\mathbf{a}, \mathbf{b})$ has nm variables but $n + m$ constraints.

Proposition 11. *\mathbf{P} is the unique solution to (26) if and only if there exists $(\mathbf{u}, \mathbf{v}) \in \mathbb{R}_+^n \times \mathbb{R}_+^m$ such that*

$$\forall (i, j) \in \llbracket n \rrbracket \times \llbracket m \rrbracket, \quad \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \quad (32)$$

and $\mathbf{P} \in \mathcal{U}(\mathbf{a}, \beta)$.

Proof. Introducing two dual variables $\mathbf{f} \in \mathbb{R}^n, \mathbf{g} \in \mathbb{R}^m$ for each marginal constraint, the Lagrangian of (26) reads

$$\mathcal{E}(\mathbf{P}, \mathbf{f}, \mathbf{g}) = \langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \mathbf{KL}(\mathbf{P} | \mathbf{a} \otimes \mathbf{b}) + \langle \mathbf{f}, \mathbf{a} - \mathbf{P} \mathbf{1}_m \rangle + \langle \mathbf{g}, \mathbf{b} - \mathbf{P}^T \mathbf{1}_n \rangle.$$

Considering first order conditions (where we ignore the positivity constraint, which can be made rigorous by showing the associated multiplier vanishes), we have

$$\frac{\partial \mathcal{E}(\mathbf{P}, \mathbf{f}, \mathbf{g})}{\partial \mathbf{P}_{i,j}} = \mathbf{C}_{i,j} + \varepsilon \log \left(\frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j} \right) - \mathbf{f}_i - \mathbf{g}_j = 0.$$

which results, for an optimal \mathbf{P} coupling to the regularized problem, in the expression $\mathbf{P}_{i,j} = \mathbf{a}_i \mathbf{b}_j e^{\frac{\mathbf{f}_i + \mathbf{g}_j - \mathbf{C}_{i,j}}{\varepsilon}}$ which can be rewritten in the form provided in the proposition using non-negative vectors $\mathbf{u} \stackrel{\text{def.}}{=} (\mathbf{a}_i e^{\mathbf{f}_i / \varepsilon})_i$ and $\mathbf{v} \stackrel{\text{def.}}{=} (\mathbf{b}_j e^{\mathbf{g}_j / \varepsilon})_j$. \square

The factorization of the optimal solution exhibited in Equation (32) can be conveniently rewritten in matrix form as $\mathbf{P} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v})$. \mathbf{u}, \mathbf{v} must therefore satisfy the following non-linear equations which correspond to the mass conservation constraints inherent to $\mathbf{U}(\mathbf{a}, \mathbf{b})$,

$$\text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \mathbf{1}_m = \mathbf{a}, \quad \text{and} \quad \text{diag}(\mathbf{v}) \mathbf{K}^T \text{diag}(\mathbf{u}) \mathbf{1}_n = \mathbf{b}, \quad (33)$$

These two equations can be further simplified, since $\text{diag}(\mathbf{v}) \mathbf{1}_m$ is \mathbf{v} , and the multiplication of $\text{diag}(\mathbf{u})$ times $\mathbf{K} \mathbf{v}$ is

$$\mathbf{u} \odot (\mathbf{K} \mathbf{v}) = \mathbf{a} \quad \text{and} \quad \mathbf{v} \odot (\mathbf{K}^T \mathbf{u}) = \mathbf{b} \quad (34)$$

where \odot corresponds to entry-wise multiplication of vectors. That problem is known in the numerical analysis community as the matrix scaling problem (see [35] and references therein). An intuitive way to try to solve these equations is to solve them iteratively, by modifying first \mathbf{u} so that it satisfies the left-hand side of Equation (34) and then \mathbf{v} to satisfy its right-hand side. These two updates define Sinkhorn's algorithm

$$\mathbf{u}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{a}}{\mathbf{K} \mathbf{v}^{(\ell)}} \quad \text{and} \quad \mathbf{v}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{b}}{\mathbf{K}^T \mathbf{u}^{(\ell+1)}}, \quad (35)$$

initialized with an arbitrary positive vector, for instance $\mathbf{v}^{(0)} = \mathbf{1}_m$. The division operator used above between two vectors is to be understood entry-wise. Note that a different initialization will likely lead to a different solution for \mathbf{u}, \mathbf{v} , since \mathbf{u}, \mathbf{v} are only defined up to a multiplicative constant (if \mathbf{u}, \mathbf{v} satisfy (33) then so do $\lambda \mathbf{u}, \mathbf{v} / \lambda$ for any $\lambda > 0$). It turns out however that these iterations converge, as we detail next.

[ToDo: Say a few word about the general probleme of scaling a matrix to a bistochastic one, and why this is non trivial for matrices with vanishing entries.]

A chief advantage, beside its simplicity, of Sinkhorn's algorithm is that the only computationnaly expensive step are matrix-vector multiplication by the Gibbs kernel, so that its complexity scales likes Knm where K is the number of Sinkhorn iteration, which can be kept polynomially in $1/\varepsilon$ if one is interested in reaching an accuracy ε on the (unregularized) transportation cost. Note however that in many situation, one is not interested in reaching high accuracy, because targeted application success is often only remotely connected to the ability to solve an optimal transport problem (but rather only being able to compare in a geometrically faithful way distribution), so that K is usually quite small. This should be contrasted with interior point methods, which also operate by introducing a barrier function of the form $-\sum_i \log(\mathbf{P}_{i,j})$. These algorithm have typically a complexity of the order $O(n^6 \log(|\varepsilon|))$ **[ToDo: check]**.

The second crucial aspect of Sinkhorn is that matrix-vector multiplication streams extremely well on GPU. Even better, if one is interested in computing many OT problem with a fixed cost matrix \mathbf{C} , one can replace many matrix-vector multiplication by matrix-matrix multiplication, so that the computation gain is enormous.

4.4 Convergence

Convergence finite dimension via alternating projections. One has

$$\langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \mathbf{KL}(\mathbf{P} | \mathbf{a} \otimes \mathbf{b}) = \varepsilon \mathbf{KL}(\mathbf{P} | \mathbf{K}) + \text{cst},$$

so that the unique solution \mathbf{P}_ε of (26) is a projection onto $\mathbf{U}(\mathbf{a}, \mathbf{b})$ of the Gibbs kernel \mathbf{K}

$$\mathbf{P}_\varepsilon = \text{Proj}_{\mathbf{U}(\mathbf{a}, \mathbf{b})}^{\mathbf{KL}}(\mathbf{K}) \stackrel{\text{def.}}{=} \underset{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})}{\text{argmin}} \mathbf{KL}(\mathbf{P} | \mathbf{K}). \quad (36)$$

Denoting

$$\mathcal{C}_{\mathbf{a}}^1 \stackrel{\text{def.}}{=} \{\mathbf{P} ; \mathbf{P} \mathbf{1}_m = \mathbf{a}\} \quad \text{and} \quad \mathcal{C}_{\mathbf{b}}^2 \stackrel{\text{def.}}{=} \left\{ \mathbf{P} ; \mathbf{P}^T \mathbf{1}_n = \mathbf{b} \right\}$$

the rows and columns constraints, one has $\mathbf{U}(\mathbf{a}, \mathbf{b}) = \mathcal{C}_{\mathbf{a}}^1 \cap \mathcal{C}_{\mathbf{b}}^2$. One can use Bregman iterative projections [11]

$$\mathbf{P}^{(\ell+1)} \stackrel{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{a}}^1}^{\mathbf{KL}}(\mathbf{P}^{(\ell)}) \quad \text{and} \quad \mathbf{P}^{(\ell+2)} \stackrel{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{b}}^2}^{\mathbf{KL}}(\mathbf{P}^{(\ell+1)}). \quad (37)$$

Since the sets $\mathcal{C}_{\mathbf{a}}^1$ and $\mathcal{C}_{\mathbf{b}}^2$ are affine, these iterations are known to converge to the solution of (36), see [11].

The two projector are simple to compute since they corresponds to scaling respectively the rows and the columns

$$\text{Proj}_{\mathcal{C}_{\mathbf{a}}^1}^{\mathbf{KL}}(\mathbf{P}) = \text{diag} \left(\frac{\mathbf{a}}{\mathbf{P} \mathbf{1}_m} \right) \mathbf{P} \quad \text{and} \quad \text{Proj}_{\mathcal{C}_{\mathbf{b}}^2}^{\mathbf{KL}}(\mathbf{P}) = \mathbf{P} \text{diag} \left(\frac{\mathbf{b}}{\mathbf{P}^T \mathbf{1}_n} \right).$$

These iterate are equivalent to Sinkhorn iterations (35) since defining

$$\mathbf{P}^{(2\ell)} \stackrel{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell)}) \mathbf{K} \text{diag}(\mathbf{v}^{(\ell)}),$$

one has

$$\begin{aligned} \mathbf{P}^{(2\ell+1)} &\stackrel{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell+1)}) \mathbf{K} \text{diag}(\mathbf{v}^{(\ell)}) \\ \text{and } \mathbf{P}^{(2\ell+2)} &\stackrel{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell+1)}) \mathbf{K} \text{diag}(\mathbf{v}^{(\ell+1)}) \end{aligned}$$

In practice however one should prefer using (35) which only requires manipulating scaling vectors and multiplication against a Gibbs kernel, which can often be accelerated (see below Remarks ?? and ??).

Such a convergence analysis using Bregman projection is however of limited interested because it only works in finite dimension. For instance, the linear convergence speed one can obtain with these analyses (because the objective is strongly convex) will degrade with the dimension (and of course also with ε). It is also possible to decay ε during the iterates to improve the speed and rely on multiscale strategies in low dimension.

Convergence for the Hilbert metric As initially explained by [26], the global convergence analysis of Sinkhorn is greatly simplified using Hilbert projective metric on $\mathbb{R}_{+,*}^n$ (positive vectors), defined as

$$\forall (\mathbf{u}, \mathbf{u}') \in (\mathbb{R}_{+,*}^n)^2, \quad d_{\mathcal{H}}(\mathbf{u}, \mathbf{u}') \stackrel{\text{def.}}{=} \|\log(\mathbf{u}) - \log(\mathbf{v})\|_V$$

where the variation semi-norm is

$$\|z\|_V = \max(z) - \min(z).$$

One can show that $d_{\mathcal{H}}$ is a distance on the projective cone $\mathbb{R}_{+,*}^n / \sim$, where $\mathbf{u} \sim \mathbf{u}'$ means that $\exists s > 0, \mathbf{u} = s\mathbf{u}'$ (the vector are equal up to rescaling, hence the naming “projective”), and that $(\mathbb{R}_{+,*}^n / \sim, d_{\mathcal{H}})$ is then a complete metric space. It was introduced independently by [8] and [39] to provide a quantitative proof of Perron-Frobenius theorem (convergence of iterations of positive matrices). Sinkhorn should be thought as a non-linear generalization of Perron-Frobenius.

Theorem 3. Let $\mathbf{K} \in \mathbb{R}_{+,*}^{n \times m}$, then for $(\mathbf{v}, \mathbf{v}') \in (\mathbb{R}_{+,*}^m)^2$

$$d_{\mathcal{H}}(\mathbf{K}\mathbf{v}, \mathbf{K}\mathbf{v}') \leq \lambda(\mathbf{K})d_{\mathcal{H}}(\mathbf{v}, \mathbf{v}') \text{ where } \begin{cases} \lambda(\mathbf{K}) \stackrel{\text{def.}}{=} \frac{\sqrt{\eta(\mathbf{K})}-1}{\sqrt{\eta(\mathbf{K})}+1} < 1 \\ \eta(\mathbf{K}) \stackrel{\text{def.}}{=} \max_{i,j,k,\ell} \frac{\mathbf{K}_{i,k}\mathbf{K}_{j,\ell}}{\mathbf{K}_{j,k}\mathbf{K}_{i,\ell}}. \end{cases}$$

The following theorem, proved by [26], makes use of this Theorem 3 to show the linear convergence of Sinkhorn's iterations.

Theorem 4. One has $(\mathbf{u}^{(\ell)}, \mathbf{v}^{(\ell)}) \rightarrow (\mathbf{u}^*, \mathbf{v}^*)$ and

$$d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) = O(\lambda(\mathbf{K})^{2\ell}), \quad d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^*) = O(\lambda(\mathbf{K})^{2\ell}). \quad (38)$$

One also has

$$d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) \leq \frac{d_{\mathcal{H}}(\mathbf{P}^{(\ell)}\mathbf{1}_m, \mathbf{a})}{1 - \lambda(\mathbf{K})} \quad \text{and} \quad d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^*) \leq \frac{d_{\mathcal{H}}(\mathbf{P}^{(\ell),\top}\mathbf{1}_n, \mathbf{b})}{1 - \lambda(\mathbf{K})}, \quad (39)$$

where we denoted $\mathbf{P}^{(\ell)} \stackrel{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell)})\mathbf{K}\text{diag}(\mathbf{v}^{(\ell)})$. Lastly, one has

$$\|\log(\mathbf{P}^{(\ell)}) - \log(\mathbf{P}^*)\|_{\infty} \leq d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) + d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^*) \quad (40)$$

where \mathbf{P}^* is the unique solution of (26).

