OPTIMAL TRANSPORT

Contents

| 1. The problems of Monge and Kantorovich | 2 |
|--|----|
| 1.1. Monge's problem | 2 |
| 1.2. Kantorovich's problem | 3 |
| 2. One-dimensional optimal transport | 5 |
| 2.1. Quantile function and one-dimensional Wasserstein spaces | 6 |
| 3. Kantorovich duality | 9 |
| 3.1. Derivation of the dual problem | 9 |
| 3.2. Strong duality | 10 |
| 3.3. Existence of solution for the dual problem | 12 |
| 3.4. Stability of optimal transport plans | 14 |
| 4. Kantorovich's functional | 15 |
| 4.1. Kantorovich's functional | 15 |
| 4.2. Solution of Monge's problem | 16 |
| 4.3. Semi-discrete optimal transport | 18 |
| 4.4. Oliker–Prussner's algorithm | 20 |
| 5. Entropy-regularized optimal transport | 21 |
| 5.1. Primal problem | 21 |
| 5.2. Dual problem | 23 |
| 5.3. Existence of a solution to the dual | 25 |
| 5.4. Sinkhorn algorithm as block-coordinate ascent | 27 |
| 5.5. Linear convergence of Sinkhorn's algorithm | 29 |
| 6. Wasserstein distances | 31 |
| 6.1. <i>p</i> -Wasserstein spaces over compact metric spaces | 31 |
| 6.2. <i>p</i> -Wasserstein geodesics on \mathbb{R}^d | 33 |
| 6.3. Geodesic convexity with respect to W_2 on \mathbb{R}^d | 35 |
| 7. Quantization and uniform quantization of measures | 38 |
| 8. Embedding of the Wasserstein space | 38 |
| 8.1. Non-embeddability results | 39 |
| 8.2. Embedding via slicing | 41 |
| 8.3. "Linearization" of the quadratic Wasserstein distance | 42 |
| 9. Stability of quadratic optimal transport maps | 44 |
| 9.1. Stability near a regular configuration | 45 |
| 9.2. Stability of potentials for entropy-regularized quadratic optimal | |
| transport | 46 |
| 10. Hölder stability of dual potentials | 49 |
| References | 51 |

OPTIMAL TRANSPORT

Why study optimal transport ? The main motivation studying optimal transport in statistics is the notion of Wasserstein distance between probability measures on a compact metric space X:

- The Wasserstein distances W_p represent faithfully the geometry of the underlying space: $x \in X \mapsto \delta_x \in \mathcal{P}(X)$ is an isometry. This means that unlike many notions of distances between functions/divergences between probability measures (E.g relative entropy),
- Application: inverse problems, Wasserstein GANs
- Application: statistics over the space of probability measures, e.g. geodesics barycenters, k-means, PCA...
- Application: PDE / particle systems

References. Introduction to optimal transport, with applications to PDE and/or calculus of variations can be found in books by Villani [42] and Santambrogio [34]. Villani's second book [43] concentrates on the application of optimal transport to geometric questions (e.g. synthetic definition of Ricci curvature), but its first chapters might be useful. We also mention Gigli, Ambrosio and Savaré [3] for the study of gradient flows with respect to the Monge-Kantorovich/Wasserstein metric.

Notation. In the following, we assume that X is a compact metric space, and we denote $\mathcal{C}^0(X)$ the space of continuous functions over X endowed with the norm of uniform convergence $\|\varphi\|_{\infty} = \sup_{x \in X} |\varphi(x)|$. We denote $\mathcal{M}(X)$ the space of Radon measures on X, which we identify with the continuous dual of $\mathcal{C}^0(X)$. We will denote $\langle \mu | \varphi \rangle = \int \varphi d\mu$. We define

$$\mathcal{M}^{+}(X) := \{ \mu \in \mathcal{M}(X) \mid \forall \varphi \in \mathcal{C}^{0}(X), \varphi \ge 0 \Longrightarrow \langle \mu | \varphi \rangle \ge 0 \}$$
$$\mathcal{P}(X) := \{ \mu \in \mathcal{M}^{+}(X) \mid \langle \mu | 1 \rangle = 1 \}$$

The support of a measure μ is denoted $spt(\mu)$.

The dual space is endowed with the total variation norm

$$\|\mu_n\|_{\mathrm{TV}} = \sup_{\varphi \in \mathcal{C}^0(X), \|\varphi\|_{\infty} \leqslant 1} \langle \mu | \varphi \rangle.$$

However, the topology that we will consider by default on $\mathcal{M}^0(X)$ is the weak* topology. We recall for instance that a sequence $(\mu_n)_{n\geq 0}$ of measures converges weak* to μ if and only if

$$\forall \varphi, \lim_{n \to +\infty} \langle \mu_n | \varphi \rangle = \langle \mu | \varphi \rangle.$$

We note that thanks to the Banach-Alaoglu theorem, any bounded sequence $(\mu_n)_{n\in\mathbb{N}}$ in $\mathcal{M}(X)$ admits a weak* converging subsequence. This applies in particular to any sequence in $\mathcal{P}(X)$: the space of probability measures is weak* sequentially compact (and even compact).

1. The problems of Monge and Kantorovich

1.1. Monge's problem.

Definition 1 (Push-forward and transport map). Let X, Y be compact metric spaces, $\mu \in \mathcal{M}(X)$ and let $T : X \to Y$ be a measurable map. The *push-forward* of μ by T is the measure $T_{\#\mu}$ on Y defined by

$$\forall B \subseteq Y, \ T_{\#}\mu(B) = \mu(T^{-1}(B)).$$

 $\mathbf{2}$

or equivalently if the following change-of-variable formula holds for all test function $\varphi \in \mathcal{C}^0(Y)$:

$$\int_{Y} \varphi(y) \mathrm{d}\nu(y) = \int_{X} \varphi(T(x)) \mathrm{d}\mu(x).$$

A measurable map $T: X \to Y$ such that $T_{\#}\mu = \nu$ is also called a *transport* map between μ and ν .

Example 1. If $Y = \{y_1, \dots, y_n\}$, then $T_{\#}\mu = \sum_{1 \le i \le n} \mu(T^{-1}(\{y_i\}))\delta_{y_i}$.

Definition 2 (Monge's problem). Consider two metric spaces X, Y, two probability measures $\mu \in \mathcal{P}(X), \nu \in \mathcal{P}(Y)$ and a *cost function* $c: X \times Y \to \mathbb{R} \cup \{+\infty\}$. Monge's problem is the following optimization problem

(MP) := inf
$$\left\{ \int_X c(x, T(x)) d\mu(x) \mid T : X \to Y \text{ and } T_\# \mu = \nu \right\}$$
 (1.1)

This problem exhibits several difficulties, one of which is that both the constraint $(T_{\#}\mu = \nu)$ and the functional are non-convex.

Example 2. There might exist no transport map between μ and ν . For instance, consider $\mu = \delta_x$ for some $x \in X$. Then, $T_{\#}\mu(B) = \mu(T^{-1}(B)) = \delta_{T(x)}$. In particular, if ν is not a Dirac mass, then there exists no transport map between μ and ν .

1.2. Kantorovich's problem.

Definition 3 (Transport plan). Let X, Y be two metric spaces and $\mu \in \mathcal{M}^+(X)$ and $\nu \in \mathcal{M}^+(Y)$ be two non-negative measures. A transport plan between μ and ν is a non-negative measure γ on the product space $X \times Y$ whose marginals are μ and ν . The set of transport plans is denoted

$$\Gamma(\mu,\nu) = \left\{ \gamma \in \mathcal{M}_+(X \times Y) \mid \Pi_{X\#}\gamma = \mu, \ \Pi_{Y\#}\gamma = \nu \right\},\,$$

where $\Pi_X : X \times Y \to X$ and $\Pi_Y : X \times Y \to Y$ are the projection maps. Note that $\Gamma(\mu, \nu)$ is a convex set, and that it is non-empty if and only if μ and ν have the same total mass, i.e. $\mu(X) = \nu(Y)$.

Definition 4 (Kantorovich's problem). Given two metric spaces X, Y, two non-negative measures $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$ and a continuous cost function $c \in \mathcal{C}^0(X \times Y)$, Kantorovich's problem is the following optimization problem

$$(KP) := \inf \left\{ \langle c | \gamma \rangle \mid \gamma \in \Gamma(\mu, \nu) \right\}$$
(1.2)

We will denote \mathcal{T}_c the associated *transport cost*

$$\mathcal{T}_{c}: \mathcal{M}^{+}(X) \times \mathcal{M}^{+}(Y) \to \mathbb{R} \cup \{+\infty\}$$

$$(\mu, \nu) \mapsto \inf\{\langle c | \gamma \rangle \mid \gamma \in \Gamma(\mu, \nu)\}.$$
(1.3)

Note that by convention, the infimum over the empty set is $+\infty$, so that $\mathcal{T}_c(\mu,\nu) = +\infty$ if $\mu(X) \neq \nu(Y)$.

Remark 1. The infimum in Kantorovich's problem is less than the infimum in Monge's problem. Indeed, consider a transport map satisfying $T_{\#}\mu = \nu$ and the associated transport plan $\gamma_T = (\mathrm{id}, T)_{\#}\mu$. Then, by the change-ofvariable formula one has

$$\langle c|\gamma_T \rangle \leqslant \int_{X \times Y} c(x, y) \mathrm{d}(id, T)_{\#} \mu(x, y) = \int_X c(x, T(x)) \mathrm{d}\mu,$$

thus proving the claim.

Example 3 (Finite support). Assume that $X = Y = \{1, ..., N\}$ and that μ, ν are the uniform probability measures over X and Y. Then, Monge's problem can be rewritten as a minimization problem over the set of bijections between the two sets X and Y:

$$\min\left\{\frac{1}{N}\sum_{1\leqslant i\leqslant N}c(i,\sigma(i))\mid \sigma\in\mathfrak{S}_N\right\}.$$

In Kantorovich's relaxation, the set of transport plans $\Gamma(\mu, \nu)$ agrees with the set of bi-stochastic matrices :

$$\gamma \in \Gamma(\mu,\nu) \Longleftrightarrow \gamma \geqslant 0, \sum_i \gamma(i,j) = 1/N = \sum_j \gamma(i,j)$$

By Birkhoff's theorem, any extremal bi-stochastic matrix is induced by a permutation. This shows that, in this case, the solution to Monge's and Kantorovich's problems agree.

Theorem 1 (Existence of solutions to (KP)). Let X, Y be compact metric spaces and let $c \in C^0(X \times Y)$. Then for any measures $(\mu, \nu) \in \mathcal{M}_+(X) \times \mathcal{M}_+(Y)$ with equal total mass, Kantorovich's problem (KP) admits a minimizer. Moreover, the transport cost \mathcal{T}_c is a convex and weak* lower semicontinuous functional on $\mathcal{M}_+(X) \times \mathcal{M}_+(Y)$.

Proposition 2. Let X, Y be compact metric spaces and let $(\mu_n)_{n\in\mathbb{N}}$ and $(\nu_n)_{n\in\mathbb{N}}$ be sequences of non-negative measures on X and Y with same total mass. Assume that these sequence weak* converge to $\mu \in \mathcal{M}^+(X)$ and $\nu \in \mathcal{M}^+(Y)$ respectively. Then, any sequence of transport plans $\gamma_n \in \Gamma(\mu_n, \nu_n)$ admits a subsequence converging to some $\gamma \in \Gamma(\mu, \nu)$.

In particular, the previous proposition implies that $\Gamma(\mu, \nu)$ is compact.

Proof. Since $\mu_n \ge 0$, one has $\|\mu_n\|_{\text{TV}} = \langle \mu_n | 1 \rangle$, which converges to $\|\mu\|_{\text{TV}}$ by weak* convergence. Thus the sequence (μ_n) is bounded. Since

$$\|\gamma_n\|_{\mathrm{TV}} = \langle \gamma_n | 1 \rangle = \langle \Pi_{\#\gamma_n} | 1 \rangle = \langle \mu_n | 1 \rangle,$$

the sequence $(\gamma_n)_{n \in \mathbb{N}}$ is also bounded. By Banach-Alaoglu's theorem, it admits a weak* converging subsequence. Relabeling if necessary, we therefore assume that γ_n converges weak* to some $\gamma \in \mathcal{M}(X \times Y)$. Then,

$$\forall \varphi \in \mathcal{C}^0(X \times Y) \text{ s.t. } \varphi \ge 0, \langle \gamma | \varphi \rangle = \lim_{n \to +\infty} \langle \gamma_n | \varphi \rangle \ge 0$$

so that γ is a non-negative measures. Given $\varphi \in \mathcal{C}^0(X)$ and $\hat{\varphi}(x, y) := \varphi(x)$, using $\prod_{X \#} \gamma_n = \mu_n$ we get $\langle \varphi | \mu_n \rangle = \langle \varphi | \Pi_{X \#} \gamma_n \rangle = \langle \hat{\varphi} | \gamma_n \rangle$. Taking the limit as $n \to +\infty$, we deduce that $\langle \varphi | \mu \rangle = \langle \hat{\varphi} | \gamma \rangle$ for all φ , implying that $\prod_{X \#} \gamma = \mu$. Similarly, we prove that $\prod_{Y \#} \gamma = \nu$, proving that $\gamma \in \Gamma(\mu, \nu)$. \Box

Proof of theorem 1. We first note that the function $\gamma \mapsto \langle c | \gamma \rangle$ is linear and continuous on $\mathcal{M}(X \times Y)$. Second, we note that if $\mu(X) = \nu(Y)$, the set $\Gamma(\mu, \nu)$ is non-empty as it contains a suitably rescaled product of μ and ν . The previous lemma shows that the set $\Gamma(\mu, \nu)$ is weak* compact, so that

 $\langle c|\gamma\rangle$ attains its minimum on this set. This shows existence of at least one solution to (KP).