Proof. One notice that for any $(\mathbf{v}, \mathbf{v}') \in (\mathbb{R}_{+,*}^m)^2$, one has

$$d_{\mathcal{H}}(\mathbf{v}, \mathbf{v}') = d_{\mathcal{H}}(\mathbf{v}/\mathbf{v}', \mathbf{1}_m) = d_{\mathcal{H}}(\mathbf{1}_m/\mathbf{v}, \mathbf{1}_m/\mathbf{v}').$$

This shows that

$$d_{\mathcal{H}}(\mathbf{u}^{(\ell+1)}, \mathbf{u}^*) = d_{\mathcal{H}}\left(\frac{\mathbf{a}}{\mathbf{K}\mathbf{v}^{(\ell)}}, \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}^*}\right) = d_{\mathcal{H}}(\mathbf{K}\mathbf{v}^{(\ell)}, \mathbf{K}\mathbf{v}^*) \leq \lambda(\mathbf{K})d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^*).$$

where we used Theorem 3. This shows (38). One also has, using the triangular inequality

$$\begin{aligned} d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) &\leq d_{\mathcal{H}}(\mathbf{u}^{(\ell+1)}, \mathbf{u}^{(\ell)}) + d_{\mathcal{H}}(\mathbf{u}^{(\ell+1)}, \mathbf{u}^*) \leq d_{\mathcal{H}}\left(\frac{\mathbf{a}}{\mathbf{K}\mathbf{v}^{(\ell)}}, \mathbf{u}^{(\ell)}\right) + \lambda(\mathbf{K})d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) \\ &= d_{\mathcal{H}}\left(\mathbf{a}, \mathbf{u}^{(\ell)} \odot (\mathbf{K}\mathbf{v}^{(\ell)})\right) + \lambda(\mathbf{K})d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*), \end{aligned}$$

which gives the first part of (39) since $\mathbf{u}^{(\ell)} \odot (\mathbf{K}\mathbf{v}^{(\ell)}) = \mathbf{P}^{(\ell)}\mathbf{1}_m$ (the second one being similar). The proof of (40) follows from [26, Lemma 3] \square

The bound (39) shows that some error measures on the marginal constraints violation, for instance $\|\mathbf{P}^{(\ell)}\mathbf{1}_m - \mathbf{a}\|_1$ and $\|\mathbf{P}^{(\ell)\top}\mathbf{1}_n - \mathbf{b}\|_1$, are useful stopping criteria to monitor the convergence. This theorem shows that Sinkhorn algorithm converges linearly, but the rates becomes exponentially bad as $\varepsilon \rightarrow 0$, since it scales like $e^{-1/\varepsilon}$. In practice, one eventually observes a linear rate after enough iteration, because the local linear rate is much better, usually of the order $1 - \varepsilon$.

5 Dual Problem

5.1 Discrete dual

The Kantorovich problem (??) is a linear program, so that one can equivalently compute its value by solving a dual linear program.

Proposition 12. *One has*

$$L_C(\mathbf{a}, \mathbf{b}) = \max_{(\mathbf{f}, \mathbf{g}) \in \mathbf{R}(\mathbf{a}, \mathbf{b})} \langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{g}, \mathbf{b} \rangle \quad (41)$$

where the set of admissible potentials is

$$\mathbf{R}(\mathbf{a}, \mathbf{b}) \stackrel{\text{def}}{=} \{(\mathbf{f}, \mathbf{g}) \in \mathbb{R}^n \times \mathbb{R}^m ; \forall (i, j) \in \llbracket n \rrbracket \times \llbracket m \rrbracket, \mathbf{f} \oplus \mathbf{g} \leq \mathbf{C}\} \quad (42)$$

Proof. For the sake of completeness, let us derive this dual problem with the use of Lagrangian duality. The Lagrangian associate to (??) reads

$$\min_{\mathbf{P} \geq 0} \max_{(\mathbf{f}, \mathbf{g}) \in \mathbb{R}^n \times \mathbb{R}^m} \langle \mathbf{C}, \mathbf{P} \rangle + \langle \mathbf{a} - \mathbf{P} \mathbf{1}_m, \mathbf{f} \rangle + \langle \mathbf{b} - \mathbf{P}^\top \mathbf{1}_n, \mathbf{g} \rangle. \quad (43)$$

For linear program, if the primal set of constraint is non-empty, one can always exchange the min and the max and get the same value of the linear program, and one thus consider

$$\max_{(\mathbf{f}, \mathbf{g}) \in \mathbb{R}^n \times \mathbb{R}^m} \langle \mathbf{a}, \mathbf{f} \rangle + \langle \mathbf{b}, \mathbf{g} \rangle + \min_{\mathbf{P} \geq 0} \langle \mathbf{C} - \mathbf{f} \mathbf{1}_m^\top - \mathbf{1}_n \mathbf{g}^\top, \mathbf{P} \rangle.$$

We conclude by remarking that

$$\min_{\mathbf{P} \geq 0} \langle \mathbf{Q}, \mathbf{P} \rangle = \begin{cases} 0 & \text{if } \mathbf{Q} \geq 0 \\ -\infty & \text{otherwise} \end{cases}$$

so that the constraint reads $\mathbf{C} - \mathbf{f} \mathbf{1}_m^\top - \mathbf{1}_n \mathbf{g}^\top = \mathbf{C} - \mathbf{f} \oplus \mathbf{g} \geq 0$. \square

The primal-dual optimality relation for the Lagrangian (43) allows to locate the support of the optimal transport plan

$$\text{Supp}(\mathbf{P}) \subset \{(i, j) \in \llbracket n \rrbracket \times \llbracket m \rrbracket ; \mathbf{f}_i + \mathbf{g}_j = \mathbf{C}_{i,j}\}. \quad (44)$$

The formulation (70) shows that $(\mathbf{a}, \mathbf{b}) \mapsto L_C(\mathbf{a}, \mathbf{b})$ is a convex function (as a supremum of linear functions). From the primal problem (??), one also sees that $\mathbf{C} \mapsto L_C(\mathbf{a}, \mathbf{b})$ is concave.

5.2 General formulation

To extend this primal-dual construction to arbitrary measures, it is important to realize that measures are naturally paired in duality with continuous functions, using the pairing $\langle f, \alpha \rangle \stackrel{\text{def}}{=} \int f d\alpha$.

Proposition 13. *One has*

$$\mathcal{L}_c(\alpha, \beta) = \max_{(f, g) \in \mathcal{R}(c)} \int_{\mathcal{X}} f(x) d\alpha(x) + \int_{\mathcal{Y}} g(y) d\beta(y), \quad (45)$$

where the set of admissible dual potentials is

$$\mathcal{R}(c) \stackrel{\text{def}}{=} \{(f, g) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y}) ; \forall (x, y), f(x) + g(y) \leq c(x, y)\}. \quad (46)$$

Here, (f, g) is a pair of continuous functions, and are often called “Kantorovich potentials”.

The discrete case (70) corresponds to the dual vectors being samples of the continuous potentials, *i.e.* $(\mathbf{f}_i, \mathbf{g}_j) = (f(x_i), g(y_j))$. The primal-dual optimality conditions allow to track the support of optimal plan, and (44) is generalized as

$$\text{Supp}(\pi) \subset \{(x, y) \in \mathcal{X} \times \mathcal{Y} ; f(x) + g(y) = c(x, y)\}. \quad (47)$$

Note that in contrast to the primal problem (21), showing the existence of solutions to (45) is non-trivial, because the constraint set $\mathcal{R}(c)$ is not compact and the function to minimize non-coercive. Using the machinery of c -transform detailed in Section ??, one can however show that optimal (f, g) are necessarily Lipschitz regular, which enable to replace the constraint by a compact one.

5.3 c -transforms

Definition. Keeping a dual potential g fixed, one can try to minimize in closed form the dual problem (45), which leads to consider

$$\sup_{g \in \mathcal{C}(\mathcal{Y})} \left\{ \int g d\beta ; \forall (x, y), g(y) \leq c(x, y) - f(x) \right\}.$$

The constraint can be replaced by

$$\forall y \in \mathcal{Y}, \quad g(y) \leq f^c(y)$$

where we define the c -transform as

$$\forall y \in \mathcal{Y}, \quad f^c(y) \stackrel{\text{def.}}{=} \inf_{x \in \mathcal{X}} c(x, y) - f(x). \quad (48)$$

Since β is positive, the maximization of $\int g d\beta$ is thus achieved at those functions such that $g = f^c$ on the support of β , which means β -almost everywhere.

Similarly, we defined the \bar{c} -transform, which a transform for the symetrized cost $\bar{c}(y, x) = c(x, y)$, i.e.

$$\forall x \in \mathcal{X}, \quad g^{\bar{c}}(x) \stackrel{\text{def.}}{=} \inf_{y \in \mathcal{Y}} c(x, y) - g(y),$$

and one checks that any function f such that $f = g^{\bar{c}}$ α -almost everywhere is solution to the dual problem for a fixed g .

The map $(f, g) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y}) \mapsto (g^{\bar{c}}, f^c) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y})$ replaces dual potentials by “better” ones (improving the dual objective \mathcal{E}). Functions that can be written in the form f^c and $g^{\bar{c}}$ are called c -concave and \bar{c} -concave functions.

Note that these partial minimizations define maximizers on the support of respectively α and β , while the definitions (48) actually define functions on the whole spaces \mathcal{X} and \mathcal{Y} . This is thus a way to extend in a canonical way solutions of (45) on the whole spaces.

Furthermore, if c is Lipschitz, then f^c and g^c are also Lipschitz functions, as we now show. This property is crucial to show existence of solution to the dual problem. Indeed, since one can impose this Lipschitz on the dual problems, the constraint set is compact via Ascoli theorem.

Proposition 14. *If c is L -Lipschitz with respect to the second variable, then f^c is L -Lipschitz.*

Proof. We apply to $F_x = c(x, \cdot) - f(x)$ the fact that if all the F_x are L -Lipschitz, then the Lipschitz constant of $F = \min_x F_x$ is L . Indeed, using the fact that $|\inf(A) - \inf(B)| \leq \sup |A - B|$ for two function A and B , then

$$|F(y) - F(y')| = |\inf_x (F_x(y)) - \inf_x (F_x(y'))| \leq \sup_x |F_x(y) - F_x(y')| \leq \sup_x Ld(y, y') = Ld(y, y').$$

□

Euclidean case. The special case $c(x, y) = -\langle x, y \rangle$ in $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ is of utmost importance because it allows one to study the W_2 problem, since for any $\pi \in \mathcal{U}(\alpha, \beta)$

$$\int \|x - y\|^2 d\pi(x, y) = \text{cst} - 2 \int \langle x, y \rangle d\pi(x, y) \quad \text{where} \quad \text{cst} = \int \|x\|^2 d\alpha(x) + \int \|y\|^2 d\beta(y).$$

For this special choice of cost, one has $f^c = -(-f)^*$ where h^* is the Fenchel-Legendre transform

$$h^*(y) \stackrel{\text{def.}}{=} \sup_x \langle x, y \rangle - h(x).$$

One has that h^* is always convex, so that f^c is always concave. For a general cost, one thus denotes functions of the form f^c as being c -concave.

The failure of alternate optimization. A crucial property of the Legendre transform is that $f^{***} = f^*$, and that f^{**} is the convex envelope of f (the largest convex function below f). These properties carries over for the more general setting of c -transforms.

Proposition 15. *The following identities, in which the inequality sign between vectors should be understood elementwise, hold, denoting $f^{c\bar{c}} \stackrel{\text{def}}{=} (f^c)^{\bar{c}}$:*

- (i) $f \leq f' \Rightarrow f^c \geq f'^c$,
- (ii) $f^{c\bar{c}} \geq f$,
- (iii) $g^{\bar{c}c} \geq g$,
- (iv) $f^{c\bar{c}c} = f^c$.

Proof. The first inequality (i) follows from the definition of c -transforms. To prove (ii), expanding the definition of $f^{c\bar{c}}$ we have

$$(f^{c\bar{c}})(x) = \min_y c(x, y) - f^c(y) = \min_y c(x, y) - \min_{x'} (c(x', y) - f(x')).$$

Now, since $-\min_{x'} c(x', y) - f(x') \geq -(c(x, y) - f(x))$, we recover

$$(f^{c\bar{c}})(x) \geq \min_y c(x, y) - c(x, y) + f(x) = f(x).$$

The relation $g^{\bar{c}c} \geq g$ is obtained in the same way. Now, to prove (iv), we first apply (ii) and then (i) with $f' = f^{c\bar{c}}$ to have $f^c \geq f^{c\bar{c}c}$. Then we apply (iii) to $g = f^c$ to obtain $f^c \leq f^{c\bar{c}c}$. \square

This invariance property shows that one can “improve” only once the dual potential this way. Indeed, starting from any pair (f, g) , one obtains the following iterates by alternating maximization

$$(f, g) \mapsto (f, f^c) \mapsto (f^{cc}, f^c) \mapsto (f^{cc}, f^{ccc}) = (f^{cc}, f^c) \dots \quad (49)$$

so that one reaches a stationary point. This failure is the classical behavior of alternating maximization on a non-smooth problem, where the non-smooth part of the functional (here the constraint) mixes the two variables. The workaround is to introduce a smoothing, which is the classical method of augmented Lagrangian, and that we will develop here using entropic regularization, and corresponds to Sinkhorn’s algorithm.

6 Semi-discrete and W_1

6.1 Semi-discrete

A case of particular interest is when $\beta = \sum_j \mathbf{b}_j \delta_{y_j}$ is discrete (of course the same construction applies if α is discrete by exchanging the role of α, β). One can adapt the definition of the \bar{c} transform (48) to this setting by restricting the minimization to the support $(y_j)_j$ of β ,

$$\forall \mathbf{g} \in \mathbb{R}^m, \forall x \in \mathcal{X}, \quad \mathbf{g}^{\bar{c}}(x) \stackrel{\text{def}}{=} \min_{j \in \llbracket m \rrbracket} c(x, y_j) - \mathbf{g}_j. \quad (50)$$

This transform maps a vector \mathbf{g} to a continuous function $\mathbf{g}^{\bar{c}} \in \mathcal{C}(\mathcal{X})$. Note that this definition coincides with (48) when imposing that the space \mathcal{X} is equal to the support of β .

Crucially, using the discrete \bar{c} -transform, when β is a discrete measure, yields a finite-dimensional optimization,

$$\mathcal{L}_c(\alpha, \beta) = \max_{\mathbf{g} \in \mathbb{R}^m} \mathcal{E}(\mathbf{g}) \stackrel{\text{def}}{=} \int_{\mathcal{X}} \mathbf{g}^{\bar{c}}(x) d\alpha(x) + \sum_j \mathbf{g}_j \mathbf{b}_j. \quad (51)$$

The Laguerre cells associated to the dual weights \mathbf{g}

$$\mathbb{L}_j(\mathbf{g}) \stackrel{\text{def.}}{=} \{x \in \mathcal{X} ; \forall j' \neq j, c(x, y_j) - \mathbf{g}_j \leq c(x, y_{j'}) - \mathbf{g}_{j'}\}$$

induce a disjoint decomposition of $\mathcal{X} = \bigcup_j \mathbb{L}_j(\mathbf{g})$. When \mathbf{g} is constant, the Laguerre cells decomposition corresponds to the Voronoi diagram partition of the space.

This allows one to conveniently rewrite the minimized energy as

$$\mathcal{E}(\mathbf{g}) = \sum_{j=1}^m \int_{\mathbb{L}_j(\mathbf{g})} (c(x, y_j) - \mathbf{g}_j) d\alpha(x) + \langle \mathbf{g}, \mathbf{b} \rangle. \quad (52)$$

The following proposition provides a formula for the gradient of this convex function.

Proposition 16. *If α has a density with respect to Lebesgue measure and if c is smooth away from the diagonal, then \mathcal{E} is differentiable and*

$$\forall j \in \llbracket m \rrbracket, \quad \nabla \mathcal{E}(\mathbf{g})_j = \mathbf{b}_j - \int_{\mathbb{L}_j(\mathbf{g})} d\alpha.$$

Proof. One has

$$\mathcal{E}(\mathbf{g} + \varepsilon \delta_j) - \mathcal{E}(\mathbf{g}) - \varepsilon \left(\mathbf{b}_j - \int_{\mathbb{L}_j(\mathbf{g})} d\alpha \right) = \sum_k \int_{\mathbb{L}_k(\mathbf{g} + \varepsilon \delta_j)} c(x, x_k) d\alpha(x) - \int_{\mathbb{L}_k(\mathbf{g})} c(x, x_k) d\alpha(x).$$

Most of the terms in the right hand side vanish (because most the Laguerre cells associated to $\mathbf{g} + \varepsilon \delta_j$ are equal to those of \mathbf{g}) and the only terms remaining correspond to neighboring cells (j, k) such that $\mathbb{L}_j(\mathbf{g}) \cap \mathbb{L}_k(\mathbf{g}) \neq \emptyset$ (for the cost $\|x - y\|^2$ and $\mathbf{g} = 0$ this forms the Delaunay triangulation). On these pairs, the right integral differs on a volume of the order of ε (since α has a density) and the function being integrated only varies on the order of ε (since the cost is smooth). So the right hand side is of the order of ε^2 . \square

The first order optimality condition shows that in order to solve the dual semi discrete problem, one needs to select the weights \mathbf{g} in order to drive the Laguerre cell in a configuration such that $\int_{\mathbb{L}_j(\mathbf{g})} d\alpha = \mathbf{b}_j$, i.e. each cell should capture the correct amount of mass. In this case, the optimal transport T such that $T_{\#}\alpha = \beta$ (which exists and is unique according to Brenier's theorem if α has a density) is piecewise constant and map $x \in \mathbb{L}_j(\mathbf{g})$ to y_j .

In the special case $c(x, y) = \|x - y\|^2$, the decomposition in Laguerre cells is also known as a “power diagram”. In this case, the cells are polyhedral and can be computed efficiently using computational geometry algorithms; see [3]. The most widely used algorithm relies on the fact that the power diagram of points in \mathbb{R}^d is equal to the projection on \mathbb{R}^d of the convex hull of the set of points $((y_j, \|y_j\|^2 - \mathbf{g}_j))_{j=1}^m \subset \mathbb{R}^{d+1}$. There are numerous algorithms to compute convex hulls; for instance, that of [18] in two and three dimensions has complexity $O(m \log(Q))$, where Q is the number of vertices of the convex hull.

Stochastic optimization. The semidiscrete formulation (52) is also appealing because the energies to be minimized are written as an expectation with respect to the probability distribution α ,

$$\mathcal{E}(\mathbf{g}) = \int_{\mathcal{X}} E(\mathbf{g}, x) d\alpha(x) = \mathbb{E}_X(E(\mathbf{g}, X)) \quad \text{where} \quad E(\mathbf{g}, x) \stackrel{\text{def.}}{=} \mathbf{g}^{\bar{c}}(x) - \langle \mathbf{g}, \mathbf{b} \rangle,$$

and X denotes a random vector distributed on \mathcal{X} according to α . Note that the gradient of each of the involved functional reads

$$\nabla_{\mathbf{g}} E(x, \mathbf{g}) = (\mathbb{1}_{\mathbb{L}_j(\mathbf{g})}(x) - \mathbf{b}_j)_{j=1}^m \in \mathbb{R}^m$$

where $\mathbb{1}_{\mathbb{L}_j(\mathbf{g})}$ is the indicator function of the Laguerre cell. One can thus use stochastic optimization methods to perform the maximization, as proposed in [27]. This allows us to obtain provably convergent algorithms

without the need to resort to an arbitrary discretization of α (either approximating α using sums of Diracs or using quadrature formula for the integrals). The measure α is used as a black box from which one can draw independent samples, which is a natural computational setup for many high-dimensional applications in statistics and machine learning.

Initializing $\mathbf{g}^{(0)} = \mathbf{0}_m$, the stochastic gradient descent algorithm (SGD; used here as a maximization method) draws at step ℓ a point $x_\ell \in \mathcal{X}$ according to distribution α (independently from all past and future samples $(x_\ell)_\ell$) to form the update

$$\mathbf{g}^{(\ell+1)} \stackrel{\text{def.}}{=} \mathbf{g}^{(\ell)} + \tau_\ell \nabla_{\mathbf{g}} E(\mathbf{g}^{(\ell)}, x_\ell). \quad (53)$$

The step size τ_ℓ should decay fast enough to zero in order to ensure that the “noise” created by using $\nabla_{\mathbf{g}} E(x_\ell, \mathbf{g})$ as a proxy for the true gradient $\nabla \mathcal{E}(\mathbf{g})$ is canceled in the limit. A typical choice of schedule is

$$\tau_\ell \stackrel{\text{def.}}{=} \frac{\tau_0}{1 + \ell/\ell_0}, \quad (54)$$

where ℓ_0 indicates roughly the number of iterations serving as a warmup phase. One can prove the convergence result

$$\mathcal{E}(\mathbf{g}^*) - \mathbb{E}(\mathcal{E}(\mathbf{g}^{(\ell)})) = O\left(\frac{1}{\sqrt{\ell}}\right),$$

where \mathbf{g}^* is a solution of (??) and where \mathbb{E} indicates an expectation with respect to the i.i.d. sampling of $(x_\ell)_\ell$ performed at each iteration.

Optimal quantization. The optimal quantization problem of some measure α corresponds to the resolution of

$$\mathcal{Q}_m(\alpha) = \min_{Y=(y_j)_{j=1}^m, (\mathbf{b}_j)_{j=1}^m} W_p(\alpha, \sum_j b_j \delta_{y_j}).$$

This problem is at the heart of the computation of efficient vector quantizer in information theory and compression, and is also the basic problem to solve for clustering in unsupervised learning. The asymptotic behavior of \mathcal{Q}_m is of fundamental importance, and its precise behavior is in general unknown. For a measure with a density in Euclidean space, it scales like $O(1/n^{1/d})$, so that quantization generally suffers from the curse of dimensionality.

This optimal quantization problem is convex with respect to \mathbf{b} , but is unfortunately non-convex with respect to $Y = (y_j)_j$. Its resolution is in general NP-hard. The only setting where this problem is simple is the 1-D case, in which case the optimal sampling is simply $y_j = \mathcal{C}_\alpha^{-1}(j/m)$. **[ToDo: see where this is proved]**

Solving explicitly for the minimization over \mathbf{b} in the formula (51) (exchanging the role of the min and the max) shows that necessarily, at optimality, one has $\mathbf{g} = 0$, so that the optimal transport maps the Voronoi cells $\mathbb{L}_j(\mathbf{g} = 0)$, which we denote $\mathbb{V}_j(Y)$ to highlight the dependency on the quantization points $Y = (y_j)_j$

$$\mathbb{V}_j(Y) = \{x ; \forall j', c(x, y_{j'}) \leq c(x, y_j)\}.$$

This also shows that the quantization energy can be rewritten in a more intuitive way, which accounts for the average quantization error induced by replacing a point x by its nearest centroid

$$\mathcal{Q}_m(\alpha) = \min_Y \mathcal{F}(Y) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \min_{1 \leq j \leq m} c(x, y_j) d\alpha(x).$$

At any local minimizer (at least if α has a density so that this function is differentiable) of this energy over Y , one sees that each y_j should be a centroid of its associated Voronoi region,

$$y_j \in \operatorname{argmin}_y \int_{\mathbb{V}_j(Y)} c(x, y) d\alpha(x).$$

For instance, when $c(x, y) = \|x - y\|^2$, one sees that any local minimizer should satisfy the fixed point equation

$$y_j = \frac{\int_{\mathbb{V}_j(Y)} x d\alpha(x)}{\int_{\mathbb{V}_j(Y)} d\alpha}.$$

The celebrated k -means algorithm, also known as Lloyd algorithm, iteratively apply this fixed point. It is not guaranteed to converge (it could in theory cycle) but in practice it always converge to a local minimum. A practical issue to obtain a good local minimizer is to seed a good initial configuration. The intuitive way to achieve this is to spread them as much as possible, and a well known algorithm to do so is the k -means++ methods, which achieve without even any iteration a quantization cost which is of the order of $\log(m)\mathcal{Q}_m(\alpha)$.

6.2 W_1

c -transform for W_1 . Here we assume that d is a distance on $\mathcal{X} = \mathcal{Y}$, and we solve the OT problem with the ground cost $c(x, y) = d(x, y)$. The following proposition highlights key properties of the c -transform (48) in this setup. In the following, we denote the Lipschitz constant of a function $f \in \mathcal{C}(\mathcal{X})$ as

$$\text{Lip}(f) \stackrel{\text{def.}}{=} \sup \left\{ \frac{|f(x) - f(y)|}{d(x, y)} ; (x, y) \in \mathcal{X}^2, x \neq y \right\}.$$

Proposition 17. *Suppose $\mathcal{X} = \mathcal{Y}$ and $c(x, y) = d(x, y)$. Then, there exists g such that $f = g^c$ if and only $\text{Lip}(f) \leq 1$. Furthermore, if $\text{Lip}(f) \leq 1$, then $f^c = -f$.*

Proof. First, suppose $f = g^c$ for some g . Then, for $x, y \in \mathcal{X}$,

$$\begin{aligned} |f(x) - f(y)| &= \left| \inf_{z \in \mathcal{X}} d(x, z) - g(z) - \inf_{z \in \mathcal{X}} d(y, z) - g(z) \right| \\ &\leq \sup_{z \in \mathcal{X}} |d(x, z) - d(y, z)| \leq d(x, y). \end{aligned}$$

The first equality follows from the definition of g^c , the next inequality from the identity $|\inf f - \inf g| \leq \sup |f - g|$, and the last from the triangle inequality. This shows that $\text{Lip}(f) \leq 1$.

If f is 1-Lipschitz, for all $x, y \in \mathcal{X}$, $f(y) - d(x, y) \leq f(x) \leq f(y) + d(x, y)$, which shows that

$$\begin{aligned} f^c(y) &= \inf_{x \in \mathcal{X}} [d(x, y) - f(x)] \geq \inf_{x \in \mathcal{X}} [d(x, y) - f(y) - d(x, y)] = -f(y), \\ f^c(y) &= \inf_{x \in \mathcal{X}} [d(x, y) - f(x)] \leq \inf_{x \in \mathcal{X}} [d(x, y) - f(y) + d(x, y)] = -f(y), \end{aligned}$$

and thus $f^c = -f$.

Applying this property to $-f$ which is also 1-Lipschitz shows that $(-f)^c = f$ so that f is indeed c -concave (i.e. it is the c -transform of a function). \square

Using the iterative c -transform scheme (49), one can replace the dual variable (f, g) by $(f^{cc}, f^c) = (-f^c, f^c)$, or equivalently by any pair $(f, -f)$ where f is 1-Lipschitz. This leads to the following alternative expression for the \mathcal{W}_1 distance

$$\mathcal{W}_1(\alpha, \beta) = \max_f \left\{ \int_{\mathcal{X}} f d(\alpha - \beta) ; \text{Lip}(f) \leq 1 \right\}. \quad (55)$$

This expression shows that \mathcal{W}_1 is actually a norm, i.e. $\mathcal{W}_1(\alpha, \beta) = \|\alpha - \beta\|_{\mathcal{W}_1}$, and that it is still valid for any measures (not necessary positive) as long as $\int_{\mathcal{X}} \alpha = \int_{\mathcal{X}} \beta$. This norm is often called the Kantorovich-Rubinstein norm [33].

For discrete measures of the form (2), writing $\alpha - \beta = \sum_k \mathbf{m}_k \delta_{z_k}$ with $z_k \in \mathcal{X}$ and $\sum_k \mathbf{m}_k = 0$, the optimization (55) can be rewritten as

$$\mathcal{W}_1(\alpha, \beta) = \max_{(\mathbf{f}_k)_k} \left\{ \sum_k \mathbf{f}_k \mathbf{m}_k ; \forall (k, \ell), |\mathbf{f}_k - \mathbf{f}_\ell| \leq d(z_k, z_\ell), \right\} \quad (56)$$

which is a finite-dimensional convex program with quadratic-cone constraints. It can be solved using interior point methods or, as we detail next for a similar problem, using proximal methods.

When using $d(x, y) = |x - y|$ with $\mathcal{X} = \mathbb{R}$, we can reduce the number of constraints by ordering the z_k 's via $z_1 \leq z_2 \leq \dots$. In this case, we only have to solve

$$\mathcal{W}_1(\alpha, \beta) = \max_{(\mathbf{f}_k)_k} \left\{ \sum_k \mathbf{f}_k \mathbf{m}_k ; \forall k, |\mathbf{f}_{k+1} - \mathbf{f}_k| \leq z_{k+1} - z_k \right\},$$

which is a linear program. Note that furthermore, in this 1-D case, a closed form expression for \mathcal{W}_1 using cumulative functions is given in (12).

\mathcal{W}_1 on Euclidean spaces In the special case of Euclidean spaces $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$, using $c(x, y) = \|x - y\|$, the global Lipschitz constraint appearing in (55) can be made local as a uniform bound on the gradient of f ,

$$\mathcal{W}_1(\alpha, \beta) = \sup_f \left\{ \int_{\mathbb{R}^d} f(d\alpha - d\beta) ; \|\nabla f\|_\infty \leq 1 \right\}. \quad (57)$$

Here the constraint $\|\nabla f\|_\infty \leq 1$ signifies that the norm of the gradient of f at any point x is upper bounded by 1, $\|\nabla f(x)\|_2 \leq 1$ for any x .