To prove that \mathcal{T}_c is lower semicontinous, we consider converging sequences $(\mu_n), (\nu_n)$ in $\mathcal{M}_+(X)$ and $\mathcal{M}_+(Y)$ respectively. with weak* limits μ and ν . Without loss of generality, we assume that μ_n and ν_n have the same total mass (if not, $\mathcal{T}_c(\mu_n, \nu_n) = +\infty$). For each n we consider $\gamma_n \in \Gamma(\mu_n, \nu_n)$ the optimal transport plan. Using the previous proposition, we assume taking a subsequence if necessary that γ_n converges to some $\gamma \in \Gamma(\mu, \nu)$. Then,

$$\mathcal{T}_{c}(\mu,\nu) \leqslant \langle c|\gamma\rangle = \lim_{n \to +\infty} \langle c|\gamma_{n}\rangle = \lim_{n \to +\infty} \mathcal{T}_{c}(\mu_{n},\nu_{n}).$$

2. One-dimensional optimal transport

Definition 5 (Monotone set). A subset S of $\mathbb{R} \times \mathbb{R}$ is called *monotone* if

$$\forall (x,y), (x',y') \in S, (x'-x) \cdot (y'-y) \ge 0.$$

Definition 6 (Submodular cost). A cost function $c : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ is called *strictly submodular* if for every $x_0 < x_1$, the function $y \mapsto c(x_1, y) - c(x_0, y)$ is decreasing.

Theorem 3. Let μ, ν be probability measures supported in $X = Y = [a, b] \subseteq \mathbb{R}$, and let c be a continuous and strictly submodular cost on $X \times Y$. Then, there exists a unique optimal transport plan $\gamma \in \Gamma(\mu, \nu)$, which is also the unique transport plan with monotone support.

Proof. Step 1. We first establish that any optimal transport plan between μ and ν must be monotone. Consider a transport plan $\gamma \in \Gamma(\mu, \nu)$ and consider (x_0, y_0) and (x_1, y_1) in $\operatorname{spt}(\gamma)$. Since we want to prove that $(x_0 - x_1)(y_0 - y_1) \leq 0$, we may assume that $x_1 \neq x_0$ and $y_1 \neq y_0$. By continuity of the cost, for any $\delta > 0$ there exists r > 0 such that:

$$\mathcal{B}((x_0, y_0), r) \cap \mathcal{B}((x_1, y_1), r) \neq \emptyset$$

 $\forall a, b \in \{x_1, x_0, y_1, y_0\}, \forall (x, y) \in \mathcal{B}((a, b), r), \ |c(x, y) - c(a, b)| \leq \delta$

Since (x_0, y_0) and (x_1, y_1) both belong to the support of γ , there must exist non-negative measures $\gamma_0 \leq \gamma$ and $\gamma_1 \leq \gamma$ with equal positive mass ε and such that $\operatorname{spt}(\gamma_i) \subseteq \operatorname{B}((x_i, y_i), r)$. Consider the marginals $\mu_i = \pi_{X\#}\gamma_i$ and $\nu_i = \pi_{Y\#}\gamma_i$, and take any coupling σ_0 (resp. σ_1) between μ_0 and ν_1 (resp. μ_1 and ν_0). Then, one can check that the measure

$$\sigma = \gamma - \gamma_0 - \gamma_1 + \sigma_0 + \sigma_1$$

is a transport plan between μ and ν (the non-negativity comes from $\gamma_i \leq \gamma$ and $\operatorname{spt}(\gamma_1) \cap \operatorname{spt}(\gamma_0) = \emptyset$). Using the optimality of γ one gets

$$0 \leq F(\sigma) - F(\gamma) = F(\sigma_0) - F(\gamma_0) + F(\sigma_1) - F(\gamma_1)$$

= $\int_{B(x_0, r) \times B(y_1, r)} cd\sigma_0 + \int_{B(x_1, r) \times B(y_0, r)} cd\sigma_1$
- $\int_{B(x_0, r) \times B(y_0, r)} cd\gamma_0 - \int_{B(x_1, r) \times B(y_1, r)} cd\gamma_1$
 $\leq \varepsilon \cdot (c(x_0, y_1) + c(x_1, y_0) - c(x_0, y_0) - c(x_1, y_1) + 4\delta)$

Since this holds for all $\delta > 0$ small enough, we deduce that

$$c(x_0, y_0) + c(x_1, y_1) \leq c(x_0, y_1) + c(x_1, y_0).$$

Assume without loss of generality that $x_0 < x_1$. Then,

 $c(x_1, y_1) - c(x_0, y_1) \leq c(x_1, y_0) - c(x_1, y_0),$

thus implying by submodularity (the function $y \mapsto c(x_1, y_1) - c(x_0, y_1)$ is decreasing) that $y_0 \leq y_1$.

Step 2. We show that there exists at most one monotone transport plan between μ and ν . Recall that a probability measure γ on \mathbb{R}^2 is uniquely defined from the values $\gamma((-\infty, a] \times (-\infty, b])$ for any $a, b \in \mathbb{R}$. This follows from the fact that such sets generate the Borel σ -algebra. Consider A = $(-\infty, a] \times (b, +\infty)$ and $B = (a, +\infty) \times (-\infty, b]$. Then, by monotonicity of $\operatorname{spt}(\gamma)$ one cannot have $\gamma(A) > 0$ and $\gamma(B) > 0$ at the same time. Hence,

$$\gamma((-\infty, a] \times (-\infty, b]) = \min(\gamma(((-\infty, a] \times (-\infty, b]) \cup A),$$
$$\gamma(((-\infty, a] \times (-\infty, b]) \cup B))$$
$$= \min(\mu((-\infty, a]), \nu((-\infty, b])).$$

This shows that $\gamma((-\infty, a] \times ((-\infty, b]))$ is uniquely defined from μ, ν , so that γ is unique.

2.1. Quantile function and one-dimensional Wasserstein spaces.

Definition 7 (Cdf and quantile function). Let μ be a probability measure on \mathbb{R} . The cumulative distribution function $F_{\mu} : \mathbb{R} \to [0, 1]$ and the inverse cumulative distribution function $T_{\mu} : [0, 1] \to \mathbb{R}$ are defined by:

 $F_{\mu}(x) = \mu((-\infty, x]) \qquad T_{\mu}(m) = \inf \left\{ x \in \mathbb{R} \mid F_{\mu}(x) \ge m \right\}.$

The function T_{μ} will also be called the *quantile function*.

In the following, we assume that X is a segment of \mathbb{R} .

Definition 8 (Wasserstein distance). The Wasserstein distance of exponent $p \ge 1$ between two probability measures $\mu, \nu \in \mathcal{P}(X)$ is defined by

$$W_p^p(\mu,\nu) = \min_{\gamma \in \Gamma(\mu,\nu)} \int \|x-y\|^p \,\mathrm{d}\gamma(x,y).$$

Proposition 4 (Quantile functions and Wasserstein distance). Let μ, ν be two probability measure on a segment $X \subseteq \mathbb{R}$. Then,

- (i) T_{μ} is a transport map between the Lebesgue measure $\lambda_{[0,1]}$ and μ
- (ii) $\gamma_{\mu \to \nu} = (T_{\mu}, T_{\nu})_{\#} \lambda_{[0,1]}$ is the unique monotone transport plan between μ and ν ;
- (*iii*) for all $p \ge 1$, $W_p(\mu, \nu) = \|T_\mu T_\nu\|_{L^p([0,1])}$.

Example 4 (Translation). If ν is obtained by translating ν by a constant $v \in \mathbb{R}$, then $T_{\nu} = T_{\mu} + v$ so that $W_p(\mu, \nu) = ||T_{\mu} - T_{\nu}||_{L^p([0,1])} = |v|$.

Example 5 (Discrete measures). If $\mu = \frac{1}{N} \sum_{i=1}^{N} \delta_{x_i}$ and the sequence $(x_i)_{1 \leq i \leq N}$ is increasing, then the quantile function satisfies

$$T_{\mu}\big|_{\left[\frac{i-1}{n},\frac{i}{n}\right]} = x_i$$

In particular, if $\nu = \frac{1}{N} \sum_{i=1}^{N} \delta_{y_i}$, where the sequence $y_{1 \leq i \leq N}$ is also increasing,

$$W_p(\mu, \nu)^p = \frac{1}{N} \sum_i ||x_i - y_i||^p$$

Proof. (i) Let $\hat{\mu} = T_{\mu \#} \lambda_{[0,1]}$. Then,

$$\begin{aligned} F_{\hat{\mu}}(x) &= \hat{\mu}((-\infty, x]) \\ &= \lambda(T_{\mu}^{-1}(-\infty, x]) \\ &= \lambda(\{m \in [0, 1], T_{\mu}(m) \leq x\} \\ &= \lambda(\{m \in [0, 1], F_{\mu}(x) \geq m\}) \\ &= F_{\mu}(x). \end{aligned}$$

were we used the equivalence $T_{\mu}(m) \leq x$ iff $F_{\mu}(x) \geq m$. This shows that $\hat{\mu} = \mu$.

(ii) Denote $\gamma := \gamma_{\mu \to \nu}$. We note first that $\Pi_{X\#}\gamma = \Pi_X \circ (T_\mu, T_\nu)_{\#}\lambda_{[0,1]} = \mu$, and similarly $\Pi_{Y\#}\gamma = \nu$. Thus, γ is a transport plan between μ and ν . In addition, γ is supported on the set $S := \{(T_\mu(m), T_\nu(m)) \mid m \in [0, 1]\}$. Given two couples $(x_i, y_i) \in S$, there exists $m_i \in [0, 1]$ such that $x_i = T_\mu(m_i)$ and $y_i = T_\nu(m_i)$. Without loss of generality, assume that $m_0 \leq m_1$. Then, $T_\mu(m_0) \leq T_\mu(m_1)$ and $T_\nu(m_0) \leq T_\nu(m_1)$ so that

$$(x_1 - x_0)(y_1 - y_0) \ge 0,$$

implying that S is monotone.

(iii) Theorem 3 proves that a solution to the optimal transport problem is given between μ and ν for the convex cost $c(x, y) = ||x - y||^p$ is given by the monotone plan, i.e.

$$\begin{split} \min_{\gamma \in \Gamma(\mu,\nu)} \int \|x - y\|^p \, \mathrm{d}\gamma(x,y) &= \int \|x - y\|^p \, \mathrm{d}\gamma_{\mu \to \nu}(x,y) \\ &= \int \|x - y\|^p \, \mathrm{d}(T_\mu, T_\nu)_{\#\lambda_{[0,1]}}(x,y) \\ &= \int_0^1 \|T_\mu(m) - T_\nu(m)\|^p \, \mathrm{d}m \\ &= \|T_\mu - T_\nu\|_{\mathrm{L}^p([0,1])}^p \end{split} \quad \Box$$

Proposition 5 (Properties of the 1D Wasserstein spaces). The following properties hold for any segment $X \subseteq \mathbb{R}$ and any $p \ge 1$:

- (i) W_p is a distance on $\mathcal{P}(X)$
- (ii) W_p metrizes weak* convergence on $\mathcal{P}(X)$, i.e. for any sequence (μ_n) in $\mathcal{P}(X)$ and any $\mu \in \mathcal{P}(X)$,

$$\lim_{n \to +\infty} W_p(\mu_n, \mu) = 0 \iff \forall \varphi \in \mathcal{C}^0(X), \lim_{n \to +\infty} \langle \mu_n | \varphi \rangle = \langle \mu | \varphi \rangle.$$

(iii) the application $\mu \mapsto T_{\mu}$ mapping a probability measure to its inverse cdf is an isometric embedding of $(\mathcal{P}(X), W_p(X))$ into $L^p([0, 1])$.

Proof. (i) We note that $W_p(\mu, \nu) = 0$ implies that $T_{\mu} = T_{\nu}$ a.e., so that $\mu = T_{\mu \#} \lambda_{[0,1]} = T_{\nu \#} \lambda_{[0,1]} = nu$. The symmetry is immediate, and the triangle inequality for W_p follows from the triangle inequality in $L^p([0,1])$.

(ii) Assume first that $W_p(\mu_n, \mu) = \|T_{\mu_n} - T_{\mu}\|_{L^p([0,1])}$ converges to zero as $n \to +\infty$. Then, $\|T_{\mu_n} - T_{\mu}\|_{L^1([0,1])}$ also converges to zero as $n \to +\infty$. Let $f: X \to \mathbb{R}$ be *L*-Lipschitz. Then,

$$\begin{split} |\langle f|\mu_n - \mu \rangle| &= \left| \int_0^1 f(T_{\mu_n}(m)) - f(T_{\mu}(m)) \mathrm{d}m \right| \\ &\leqslant L \int_0^1 \|T_{\mu_n}(m) - T_{\mu}(m)\| \mathrm{d}m \\ &= L \operatorname{W}_1(\mu_n, \mu) \xrightarrow{n \to +\infty} 0 \end{split}$$

Since continuous functions on X can be uniformly approximated by Lipschitz functions, we get weak* convergence.