Considering the dual problem to (57), denoting $\xi \stackrel{\text{def.}}{=} \alpha - \beta$, and using that

$$\iota_{\|\cdot\|_{\mathbb{R}^d} \leq 1}(u) = \max_v \langle u, v \rangle - \|v\|_{\mathbb{R}^d}$$

one has a maximization on flow vector fields $s : \mathbb{R}^d \rightarrow \mathbb{R}^d$

$$\begin{aligned} \mathcal{W}_1(\alpha, \beta) &= \sup_f \inf_{s(x) \in \mathbb{R}^d} \int_{\mathbb{R}^d} f d\xi - \int \langle \nabla f(x), s(x) \rangle dx + \int \|s(x)\|_{\mathbb{R}^d} dx \\ &= \inf_{s(x) \in \mathbb{R}^d} \int \|s(x)\|_{\mathbb{R}^d} dx + \sup_f \int f(x) (d\xi - \text{div}(s) dx) \end{aligned}$$

one obtains an optimization problem under fixed divergence constraint

$$\mathcal{W}_1(\alpha, \beta) = \inf_s \left\{ \int_{\mathbb{R}^d} \|s(x)\|_{\mathbb{R}^d} dx ; \text{div}(s) = \alpha - \beta \right\}, \quad (58)$$

which is often called the Beckmann formulation [5]. Here the vectorial function $s(x) \in \mathbb{R}^2$ can be interpreted as a flow field, describing locally the movement of mass. Outside the support of the two input measures, $\text{div}(s) = 0$, which is the conservation of mass constraint. Once properly discretized using finite elements, Problems (57) and (58) become a nonsmooth convex optimization problems.

The previous formulations (57) and (58) of \mathcal{W}_1 can be generalized to the setting where \mathcal{X} is a Riemannian manifold, i.e. $c(x, y) = d(x, y)$ where d is the associated geodesic distance (and then for smooth manifolds, the gradient and divergence should be understood as the differential operators on manifold). In a similar way it can be extended on a graph (where the geodesic distance is the length of the shortest path), in this case, the gradient and divergence are the corresponding finite difference operations operating along the edges of the graph. In this setting, the corresponding linear program can be solved using a min-cost flow simplex in complexity $O(n^2 \log(n))$ for sparse graph (e.g. grids).

6.3 Dual norms (Integral Probability Metrics)

Formulation (57) is a special case of a dual norm. A dual norm is a convenient way to design “weak” norms that can deal with arbitrary measures. For a symmetric convex set B of measurable functions, one defines

$$\|\alpha\|_B \stackrel{\text{def.}}{=} \sup_f \left\{ \int_{\mathcal{X}} f(x) d\alpha(x) ; f \in B \right\}. \quad (59)$$

These dual norms are often called “integral probability metrics”; see [43].

Example 1 (Total variation). The total variation norm (Example 5) is a dual norm associated to the whole space of continuous functions

$$B = \{f \in \mathcal{C}(\mathcal{X}) ; \|f\|_{\infty} \leq 1\}.$$

The total variation distance is the only nontrivial divergence that is also a dual norm; see [42].

Example 2 (\mathcal{W}_1 norm). \mathcal{W}_1 as defined in (57), is a special case of dual norm (59), using

$$B = \{f ; \text{Lip}(f) \leq 1\}$$

the set of 1-Lipschitz functions.

Example 3 (Flat norm and Dudley metric). If the set B is bounded, then $\|\cdot\|_B$ is a norm on the whole space $\mathcal{M}(\mathcal{X})$ of measures. This is not the case of \mathcal{W}_1 , which is only defined for α such that $\int_{\mathcal{X}} d\alpha = 0$ (otherwise $\|\alpha\|_B = +\infty$). This can be alleviated by imposing a bound on the value of the potential f , in order to define for instance the flat norm,

$$B = \{f ; \text{Lip}(f) \leq 1 \quad \text{and} \quad \|f\|_{\infty} \leq 1\}. \quad (60)$$

It metrizes the weak convergence on the whole space $\mathcal{M}(\mathcal{X})$. Formula (56) is extended to compute the flat norm by adding the constraint $|\mathbf{f}_k| \leq 1$. The flat norm is sometimes called the “Kantorovich–Rubinstein” norm [30] and has been used as a fidelity term for inverse problems in imaging [34]. The flat norm is similar to the Dudley metric, which uses

$$B = \{f ; \|\nabla f\|_{\infty} + \|f\|_{\infty} \leq 1\}.$$

The following proposition shows that to metrize the weak convergence, the dual ball B should not be too large (because otherwise one obtain a strong norm), namely one needs $\mathcal{C}(\mathcal{X}) \subset \overline{\text{Span}(B)}$.

Proposition 18. (i) If $\mathcal{C}(\mathcal{X}) \subset \overline{\text{Span}(B)}$ (i.e. if the span of B is dense in continuous functions for the sup-norm $\|\cdot\|_{\infty}$), then $\|\alpha_k - \alpha\|_B \rightarrow 0$ implies $\alpha_k \rightarrow \alpha$.

(ii) If $B \subset \mathcal{C}(\mathcal{X})$ is compact for $\|\cdot\|_{\infty}$ then $\alpha_k \rightarrow \alpha$ implies $\|\alpha_k - \alpha\|_B \rightarrow 0$.

Proof. (i) If $\|\alpha_k - \alpha\|_B \rightarrow 0$, then by duality, for any $f \in B$, since $\langle f, \alpha_k - \alpha \rangle \leq \|\alpha_k - \alpha\|_B$ then $\langle f, \alpha_k \rangle \rightarrow \langle f, \alpha \rangle$. By linearity, this property extends to $\text{Span}(B)$. By density, this extends to $\overline{\text{Span}(B)}$, indeed $|\langle f, \alpha_k \rangle - \langle f', \alpha_k \rangle| \leq \|f - f'\|_{\infty}$.

(ii) We assume $\alpha_k \rightarrow \alpha$ and we consider a sub-sequence α_{n_k} such that

$$\|\alpha_{n_k} - \alpha\|_B \rightarrow \limsup_k \|\alpha_k - \alpha\|_B$$

Since B is compact, the maximum appearing in the definition of $\|\alpha_{n_k} - \alpha\|_B$ is reached, so that there exists some 1-Lipschitz function f_{n_k} so that $\langle \alpha_{n_k} - \alpha, f_{n_k} \rangle = \|\alpha_{n_k} - \alpha\|_B$. Once again, by Ascoli-Arzelà theorem, we can extract from $(f_{n_k})_k$ a (not relabelled for simplicity) subsequence converging to some $f \in B$. One has $\|\alpha_{n_k} - \alpha\|_B = \langle \alpha_{n_k} - \alpha, f_{n_k} \rangle$, and this quantity converges to 0 because one can decompose it as

$$\langle \alpha_{n_k} - \alpha, f_{n_k} \rangle = \langle \alpha_{n_k} - \alpha, f \rangle + \langle \alpha_{n_k}, f_{n_k} - f \rangle - \langle \alpha, f_{n_k} - f \rangle$$

and these three terms goes to zero because $\alpha_{n_k} - \alpha \rightarrow 0$ (first term) and $\|f_{n_k} - f\|_{\infty} \rightarrow 0$ (two others, recall that $|\langle \alpha_{n_k}, f_{n_k} - f \rangle| \leq \|f_{n_k} - f\|_{\infty}$). \square

Corollary 2. *On a compact space, the Wasserstein- p distance metrizes the weak convergence.*

Proof. Denoting $B = \{f ; \text{Lip}(f) \leq 1\}$.

For (i), one has that then $\text{Span}(B)$ is the space of Lipschitz functions. The adherence of Lipschitz functions for $\|\cdot\|_\infty$ is the space of continuous functions. For (ii), for probability distributions, without loss of generality, functions f in B can be taken up to an additive constant, so that we can impose $f(x_0) = 0$ for some fixed $x_0 \in \mathcal{X}$, and since \mathcal{X} is compact, $\|f\|_\infty \leq \text{diam}(\mathcal{X})$ so that we can consider in place of B another ball of equicontinuous bounded functions. By Ascoli-Arzelà theorem, it is hence compact. Proposition 8 shows that W_p has the same topology as W_1 so that it is also the topology of convergence in law. \square

Dual RKHS Norms and Maximum Mean Discrepancies. It is also possible to define “Euclidean” norms (built using quadratic functionals) on measures using the machinery of kernel methods and more specifically reproducing kernel Hilbert spaces (RKHS; see [40] for a survey of their applications in data sciences), of which we recall first some basic definitions.

Definition 4. *A symmetric function k defined on $\mathcal{X} \times \mathcal{X}$ is said to be positive definite if for any $n \geq 0$, for any family $x_1, \dots, x_n \in \mathcal{X}$ the matrix $(k(x_i, x_j))_{i,j}$ is positive (i.e. has positive eigenvalues), i.e. for all $r \in \mathbb{R}^n$*

$$\sum_{i,j=1}^n r_i r_j k(x_i, x_j) \geq 0, \quad (61)$$

The kernel is said to be conditionally positive if positivity only holds in (61) for zero mean vectors r (i.e. such that $\langle r, \mathbb{1}_n \rangle = 0$).

One of the most popular kernels is the Gaussian one $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$, which is a positive universal kernel on $\mathcal{X} = \mathbb{R}^d$. Another type of kernels are energy distances, which are more global (scale free) and are studied in Section 7.2.

If k is conditionally positive, one defines the following norm for $\xi = \alpha - \beta$ being a signed measure

$$\|\xi\|_k^2 \stackrel{\text{def.}}{=} \int_{\mathcal{X} \times \mathcal{X}} k(x, y) d\xi(x) d\xi(y). \quad (62)$$

These norms are often referred to as “maximum mean discrepancy” (MMD) (see [29]) and have also been called “kernel norms” in shape analysis [28]. This expression (62) can be rephrased, introducing two independent random vectors (X, X') on \mathcal{X} distributed with law α , as

$$\|\alpha\|_k^2 = \mathbb{E}_{X, X'}(k(X, X')).$$

One can show that $\|\cdot\|_k^2$ is the dual norm in the sense of (59) associated to the unit ball B of the RKHS associated to k . We refer to [7, 31, 40] for more details on RKHS functional spaces.

Remark 4 (Universal kernels). According to Proposition 18, the MMD norm $\|\cdot\|_k$ metrizes the weak convergence if the span of the dual ball B is dense in the space of continuous functions $\mathcal{C}(\mathcal{X})$. This means that finite sums of the form $\sum_{i=1}^n a_i k(x_i, \cdot)$ (for arbitrary choice of n and points $(x_i)_i$) are dense in $\mathcal{C}(\mathcal{X})$ for the uniform norm $\|\cdot\|_\infty$. For translation-invariant kernels over $\mathcal{X} = \mathbb{R}^d$, $k(x, y) = k_0(x - y)$, this is equivalent to having a nonvanishing Fourier transform, $\hat{k}_0(\omega) > 0$.

In the special case where α is a discrete measure of the form (??), one thus has the simple expression

$$\|\alpha\|_k^2 = \sum_{i=1}^n \sum_{i'=1}^n \mathbf{a}_i \mathbf{a}_{i'} \mathbf{k}_{i,i'} = \langle \mathbf{k} \mathbf{a}, \mathbf{a} \rangle \quad \text{where} \quad \mathbf{k}_{i,i'} \stackrel{\text{def.}}{=} k(x_i, x_{i'}).$$

In particular, when $\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i}$ and $\beta = \sum_{i=1}^n \mathbf{b}_i \delta_{x_i}$ are supported on the same set of points, $\|\alpha - \beta\|_k^2 = \langle \mathbf{k}(\mathbf{a} - \mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$, so that $\|\cdot\|_k$ is a Euclidean norm (proper if \mathbf{k} is positive definite, degenerate otherwise if \mathbf{k} is semidefinite) on the simplex Σ_n . To compute the discrepancy between two discrete measures of the form (??), one can use

$$\|\alpha - \beta\|_k^2 = \sum_{i,i'} \mathbf{a}_i \mathbf{a}_{i'} k(x_i, x_{i'}) + \sum_{j,j'} \mathbf{b}_j \mathbf{b}_{j'} k(y_j, y_{j'}) - 2 \sum_{i,j} \mathbf{a}_i \mathbf{b}_j k(x_i, y_j). \quad (63)$$

6.4 φ -divergences

We now consider a radically different class of methods to compare distributions, which are simpler to compute ($O(n)$ for discrete distributions) but never metrize the weak* convergence. Note that yet another way is possible, using Bregman divergence, which might metrize the weak* convergence in the case where the associated entropy function is weak* regular.

Definition 5 (Entropy function). *A function $\varphi : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$ is an entropy function if it is lower semicontinuous, convex, $\text{dom } \varphi \subset [0, \infty[$, and satisfies the following feasibility condition: $\text{dom } \varphi \cap]0, \infty[\neq \emptyset$. The speed of growth of φ at ∞ is described by*

$$\varphi'_\infty = \lim_{x \rightarrow +\infty} \varphi(x)/x \in \mathbb{R} \cup \{\infty\}.$$

If $\varphi'_\infty = \infty$, then φ grows faster than any linear function and φ is said *superlinear*. Any entropy function φ induces a φ -divergence (also known as Ciszár divergence [19, 2] or f -divergence) as follows.

Definition 6 (φ -Divergences). *Let φ be an entropy function. For $\alpha, \beta \in \mathcal{M}(\mathcal{X})$, let $\frac{d\alpha}{d\beta}\beta + \alpha^\perp$ be the Lebesgue decomposition¹ of α with respect to β . The divergence \mathcal{D}_φ is defined by*

$$\mathcal{D}_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi \left(\frac{d\alpha}{d\beta} \right) d\beta + \varphi'_\infty \alpha^\perp(\mathcal{X}) \quad (64)$$

if α, β are nonnegative and ∞ otherwise.

The additional term $\varphi'_\infty \alpha^\perp(\mathcal{X})$ in (64) is important to ensure that \mathcal{D}_φ defines a continuous functional (for the weak topology of measures) even if φ has a linear growth at infinity, as this is, for instance, the case for the absolute value (68) defining the TV norm. If φ has a superlinear growth, e.g. the usual entropy (67), then $\varphi'_\infty = +\infty$ so that $\mathcal{D}_\varphi(\alpha|\beta) = +\infty$ if α does not have a density with respect to β .