Conversely, assume that μ_n converges weakly to μ . The non-decreasing map T_{μ} is continuous on $[0,1] \setminus Z$, where Z is at most countable. It is standard that for any $x \notin Z$, $T_{\mu_n}(x)$ converges to $T_{\mu}(x)$ as $n \to +\infty$, i.e. T_{μ_n} converges a.e. to T_{μ} . Since in addition T_{μ_n} is bounded, we deduce that convergence holds in $L^p([0,1])$ for any $p \ge 1$.

Definition 9 (Geodesic). Let (E, d) be a metric space. A constant speed geodesic between two points $x_0, x_1 \in E$ is a continuous curve $x : [0, 1] \to E$ such that for every $s, t \in [0, 1], d(x_s, x_t) = |s - t| d(x_0, x_1)$.

Proposition 6. Let X be a segment of \mathbb{R} and let $\mu_0, \mu_1 \in \mathcal{P}(X)$. Define

$$\mu_t := T_{t\#} \lambda_{[0,1]}, \text{ where } T_t = (1-t)T_{\mu_0} + tT_{\mu_1}$$

Then, the curve μ_t is a constant speed geodesic between μ_0 and μ_1 in the space $(\mathcal{P}(X), W_p)$, for any exponent $p \ge 1$. In particular, this space is a geodesic space, meaning that any $\mu_0, \mu_1 \in \mathcal{P}_p(X)$ can be joined by (at least one) constant speed geodesic.

Proof. First note that if $0 \leq s \leq t \leq 1$,

$$W_p(\mu_0, \mu_1) \leq W_p(\mu_0, \mu_s) + W_p(\mu_s, \mu_t) + W_p(\mu_t, \mu_1),$$

so that it suffices to prove the inequality $W_p(\mu_s, \mu_t) \leq |t-s| W_p(\mu_0, \mu_1)$ for all $0 \leq s \leq t \leq 1$ to get equality. The inequality is easily checked by taking $\gamma_{st} := (T_s, T_t)_{\#} \lambda_{[0,1]} \in \Gamma(\mu_s, \mu_t)$, so that

$$W_{p}(\mu_{s},\mu_{t})^{p} \leq \int ||T_{s}(m) - T_{t}(m)||^{p} dm$$

= $\int ||(1-s)T_{0}(m) + sT_{1}(m) - ((1-t)T_{0}(m) + tT_{1}(m))||^{p} dm$
= $\int ||(t-s)(T_{0}(m) - T_{t}(m))||^{p} dm = (t-s)^{p} W_{p}(\mu,\nu)^{p}$

Remark 2 (Barycenters). We can also consider barycenters in the Wasserstein, at least in the case p = 2 and on a segment X. The weighted barycenter of probability measures $\mu_0, \ldots, \mu_k \in \mathcal{P}(X)$ with weights $\alpha_1, \ldots, \alpha_k > 0$ is the unique minimizer of

$$\min_{\mu \in \mathcal{P}(X)} \sum_{1 \leqslant i \leqslant k} \alpha_k \operatorname{W}_2^2(\mu_k, \mu)$$

The quantile function of the barycenter μ therefore solves the following minimization problem

$$T_{\mu} \in \arg\min_{T} \sum_{1 \leq i \leq k} \alpha_k \|T_{\mu_k} - T\|^2_{\mathrm{L}^2([0,1])},$$

so that T_{μ} is simply a weighted average of the T_{μ_k} :

$$T_{\mu} = \frac{1}{\sum_{k} \alpha_{k}} \sum_{1 \leqslant i \leqslant k} \alpha_{k} T_{\mu_{k}}$$

The barycenter is finally recovered thanks to the formula $\mu = T_{\mu \#} \lambda_{[0,1]}$, i.e.

$$\mu = \left(\frac{1}{\sum_k \alpha_k} \sum_{1 \leq i \leq k} \alpha_k T_{\mu_k}\right)_{\#} \lambda_{[0,1]}.$$

3. KANTOROVICH DUALITY

3.1. Derivation of the dual problem. The primal Kantorovich problem (KP) can be reformulated by introducing Lagrange multipliers for the constraints. Namely, we use that for any $\gamma \in \mathcal{M}^+(X \times Y)$,

$$\sup_{\varphi \in \mathcal{C}^{0}(X)} -\langle \varphi \otimes 1 | \gamma \rangle + \langle \varphi | \mu \rangle = \begin{cases} 0 & \text{if } \Pi_{X \#} \gamma = \mu \\ +\infty & \text{if not} \end{cases}$$
$$\sup_{\varphi \in \mathcal{C}^{0}(X)} -\langle 1 \otimes \psi | \gamma \rangle + \langle \psi | \mu \rangle = \begin{cases} 0 & \text{if } \Pi_{X \#} \gamma = \mu \\ +\infty & \text{if not} \end{cases}$$

to deduce that for any $\gamma \in \mathcal{M}_+(X \times Y)$,

 φ

$$\sup_{\in \mathcal{C}^{0}(X), \psi \in \mathcal{C}^{0}(Y)} \langle \varphi | \mu \rangle + \langle \psi | \nu \rangle - \langle \varphi \oplus \psi | \gamma \rangle = \begin{cases} 0 & \text{if } \gamma \in \Gamma(\mu, \nu) \\ +\infty & \text{if not.} \end{cases}$$

This leads to the following formulation of the Kantorovich problem

$$(\mathrm{KP}) = \inf_{\gamma \in \mathcal{M}^+(X \times Y)} \sup_{(\varphi, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y)} \langle c - (\varphi \oplus \psi) | \gamma \rangle + \langle \varphi | \mu \rangle - \langle \psi | \nu \rangle$$

Kantorovich dual problem is simply obtained by inverting the infimum and the supremum:

(KD) :=
$$\sup_{\varphi,\psi} \inf_{\gamma \ge 0} \langle c - (\varphi \oplus \psi) | \gamma \rangle + \langle \varphi | \mu \rangle - \langle \psi | \nu \rangle.$$

Note that we will often omit the assumptions that $\gamma \in \mathcal{M}(X \times Y)$ and φ, ψ are continuous, when the context is clear. The dual problem can further be simplified by remarking that

$$\inf_{\gamma \geqslant 0} \langle c - \varphi \oplus \psi | \gamma \rangle = \begin{cases} 0 & \text{if } \varphi \oplus \psi \leqslant c \\ -\infty & \text{if not.} \end{cases}$$

Definition 10 (Kantorovich's dual problem). Given $\mu \in \mathcal{P}(X)$ and $\nu \in \mathcal{P}(Y)$ with X, Y compact metric spaces and $c \in \mathcal{C}^0(X \times Y)$, we define Kantorovich's dual problem by

$$(\mathrm{KD}) = \sup\left\{\int_X \varphi \mathrm{d}\mu + \int_Y \psi \mathrm{d}\nu \mid (\varphi, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y), \varphi \oplus \psi \leqslant c\right\}$$
(3.4)

Proposition 7. Weak duality holds, i.e. $(KP) \ge (KD)$.

Proof. Given $(\varphi, \psi, \gamma) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y) \times \Gamma(\mu, \nu)$ satisfying the constraint $\varphi \oplus \psi \leq c$, one has

$$\langle \varphi | \mu \rangle + \langle \psi | \nu \rangle = \langle \varphi \oplus \psi | \gamma \rangle \leqslant \langle c | \gamma \rangle,$$

where we used $\gamma \in \Gamma(\mu, \nu)$ to get the equality and $\varphi \oplus \psi \leq c$ to get the inequality. As a conclusion,

$$(\mathrm{KD}) = \sup_{\varphi \oplus \psi \leqslant c} \langle \varphi | \mu \rangle - \langle \psi | \nu \rangle \leqslant \min_{\gamma \in \Gamma(\mu, \nu)} \langle c | \gamma \rangle = (\mathrm{KP}) \qquad \Box$$

Remark 3. As often, the Lagrange multipliers (or Kantorovich potentials) φ, ψ have an economic interpretation as prices. For instance, imagine that μ is the distribution of sand available at quarries, and ν describes the amount of sand required by construction work. Then, (KP) can be interpreted as finding the cheapest way of transporting the sand from μ to ν for a construction company. Imagine that this company wants to externalize the transport, by paying a loading coast $\varphi(x)$ at a point x (in a quarry) and an unloading coast $\psi(y)$ at a point y (at a construction place). Then, the constraint $\varphi(x)+\psi(y) \leq c(x,y)$ translates the fact that the construction company would not externalize if its cost is higher than the cost of transporting the sand by itself. Then, Kantorovich's dual problem (KD) describes the problem of a transporting company: maximizing its revenue $\int \varphi d\mu + \int \psi d\nu$ under the constraint $\varphi \oplus \psi \leq c$ imposed by the construction company. The economic interpretation of the strong duality (KP) = (KD) is that in this setting, externalization has exactly the same cost as doing the transport by oneself.

The questions that we will address now are the following:

- When does strong duality ((KP) = (KD)) hold ?
- When is the supremum in Kantorovich's dual problem attained ?
- What does Kantorovich's duality imply about Monge's problem, stability of optimal transport maps/plans, numerics, etc ?

3.2. Strong duality. We prove strong duality using a strategy recently proposed by Savaré and Sodini [35], which relies only the Fenchel-Moreau theorem from convex analysis. In addition to the transport cost functional,

$$\mathcal{T}_{c}: \mathcal{M}(X) \times \mathcal{M}(Y) \to \mathbb{R} \cup \{+\infty\}$$

$$(\mu, \nu) \mapsto \begin{cases} \inf\{\langle c|\gamma\rangle \mid \gamma \in \Gamma(\mu, \nu)\} & \text{if } \mu \ge 0, \nu \ge 0, \text{ and } \mu(X) = \nu(Y) \\ +\infty & \text{otherwise} \end{cases}$$

$$(3.5)$$

we will consider the following, non-convex and very singular functional, which encodes the cost of transport between Dirac masses with the same weight:

$$F_c: \mathcal{M}(X) \times \mathcal{M}(Y) \to \mathbb{R} \cup \{+\infty\}$$

$$(\mu, \nu) \mapsto \begin{cases} mc(x, y) & \text{if } \mu = m\delta_x, \nu = m\delta_y \text{ and } m \ge 0 \\ +\infty & \text{otherwise} \end{cases}$$
(3.6)

Theorem 8 (Savaré and Sodini). $\mathcal{T}_c = F_c^{**}$

Corollary 9 (Strong duality in Kantorovich's problem). (KP) = (KD).

The proof of these results rely on the Fenchel-Moreau theorem from convex analysis. To state this theorem, we need to define the convex and convex biconjugate of a function on a topological vector space.

Definition 11 (Convex conjugate). Let E be a topological vector space. The *convex conjugate* of a function $F: E \to \mathbb{R} \cup \{+\infty\}$ is the function F^* on the dual space E^* defined by

$$F^*(x^*) = \sup_{x \in E} \langle x^* | x \rangle - F(x).$$

The *biconjugate* of F is then defined as $F^{**}: E \to \mathbb{R} \cup \{+\infty\}$ by

$$F^{**}(x) = \sup_{x^* \in E^*} \langle x^* | x \rangle - F^*(x^*)$$

It is quite easy to see that F^* and F^{**} are convex and lower semicontinuous, as suprema of continuous affine functions. Fenchel-Moreau's theorem show that F^{**} is in fact the *lower semicontinuous convex envelope* of F, i.e. the largest lsc convex function that lies below F.

Theorem 10 (Fenchel-Moreau). Let E be a locally convex and separated topological vector space and let $F : E \to \mathbb{R} \cup \{+\infty\}$. Then F^{**} is the lsc convex envelope of F, i.e. the largest lsc convex function that lies below F. In particular, $F = F^{**}$ if and only if F is convex and lower semicontinuous.

Proof. Let G be the lsc convex envelope of F. We first prove that $F^{**} \leq G$. Given any point $x \in E$, the definition of F^* as a supremum gives $F^*(x^*) \geq \langle x^* | x \rangle - F(x)$. Thus,

$$F^{**}(x) = \sup_{x^* \in E^*} \langle x^* | x \rangle - F^*(x^*) \leqslant \sup_{x^* \in E^*} \langle x^* | x \rangle - (\langle x^* | x \rangle - F(x)) = F(x).$$

This shows that the lsc convex function F^{**} lies below F, so that F^{**} lies below the lsc convex envelope of F.

To prove that $F^{**} \ge G$, we use the following representation of G as the maximum of continuous affine functions that lie below F:

$$G(x) = \sup \left\{ \langle x^* | x \rangle + \alpha \mid (x^*, \alpha) \in X^* \times \mathbb{R} \text{ s.t. } \langle x^* | \cdot \rangle + \alpha \leqslant F \right\}.$$

We now choose some affine function defined by $(x^*, \alpha) \in E^* \times \mathbb{R}$ and lying below F, i.e. such that $F \ge \langle x^* | \cdot \rangle + \alpha$. Then,

$$F^*(x^*) \leqslant \sup_{x \in X} \langle x^* | x \rangle - F(x) \leqslant \sup_{x \in X} \langle x^* | x \rangle - (\langle x^* | x \rangle + \alpha) = -\alpha.$$

This implies that $F^{**}(x) \ge \langle x^* | x \rangle - F^*(x) \ge \langle x^* | x \rangle + \alpha$. In other words, F^{**} is larger than any affine function that lies below F, i.e. $F^{**} \ge G$. \Box

Proof of Theorem 8. We need to compute the convex conjugate and biconjucate of the functional F_c . This functional is defined on the space $\mathcal{M}(X) \times \mathcal{M}(Y)$ endowed with the product of the weak*-topologies, making it a locally convex and separated topological vector space. By definition of the weak* topology, $\mathcal{M}(X)^* = \mathcal{C}^0(X)$, so that we may identify $(\mathcal{M}(X) \times \mathcal{M}(X))^*$ with $\mathcal{C}^0(X) \times \mathcal{C}^0(Y)$. We have

$$\begin{split} F_c^*(\varphi,\psi) &= \sup_{\mu,\nu} \langle \mu | \varphi \rangle + \langle \nu | \psi \rangle - F(\mu,\nu) \\ &= \sup_{x,y \in X, m \geqslant 0} m(\langle \delta_x | \varphi \rangle + \langle \delta_y | \psi \rangle - c(x,y)) \\ &= \sup_{x,y \in X, m \geqslant 0} m(\varphi(x) + \psi(y) - c(x,y)) \\ &= \begin{cases} 0 & \text{if } \varphi \oplus \psi \leqslant c \\ +\infty & \text{otherwise} \end{cases} \end{split}$$

Therefore, the biconjugate of F_c is given by

$$\begin{split} F_c^{**}(\mu,\nu) &= \sup_{\varphi,\psi} \langle \mu | \varphi \rangle + \langle \nu | \psi \rangle - F_c^*(\varphi,\psi) \\ &= \sup_{\varphi \oplus \psi \leqslant c} \langle \mu | \varphi \rangle + \langle \nu | \psi \rangle = (\mathrm{KD}). \end{split}$$

Recall that F_c^{**} is the largest lsc convex function that lie below F_c . Since \mathcal{T}_c is lsc convex and also lies below F_c , we deduce that $(\text{KD}) = F_c^{**} \ge \mathcal{T}_c = (\text{KP})$. Since we already know (by weak duality) that $(\text{KP}) \ge (\text{KD})$, we deduce strong duality ((KP) = (KD)) and $F_c^{**} = \mathcal{T}_c$.

3.3. Existence of solution for the dual problem. Kantorovich's dual problem (KD) consists in maximizing a concave (actually linear) functional under linear inequality constraints. It can also also easily be turned into an unconstrained minimization problem. The idea is quite simple: given a certain $\psi \in C^0(Y)$, one wishes to select φ on X which is as large as possible (to maximize the term $\langle \varphi | \mu \rangle$ in (KD)) while satisfying the constraint $\varphi \oplus \psi \leq c$. This constraint can be rewritten as

$$\forall x \in X, \ \varphi(x) \leq \min_{y \in Y} c(x, y) - \psi(y).$$

The largest function φ satisfying it is $\varphi(x) = \min_{y \in Y} c(x, y) - \psi(y)$. Thus,

$$\begin{split} (\mathrm{KP}) &= \sup_{\varphi \oplus \psi \leqslant c} \langle \varphi | \mu \rangle + \langle \psi | \nu \rangle \\ &= \sup_{\psi \in \mathcal{C}^0(Y)} \int_X \left(\min_{y \in Y} c(x, y) - \psi(y) \right) \mathrm{d}\mu(x) + \int \psi(y) \mathrm{d}\nu(y). \end{split}$$

This idea is at the basis of many algorithms to solve discrete instances of optimal transport, but also useful in theory. It also suggests to introduce the notion of c-transform.

Definition 12 (*c*-Transform, *c*-Concavity). The *c*-transform (resp. \overline{c} -transform) of a function $\psi: Y \to \mathbb{R} \cup \{+\infty\}$ (resp. $\varphi: X \to \mathbb{R} \cup \{+\infty\}$) is

$$\psi^c : x \in X \mapsto \min_{y \in Y} c(x, y) - \psi(y) \tag{3.7}$$

$$\varphi^{\overline{c}} : y \in Y \mapsto \min_{x \in X} c(x, y) - \varphi(x)$$
(3.8)

A function φ on X is called *c*-concave if $\varphi = \psi^c$ for some $\psi \in \mathcal{C}^0(Y)$. Similarly, a function ψ on Y is called \overline{c} -concave if $\psi = \varphi^{\overline{c}}$ for some $\varphi \in \mathcal{C}^0(X)$. Thanks to this notion of c-transform, one can reformulate the dual problem (KD) as an unconstrained maximization problem:

$$(\mathrm{KD}) = \sup_{\psi \in \mathcal{C}^0(Y)} \int_X \psi^c \mathrm{d}\mu + \int_Y \psi \mathrm{d}\nu.$$
(3.9)

Theorem 11 (Existence of dual potentials). The dual Kantorovich problem (KD) admits a maximizer. Moreover, for any $x_0 \in X$ there exists a maximizer of the form (φ, ψ) , such that $\varphi = \psi^c$ and $\psi = \varphi^{\overline{c}}$, and satisfying $\varphi(x_0) = 0$.

The existence of maximizers follows from the fact that a c-concave/ \bar{c} -convex function has the same modulus of continuity as c.

Definition 13 (Modulus of continuity). A real-valued function f on a metric space (Z, d_Z) has modulus of continuity $\omega : \mathbb{R}^+ \to \mathbb{R}$ if ω satisfies $\lim_{t\to 0} \omega(t) = 0$ and if for every $z, z' \in Z$, $|f(z) - f(z')| \leq \omega(\mathrm{d}_Z(z, z'))$.

Lemma 12 (Properties of c-transforms). Let $\omega : \mathbb{R}^+ \to \mathbb{R}^+$ be a modulus of continuity for $c \in \mathcal{C}^0(X \times Y)$ for the distance

$$d_{X \times Y}((x, y), (x', y')) = d_X(x, x') + d_Y(y, y').$$

Then for every $\varphi \in \mathcal{C}^0(X)$ and every $\psi \in \mathcal{C}^0(Y)$,

- $\varphi^{\overline{c}}$ and ψ^c also admits ω as modulus of continuity.
- $\psi^{c\overline{c}} \ge \psi$ and $\psi^{c\overline{c}c} = \psi^c$.
- $\varphi^{\overline{c}c} \geqslant \varphi$ and $\varphi^{\overline{c}c\overline{c}} = \varphi^{\overline{c}}$.

Proof. (i) Let $\psi \in \mathcal{C}^0(Y)$ and let x be a point in X. By compactness, there exists a point y_x in Y realizing the minimum in the definition of ψ^c . Then, for every $x' \in X$,

$$\psi^{c}(x') = \min_{y \in Y} c(x', y) - \psi(y)$$

$$\leqslant c(x', y_{x}) - \psi(y_{x}) = \psi^{c}(x) + c(x', y_{x}) - c(x, y_{x})$$

$$\leqslant \psi^{c}(x) + \omega(d_{X}(x, x')).$$

Exchanging the role of x and x' we get $|\psi^c(x') - \psi^c(x)| \leq \omega(d_X(x, x'))$ as desired. The proof that $\varphi^{\overline{c}}$ has the ω as modulus of continuity is similar. (ii) By definition, of the c and \overline{c} -transforms, one has

$$\psi^{c\overline{c}}(y) = \min_{x \in X} \left(c(x,y) - \min_{\tilde{y} \in Y} c(x,\tilde{y}) - \psi(\tilde{y}) \right).$$

Taking $\tilde{y} = y$, one gets $\psi^{c\bar{c}}(y) \ge \psi(y)$. Again, by definition, we have

$$\psi^{c\overline{c}c}(x) = \min_{y \in Y} \left(c(x,y) - \min_{\tilde{x} \in X} \left(c(\tilde{x},y) - \min_{\tilde{y} \in Y} c(\tilde{x},\tilde{y}) - \psi(\tilde{y}) \right) \right).$$

By taking $\tilde{x} = x$, one gets $\psi^{c\bar{c}c}(x) \ge \psi^c(x)$, while taking $\tilde{y} = y$ gives us $\psi^{c\bar{c}c}(x) \le \psi^c(x)$. The claim (iii) is proven similarly.

Proof of Theorem 11. Let $(\varphi_n, \psi_n)_{n \in \mathbb{N}}$ be a maximizing sequence for (KD), i.e. $\varphi_n \oplus \psi_n \leqslant c$ and $\lim_{n \to +\infty} \langle \varphi_n | \mu \rangle + \langle \psi_n | \nu \rangle =$ (KD). Define $\hat{\varphi}_n = \psi_n^c$ and $\hat{\psi}_n = \hat{\varphi}_n^{\overline{c}}$. Then $\hat{\varphi}_n \oplus \hat{\psi}_n \leqslant c$, $\varphi_n \leqslant \hat{\varphi}_n$ and $\psi_n \leqslant \hat{\psi}_n$, which implies

$$\langle \varphi_n | \mu \rangle + \langle \psi_n | \nu \rangle \leqslant \langle \hat{\varphi}_n | \mu \rangle + \langle \psi_n | \nu \rangle.$$

Thus, the sequence $(\hat{\varphi}_n, \hat{\psi}_n)_{n \in \mathbb{N}}$ is also a maximizing sequence for (KD). We note at this point that it is possible to assume that $\hat{\varphi}_n(x_0) = 0$ for all n, where x_0 is a given point in X. Indeed, if this is not the case, we may replace the original sequence $(\hat{\varphi}_n, \hat{\psi}_n)_{n \in \mathbb{N}}$ by $(\hat{\varphi}_n - \hat{\varphi}_n(x_0), \hat{\psi}_n + \hat{\varphi}_n(x_0))_{n \in \mathbb{N}}$, which is also admissible and has the same dual value.

We now prove that the sequence $(\hat{\varphi}_n, \psi_n)$ admits a converging subsequence. By Lemma 12, the sequences $(\hat{\varphi}_n)_n$ and $(\hat{\psi}_n)_n$ are equicontinuous. Since $\hat{\varphi}_n(x_0) = 0$, we deduce from uniform continuity that the sequence $(\hat{\varphi}_n)_{n \in \mathbb{N}}$ is uniformly bounded. Then, using

$$\hat{\psi}_n(y) = \hat{\varphi}_n^{\overline{c}}(y) = \max_{x \in X} c(x, y) - \hat{\varphi}_n(x),$$

we deduce that $\|\hat{\psi}_n\|_{\infty} \leq \|c\|_{\infty} + \|\hat{\varphi}_n\|_{\infty}$ so that $(\hat{\varphi}_n)_{n \in \mathbb{N}}$ is also uniformly bounded. By Arzelà-Ascoli's theorem, both sequences therefore admit converging subsequences. The limit potentials are then maximizers for (KD) because the functional which is maximized in (KD) is continuous.

3.4. Stability of optimal transport plans.

Proposition 13 (Support of OT plans). Let $(\varphi, \psi) \in C^0(X) \times C^0(Y)$ be admissible for the problem (KD), i.e. $\varphi \oplus \psi \leq c$, and let $\gamma \in \Gamma(\mu, \nu)$ be a transport plan. Then the two assertions are equivalent

- γ is an optimal transport plan and (φ, ψ) is a maximizer in (KD)
- $\operatorname{spt}(\gamma) \subseteq \{(x,y) \in X \times Y \mid \varphi(x) \oplus \psi(y) = c(x,y)\}.$

Proof. Using first the admissibility of (φ, ψ) and then $\gamma \in \Gamma(\mu, \nu)$,

$$0 \leqslant \langle c | \gamma \rangle - \langle \varphi \oplus \psi | \gamma \rangle = \langle c | \gamma \rangle - (\langle \varphi | \mu \rangle + \langle \psi | \nu \rangle$$

We see that the last term vanishes if and only if γ minimizes (KP) and (φ, ψ) maximizes (KD) (and if strong duality, (KP) = (KD), holds). But this term also vanishes if and only if the first inequality is an equality. Since $\varphi \oplus \psi \leq c$, this is equivalent to $c - \varphi \oplus \psi = 0$ γ -almost everywhere.

Because of this proposition, one can think of the dual Kantorovich potentials, the prices in the economic interpretation of OT, as an "optimality certificate" for an optimal transport plan (i.e. a way to convince someone that you actually found the optimum). This leads to the following stability theorem for optimal transport maps.

Theorem 14 (Stability of OT plans). Let X, Y be compact metric spaces and let $c \in C^0(X \times Y)$. Consider $(\mu_k)_{k \in \mathbb{N}}$ and $(\nu_k)_{k \in \mathbb{N}}$ in $\mathcal{P}(X)$ and $\mathcal{P}(Y)$ converging weakly to μ and ν respectively.

- If $\gamma_k \in \Gamma(\mu_k, \nu_k)$ is optimal then, up to subsequences, (γ_k) converges weakly to an optimal transport plan $\gamma \in \Gamma(\mu, \nu)$.
- Let (φ_k, ψ_k) be optimal Kantorovich potentials in the dual problem between μ_k and ν_k , satisfying $\psi_k = \varphi_k^{\overline{c}}$, $\varphi_k = \psi_k^c$ and $\varphi_k(x_0) = 0$ for some $x_0 \in X$. Then, up to subsequences, the sequence (φ_k, ψ_k) converges uniformly to a maximizing pair (φ, ψ) for (KD) also satisfying $\varphi = \psi^c$ and $\psi = \varphi^{\overline{c}}$.

We will use the following lemma about the convergence of the supports of weak* converging measures.