In the discrete setting, assuming

$$\alpha = \sum_i \mathbf{a}_i \delta_{x_i} \quad \text{and} \quad \beta = \sum_i \mathbf{b}_i \delta_{x_i} \quad (65)$$

are supported on the same set of n points $(x_i)_{i=1}^n \subset \mathcal{X}$, (64) defines a divergence on Σ_n

$$\mathbf{D}_\varphi(\mathbf{a}|\mathbf{b}) = \sum_{i \in \text{Supp}(\mathbf{b})} \varphi \left(\frac{\mathbf{a}_i}{\mathbf{b}_i} \right) \mathbf{b}_i + \varphi'_\infty \sum_{i \notin \text{Supp}(\mathbf{b})} \mathbf{a}_i, \quad (66)$$

where $\text{Supp}(\mathbf{b}) \stackrel{\text{def.}}{=} \{i \in \llbracket n \rrbracket ; b_i \neq 0\}$.

Proposition 1. *If φ is an entropy function, then \mathcal{D}_φ is jointly 1-homogeneous, convex and weakly* lower semicontinuous in (α, β) .*

Proof. **[ToDo: write me, perspective function]** □

Example 4 (Kullback–Leibler divergence). The Kullback–Leibler divergence $\text{KL} \stackrel{\text{def.}}{=} \mathcal{D}_{\varphi_{\text{KL}}}$, also known as the relative entropy, was already introduced in (31) and (25). It is the divergence associated to the Shannon–Boltzman entropy function φ_{KL} , given by

$$\varphi_{\text{KL}}(s) = \begin{cases} s \log(s) - s + 1 & \text{for } s > 0, \\ 1 & \text{for } s = 0, \\ +\infty & \text{otherwise.} \end{cases} \quad (67)$$

¹The Lebesgue decomposition theorem asserts that, given β , α admits a unique decomposition as the sum of two measures $\alpha^s + \alpha^\perp$ such that α^s is absolutely continuous with respect to β and α^\perp and β are singular.

Example 5 (Total variation). The total variation distance $\text{TV} \stackrel{\text{def}}{=} \mathcal{D}_{\varphi_{\text{TV}}}$ is the divergence associated to

$$\varphi_{\text{TV}}(s) = \begin{cases} |s - 1| & \text{for } s \geq 0, \\ +\infty & \text{otherwise.} \end{cases} \quad (68)$$

It actually defines a norm on the full space of measure $\mathcal{M}(\mathcal{X})$ where

$$\text{TV}(\alpha|\beta) = \|\alpha - \beta\|_{\text{TV}}, \quad \text{where} \quad \|\alpha\|_{\text{TV}} = |\alpha|(\mathcal{X}) = \int_{\mathcal{X}} d|\alpha|(x). \quad (69)$$

If α has a density ρ_α on $\mathcal{X} = \mathbb{R}^d$, then the TV norm is the L^1 norm on functions, $\|\alpha\|_{\text{TV}} = \int_{\mathcal{X}} |\rho_\alpha(x)| dx = \|\rho_\alpha\|_{L^1}$. If α is discrete as in (65), then the TV norm is the ℓ^1 norm of vectors in \mathbb{R}^n , $\|\alpha\|_{\text{TV}} = \sum_i |\mathbf{a}_i| = \|\mathbf{a}\|_{\ell^1}$.

Remark 5 (Strong vs. weak topology). The total variation norm (69) defines the so-called “strong” topology on the space of measure. On a compact domain \mathcal{X} of radius R , one has

$$\mathcal{W}_1(\alpha, \beta) \leq R \|\alpha - \beta\|_{\text{TV}}$$

so that this strong notion of convergence implies the weak convergence metrized by Wasserstein distances. The converse is, however, not true, since δ_x does not converge strongly to δ_y if $x \rightarrow y$ (note that $\|\delta_x - \delta_y\|_{\text{TV}} = 2$ if $x \neq y$). A chief advantage is that $\mathcal{M}_+^1(\mathcal{X})$ (once again on a compact ground space \mathcal{X}) is compact for the weak topology, so that from any sequence of probability measures $(\alpha_k)_k$, one can always extract a converging subsequence, which makes it a suitable space for several optimization problems, such as those considered in Chapter ??.

Proposition 2 (Dual expression). *A φ -divergence can be expressed using the Legendre transform*

$$\varphi^*(s) \stackrel{\text{def}}{=} \sup_{t \in \mathbb{R}} st - \varphi(t)$$

of φ as

$$\mathcal{D}_\varphi(\alpha|\beta) = \sup_{f: \mathcal{X} \rightarrow \mathbb{R}} \int_{\mathcal{X}} f(x) d\alpha(x) - \int_{\mathcal{X}} \varphi^*(f(x)) d\beta(x). \quad (70)$$

which equivalently reads that the Legendre transform of $\mathcal{D}_\varphi(\cdot|\beta)$ reads

$$\forall f \in \mathcal{C}(\mathcal{X}), \quad \mathcal{D}_\varphi^*(f|\beta) = \int_{\mathcal{X}} \varphi^*(f(x)) d\beta(x). \quad (71)$$

Proof. [ToDo: write me] □

GANs via duality. The goal is to fit a generative parametric model $\alpha_\theta = g_{\theta, \#} \zeta$ to empirical data $\beta = \frac{1}{m} \sum_j \delta_{y_j}$, where $\zeta \in \mathcal{M}_+^1(\mathcal{Z})$ is a fixed density over the latent space and $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$ is the generator, often a neural network. We consider first a dual norm (59) minimization, in which case one aim at solving a min-max saddle point problem

$$\min_{\theta} \|\alpha_\theta - \beta\|_B = \min_{\theta} \sup_{f \in B} \int_{\mathcal{X}} f(x) d(\alpha_\theta - \beta)(x) = \min_{\theta} \sup_{f \in B} \int_{\mathcal{Z}} f(g_\theta(z)) d\zeta - \frac{1}{m} \sum_j f(y_j).$$

Instead of a dual norm, one can consider any convex function and represent is as a maximization, for instance a φ -divergence, which, thanks to the dual formula (70), leads to

$$\min_{\theta} \mathcal{D}_\varphi(\alpha_\theta|\beta) = \min_{\theta} \sup_f \int_{\mathcal{X}} f(x) d\alpha_\theta(x) - \mathcal{D}_\varphi^*(f|\beta) = \min_{\theta} \sup_f \int_{\mathcal{Z}} f(g_\theta(z)) d\zeta - \frac{1}{m} \sum_j \varphi^*(f(y_j)).$$

The GAN’s idea corresponds to replacing f by a parameterized network $f = f_\xi$ and doing the maximization over the parameter ξ . For instance, Wasserstein GAN consider weight clipping by constraining $\|\xi\|_\infty \leq 1$ in order to ensure $f_\xi \in B = \{f ; \text{Lip}(f) \leq 1\}$. This set of network is both in practice smaller and non-convex so that no theoretical analysis of this method currently exists.

7 Sinkhorn Divergences

7.1 Dual of Sinkhorn

Discrete dual. The following proposition details the dual problem associated to (26).

Proposition 19. *One has*

$$L_{\mathbf{C}}^{\varepsilon}(\mathbf{a}, \mathbf{b}) = \max_{\mathbf{f} \in \mathbb{R}^n, \mathbf{g} \in \mathbb{R}^m} \langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{g}, \mathbf{b} \rangle - \varepsilon \sum_{i,j} \exp\left(\frac{\mathbf{f}_i + \mathbf{g}_j - \mathbf{C}_{i,j}}{\varepsilon}\right) \mathbf{a}_i \mathbf{b}_j + \varepsilon. \quad (72)$$

The optimal (\mathbf{f}, \mathbf{g}) are linked to scalings (\mathbf{u}, \mathbf{v}) appearing in (32) through

$$(\mathbf{u}, \mathbf{v}) = (\mathbf{a}_i e^{\mathbf{f}/\varepsilon}, \mathbf{b}_j e^{\mathbf{g}/\varepsilon}). \quad (73)$$

Proof. We introduce Lagrange multiplier and consider

$$\min_{\mathbf{P} \geq 0} \max_{\mathbf{f}, \mathbf{g}} \langle \mathbf{C}, \mathbf{P} \rangle + \varepsilon \mathbf{KL}(\mathbf{P} | \mathbf{a} \otimes \mathbf{b}) + \langle \mathbf{a} - \mathbf{P} \mathbf{1}, \mathbf{f} \rangle + \langle \mathbf{b} - \mathbf{P}^{\top} \mathbf{1}, \mathbf{g} \rangle.$$

One can check that strong duality holds since the function is continuous and that one can exchange the min with the max to get

$$\max_{\mathbf{f}, \mathbf{g}} \langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{g}, \mathbf{b} \rangle - \varepsilon \min_{\mathbf{P} \geq 0} \left\langle \frac{\mathbf{f} \oplus \mathbf{g} - \mathbf{C}}{\varepsilon}, \mathbf{P} \right\rangle - \mathbf{KL}(\mathbf{P} | \mathbf{a} \otimes \mathbf{b}) = \langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{g}, \mathbf{b} \rangle - \varepsilon \mathbf{KL}^* \left(\frac{\mathbf{f} \oplus \mathbf{g} - \mathbf{C}}{\varepsilon} | \mathbf{a} \otimes \mathbf{b} \right).$$

One concludes by verifying that **[ToDo: write the proof]** (see also formula (71))

$$\mathbf{KL}^*(H | \mathbf{a} \otimes \mathbf{b}) = \sum_{i,j} e^{H_{i,j}} \mathbf{a}_i \mathbf{b}_j - 1.$$

□

Discret dual. Since the dual problem (72) is smooth, one can consider an alternating minimization. For a fixed \mathbf{g} , one can minimize with respect to \mathbf{f} , which leads to the following equation to be solved when zeroing the derivative with respect to \mathbf{f}

$$\mathbf{a}_i - e^{\frac{\mathbf{f}_i}{\varepsilon}} \mathbf{a}_i \sum_j \exp\left(\frac{\mathbf{g}_j - \mathbf{C}_{i,j}}{\varepsilon}\right) \mathbf{b}_j = 0$$

which leads to the explicit solution

$$\mathbf{f}_i = -\varepsilon \log \sum_j \exp\left(\frac{\mathbf{g}_j - \mathbf{C}_{i,j}}{\varepsilon}\right) \mathbf{b}_j.$$

We conveniently introduce the soft-min operator of some vector $h \in \mathbb{R}^m$

$$\min_{\mathbf{b}}^{\varepsilon}(h) \stackrel{\text{def.}}{=} -\varepsilon \log \sum_j e^{-h_j/\varepsilon} \mathbf{b}_j$$

which is a smooth approximation of the minimum operator, and the optimal \mathbf{f} for a fixed \mathbf{g} is computed by a soft version of the c -transform

$$\mathbf{f}_i = \min_{\mathbf{b}}^{\varepsilon}(\mathbf{C}_{i,\cdot} - \mathbf{g}). \quad (74)$$

In a similar way, the optimal \mathbf{g} for a fixed \mathbf{f} is

$$\mathbf{g}_j = \min_{\mathbf{a}}^{\varepsilon}(\mathbf{C}_{\cdot,j} - \mathbf{f}). \quad (75)$$

Exponentiating these iterations, one retrieves exactly Sinkhorn algorithm. These iterations are however unstable for small ε . To be able to apply the algorithm in this regime, one needs to stabilize it using the celebrated log-sum-exp trick. It follows from noticing that similarly to the minimum operator, one has

$$\min_{\mathbf{b}}^{\varepsilon}(h - \text{cst}) = \min_{\mathbf{b}}^{\varepsilon}(h) - \text{cst}$$

and to replace the computation of $\min_{\mathbf{b}}^{\varepsilon}(h)$ by its stabilized version (equal when using infinite precision computation) $\min_{\mathbf{b}}^{\varepsilon}(h - \min(h)) + \min(h)$.

Continuous dual and soft-transforms. For generic (non-necessarily discrete) input measures (α, β) , the dual problem (72) reads

$$\sup_{f, g \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y})} \int_{\mathcal{X}} f(x) d\alpha(x) + \int_{\mathcal{Y}} g(y) d\beta(y) - \varepsilon \int_{\mathcal{X} \times \mathcal{Y}} \left(e^{\frac{f \oplus g - c}{\varepsilon}} - 1 \right) d\alpha \otimes d\beta \quad (76)$$

This corresponds to a smoothing of the constraint $\mathcal{R}(c)$ appearing in the original problem (45), which is retrieved in the limit $\varepsilon \rightarrow 0$.

The corresponding soft c -transform, which minimize this dual problem with respect to either f or g reads

$$\begin{aligned} \forall y \in \mathcal{Y}, \quad f^{c, \varepsilon}(y) &\stackrel{\text{def.}}{=} -\varepsilon \log \left(\int_{\mathcal{X}} e^{\frac{-c(x, y) + f(x)}{\varepsilon}} d\alpha(x) \right), \\ \forall x \in \mathcal{X}, \quad g^{c, \varepsilon}(x) &\stackrel{\text{def.}}{=} -\varepsilon \log \left(\int_{\mathcal{Y}} e^{\frac{-c(x, y) + g(y)}{\varepsilon}} d\beta(y) \right). \end{aligned}$$

In the case of discrete measures, one retrieves the formula (74) and (75).

We omit the details, but similarly to the unregularized case, one can define an entropic semi-discrete problem and develop stochastic optimization method to solve it.

[ToDo: dual potentials are convex functions for $W2_{\varepsilon}$]

[ToDo: sinkhorn between Gaussians is a Gaussian]

[ToDo: Sinkhorn is smooth, computes is Eulerian gradient Lagrangian gradient Fitting, auto diff]

7.2 Sinkhorn Divergences

Entropic bias. A major issue of the value of Sinkhorn problem (30) is that $\mathcal{L}_c^{\varepsilon}(\alpha, \beta) > 0$. So in particular,

$$\alpha_{\varepsilon} = \underset{\beta}{\operatorname{argmin}} \mathcal{L}_c^{\varepsilon}(\alpha, \beta)$$

does not satisfy $\alpha_{\varepsilon} = \alpha$ unless $\varepsilon = 0$. The following proposition shows that the bias induced by this entropic regularization has a catastrophic influence in the large ε limit.