Lemma 15. If a sequence of non-negative measures $(\mu_n)_{n\in\mathbb{N}}$ weak*-converges to μ , then any point x in spt (μ) is the limit as $n \to +\infty$ of points x_n in spt (μ_n) .

Proof of Theorem 14. As c-concave functions, φ_k and ψ_k have the same modulus of continuity as the cost function c (see Lemma 12), and they are uniformly bounded (using $\varphi_k(x_0) = 0$). Using Arzelà-Ascoli theorem, we can therefore assume that up to subsequences, (φ_k) (resp. (ψ_k)) converges to some φ (resp ψ) uniformly. Then, one easily sees that $\varphi \oplus \psi \leq c$ so that (φ, ψ) are admissible for the limit dual problem (KD). By Proposition 2, we can assume, taking subsequences if necessary, that the sequence $\gamma_k \in \Gamma(\mu_k, \nu_k)$ converges to some $\gamma \in \Gamma(\mu, \nu)$.

By Proposition 13, we see that γ_k is supported on the set

$$S_k = \{(x, y) \in X \times Y \mid \varphi_k(x) + \psi_k(y) = c(x, y)\}$$

Moreover, by Lemma 15, every pair $(x, y) \in \operatorname{spt}(\gamma)$ can be approximated by a sequence of pairs $(x_k, y_k) \in \operatorname{spt}(\gamma_k)$ i.e. $\lim_{k\to\infty} (x_k, y_k) = (x, y)$. Since γ_k is supported on S_k one has $c(x_k, y_k) = \varphi_k(x_k) + \psi_k(x_k)$. This gives at the limit $c(x, y) = \varphi(x) + \psi(y)$. We have just shown that for every point pair (x, y) in $\operatorname{spt}(\gamma), c(x, y) = \varphi(x) + \psi(y)$ where φ, ψ is admissible. Applying Proposition 13 again, this shows that γ and (φ, ψ) are optimal for their respective problems.

4. KANTOROVICH'S FUNCTIONAL

4.1. Kantorovich's functional. As already mentioned in (3.9), the Kantorovich's dual problem (KD) can be expressed as an unconstrained maximization problem, involving the *c*-transform.

Definition 14. The Kantorovitch functional is defined on $\mathcal{C}^0(Y)$ by

$$\mathcal{K}_{\mu}(\psi) = \int_{X} \psi^{c} \mathrm{d}\mu \qquad (4.10)$$

The Kantorovitch dual problem therefore amounts to maximizing the Kantorovitch functional plus a linear term:

(KD) =
$$\max_{\psi \in \mathcal{C}^0(Y)} \mathcal{K}_{\mu}(\psi) + \langle \psi | \nu \rangle.$$

It is quite easy to see that \mathcal{K}_{μ} is concave, recalling the definition of the *c*-transform as a minimum. If (φ, ψ) are maximizers in the Kantorovich's dual problem (KD) between μ and ν , then ψ is a maximizer of $\mathcal{K}_{\mu} + \langle \cdot | \nu \rangle$.

This subsection is devoted to the computation of the superdifferential of Kantorovich's functional, in particular when the source measure μ is absolutely continuous. This computation will be used to establish existence of solutions to Monge's problem (following Brenier and Gangbo-McCann) and to construct and study algorithms for (semi-)discretized optimal transport.

Definition 15 (Response map). Given a potential $\psi \in \mathcal{C}^0(Y)$, we call *response map* the set-valued map \hat{T}_{ψ} defined by

$$\hat{T}_{\psi}(x) = \arg\min_{y \in Y} c(x, y) - \psi(y) = \{ y \in Y \mid c(x, y) - \psi(y) = \psi^{c}(x) \}.$$

Remark 4 (Construction of optimal transports). One can easily sees that the graph of \hat{T}_{ψ} is

$$\operatorname{Graph}(\hat{T}_{\psi}) = \{(x, y) \in X \times Y \mid \psi^c(x) + \psi(y) = c(x, y)\}.$$

We note that if ψ is a maximizer of $\mathcal{K}_{\mu} + \langle \cdot | \nu \rangle$, then (ψ^c, ψ) is a maximizer of (KD). By proposition Proposition 13, we see that the set of optimal transport plans between μ and ν is equal to

$$\{\gamma \in \Gamma(\mu, \nu) \mid \operatorname{spt}(\gamma) \subseteq \operatorname{Graph}(\hat{T}_{\psi})\},\tag{4.11}$$

making it a priori possible to recover a solution to the primal problem from a maximizer of the $\mathcal{K}_{\mu} + \langle \cdot | \nu \rangle$.

Proposition 16. Let X, Y be compact metric spaces and let $c \in C^0(X \times Y)$. Then, for all measure $\mu \in \mathcal{P}(X)$ and any $\psi \in C^0(Y)$, one has

$$\partial^{+} \mathcal{K}_{\mu}(\psi) = \left\{ -\nu \mid \exists \gamma \in \Gamma(\mu, \nu) \ s.t. \ \operatorname{spt}(\gamma) \subseteq \operatorname{Graph}(\hat{T}_{\psi}) \right\}.$$

Proof. Let $\psi \in \mathcal{C}^0(Y)$ and let $\nu \in (\mathcal{C}^0(Y))^* = \mathcal{M}(Y)$. Assume that $-\nu$ belongs to $\partial^+ \mathcal{K}_{\mu}(\psi)$. Then,

$$\forall \psi' \in \mathcal{C}^0(Y), \mathcal{K}_\mu(\psi') \leqslant \mathcal{K}_\nu(\psi) - \langle \psi' - \psi | \nu \rangle,$$

which is equivalent to

$$\forall \psi' \in \mathcal{C}^0(Y), \langle (\psi')^c | \mu \rangle + \langle \psi' | \nu \rangle \leqslant \langle \psi^c | \mu \rangle + \langle \psi | \nu \rangle,$$

so that (ψ^c, ψ) is a maximizer of the dual Kantorovich problem between μ and ν . By strong Kantorovich duality $(\mathcal{T}_c(\mu, \nu) = (\text{KD}))$, this implies that ν is non-negative, with same mass as μ , and that $\langle \psi^c | \mu \rangle + \langle \psi | \nu \rangle = \mathcal{T}_c(\mu, \nu)$. Let $\gamma \in \Gamma(\mu, \nu)$ be an optimal transport plan between μ and ν for the cost c. Then, by Proposition 13, we see that $\psi^c \oplus \psi = c$ on $\operatorname{spt}(\gamma)$ as desired.

Conversely, if a measure ν is such that there exists $\gamma \in \Gamma(\mu, \nu)$ supported on $\psi^c \oplus \psi = c$, we get using $(\psi')^c \oplus \psi \leq c$

$$\begin{aligned} \mathcal{K}_{\mu}(\psi') &= \langle (\psi')^{c} | \mu \rangle = \langle (\psi')^{c} \oplus \psi' | \gamma \rangle - \langle \psi' | \nu \rangle \\ &\leq \langle c | \gamma \rangle - \langle \psi | \nu \rangle \\ &= \langle \psi^{c} \oplus \psi | \gamma \rangle - \langle \psi' | \nu \rangle \\ &= \mathcal{K}_{\mu}(\psi) + \langle \psi' - \psi | - \nu \rangle, \end{aligned}$$

thus proving that $-\nu \in \partial^+ \mathcal{K}_{\mu}(\psi)$.

4.2. Solution of Monge's problem. We now use Proposition 16 to prove the existence of optimal transport maps when the source measure is absolutely continuous on a compact subset of \mathbb{R}^d and when the cost function satisfies a *twist condition*. This result is due to Brenier [12] in the case of the quadratic cost, that is $c(x, y) = ||x - y||^2$ on \mathbb{R}^d , and Gangbo-McCann in the general case of twisted costs [20]. The question is to determine conditions under which the response map \hat{T}_{ψ} is single-valued μ -almost everywhere.

Definition 16 (Twisted cost). Let Ω_X, Ω_Y be open subsets of \mathbb{R}^d , and let $c \in \mathcal{C}^1(\Omega_X \times \Omega_Y)$. The cost function c satisfies the *twist condition* if

$$\forall x_0 \in \Omega_X, \text{ the map } y \in \Omega_Y \mapsto \nabla_x c(x_0, y) \in \mathbb{R}^d \text{ is injective}, \qquad (4.12)$$

where $\nabla_x c(x_0, y)$ is the gradient of $x \mapsto c(\cdot, y)$ at $x = x_0$.

Proposition 17. Let Ω_X , Ω_Y be open subsets of \mathbb{R}^d , let $c \in \mathcal{C}^1(\Omega_X \times \Omega_Y)$ be a cost satisfying the twist condition (4.12), and let X, Y be compact subsets of Ω_X and Ω_Y . Then, for Lebesgue-almost every $x \in X$, the response map is a singleton:

$$\hat{T}_{\psi}(x) = \arg\min_{y \in Y} c(x, y) - \psi(y) =: \{T_{\psi}(x)\}.$$

In particular, if $\mu \in \mathcal{P}(X)$ is absolutely continuous, then

$$\nabla \mathcal{K}_{\mu}(\psi) = -T_{\psi \#}\mu.$$

Proof. Define $\varphi = \psi^c$, i.e. $\varphi(x) = \min_{y \in Y} c(x, y) - \psi(y)$. If the minimum in the definition of the response map is not unique, there exists two distinct points y_0, y_1 in $\hat{T}_{\psi}(x)$. For any $i \in \{0, 1\}$, we have

$$\varphi(x') = \min_{y \in Y} c(x, y) - \psi(y) \leqslant c(x', y_i) - \psi(y_i),$$

with equality at x' = x. Since $\nabla c(x', y_1) \neq \nabla c(x', y_0)$ by injectivity of $y \mapsto \nabla c(x', y)$, we see that φ is not differentiable at x.

Using $c \in \operatorname{Lip}(X \times Y)$, we get that φ is Lipschitz. Rademacher's theorem then implies that φ is differentiable on a set B with full Lebesgue measure in X. By the previous paragraph, we obtain that \hat{T}_{ψ} is a singleton at any point of B. We conclude with the next lemma. \Box

Lemma 18. Let $\mu \in \mathcal{P}(X)$ and let $\hat{T} : X \to Y$ be a set-valued map such that $\hat{T}(x) = \{T(x)\}$ for μ -almost every x. Then, there exists only one transport plan $\gamma \in \Gamma(\mu, \nu)$ satisfying $\operatorname{spt}(\gamma) \subseteq \gamma(\operatorname{Graph}(\hat{T}))$. This transport plan is induced by the map T, i.e. $\gamma = (\operatorname{id}, T)_{\#}\gamma$.

Proof. By definition of $\gamma_T = (\mathrm{id}, T)_{\#}\gamma$ one has $\gamma_T(A \times B) = \mu(T^{-1}(B) \cap A)$ for all Borel sets $A \subseteq X$ and $B \subseteq Y$. On the other hand, consider the set $X' \subseteq X$ of points such that $\hat{T}(x) = \{T(x)\}$, so that $X \setminus X'$ is μ -negligible by assumption. Then,

$$\begin{split} \gamma(A \times B) &= \gamma((A \cap X') \times B) \\ &= \gamma(\{(x, y) \mid x \in A \cap X', \text{ and } y \in B\}) \\ &= \gamma(\{(x, y) \mid x \in A \cap X', y \in B \text{ and } y = T(x)\}) \\ &= \gamma(\{(x, y) \mid x \in A \cap X' \cap T^{-1}(B), y = T(x)\} \\ &= \mu(A \cap X' \cap T^{-1}(B)) \\ &= \mu(A \cap T^{-1}(X)) \end{split}$$

thus proving the claim.

Theorem 19 (Gangbo-McCann [20]). Let Ω_X, Ω_Y be open subsets of \mathbb{R}^d and let $c \in \mathcal{C}^1(\Omega_X \times \Omega_Y)$ be a cost satisfying the twist condition (4.12). Given compact subsets X and Y of Ω_X and Ω_Y and two probability measures $(\mu, \nu) \in \mathcal{P}^{\mathrm{ac}}(X) \times \mathcal{P}(Y)$. Then, there exists $\psi \in \mathcal{C}^0(Y)$ such that the unique optimal transport map between μ and ν is induced by T_{ψ} .

Proof of Theorem 19. Let ψ be a maximizer of $\mathcal{K}_{\mu} + \langle \nu | \cdot \rangle$. By equation (4.11), the set of optimal transport plans is $\{\gamma \in \Gamma(\mu, \nu) \mid \operatorname{spt}(\gamma) \subseteq \operatorname{Graph}(\hat{T}_{\psi})\}$.

Combining Proposition 17 and Lemma 18, we deduce that the unique element of this set is $\gamma = (\mathrm{id}, T_{\psi})_{\#} \mu$.

Here, we obtain Brenier's theorem as a corollary of Gangbo-McCann's result — even though historically Brenier's theorem has been proven first.

Corollary 20 (Brenier [12]). Let X, Y be two compact subsets of \mathbb{R}^d , let $c(x,y) = ||x-y||^2$ and let $(\mu,\nu) \in \mathcal{P}^{\mathrm{ac}}(X) \times \mathcal{P}(Y)$. Then, there exists $\varphi : \mathbb{R}^d \to \mathbb{R}$ convex such that $\nabla \varphi_{\#} \mu = \nu$ and the unique optimal transport plan between μ and ν is induced by the map $T = \nabla \varphi$.

Proof. We need to compute the response map associated to the maximizer ψ of $\mathcal{K}_{\mu} + \langle \cdot | \nu \rangle$ for the quadratic cost:

$$T_{\psi}(x) = \arg\min_{y} ||x - y||^2 - \psi(y)$$

= $\arg\min_{y} ||y||^2 - 2\langle x|y \rangle - \psi(y)$
= $\arg\max_{y} \langle x|y \rangle - \frac{1}{2}(||y||^2 - \psi(y)).$

Recalling the definition of the convex conjugate, one can see at once that $T_{\psi} = \nabla u$ where $u = \left(\frac{1}{2}(\|\cdot\|^2 - \psi)\right)^*$.

Remark 5 (Monge-Kantorovich quantiles). Given a fixed probability density ρ on a compact domain of \mathbb{R}^d , e.g. $\rho \equiv 1$ on $[0, 1]^d$, and any compactly supported $\nu \in \mathcal{P}(\mathbb{R}^d)$, one can denote T_{ν} the quadratic optimal transport map between ρ and ν . In dimension d = 1, one recovers the quantile function. In higher dimension, there is no canonical definition of a quantile function, but T_{ν} was proposed as a challenger under the name "Monge-Kantorovich quantile" by Chernozhukov, Galichon, Hallin, Henry in [15]. Being the gradient of a convex function, the Monge-Kantorovich quantile is monotone, i.e.

for a.e.
$$x, y \in \operatorname{spt}(\nu), \langle T_{\nu}(x) - T_{\nu}(y) | x - y \rangle \ge 0.$$

This notion can be used to define multivariate notions of ranks and depth.

4.3. Semi-discrete optimal transport. Our working assumptions for the remainder of this section are the following:

- Ω_X, Ω_Y are two open subsets of \mathbb{R}^d . The cost function c belongs to $\mathcal{C}^1(\Omega_X \times \Omega_Y)$ and satisfies the twist condition (4.12).
- the source measure ρ is absolutely continuous with respect to the Lebesgue measure and is supported in a compact subset X of Ω_X .
- the target space Y is finite so that $\nu \in \mathcal{P}(Y)$ can be written under the form $\nu = \sum_{y \in Y} \nu_y \delta_y$. For simplicity, we assume that $\min_y \nu_y > 0$.

Note that by an abuse of notation, we will often conflate ρ with its density with respect to the Lebesgue measure.

Definition 17 (Laguerre tessellation). The Laguerre tessellation associated to a set of prices $\psi : Y \to \mathbb{R}$ is a decomposition of the space into *Laguerre cells* defined by

$$\operatorname{Lag}_{y}(\psi) := \{ x \in \Omega_{X} \mid \forall z \in Y, c(x, y) - \psi(y) \leq c(x, z) - \psi(z) \}.$$
(4.13)



FIGURE 1. (Left) The domain X (with boundary in blue) is endowed with a probability density pictured in grayscale representing the density of population in a city. The set Y (in red) represents the location of bakeries. Here, $X, Y \subseteq \mathbb{R}^2$ and $c(x,y) = |x-y|^2$ (Middle) The Voronoi tessellation induced by the bakeries (Right) The Laguerre tessellation: the price of bread the bakery near the center of X is higher than at the other bakeries, effectively shrinking its Laguerre cell.

When $\psi \equiv 0$, the Laguerre cells are called Voronoi cells. The Voronoi cell of the point $y \in Y$ is denoted $\operatorname{Vor}_{y}(\psi)$.

Remark 6 (Response map). Let $\psi \in \mathbb{R}^{Y}$. The response map T_{ψ} is constant on the interior of the Laguerre cells (and undefined on their boundary) by:

$$\forall y \in Y, \ T_{\psi} | \operatorname{Lag}_y = y.$$

In particular,

$$T_{\psi \#} \rho = \sum_{y \in Y} G_y(\psi) \delta_y, \text{ where } G_y(\psi) = \rho(\text{Lag}_y).$$
(4.14)

Theorem 21 (Aurenhammer, Hoffman, Aronov). Under the assumptions of this paragraph, the Kantorovich functional \mathcal{K}_{μ} is \mathcal{C}^1 -smooth on \mathbb{R}^Y . Its qradient is given by

$$\nabla \mathcal{K}_{\rho}(\psi) = -\sum_{y \in Y} \rho(\operatorname{Lag}_{y}(\psi)) \delta_{y}$$
(4.15)

In particular $\psi \in \mathbb{R}^Y$ maximizes $\mathcal{K}_{\rho} + \langle \cdot | \nu \rangle$, where $\nu \in \mathcal{P}(Y)$, if and only if $\forall y \in Y, \ \rho(\operatorname{Lag}_{u}(\psi)) = \nu(\{y\}).$

The only new statement in this theorem, compared to Proposition 17 is that \mathcal{K}_{μ} is \mathcal{C}^1 . This is proven as point (iv) of the following lemma. In what follows, we will denote R the oscillation of the cost function:

$$R := \max_{X \times Y} c - \min_{X \times Y} c, \tag{4.16}$$

Lemma 22. Assume c is twisted (Def. 16) and $\rho \in \mathcal{P}^{\mathrm{ac}}(X)$. Then,

- (i) $\forall y \in Y$, the map $t \mapsto G_y(\psi + t\mathbf{1}_y)$ is non-decreasing,
- (ii) $\forall y \neq z \in Y$, the map $t \mapsto G_y(\psi + t\mathbf{1}_z)$ is non-increasing,
- (iii) if $\psi \in \mathbb{R}^Y$ is such that $G_{y_0}(\psi) > 0$, then $\psi(y_0) \leq \min_Y \psi + R$, (iv) for all $y \in Y$, the function G_y is continuous.

Proof. The properties (i), (ii) are straightforward consequences of the definition of Laguerre cells. To prove (iii), take ψ such that $G_{y_0}(\psi) > 0$, implying in particular that the Laguerre cell $\operatorname{Lag}_{y_0}(\psi)$ is non-empty and contains a point $x \in X$. Then, by definition of the cell one has for all $y \in Y \setminus \{y_0\}$, $c(x, y_0) + \psi(y_0) \leq c(x, y) + \psi(y)$, thus showing that $\psi(y_0) \leq \min_Y \psi + R$.

It remains to establish that each of the maps G_y is continuous. For this purpose, we consider a sequence $(\psi_n)_{n\in\mathbb{N}}$ converging to some ψ_{∞} . We first note that thanks to the Twist hypothesis, the set S defined by

$$S = \{x \in X \mid \exists y \neq z \in Y \text{ s.t. } c(x,y) - \psi(y) = c(x,y) - \psi(z)\}$$
$$\subseteq \bigcup_{y \in Y, z \in Y \setminus \{y\}} \{x \in X \mid c(x,y) - \psi(y) = c(x,y) - \psi(z)\}.$$

is included in a finite union of (d-1)-dimensional submanifolds, which are all Lebesgue-negligible. Thus, S is also ρ -negligible. Defining $\chi = \mathbf{1}_{\text{Lag}_y(\psi)}$ and $\chi_n = \mathbf{1}_{\text{Lag}_y(\psi_n)}$, we have

$$G_y(\psi_n) = \int \chi_n d\rho$$
, and $G(\psi) = \int \chi d\rho$.

To prove that $\lim_{n\to+\infty} G_y(\psi_n) = G_y(\psi)$ it suffices to establish that χ_n converges to χ on $X \setminus S$, which is straightforward (because the inequalities defining the set $X \setminus S$ are strict), and to apply Lebesgue's dominated convergence theorem.

4.4. **Oliker-Prussner's algorithm.** Oliker-Prussner's algorithm for solving $G(\psi) = \nu$ is described in Algorithm 1, and bears strong resemblance with Bertsekas' auction algorithm for the assignment problem [8, 9]. In particular, the values of ψ are evolved in a monotonic way.

Algorithm 1 Oliker-Prussner algorithm

Input: A tolerence parameter $\delta > 0$. **Initialization:** Fix some $y_0 \in Y$ once for all. Set

$$\psi^{(0)}(y) := \begin{cases} 0 & \text{if } y = y_0 \\ R & \text{if not.} \end{cases}$$

While: $\exists y \in Y \setminus \{y_0\}$ such that $G_y(\psi^{(k)}) \leq \nu_y - \frac{\delta}{N}$ Step 1: Compute

$$t_y = \min\{t \ge 0 \mid G_y(\psi^{(k)} + t\mathbf{1}_y) \ge \nu_y\}.$$
(4.17)

Step 2: Set $\psi^{(k+1)} = \psi^{(k)} + t\mathbf{1}_y$. Output: A vector $\psi^{(k)}$ that satisfies $\max_y \left\| G_y(\psi^{(k)}) - \nu(\{y\}) \right\|_{\infty} \leq \delta$.

Theorem 23 (Oliker-Prussner). Assume that the cost $c \in C^2(\Omega_X \times \Omega_Y)$ is twisted (Def. 16) and that $\rho \in \mathcal{P}^{\mathrm{ac}}(X) \cap \mathrm{L}^{\infty}(X)$. Then,

- Oliker-Prussner's algorithm terminates in a finite number of steps.
- Furthermore, at the final step k, one has

$$\max_{y \in Y} \left| G_i(\psi^{(k)}) - \nu_i \right| \leqslant \delta$$

Proof of Theorem 23.

Step 1 (Correctness) When Algorithm 1 terminates with $\psi := \psi^{(k)}$, one has for any $y \neq y_0$, $\rho(\operatorname{Lag}_y(\psi)) \leq \nu_y$. When it stops, it also means that one has $\rho(\operatorname{Lag}_y(\psi)) \geq \nu_y - \frac{\delta}{N}$. Then, as desired, we get

$$\rho(\operatorname{Lag}_{y_0}(\psi)) = 1 - \sum_{y \neq y_0} \rho(\operatorname{Lag}_{y_0}(\psi)) \in [\nu_{y_0}, \nu_{y_0} + \delta].$$

Step 2 (A priori bound on ψ_k) By construction one has $\rho(\text{Lag}_y(\psi^{(k)})) \leq \nu_y$, which also imply that

$$\rho(\operatorname{Lag}_{y_0}(\psi^{(k)})) = 1 - \sum_{y \in Y \setminus \{y_0\}} \rho(\operatorname{Lag}_y(\psi^{(k)})) \ge \nu_{y_0} > 0.$$

By Proposition 22-(iii), we get $0 = \psi^k(y_0) \leq \min_Y \psi^{(k)} + R$. Since the price of y_0 is never changed, $\psi^{(k)}(y_0) = 0$ and $R \geq \psi^{(k)} \geq -R$.

Step 3 (Minimum decrease and termination) Since by Lemma 22–(iv) G_y is continuous, it admits a continuity modulus on the compact set $[-R, R]^Y$, i.e. a function $\omega_y : \mathbb{R} \to \mathbb{R}$ such that $\lim_{t\to 0} \omega_y(t) = 0$ and such that

$$\forall \psi, \psi' \in [-R, R]^Y, \left| G_y(\psi) - G_y(\psi') \right| \le \left\| \psi - \psi' \right\|_{\infty}$$

In the second step of the algorithm, when $\psi^{(k)}$ is updated one has $G_y(\psi^{(k)} - t_y \mathbf{1}_y) \ge G_y(\psi^{(k)}) + \frac{\delta}{N}$. Using the uniform continuity of G_y , we have

$$\frac{\delta}{N} \leqslant \left| G_y(\psi^{(k)} - t_y \mathbf{1}_y) - G_y(\psi^{(k)}) \right| \leqslant \omega(t_y),$$

implying that there exists $\tau > 0$ such that $t_y \ge \tau$. Since for any $k, \psi_k(y) \in [-R, R]$, the number of times k_y the price of a point $y \in Y$ has been updated is bounded: $k_y \le 2R/\tau$. Thus, the algorithm terminates in finite time. \Box

Remark 7 (Quadratic cost). For the cost $c(x, y) = ||x - y||^2$, but also in more general cases (see e.g. [27]), one can show that G is Lipschitz, with constant larger than CN. In this case, the number of iterations is of the algorithm is bounded by $O(N^3)$.

5. ENTROPY-REGULARIZED OPTIMAL TRANSPORT

5.1. **Primal problem.** We start from the primal formulation of the optimal transport problem, but instead of imposing the non-negativity constraints $\gamma \ge 0$, we add a term to the transport cost, which promotes (minus) the entropy of the transport plan and acts as a barrier for the non-negativity constraint. The entropy of a measure $\mu \in \mathcal{M}(X)$ on a compact metric space X with respect to a probability measure ω on X is defined by

$$H(\mu \mid \omega) = \begin{cases} \int h(\rho) d\omega & \text{if } d\mu = \rho d\omega \\ +\infty & \text{otherwise }, \end{cases}$$

where $h(r) = \begin{cases} r(\log r - 1) & \text{if } r > 0, \\ 0 & \text{if } r = 0, \\ +\infty & \text{if } r < 0. \end{cases}$ (5.18)

The regularized optimal transport problem is then defined as

$$(\mathrm{KP}^{\varepsilon}) := \inf_{\gamma \in \Gamma(\mu,\nu)} \langle c | \gamma \rangle + \varepsilon H(\gamma \mid \mu \otimes \nu).$$
(5.19)

We will rely on the following dual representation of entropy.