Proposition 20. *One has $\mathcal{L}_c^{\varepsilon}(\alpha, \beta) \rightarrow \int c d\alpha \otimes \beta$ as $\varepsilon \rightarrow +\infty$.*

Proof. The intuition of the proof follows from the fact that the optimal coupling converges to $\alpha \otimes \beta$. □

So in the large ε limit, $\mathcal{L}_c^{\varepsilon}$ behaves like an inner product and not like a norm. For instance, in the case

$$\alpha_{\varepsilon} \rightarrow \min_{\beta} \left\langle \int c(x, \cdot) d\alpha(x), \beta \right\rangle = \delta_{y^*(\alpha)} \quad \text{where} \quad y^*(\alpha) = \underset{y}{\operatorname{argmin}} \int c(x, y) d\alpha(x).$$

For instance, when $c(x, y) = \|x - y\|^2$ then α_{ε} collapses towards a Dirac located at the mean $\int x d\alpha(x)$ of α .

Sinkhorn divergences. The usual way to go from an inner product to a norm is to use the polarization formula, we thus also consider for the Sinkhorn cost, in order to define the debiased Sinkhorn divergence

$$\bar{\mathcal{L}}_c^\varepsilon(\alpha, \beta) \stackrel{\text{def.}}{=} \mathcal{L}_c^\varepsilon(\alpha, \beta) - \frac{1}{2}\mathcal{L}_c^\varepsilon(\alpha, \alpha) - \frac{1}{2}\mathcal{L}_c^\varepsilon(\alpha, \beta).$$

It is not (yet) at all clear why this quantity should be positive.

Before going on, we prove a fundamental lemma which states that the dual cost has a simple form where the regularization vanish at a solution (and actually it vanishes also during Sinkhorn's iteration by the same proof).

Lemma 1. *Denoting $(f_{\alpha, \beta}, g_{\alpha, \beta})$ optimal dual potentials (which can be shown to be unique up to an additive constant), one has*

$$\mathcal{L}_c^\varepsilon(\alpha, \beta) = \langle f_{\alpha, \beta}, \alpha \rangle + \langle g_{\alpha, \beta}, \beta \rangle. \quad (77)$$

Proof. We first notice that at optimality, the relation

$$f_{\alpha, \beta} = -\varepsilon \log \int_{\mathcal{Y}} e^{\frac{g_{\alpha, \beta}(y) - c(x, y)}{\varepsilon}} d\beta(y)$$

after taking the exponential, equivalently reads

$$1 = \int_{\mathcal{Y}} e^{\frac{f_{\alpha, \beta}(x) + g_{\alpha, \beta}(y) - c(x, y)}{\varepsilon}} d\beta(y) \implies \int_{\mathcal{X} \times \mathcal{Y}} \left(e^{\frac{f_{\alpha, \beta} \oplus g_{\alpha, \beta} - c}{\varepsilon}} - 1 \right) d\alpha \otimes \beta = 0.$$

Plugging this in formula (76), one obtains the result. \square

Let us first show that its asymptotic makes sense.

Proposition 21. *One has $\bar{\mathcal{L}}_c^\varepsilon(\alpha, \beta) \rightarrow \mathcal{L}_c(\alpha, \beta)$ when $\varepsilon \rightarrow 0$ and*

$$\bar{\mathcal{L}}_c^\varepsilon(\alpha, \beta) \rightarrow \frac{1}{2} \int -cd(\alpha - \beta) \otimes d(\alpha - \beta) \quad \text{when } \varepsilon \rightarrow +\infty.$$

Proof. For discrete measures, the convergence is already proved in Proposition (28), we now give a general treatment. **Case $\varepsilon \rightarrow 0$.** **[ToDo: Red the proof of $\varepsilon \rightarrow 0$ on non-discrete spaces]** **Case $\varepsilon \rightarrow +\infty$.** We denote $(f_\varepsilon, g_\varepsilon)$ optimal dual potential. Optimality condition on f_ε (equivalently Sinkhorn fixed point on f_ε) reads

$$\begin{aligned} f_\varepsilon &= -\varepsilon \log \int \exp \left(\frac{g_\varepsilon(y) - c(\cdot, y)}{\varepsilon} \right) d\beta(y) = -\varepsilon \log \int \left(1 + \frac{g_\varepsilon(y) - c(\cdot, y)}{\varepsilon} + o(1/\varepsilon) \right) d\beta(y) \\ &= -\varepsilon \int \left(\frac{g_\varepsilon(y) - c(\cdot, y)}{\varepsilon} + o(1/\varepsilon) \right) d\beta(y) = - \int g_\varepsilon d\beta + \int c(\cdot, y) d\beta(y) + o(1). \end{aligned}$$

Plugging this relation in the dual expression (77)

$$\mathcal{L}_c^\varepsilon(\alpha, \beta) = \int f_\varepsilon d\alpha + \int g_\varepsilon d\beta = - \iint c(x, y) d\alpha(x) d\beta(y) + o(1).$$

\square

In the case where $-c$ defines a conditionnaly positive definite kernel, then $\bar{\mathcal{L}}_c^\varepsilon(\alpha, \beta)$ converges to the square of an Hilbertian kernel norm. A typical example is when $c(x, y) = \|x - y\|^p$ for $0 < p < 2$, which corresponds to the so-called Energy distance kernel. This kernel norm is the dual of a homogeneous Sobolev norm.

We now show that this debiased Sinkhorn divergence is positive.

Proposition 22. *If $k(x, y) = e^{-c(x, y)/\varepsilon}$ is positive definite, then $\bar{\mathcal{L}}_c^\varepsilon(\alpha, \beta) \geq 0$ and is zero if and only if $\alpha = \beta$.*

Proof. In the following, we denote $(f_{\alpha, \beta}, g_{\alpha, \beta})$ optimal dual potential for the dual Schrodinger problem between α and β . We denote $f_{\alpha, \alpha} = g_{\alpha, \alpha}$ (one can assume they are equal by symmetry) the solution for the problem between α and itself. Using the suboptimal function $(f_{\alpha, \alpha}, g_{\beta, \beta})$ in the dual maximization problem, and using relation (77) for the simplified expression of the dual cost, one obtains

$$\mathcal{L}_c^\varepsilon(\alpha, \beta) \geq \langle f_{\alpha, \alpha}, \alpha \rangle + \langle g_{\beta, \beta}, \beta \rangle - \varepsilon \langle e^{\frac{f_{\alpha, \alpha} + g_{\beta, \beta} - c}{\varepsilon}} - 1, \alpha \otimes \beta \rangle$$

But one has $\langle f_{\alpha, \alpha}, \alpha \rangle = \frac{1}{2} \mathcal{L}_c^\varepsilon(\alpha, \alpha)$ and same for β , so that the previous inequality equivalently reads

$$\frac{1}{\varepsilon} \bar{\mathcal{L}}_c^\varepsilon(\alpha, \beta) \geq 1 - \langle e^{\frac{f_{\alpha, \alpha} + g_{\beta, \beta} - c}{\varepsilon}} - 1, \alpha \otimes \beta \rangle = 1 - \langle \tilde{\alpha}, \tilde{\beta} \rangle_k$$

where $\tilde{\alpha} = e^{f_{\alpha, \alpha}} \alpha$, $\tilde{\beta} = e^{f_{\beta, \beta}} \beta$ and we introduced the inner product (which is a valid one because k is positive) $\langle \tilde{\alpha}, \tilde{\beta} \rangle_k \stackrel{\text{def}}{=} \int k(x, y) d\tilde{\alpha}(x) d\tilde{\beta}(y)$. One note that Sinkhorn fixed point equation, once exponentiated, reads $e^{f_{\alpha, \alpha}} \odot [k(\tilde{\alpha})] = 1$ and hence

$$\|\tilde{\alpha}\|_k^2 = \langle k(\tilde{\alpha}), \tilde{\alpha} \rangle = \langle e^{f_{\alpha, \alpha}} \odot k(\tilde{\alpha}), \alpha \rangle = \langle 1, \alpha \rangle = 1$$

and similarly $\|\tilde{\beta}\|_k^2 = 1$. So by Cauchy-Schwartz, one has $1 - \langle \tilde{\alpha}, \tilde{\beta} \rangle_k \geq 0$. Showing strict positivity is more involved, and is not proved here. \square

One can furthermore show that this debiased divergence metrizes the convergence in law.

8 Barycenters

8.1 Frechet Mean over the Wasserstein Space

This barycenter problem (79) was originally introduced by [1] following earlier ideas of [16]. They proved in particular uniqueness of the barycenter for $c(x, y) = \|x - y\|^2$ over $\mathcal{X} = \mathbb{R}^d$, if one of the input measure has a density with respect to the Lebesgue measure (and more generally under the same hypothesis as the one guaranteeing the existence of a Monge map, see Remark ??).

Given a set of input measure $(\beta_s)_s$ defined on some space \mathcal{X} , the barycenter problem becomes

$$\min_{\alpha \in \mathcal{M}_+^1(\mathcal{X})} \sum_{s=1}^S \lambda_s \mathcal{L}_c(\alpha, \beta_s). \quad (78)$$

In the case where $\mathcal{X} = \mathbb{R}^d$ and $c(x, y) = \|x - y\|^2$, [1] shows that if one of the input measures has a density, then this barycenter is unique. Problem (78) can be viewed as a generalization of the problem of computing barycenters of points $(x_s)_{s=1}^S \in \mathcal{X}^S$ to arbitrary measures. Indeed, if $\beta_s = \delta_{x_s}$ is a single Dirac mass, then a solution to (78) is δ_{x^*} where x^* is a Fréchet mean solving (?). Note that for $c(x, y) = \|x - y\|^2$, the mean of the barycenter α^* is necessarily the barycenter of the mean, *i.e.*

$$\int_{\mathcal{X}} x d\alpha^*(x) = \sum_s \lambda_s \int_{\mathcal{X}} x d\alpha_s(x),$$

and the support of α^* is located in the convex hull of the supports of the $(\alpha_s)_s$. The consistency of the approximation of the infinite dimensional optimization (78) when approximating the input distribution using discrete ones (and thus solving (79) in place) is studied in [17]. Let us also note that it is possible to recast (78) as a multi-marginal OT problem, see Remark ??.

[ToDo: Write me: existence, dual, Monge map]

8.2 1-D Case

[ToDo: Write me]

8.3 Gaussians Case

[ToDo: Write me]

8.4 Discrete Barycenters

Given input histogram $\{\mathbf{b}_s\}_{s=1}^S$, where $\mathbf{b}_s \in \Sigma_{n_s}$, and weights $\lambda \in \Sigma_S$, a Wasserstein barycenter is computed by minimizing

$$\min_{\mathbf{a} \in \Sigma_n} \sum_{s=1}^S \lambda_s L_{\mathbf{C}_s}(\mathbf{a}, \mathbf{b}_s) \quad (79)$$

where the cost matrices $\mathbf{C}_s \in \mathbb{R}^{n \times n_s}$ need to be specified. A typical setup is “Eulerian”, so that all the barycenters are defined on the same grid, $n_s = n$, $\mathbf{C}_s = \mathbf{C} = \mathbf{D}^p$ is set to be a distance matrix, so that one solves

$$\min_{\mathbf{a} \in \Sigma_n} \sum_{s=1}^S \lambda_s W_p^p(\mathbf{a}, \mathbf{b}_s).$$

The barycenter problem for histograms (79) is in fact a linear program, since one can look for the S couplings $(\mathbf{P}_s)_s$ between each input and the barycenter itself

$$\min_{\mathbf{a} \in \Sigma_n, (\mathbf{P}_s \in \mathbb{R}^{n \times n_s})_s} \left\{ \sum_{s=1}^S \lambda_s \langle \mathbf{P}_s, \mathbf{C}_s \rangle ; \forall s, \mathbf{P}_s^\top \mathbf{1}_{n_s} = \mathbf{a}, \mathbf{P}_s^\top \mathbf{1}_n = \mathbf{b}_s \right\}.$$

Although this problem is an LP, its scale forbids the use generic solvers for medium scale problems. One can therefore resort to using first order methods such as subgradient descent on the dual [17].

8.5 Sinkhorn for barycenters

[ToDo: Explain the key difference with the regularized OT problem: here there is no more a “canonical” reference measure $\alpha \otimes \beta$ since the barycenter is unknown.]

One can use entropic smoothing and approximate the solution of (79) using

$$\min_{\mathbf{a} \in \Sigma_n} \sum_{s=1}^S \lambda_s L_{\mathbf{C}_s}^\varepsilon(\mathbf{a}, \mathbf{b}_s) \quad (80)$$

for some $\varepsilon > 0$. This is a smooth convex minimization problem, which can be tackled using gradient descent [20]. An alternative is to use descent method (typically quasi-Newton) on the semi-dual [21], which is useful to integrate additional regularizations on the barycenter (e.g. to impose some smoothness). A simple but effective approach, as remarked in [6] is to rewrite (80) as a (weighted) KL projection problem

$$\min_{(\mathbf{P}_s)_s} \left\{ \sum_s \lambda_s \text{KL}(\mathbf{P}_s | \mathbf{K}_s) ; \forall s, \mathbf{P}_s^\top \mathbf{1}_m = \mathbf{b}_s, \mathbf{P}_1 \mathbf{1}_1 = \dots = \mathbf{P}_S \mathbf{1}_S \right\} \quad (81)$$

where we denoted $\mathbf{K}_s \stackrel{\text{def.}}{=} e^{-\mathbf{C}_s/\varepsilon}$. Here, the barycenter \mathbf{a} is implicitly encoded in the row marginals of all the couplings $\mathbf{P}_s \in \mathbb{R}^{n \times n_s}$ as $\mathbf{a} = \mathbf{P}_1 \mathbf{1}_1 = \dots = \mathbf{P}_S \mathbf{1}_S$. As detailed in [6], one can generalize Sinkhorn to this problem, which also corresponds to iterative projection. This can also be seen as a special case of the generalized Sinkhorn detailed in §??. The optimal couplings $(\mathbf{P}_s)_s$ solving (81) are computed in scaling form as

$$\mathbf{P}_s = \text{diag}(\mathbf{u}_s) \mathbf{K} \text{diag}(\mathbf{v}_s), \quad (82)$$

and the scalings are sequentially updated as

$$\forall s \in \llbracket 1, S \rrbracket, \quad \mathbf{v}_s^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{b}_s}{\mathbf{K}_s^T \mathbf{u}_s^{(\ell)}}, \quad (83)$$

$$\forall s \in \llbracket 1, S \rrbracket, \quad \mathbf{u}_s^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{a}^{(\ell+1)}}{\mathbf{K}_s \mathbf{v}_s^{(\ell+1)}}, \quad (84)$$

$$\text{where } \mathbf{a}^{(\ell+1)} \stackrel{\text{def.}}{=} \prod_s (\mathbf{K}_s \mathbf{v}_s^{(\ell+1)})^{\lambda_s}. \quad (85)$$

An alternative way to derive these iterations is to perform alternate minimization on the variables of a dual problem, which detailed in the following proposition.

Proposition 23. *The optimal $(\mathbf{u}_s, \mathbf{v}_s)$ appearing in (82) can be written as $(\mathbf{u}_s, \mathbf{v}_s) = (e^{\mathbf{f}_s/\varepsilon}, e^{\mathbf{g}_s/\varepsilon})$ where $(\mathbf{f}_s, \mathbf{g}_s)_s$ are the solutions of the following program (whose value matches the one of (80))*

$$\max_{(\mathbf{f}_s, \mathbf{g}_s)_s} \left\{ \sum_s \lambda_s \left(\langle \mathbf{g}_s, \mathbf{b}_s \rangle - \varepsilon \langle \mathbf{K}_s e^{\mathbf{g}_s/\varepsilon}, e^{\mathbf{f}_s/\varepsilon} \rangle \right) ; \sum_s \lambda_s \mathbf{f}_s = 0 \right\}. \quad (86)$$

Proof. Introducing Lagrange multipliers in (81) leads to

$$\min_{(\mathbf{P}_s)_s, \mathbf{a}} \max_{(\mathbf{f}_s, \mathbf{g}_s)_s} \sum_s \lambda_s \left(\varepsilon \mathbf{KL}(\mathbf{P}_s | \mathbf{K}_s) + \langle \mathbf{a} - \mathbf{P}_s \mathbf{1}_m, \mathbf{f}_s \rangle \right. \\ \left. + \langle \mathbf{b}_s - \mathbf{P}_s^T \mathbf{1}_m, \mathbf{g}_s \rangle \right).$$

Strong duality holds, so that one can exchange the min and the max, and gets

$$\max_{(\mathbf{f}_s, \mathbf{g}_s)_s} \sum_s \lambda_s \left(\langle \mathbf{g}_s, \mathbf{b}_s \rangle + \min_{\mathbf{P}_s} \varepsilon \mathbf{KL}(\mathbf{P}_s | \mathbf{K}_s) - \langle \mathbf{P}_s, \mathbf{f}_s \oplus \mathbf{g}_s \rangle \right) \\ + \min_{\mathbf{a}} \left\langle \sum_s \lambda_s \mathbf{f}_s, \mathbf{a} \right\rangle.$$

The explicit minimization on \mathbf{a} gives the constraint $\sum_s \lambda_s \mathbf{f}_s = 0$ together with

$$\max_{(\mathbf{f}_s, \mathbf{g}_s)_s} \sum_s \lambda_s \langle \mathbf{g}_s, \mathbf{b}_s \rangle - \varepsilon \mathbf{KL}^* \left(\frac{\mathbf{f}_s \oplus \mathbf{g}_s}{\varepsilon} | \mathbf{K}_s \right)$$

where $\mathbf{KL}^*(\cdot | \mathbf{K}_s)$ is the Legendre transform (71) of the function $\mathbf{KL}^*(\cdot | \mathbf{K}_s)$. This Legendre transform reads

$$\mathbf{KL}^*(\mathbf{U} | \mathbf{K}) = \sum_{i,j} \mathbf{K}_{i,j} (e^{\mathbf{U}_{i,j}} - 1), \quad (87)$$

which shows the desired formula. To show (87), since this function is separable, one needs to compute

$$\forall (u, k) \in \mathbb{R}_+^2, \quad \mathbf{KL}^*(u | k) \stackrel{\text{def.}}{=} \max_r ur - (r \log(r/k) - r + k)$$

whose optimality condition reads $u = \log(r/k)$, i.e. $r = ke^u$, hence the result. \square

Minimizing (86) with respect to each \mathbf{g}_s , while keeping all the other variable fixed, is obtained in closed form by (83). Minimizing (86) with respect to all the $(\mathbf{f}_s)_s$ requires to solve for \mathbf{a} using (85) and leads to the expression (84).

Figures ?? and ?? show applications to 2-D and 3-D shapes interpolation. Figure ?? shows a computation of barycenters on a surface, where the ground cost is the square of the geodesic distance. For this figure, the computations are performed using the geodesic in heat approximation detailed in Remark ?. We refer to [41] for more details and other applications to computer graphics and imaging sciences.

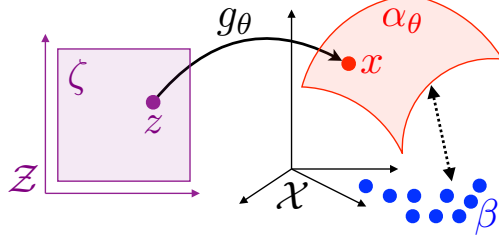


Figure 3: Schematic display of the density fitting problem 89.

9 Wasserstein Estimation

9.1 Wasserstein Loss

In statistics, text processing or imaging, one must usually compare a probability distribution β arising from measurements to a model, namely a parameterized family of distributions $\{\alpha_\theta, \theta \in \Theta\}$ where Θ is a subset of an Euclidean space. Such a comparison is done through a “loss” or a “fidelity” term, which, in this section, is the Wasserstein distance. In the simplest scenario, the computation of a suitable parameter θ is obtained by minimizing directly

$$\min_{\theta \in \Theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \mathcal{L}_c(\alpha_\theta, \beta). \quad (88)$$

Of course, one can consider more complicated problems: for instance, the barycenter problem described in §?? consists in a sum of such terms. However, most of these more advanced problems can be usually solved by adapting tools defined for basic case: either using the chain rule to compute explicitly derivatives, or using automatic differentiation.

The Wasserstein distance between two histograms or two densities is convex with respect to these inputs, as shown by (70) and (45) respectively. Therefore, when the parameter θ is itself a histogram, namely $\Theta = \Sigma_n$ and $\alpha_\theta = \theta$, or more generally when θ describes K weights in the simplex, $\Theta = \Sigma_K$, and $\alpha_\theta = \sum_{i=1}^K \theta_i \alpha_i$ is a convex combination of known atoms $\alpha_1, \dots, \alpha_K$ in Σ_N , Problem (88) remains convex (the first case corresponds to the barycenter problem, the second to one iteration of the dictionary learning problem with a Wasserstein loss [38]). However, for more general parameterizations $\theta \mapsto \alpha_\theta$, Problem (88) is in general not convex.

A practical problem of paramount importance in statistic and machine learning is density fitting. Given some discrete samples $(x_i)_{i=1}^n \subset \mathcal{X}$ from some unknown distribution, the goal is to fit a parametric model $\theta \mapsto \alpha_\theta \in \mathcal{M}(\mathcal{X})$ to the observed empirical input measure β

$$\min_{\theta \in \Theta} \mathcal{L}(\alpha_\theta, \beta) \quad \text{where} \quad \beta = \frac{1}{n} \sum_i \delta_{x_i}, \quad (89)$$

where \mathcal{L} is some “loss” function between a discrete and a “continuous” (arbitrary) distribution (see Figure 3).

In the case where α_θ as a density $\rho_\theta \stackrel{\text{def.}}{=} \rho_{\alpha_\theta}$ with respect to the Lebesgue measure (or any other fixed reference measure), the maximum likelihood estimator (MLE) is obtained by solving

$$\min_{\theta} \mathcal{L}_{\text{MLE}}(\alpha_\theta, \beta) \stackrel{\text{def.}}{=} - \sum_i \log(\rho_\theta(x_i)).$$

This corresponds to using an empirical counterpart of a Kullback-Leibler loss since, assuming the x_i are i.i.d. samples of some $\bar{\beta}$, then

$$\mathcal{L}_{\text{MLE}}(\alpha, \beta) \xrightarrow{n \rightarrow +\infty} \text{KL}(\alpha | \bar{\beta})$$

This MLE approach is known to lead to optimal estimation procedures in many cases (see for instance [36]). However, it fails to work when estimating singular distributions, typically when the α_θ does not

has a density (so that $\mathcal{L}_{\text{MLE}}(\alpha_\theta, \beta) = +\infty$) or when $(x_i)_i$ are samples from some singular $\bar{\beta}$ (so that the α_θ should share the same support as β for $\text{KL}(\alpha|\bar{\beta})$ to be finite, but this support is usually unknown). Another issue is that in several cases of practical interest, the density ρ_θ is inaccessible (or too hard to compute).

A typical setup where both problems (singular and unknown densities) occur is for so-called generative models, where the parametric measure is written as a push-forward of a fixed reference measure $\zeta \in \mathcal{M}(\mathcal{Z})$

$$\alpha_\theta = h_{\theta,\#}\zeta \quad \text{where} \quad h_\theta : \mathcal{Z} \rightarrow \mathcal{X}$$

where the push-forward operator is introduced in Definition 1. The space \mathcal{Z} is usually low-dimensional, so that the support of α_θ is localized along a low-dimensional “manifold” and the resulting density is highly singular (it does not have a density with respect to Lebesgue measure). Furthermore, computing this density is usually intractable, while generating i.i.d. samples from α_θ is achieved by computing $x_i = h_\theta(z_i)$ where $(z_i)_i$ are i.i.d. samples from ζ .

In order to cope with such difficult scenario, one has to use weak metrics in place of the MLE functional \mathcal{L}_{MLE} , which needs to be written in dual form as

$$\mathcal{L}(\alpha, \beta) \stackrel{\text{def.}}{=} \max_{(f,g) \in \mathcal{C}(\mathcal{X})^2} \left\{ \int_{\mathcal{X}} f(x) d\alpha(x) + \int_{\mathcal{X}} g(x) d\beta(x) ; (f, g) \in \mathcal{R} \right\}. \quad (90)$$

Dual norms exposed in §6.3 correspond to imposing $\mathcal{R} = \{(f, -f) ; f \in B\}$, while optimal transport (45) sets $\mathcal{R} = \mathcal{R}(c)$ as defined in (46).

For a fixed θ , evaluating the energy to be minimized in (89) using such a loss function corresponds to solving a semi-discrete optimal transport, which is the focus of Chapter ???. Minimizing the energy with respect to θ is much more involved, and is typically highly non-convex.

The class of estimators obtained using $\mathcal{L} = \mathcal{L}_c$, often called “Minimum Kantorovitch Estimators” (MKE), was initially introduced in [4], see also [14].

9.2 Wasserstein Derivatives

[ToDo: Write me.]

Eulerian vs Lagrangian.

Derivatives.

Sinkhorn smoothing.

9.3 Sample Complexity

In an applied setting, given two input measures $(\alpha, \beta) \in \mathcal{M}_+^1(\mathcal{X})^2$, an important statistical problem is to approximate the (usually unknown) divergence $D(\alpha, \beta)$ using only samples $(x_i)_{i=1}^n$ from α and $(y_j)_{j=1}^m$ from β . These samples are assumed to be independently identically distributed from their respective distributions. For both Wasserstein distances \mathcal{W}_p (see 24) and MMD norms (see §6.3), a straightforward estimator of the unknown distance between distributions is compute it directly between the empirical measures, hoping ideally that one can control the rate of convergence of the latter to the former,

$$D(\alpha, \beta) \approx D(\alpha_n, \beta_m) \quad \text{where} \quad \begin{cases} \alpha_n \stackrel{\text{def.}}{=} \frac{1}{n} \sum_i \delta_{x_i}, \\ \beta_m \stackrel{\text{def.}}{=} \frac{1}{m} \sum_j \delta_{y_j}. \end{cases}$$

Note that here both α_n and β_m are random measures, so $D(\alpha_n, \beta_m)$ is a random number. For simplicity, we assume that \mathcal{X} is compact (handling unbounded domain requires extra constraints on the moments of the input measures).

For such a dual distance that metrizes the weak convergence (see Definition 2), since there is the weak convergence $\hat{\alpha}_n \rightarrow \alpha$, one has $D(\alpha_n, \beta_n) \rightarrow D(\alpha, \beta)$ as $n \rightarrow +\infty$. But an important question is the speed of convergence of $D(\alpha_n, \beta_n)$ toward $D(\alpha, \beta)$, and this rate is often called the “sample complexity” of D .

Rates for OT. For $\mathcal{X} = \mathbb{R}^d$ and measure supported on bounded domain, it is shown by [23] that for $d > 2$, and $1 \leq p < +\infty$,

$$\mathbb{E}(|\mathcal{W}_p(\alpha_n, \beta_n) - \mathcal{W}_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}}),$$

where the expectation \mathbb{E} is taken with respect to the random samples $(x_i, y_i)_i$. This rate is tight in \mathbb{R}^d if one of the two measures has a density with respect to the Lebesgue measure. This result was proved for general metric spaces [23] using the notion of covering numbers and was later refined, in particular for $\mathcal{X} = \mathbb{R}^d$ in [22, 25]. This rate can be refined when the measures are supported on low-dimensional subdomains: [45] show that, indeed, the rate depends on the intrinsic dimensionality of the support. [45] also study the nonasymptotic behavior of that convergence, such as for measures which are discretely approximated (*e.g.* mixture of Gaussians with small variances). It is also possible to prove concentration of $\mathcal{W}_p(\alpha_n, \beta_n)$ around its mean $\mathcal{W}_p(\alpha, \beta)$; see [10, 9, 45].

Rates for MMD. For weak norms $\|\cdot\|_k^2$ which are dual of RKHS norms (also called MMD), as defined in (62), and contrary to Wasserstein distances, the sample complexity rate does not depend on the ambient dimension

$$\mathbb{E}(|\|\alpha_n - \beta_n\|_k^2 - \|\alpha - \beta\|_k^2|) = O(n^{-\frac{1}{2}}),$$

see [43]. Note however that the constant appearing in this rate might depend on the dimension. This corresponds to the classical rate when using a Monte-Carlo method to estimate an integral using random samples. For instance, one has, denoting $\xi = \alpha - \beta$ and $\xi_n = \alpha_n - \beta_n$

$$\mathbb{E}(|\|\xi_n\|_k - \|\alpha - \beta\|_k|) = \mathbb{E}(|\int kd(\xi \otimes \xi - \xi_n \otimes \xi_n)|^2).$$

[ToDo: Explain that this corresponds to the Monte-Carlo approximation of integrals using sums. Explain that the constant might depends on the dimension. Give a proof.]