Proposition 24 (Donsker-Varadhan). Let Z be a compact space, and let $\omega \in \mathcal{M}_+(Z)$ be finite. Then, for any measure $\mu \in \mathcal{M}(Z)$,

$$H(\mu \mid \omega) = \sup_{f \in \mathcal{C}^0(Z)} \langle f \mid \mu \rangle - \langle e^f \mid \omega \rangle.$$
(5.20)

In particular, $\mu \mapsto H(\mu \mid \omega)$ is convex and weak* lsc. In addition:

- (i) the supremum in (5.20) is attained at $f \in C^0(Z)$ if and only if e^f is the density of μ with respect to ω .
- (ii) the restriction of $\mu \mapsto H(\mu \mid \omega)$ to the set of absolutely continuous measures with respect to ω is strictly convex.

Remark 8 (Finite entropy implies non-negativity). We can prove thanks to (5.20) that if $\mu \notin \mathcal{M}^+(Z)$, then $H(\mu \mid \omega) = +\infty$. Indeed, if $\langle \mu | g \rangle < 0$ for some continuous function $g \ge 0$, one can check by taking $f = -\lambda g$ that

$$H(\mu \mid \gamma) \ge \lambda \underbrace{\langle \mu \mid -g \rangle}_{>0} - \langle \underbrace{e^{\lambda g}}_{\leqslant 1} \mid \omega \rangle \xrightarrow{\lambda \to +\infty} +\infty.$$

This means that the regularized optimal transport problem can be equivalently written by removing non-negativity constraint $\gamma \ge 0$:

$$(\mathrm{KP}^{\varepsilon}) = \inf_{\gamma \in \mathcal{M}(X \times Y) \mid \Pi_{X \#} \gamma = \mu, \Pi_{Y \#} \gamma = \nu} \langle c | \gamma \rangle + \varepsilon H(\gamma \mid \mu \otimes \nu).$$

Proof. Note that for r > 0, $h'(r) = \ln(r)$ for r > 0. The convex conjugate of h is therefore given by

$$h^*(s) = \sup_{r>0} rs - h(r) = e^r.$$

The Fenchel-Young inequality reads $h^*(s) + h(r) \ge rs$ with equality if and only if $r = e^s$. Assume that μ has density ρ with respect to ω . Then,

$$H(\mu \mid \omega) = \int h(\rho(x)) d\omega(x)$$
$$= \int h^{**}(\rho(x)) d\omega(x)$$
$$= \int \sup_{s} s\rho(x) - h^{*}(\rho(x)) d\omega(x)$$

In particular, for any bounded measurable function f we have

$$H(\mu \mid \omega) \ge \langle f \mid \rho \omega \rangle - \langle e^f \mid \omega \rangle = \langle f \mid \mu \rangle - \langle e^f \mid \omega \rangle$$

with equality if $f = e^{\rho}$ a.e.

Proposition 25. The regularized optimal transport problem admits a unique solution. Moreover, the density of γ with respect to $\mu \otimes \nu$ is positive a.e.

Remark 9 (No transport maps). In this entropy regularized setting, one cannot expect to find an optimal transport map, since minimizers of the regularized optimal transport problem are supported on the whole support of the product $\mu \otimes \nu$.

22

Remark 10 (Barrier). The main ingredient of the previous proposition is that the slope of $h: r \mapsto r \ln r$ is $+\infty$ at r = 0, which forbids the density of γ with respect to $\mu \otimes \nu$ to vanish on sets of positive measure. A stronger effect could be obtained by using a penalization of the form $\varepsilon G(\gamma \mid \mu \otimes \nu)$ instead of $\varepsilon H(\gamma \mid \mu \otimes \nu)$ where

$$G(\mu \mid \omega) = \begin{cases} \int g(\rho) d\mu \otimes \nu & \text{if } d\mu = \rho d\omega \\ +\infty & \text{otherwise }, \end{cases}$$
(5.21)

where

$$g(r) = \begin{cases} -\log r & \text{if } r > 0, \\ +\infty & \text{if } r \leqslant 0. \end{cases}$$

This barrier is stronger, as it forbis r = 0. When X and Y are finite, this choice is related to the interior point method for solving the optimal transport problem, where one would solve subsequent problems of the form

$$\min_{\gamma \in \Gamma(\mu,\nu)} \langle c | \gamma \rangle + \varepsilon_k H(\gamma \mid \mu \otimes \nu)$$

for a sequence of parameters ε_k converging to zero.

Proof. Existence follows from lower semi-continuity of the functional and compactness of $\Gamma(\mu, \nu)$, while uniqueness follows from the strict convexity.

Let γ^* be the optimizer of $(\mathrm{KP}^{\varepsilon})$, and let ρ be the density of γ^* with respect to $\mu \otimes \nu$. We will prove by contradiction that the set $Z := \{(x, y) \mid \rho = 0\}$ satisfies $\rho(Z) = 0$. For this purpose, we define a new transport plan γ^t between μ and ν by setting $\gamma^t = (1 - t)\gamma^* + t\mu \otimes \nu$. The density of γ^t with respect to $\mu \otimes \nu$ is $\rho^t = (1 - t)\rho + t$. We give an upper bound on the energy of γ^t . We first observe that by convexity of $h(r) = r(\ln r - 1)$, we have

$$\int_{X \times Y \setminus Z} h(\rho^t) d\mu \otimes \nu \leqslant (1-t) \int_{X \times Y \setminus Z} h(\rho^t) d\mu \otimes \nu + t \int_{X \times Y \setminus Z} h(1) d\mu \otimes \nu$$
$$= (1-t) H(\gamma^t \mid \mu \otimes \nu) - t \cdot \mu \otimes \nu(X \times Y \setminus Z).$$

On the other hand, on Z we have $\rho^t = t$, so that

$$\int_{X \times Y \setminus Z} h(\rho^t) \mathrm{d}\mu \otimes \nu = t(\ln(t) - 1) \cdot \mu \otimes \nu(Z).$$

Finally, we note that $\langle c|\gamma^t\rangle = \langle c|\gamma^*\rangle + t(\langle \mu \otimes \nu - \gamma^*|c\rangle)$. Summing these equalities and inequalities, we get

 $\langle c|\gamma^t \rangle + \varepsilon H(\gamma^t \mid \mu \otimes \nu) \leq \langle c|\gamma^* \rangle + \varepsilon H(\gamma^* \mid \mu \otimes \nu) + t(C + \ln(t) \cdot \mu \otimes \nu(Z)).$ Taking t small enough, one get a contradiction on the optimality of γ^* , unless the set Z has zero $\mu \otimes \nu$ -measure.

5.2. **Dual problem.** The dual problem is constructed, as before, by introducing Lagrange multipliers $\varphi \in C^0(X)$ and $\psi \in C^0(Y)$ for the constraints $\Pi_{X\#}\gamma = \mu$ and $\Pi_{Y\#}\gamma = \nu$, and also dualizing the entropy using the Donsker-Varadhan formula. We have

$$(\mathrm{KP}^{\varepsilon}) = \inf_{\substack{\gamma \mid \Pi_{X \#} \gamma = \mu \text{ and } \Pi_{X \#} \gamma}} \langle c \mid \gamma \rangle + \varepsilon H(\gamma \mid \mu \otimes \nu)$$
$$= \inf_{\substack{\gamma \mid \varphi, \psi, f}} \sup \langle c - \varphi \oplus \psi \mid \gamma \rangle + \langle \varphi \mid \mu \rangle + \langle \psi \mid \nu \rangle + \varepsilon (\langle f \mid \gamma \rangle - \langle e^{f} \mid \mu \otimes \nu \rangle)$$

The dual problem is constructed by inverting the infimum and the supremum:

$$(\mathrm{KD}^{\varepsilon}) = \sup_{\varphi,\psi,f} \inf_{\gamma} \langle c - \varphi \oplus \psi + \varepsilon f | \gamma \rangle + \langle \varphi | \mu \rangle + \langle \psi | \nu \rangle - \varepsilon \langle e^{f} | \mu \otimes \nu \rangle)$$

One notices that the infimum is $-\infty$ unless $c - \varphi \oplus \psi + \varepsilon f = 0$, i.e. $f = \frac{\varphi \oplus \psi - c}{\varepsilon}$. This gives us the following dual formulation

$$(\mathrm{KD}^{\varepsilon}) = \sup_{\varphi \in \mathcal{C}^0(X), \psi \in \mathcal{C}^0(Y)} \mathcal{K}^{\varepsilon}(\varphi, \psi)$$

with

$$\mathcal{K}^{\varepsilon}(\varphi,\psi) = \langle \varphi | \mu \rangle + \langle \psi | \nu \rangle - \varepsilon \langle e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} | \mu \otimes \nu \rangle,$$

which is a concave maximization problem.

Remark 11 (Penalization of $\varphi \oplus \psi \leq c$). The dual of the entropy-regularized $(\mathrm{KD}^{\varepsilon})$ resembles the dual of the standard optimal transport problem, but where the hard constraint $\varphi \oplus \psi \leq c$ is replaced by a soft penalization: for small values of ε , $e^{\frac{\varphi \oplus \psi - c}{\varepsilon}}$ is small only $\varphi \oplus \psi - c$ is not much larger than zero. Lemma 26 (Weak duality). For any potentials $(\varphi, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y)$ and any transport plan $\gamma \in \Gamma(\mu, \nu)$, one has

$$\mathcal{K}^{\varepsilon}(\varphi,\psi) \geqslant \langle c | \gamma \rangle + \varepsilon H(\gamma \mid \mu \otimes \nu),$$

with equality if $\gamma = e^{\frac{\varphi + \psi - c}{\varepsilon}} \mu \otimes \nu$. In particular, weak duality $(KP^{\varepsilon}) \ge (KD^{\varepsilon})$ holds.

Proof. Denote $f = \frac{\varphi + \psi - c}{\varepsilon}$. Then,

$$\begin{split} \langle \varphi | \mu \rangle + \langle \psi | \nu \rangle - \varepsilon \langle e^{\frac{\varphi + \psi - c}{\varepsilon}} | \mu \otimes \nu \rangle &= \langle c | \gamma \rangle + \varepsilon \langle f | \gamma \rangle - \varepsilon \langle e^{f} | \mu \otimes \nu \rangle \\ &\geqslant \langle c | \gamma \rangle + \varepsilon H(\gamma \mid \mu \otimes \nu), \end{split}$$

with equality if and only if the density of γ with respect to $\mu \otimes \nu$ is e^f . \Box Lemma 27 (Optimality condition). The gradients of $\mathcal{K}^{\varepsilon}$ are given by:

$$\nabla_{\varphi} \mathcal{K}^{\varepsilon}(\varphi, \psi) = \mu - \prod_{X \#} e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} \mu \otimes \nu$$
$$\nabla_{\psi} \mathcal{K}^{\varepsilon}(\varphi, \psi) = \nu - \prod_{Y \#} e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} \mu \otimes \nu$$

Proof. We compute the first gradient, the second being similar. Let $(\varphi, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y)$ and let $v \in \mathcal{C}^0(X)$. Then,

$$\frac{1}{t}(\mathcal{K}^{\varepsilon}(\varphi+tv,\psi)-\mathcal{K}^{\varepsilon}(\varphi,\psi))=\langle v|\mu\rangle-\frac{\varepsilon}{t}\langle e^{\frac{(\varphi+tv)\oplus\psi-c}{\varepsilon}}-e^{\frac{(\varphi)\oplus\psi-c}{\varepsilon}}|\mu\otimes\nu\rangle.$$

Taking the limit as $t \to 0$, we get

$$\begin{split} \langle \nabla \mathcal{K}^{\varepsilon}(\varphi,\psi) | \rangle &= \langle v | \mu \rangle - \langle v e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} | \mu \otimes \nu \rangle \\ &= \langle v | \mu - \Pi_{X\#} e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} \mu \otimes \nu \rangle. \end{split}$$

Remark 12 (Existence of a maximizer to $(\mathrm{KD}^{\varepsilon})$ implies strong duality). If the dual problem admits a maximizer $(\varphi, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y)$, then the optimality conditions read $\Pi_{X\#}\gamma = \mu$ and $\Pi_{Y\#}\gamma = \nu$, where

$$\gamma = e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} \mu \otimes \nu.$$

Thus, by Lemma 26, we see that γ is a minimizer for the primal problem, and that strong duality holds.