Rates for Sinkhorn. **[ToDo: Give the intuition : smoothness of the potentials, Sobolev ball.]**

10 Gradient Flows

10.1 Optimization over Measures

Example : neural net training, super-resolution, and other functional over measures.
Eulerian vs Lagrangian derivative

10.2 Particle System and Lagrangian Flows

The intuition : at a Lagrangian OT as ℓ_2 metric on points. OT flow is a flow on particle locations.
Gradient descent schemes.
Study of mean field limits.

10.3 Wasserstein Gradient Flows

Implicit stepping.
Fokker planck,
Unbalanced gradient flows.

10.4 Langevin Flows

1 random particles as opposed to many deterministic particles. Crucial in high dimension.

11 Extensions

11.1 Dynamical formulation

[ToDo: Write me.]

11.2 Unbalanced OT

[ToDo: Write me.]

11.3 Gromov Wasserstein

Optimal transport needs a ground cost \mathbf{C} to compare histograms (\mathbf{a}, \mathbf{b}) , it can thus not be used if the histograms are not defined on the same underlying space, or if one cannot pre-register these spaces to define a ground cost. To address this issue, one can instead only assume a weaker assumption, namely that one has at its disposal two matrices $\mathbf{D} \in \mathbb{R}^{n \times n}$ and $\mathbf{D}' \in \mathbb{R}^{m \times m}$ that represent some relationship between the points on which the histograms are defined. A typical scenario is when these matrices are (power of) distance matrices. The Gromov-Wasserstein problem reads

$$\text{GW}((\mathbf{a}, \mathbf{D}), (\mathbf{b}, \mathbf{D}'))^2 \stackrel{\text{def.}}{=} \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \mathcal{E}_{\mathbf{D}, \mathbf{D}'}(\mathbf{P}) \stackrel{\text{def.}}{=} \sum_{i,j,i',j'} |\mathbf{D}_{i,i'} - \mathbf{D}'_{j,j'}|^2 \mathbf{P}_{i,j} \mathbf{P}_{i',j'}. \quad (91)$$

This is a non-convex problem, which can be recast as a Quadratic Assignment Problem (QAP) [?] and is in full generality NP-hard to solve for arbitrary inputs. It is in fact equivalent to a graph matching problem [?] for a particular cost.

One can show that GW satisfies the triangular inequality, and in fact it defines a distance between metric spaces equipped with a probability distribution (here assumed to be discrete in definition (91)) up to isometries preserving the measures. This distance was introduced and studied in details by Memoli in [?]. An in-depth mathematical exposition (in particular, its geodesic structure and gradient flows) is given in [?]. See also [?] for applications in computer vision. This distance is also tightly connected with the Gromov-Hausdorff distance [?] between metric spaces, which have been used for shape matching [?, ?].

Remark 6. Gromov-Wasserstein distance The general setting corresponds to computing couplings between metric measure spaces $(\mathcal{X}, d_{\mathcal{X}}, \alpha_{\mathcal{X}})$ and $(\mathcal{Y}, d_{\mathcal{Y}}, \alpha_{\mathcal{Y}})$ where $(d_{\mathcal{X}}, d_{\mathcal{Y}})$ are distances and $(\alpha_{\mathcal{X}}, \alpha_{\mathcal{Y}})$ are measures on their respective spaces. One defines

$$\mathcal{GW}((\alpha_{\mathcal{X}}, d_{\mathcal{X}}), (\alpha_{\mathcal{Y}}, d_{\mathcal{Y}}))^2 \stackrel{\text{def.}}{=} \min_{\pi \in \mathbf{U}(\alpha_{\mathcal{X}}, \alpha_{\mathcal{Y}})} \int_{\mathcal{X}^2 \times \mathcal{Y}^2} |d_{\mathcal{X}}(x, x') - d_{\mathcal{Y}}(y, y')|^2 d\pi(x, y) d\pi(x', y'). \quad (92)$$

\mathcal{GW} defines a distance between metric measure spaces up to isometries, where one says that $(\alpha_{\mathcal{X}}, d_{\mathcal{X}})$ and $(\alpha_{\mathcal{Y}}, d_{\mathcal{Y}})$ are isometric if there exists $\varphi : \mathcal{X} \rightarrow \mathcal{Y}$ such that $\varphi_{\#} \alpha_{\mathcal{X}} = \alpha_{\mathcal{Y}}$ and $d_{\mathcal{Y}}(\varphi(x), \varphi(x')) = d_{\mathcal{X}}(x, x')$.

Remark 7. Gromov-Wasserstein geodesics The space of metric spaces (up to isometries) endowed with this \mathcal{GW} distance (92) has a geodesic structure. [?] shows that the geodesic between $(\mathcal{X}_0, d_{\mathcal{X}_0}, \alpha_0)$ and $(\mathcal{X}_1, d_{\mathcal{X}_1}, \alpha_1)$ can be chosen to be $t \in [0, 1] \mapsto (\mathcal{X}_0 \times \mathcal{X}_1, d_t, \pi^*)$ where π^* is a solution of (92) and for all $((x_0, x_1), (x'_0, x'_1)) \in (\mathcal{X}_0 \times \mathcal{X}_1)^2$,

$$d_t((x_0, x_1), (x'_0, x'_1)) \stackrel{\text{def.}}{=} (1-t)d_{\mathcal{X}_0}(x_0, x'_0) + td_{\mathcal{X}_1}(x_1, x'_1).$$

This formula allows one to define and analyze gradient flows which minimize functionals involving metric spaces, see [?]. It is however difficult to handle numerically, because it involves computations over the product space $\mathcal{X}_0 \times \mathcal{X}_1$. A heuristic approach is used in [?] to define geodesics and barycenters of metric measure spaces while imposing the cardinality of the involved spaces and making use of the entropic smoothing (93) detailed below.

To approximate the computation of GW, and to help convergence of minimization schemes to better minima, one can consider the entropic regularized variant

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \mathcal{E}_{\mathbf{D}, \mathbf{D}'}(\mathbf{P}) - \varepsilon \mathbf{H}(\mathbf{P}). \quad (93)$$

As proposed initially in [?, ?], and later revisited in [?] for applications in graphics, one can use iteratively Sinkhorn’s algorithm to progressively compute a stationary point of (93). Indeed, successive linearizations of the objective function lead to consider the succession of updates

$$\mathbf{P}^{(\ell+1)} \stackrel{\text{def.}}{=} \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C}^{(\ell)} \rangle - \varepsilon H(\mathbf{P}) \quad \text{where} \quad (94)$$

$$\mathbf{C}^{(\ell)} \stackrel{\text{def.}}{=} \nabla \mathcal{E}_{\mathbf{D}, \mathbf{D}'}(\mathbf{P}^{(\ell)}) = -\mathbf{D}'^T \mathbf{P}^{(\ell)} \mathbf{D},$$

which can be interpreted as a mirror-descent scheme [?]. Each update can thus be solved using Sinkhorn iterations (35) with cost $\mathbf{C}^{(\ell)}$. Figure (??) illustrates the use of this entropic Gromov-Wasserstein to compute soft maps between domains.

11.4 Quantum OT

Static formulation.

Gurvits algorithm, Q-sinkhorn.

References

- [1] Martial Agueh and Guillaume Carlier. Barycenters in the Wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011.
- [2] Syed Mumtaz Ali and Samuel D Silvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society. Series B (Methodological)*, 28(1):131–142, 1966.
- [3] Franz Aurenhammer. Power diagrams: properties, algorithms and applications. *SIAM Journal on Computing*, 16(1):78–96, 1987.
- [4] Federico Bassetti, Antonella Bodini, and Eugenio Regazzini. On minimum kantorovich distance estimators. *Statistics & Probability Letters*, 76(12):1298–1302, 2006.
- [5] Martin Beckmann. A continuous model of transportation. *Econometrica*, 20:643–660, 1952.
- [6] Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative Bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.
- [7] Alain Berlinet and Christine Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer Academic Publishers, 2003.
- [8] Garrett Birkhoff. Extensions of jentzsch’s theorem. *Transactions of the American Mathematical Society*, 85(1):219–227, 1957.
- [9] Emmanuel Boissard. Simple bounds for the convergence of empirical and occupation measures in 1-Wasserstein distance. *Electronic Journal of Probability*, 16:2296–2333, 2011.
- [10] Francois Bolley, Arnaud Guillin, and Cédric Villani. Quantitative concentration inequalities for empirical measures on non-compact spaces. *Probability Theory and Related Fields*, 137(3):541–593, 2007.

- [11] Lev M Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):200–217, 1967.
- [12] Yann Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Communications on Pure and Applied Mathematics*, 44(4):375–417, 1991.
- [13] Donald Bures. An extension of Kakutani’s theorem on infinite product measures to the tensor product of semifinite w^* -algebras. *Transactions of the American Mathematical Society*, 135:199–212, 1969.
- [14] Guillermo Canas and Lorenzo Rosasco. Learning probability measures with respect to optimal transport metrics. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 2492–2500. 2012.
- [15] Guillaume Carlier, Victor Chernozhukov, and Alfred Galichon. Vector quantile regression beyond correct specification. *arXiv preprint arXiv:1610.06833*, 2016.
- [16] Guillaume Carlier and Ivar Ekeland. Matching for teams. *Economic Theory*, 42(2):397–418, 2010.
- [17] Guillaume Carlier, Adam Oberman, and Edouard Oudet. Numerical methods for matching for teams and Wasserstein barycenters. *ESAIM: Mathematical Modelling and Numerical Analysis*, 49(6):1621–1642, 2015.
- [18] Timothy M Chan. Optimal output-sensitive convex hull algorithms in two and three dimensions. *Discrete & Computational Geometry*, 16(4):361–368, 1996.
- [19] Imre Ciszár. Information-type measures of difference of probability distributions and indirect observations. *Studia Scientiarum Mathematicarum Hungarica*, 2:299–318, 1967.
- [20] Marco Cuturi and Arnaud Doucet. Fast computation of Wasserstein barycenters. In *Proceedings of ICML*, volume 32, pages 685–693, 2014.
- [21] Marco Cuturi and Gabriel Peyré. A smoothed dual approach for variational Wasserstein problems. *SIAM Journal on Imaging Sciences*, 9(1):320–343, 2016.
- [22] Steffen Dereich, Michael Scheutzw, and Reik Schottstedt. Constructive quantization: Approximation by empirical measures. In *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, volume 49, pages 1183–1203, 2013.
- [23] Richard M. Dudley. The speed of mean Glivenko-Cantelli convergence. *Annals of Mathematical Statistics*, 40(1):40–50, 1969.
- [24] Peter J Forrester and Mario Kieburg. Relating the Bures measure to the Cauchy two-matrix model. *Communications in Mathematical Physics*, 342(1):151–187, 2016.
- [25] Nicolas Fournier and Arnaud Guillin. On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3-4):707–738, 2015.
- [26] Joel Franklin and Jens Lorenz. On the scaling of multidimensional matrices. *Linear Algebra and its Applications*, 114:717–735, 1989.
- [27] Aude Genevay, Marco Cuturi, Gabriel Peyré, and Francis Bach. Stochastic optimization for large-scale optimal transport. In *Advances in Neural Information Processing Systems*, pages 3440–3448, 2016.
- [28] Joan Glaunes, Alain Trounev, and Laurent Younes. Diffeomorphic matching of distributions: a new approach for unlabelled point-sets and sub-manifolds matching. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, 2004.

- [29] Arthur Gretton, Karsten M Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J Smola. A kernel method for the two-sample-problem. In *Advances in Neural Information Processing Systems*, pages 513–520, 2007.
- [30] Leonid G Hanin. Kantorovich-Rubinstein norm and its application in the theory of Lipschitz spaces. *Proceedings of the American Mathematical Society*, 115(2):345–352, 1992.
- [31] Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. Kernel methods in machine learning. *Annals of Statistics*, 36(3):1171–1220, 2008.
- [32] Leonid Kantorovich. On the transfer of masses (in russian). *Doklady Akademii Nauk*, 37(2):227–229, 1942.
- [33] LV Kantorovich and G.S. Rubinstein. On a space of totally additive functions. *Vestn Leningrad Universitet*, 13:52–59, 1958.
- [34] Jan Lellmann, Dirk A Lorenz, Carola Schönlieb, and Tuomo Valkonen. Imaging with Kantorovich–Rubinstein discrepancy. *SIAM Journal on Imaging Sciences*, 7(4):2833–2859, 2014.
- [35] Arkadi Nemirovski and Uriel Rothblum. On complexity of matrix scaling. *Linear Algebra and its Applications*, 302:435–460, 1999.
- [36] Art B Owen. *Empirical Likelihood*. Wiley Online Library, 2001.
- [37] Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- [38] Antoine Rolet, Marco Cuturi, and Gabriel Peyré. Fast dictionary learning with a smoothed Wasserstein loss. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 630–638, 2016.
- [39] Hans Samelson et al. On the perron-frobenius theorem. *Michigan Mathematical Journal*, 4(1):57–59, 1957.
- [40] Bernhard Schölkopf and Alexander J Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- [41] Justin Solomon, Fernando De Goes, Gabriel Peyré, Marco Cuturi, Adrian Butscher, Andy Nguyen, Tao Du, and Leonidas Guibas. Convolutional Wasserstein distances: efficient optimal transportation on geometric domains. *ACM Transactions on Graphics*, 34(4):66:1–66:11, 2015.
- [42] Bharath K Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, and Gert RG Lanckriet. On integral probability metrics, φ -divergences and binary classification. *arXiv preprint arXiv:0901.2698*, 2009.
- [43] Bharath K Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, and Gert RG Lanckriet. On the empirical estimation of integral probability metrics. *Electronic Journal of Statistics*, 6:1550–1599, 2012.
- [44] Cedric Villani. *Topics in Optimal Transportation*. Graduate Studies in Mathematics Series. American Mathematical Society, 2003.
- [45] Jonathan Weed and Francis Bach. Sharp asymptotic and finite-sample rates of convergence of empirical measures in Wasserstein distance. *arXiv preprint arXiv:1707.00087*, 2017.