Lemma 28 (Uniqueness of maximizer up to a constant). If (φ^*, ψ^*) is a maximizer of $(\mathrm{KD}^{\varepsilon})$, then for any other maximizer (φ, ψ) of $(\mathrm{KD}^{\varepsilon})$, there exists a constant C such that

 $\varphi = \varphi^* + C \ \mu\text{-a.e.}, \qquad \psi = \psi^* - C \ \nu\text{-a.e.}.$

Proof. Let φ, ψ be another maximizer of (KD^{ε}) , and let

$$\varphi' = \frac{1}{2}\varphi + \frac{1}{2}\varphi^*, \quad \psi' = \frac{1}{2}\psi + \frac{1}{2}\psi^*.$$

Then, by optimality of (φ, ψ) and (φ^*, ψ^*) , we have

$$0 \ge \mathcal{K}^{\varepsilon}(\varphi',\psi') - \frac{1}{2}\mathcal{K}^{\varepsilon}(\varphi',\psi') - \frac{1}{2}\mathcal{K}^{\varepsilon}(\varphi^{*},\psi^{*})$$
$$= -\int \left(e^{\frac{\varphi'\oplus\psi'-c}{\varepsilon}} - \frac{1}{2}e^{\frac{\varphi^{*}\oplus\psi^{*}-c}{\varepsilon}} - e^{\frac{1}{2}\frac{\varphi\oplus\psi-c}{\varepsilon}}\right) \mathrm{d}\mu \otimes \nu$$

By strong convexity of $t \mapsto e^t$, this is possible if and only if $\varphi' \oplus \psi' = \varphi^* \oplus \psi^*$ $\mu \otimes \nu$ -almost everywhere. Now, choose $x^* \in \operatorname{spt}\mu$, and define $C = \langle \varphi - \varphi^* | \mu \rangle$. Then, expanding the square in the following expression and using Fubini's theorem, we obtain

$$0 = \int (\varphi^* \oplus \psi^* - \varphi \oplus \psi)^2 d\mu \otimes \nu$$

=
$$\int (\varphi^*(x) - \varphi(x) - C + \psi^*(y) - \psi(y) + C)^2 d(\mu \otimes \nu)$$

=
$$\int (\varphi^*(x) - \varphi(x) - C)^2 d\mu(x) + \int (\psi^*(y) - \psi(y) + C)^2 d\nu(y) \qquad \Box$$

5.3. Existence of a solution to the dual. We now prove the existence of a solution to the dual problem. As in optimal transport the trick is to prove that the maximum can be taken over a compact subset of $\mathcal{C}^0(X) \times \mathcal{C}^0(Y)$, where the potentials are uniformly continuous. This is obtained by taking the maximum with respect to one of the two variables only. For instance, let $\psi \in \mathcal{C}^0(Y)$. Then, the maximum of $\mathcal{K}^{\varepsilon}(\cdot, \psi)$ is attained for some φ satisfying

$$\nabla_{\varphi} \mathcal{K}^{\varepsilon}(\varphi, \psi) = 0 = \mu - \prod_{X \#} e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} \mu \otimes \nu.$$

A sufficient condition is that for μ -almost every $x \in X$,

$$1 = \int_{Y} e^{\frac{\varphi(x) + \psi(y) - c(x,y)}{\varepsilon}} \mathrm{d}\nu(y) = e^{\frac{\varphi(x)}{\varepsilon}} \langle e^{\frac{\psi - c(x,\cdot)}{\varepsilon}} | \nu \rangle.$$

Definition 18 ((c, ε)-Transform). We define the (c, ε)-transform of $\psi \in C^0(Y)$ and the ($\overline{c}, \varepsilon$)-transform of $\varphi \in C^0(X)$ by

$$\psi^{c,\varepsilon}(x) = -\varepsilon \log\left(\langle e^{\frac{\psi-c(x,\cdot)}{\varepsilon}} | \nu \rangle\right)$$

$$\varphi^{\overline{c},\varepsilon}(y) = -\varepsilon \log\left(\langle e^{\frac{\varphi-c(\cdot,y)}{\varepsilon}} | \mu \rangle\right)$$

(5.22)

Remark 13 (Convergence to the c-transform as $\varepsilon \to 0$). Bounding the term in the exponential in the integral defining $\psi^{c,\varepsilon}$ from below, one clearly sees

$$\psi^{c,\varepsilon}(x) \leq \min_{y \in \operatorname{spt}(\nu)} c(x,y) - \psi(y).$$
(5.23)

On the other hand, by definition of the support of ν and by continuity of $c(x, y) - \psi(y)$, for any $\eta > 0$ there exists a measurable set $A \subseteq \operatorname{spt}(A)$ with $\nu(A) > 0$ and such that

$$\forall z \in A, \ c(x,z) - \psi(z) \leq \min_{y \in \operatorname{spt}(\nu)} c(x,y) - \psi(y) + \eta = \eta$$

Then,

$$\begin{split} \psi^{c,\varepsilon}(x) &\ge -\varepsilon \log \left(\int_A e^{\frac{\psi(z) - c(x,z)}{\varepsilon}} \mathrm{d}\nu(z) \right) \\ &\ge -\varepsilon \log \left(\int_A e^{\frac{\min_{y \in \operatorname{spt}(\nu)} c(x,y) - \psi(y) + \eta}{\varepsilon}} \mathrm{d}\nu(z) \right) \\ &\ge \min_{y \in \operatorname{spt}(\nu)} c(x,y) - \psi(y) + \eta - \varepsilon \log \nu(A) \end{split}$$

Thus, $\liminf_{\varepsilon \to 0} \psi^{c,\varepsilon}(x) \ge \min_{y \in \operatorname{spt}(\nu)} c(x, y) - \psi(y) + \eta$. Since this holds for all $\eta > 0$, we deduce with (5.23) that if $\operatorname{spt}(\nu) = Y$, then

$$\lim_{\varepsilon \to 0} \psi^{c,\varepsilon}(x) = \psi^c(x)$$

Lemma 29 (Modulus of continuity). For any $(\varphi, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y)$, the transforms $\psi^{c,\varepsilon}$ and $\varphi^{\overline{c},\varepsilon}$ have the same modulus of continuity as the cost c.

Proof. We only prove this property for $\psi^{c,\varepsilon}$, denoting ω_c the continuity modulus of the cost c:

$$\begin{split} \psi^{c,\varepsilon}(x') - \psi^{c,\varepsilon}(x) &= \varepsilon \left(\log \left(\langle e^{\frac{\psi - c(x,\cdot)}{\varepsilon}} | \nu \rangle \right) - \log \left(\langle e^{\frac{\psi - c(x',\cdot)}{\varepsilon}} | \nu \rangle \right) \right) \\ &= \varepsilon \left(\log \left(\langle e^{\frac{\psi - c(x',\cdot)}{\varepsilon}} e^{\frac{c(x',\cdot) - c(x,\cdot)}{\varepsilon}} | \nu \rangle \right) - \log \left(\langle e^{\frac{\psi - c(x',\cdot)}{\varepsilon}} | \nu \rangle \right) \right) \\ &\leqslant \varepsilon \left(\log \left(\langle e^{\frac{\psi - c(x',\cdot)}{\varepsilon}} e^{\frac{\omega c(d_X(x,x'))}{\varepsilon}} | \nu \rangle \right) - \log \left(\langle e^{\frac{\psi - c(x',\cdot)}{\varepsilon}} | \nu \rangle \right) \right) \\ &\leqslant \omega_c (d_X(x,x')). \end{split}$$

Corollary 30 (Existence of solution to $(\mathrm{KD}^{\varepsilon})$). The supremum in the definition of $(\mathrm{KD}^{\varepsilon})$ is attained for a couple $(\varphi, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(Y)$ such that

φ, ψ have the same continuity modulus as c,
⟨ψ|ν⟩ = 0

• $\langle \psi | \nu \rangle = 0$

Then, $(KP^{\varepsilon}) = (KD^{\varepsilon})$ and the unique solution to (KP^{ε}) is given by

$$\gamma = e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} \mu \otimes \nu$$

Proof. We note that by definition of the (c, ε) and $(\overline{c}, \varepsilon)$ -transforms,

$$\begin{split} \sup_{\varphi,\psi} \mathcal{K}^{\varepsilon}(\varphi,\psi) &= \sup_{\psi} \mathcal{K}^{\varepsilon}(\psi^{\overline{c},\varepsilon},\psi) \\ &= \sup_{\psi} \mathcal{K}^{\varepsilon}(\psi^{\overline{c},\varepsilon},(\psi^{\overline{c},\varepsilon})^{c,\varepsilon}) \\ &= \sup_{\psi} \mathcal{K}^{\varepsilon}(((\psi^{\overline{c},\varepsilon})^{c,\varepsilon})^{\overline{c},\varepsilon},(\psi^{\overline{c},\varepsilon})^{c,\varepsilon}) \\ &= \sup_{\psi \in \mathcal{C}^{0,\omega_c}(X)} \mathcal{K}^{\varepsilon}(\psi^{\overline{c},\varepsilon},\psi), \end{split}$$

where $\mathcal{C}^{0,\omega}(X)$ denotes the space of continuous functions with continuity modulus ω . Since for any constant $C \in \mathbb{R}$, one has $\mathcal{K}^{\varepsilon}(\varphi + C, \psi - C) = \mathcal{K}^{\varepsilon}(\varphi, \psi)$, we may impose without loss of generality that $\langle \psi | \nu \rangle = 0$ in the optimization problem. Thus,

$$(\mathrm{KD}^{\varepsilon}) = \sup_{\psi \in \mathcal{C}^{0,\omega_c}(Y) | \langle \psi | \nu \rangle = 0} \mathcal{K}^{\varepsilon}(\psi^{\overline{c},\varepsilon},\psi).$$

Since ψ belongs to $\mathcal{C}^{0,\omega_c}(Y)$, we have

$$\operatorname{osc}(\psi) := \max_{Y} \psi - \min_{Y} \psi \leqslant \operatorname{osc}(c) \leqslant 2 \, \|c\|_{\infty} \, .$$

Using in addition that $\langle \psi | \nu \rangle = 0$, we get $\|\psi\|_{\infty} \leq 2 \|c\|_{\infty}$. This shows that the set

$$\{\psi \in \mathcal{C}^{0,\omega_c}(Y) \mid \langle \psi | \nu \rangle = 0\}$$

is a compact subset of $\mathcal{C}^0(Y)$. Finally, we check that $\psi \mapsto \mathcal{K}^{\varepsilon}(\psi^{\overline{c},\varepsilon},\psi)$ is continuous on this set, and we conclude by Arzelà-Ascoli's theorem that the maximum in $(\mathrm{KD}^{\varepsilon})$ is attained.

5.4. Sinkhorn algorithm as block-coordinate ascent. We study in this section the algorithm that consists in computing a maximizer to the dual problem (KD^{ε}) by optimizing the functional $\mathcal{K}^{\varepsilon}$ alternatively in φ and ψ . The iterations are defined by

$$\begin{cases} \varphi^{(k+1)} = (\psi^{(k)})^{c,\varepsilon} \\ \psi^{(k+1)} = (\varphi^{(k+1)})^{\overline{c},\varepsilon}. \end{cases}$$

$$(5.24)$$

or equivalently $\psi^{(k+1)} = S(\psi^{(k)})$ where

$$S(\psi) = (\psi^{c,\varepsilon})^{\overline{c},\varepsilon}.$$
(5.25)

Remark 14 (Fixed point). Assume that (φ, ψ) is a fixed point of the algorithm, i.e. $\varphi = \psi^{c,\varepsilon}$ and $\psi = \varphi^{\overline{c},\varepsilon}$, and denote $\gamma = e^{\frac{\varphi \oplus \psi - c}{\varepsilon}} \mu \otimes \nu$. Thus,

$$\max_{\hat{\varphi}} \mathcal{K}^{\varepsilon}(\hat{\varphi}, \psi) = \mathcal{K}^{\varepsilon}(\varphi, \psi).$$

The first-order optimality condition for this problem, $\nabla_{\varphi} \mathcal{K}^{\varepsilon}(\varphi, \psi) = 0$, implies that $\Pi_{X\#} \gamma = \mu$. Similarly, we get $\Pi_{Y\#} \gamma = \nu$, showing by Lemma 26 that (φ, ψ) maximizes $(\mathrm{KD}^{\varepsilon})$ and γ minimizes $(\mathrm{KP}^{\varepsilon})$.

Remark 15 (Relation to matrix factorization). Algorithm (5.24) is in fact a reformulation, using a logarithmic change of variable, of Sinkhorn's algorithm for finding a factorization of non-negative matrices [38]. Let X = $\{x_1, \ldots, x_N\}, Y = \{y_1, \ldots, y_M\}, c_{ij} = c(x_i, y_j), \mu = \sum_i \mu_i \delta_{x_i} \text{ and } \nu =$ $\sum_j \nu_j \delta_{y_j}$. Then, by the discussion of the previous paragraph, $\gamma = \sum_{i,j} \gamma_{ij} \delta_{ij}$ is a solution to the entropy-regularized optimal transport problem between μ and ν if there exists $\varphi \in \mathbb{R}^N$ and $\psi \in \mathbb{R}^N$ such that

$$\begin{split} \gamma_{ij} &= e^{\frac{\varphi_i + \psi_j - c_{ij}}{\varepsilon}} \\ \text{s.t.} & \begin{cases} \forall i \in \{1, \dots, N\}, \ \sum_{1 \leq j \leq N} \gamma_{ij} = \mu_i \\ \forall j \in \{1, \dots, N\}, \ \sum_{1 \leq i \leq N} \gamma_{ij} = \nu_j. \end{cases} \end{split}$$

Denote $K_{ij} = e^{-\frac{c_{ij}}{\varepsilon}}$. The iterates of Sinkhorn's algorithm are

$$\begin{cases} \varphi_i^{k+1} = -\varepsilon \log \left(\sum_j e^{\frac{\psi_j^k - c_{ij}}{\varepsilon}} \nu_j \right) \\ \psi_j^{k+1} = -\varepsilon \log \left(\sum_i e^{\frac{\varphi_i^{k+1} - c_{ij}}{\varepsilon}} \mu_i \right) \end{cases}$$
(5.26)

One may also record the transport plan γ^k induced by φ^k and ψ^k :

$$\gamma_{ij}^k = e^{\frac{\varphi_i^k + \psi_j^k - c_{ij}}{\varepsilon}} \mu_i \nu_j$$

Denoting $u_i^k = e^{\frac{\varphi^k}{\varepsilon}} \mu_i$, $v_i^k = e^{\frac{\varphi^k}{\varepsilon}} \nu_i$ and $K_{ij} = e^{\frac{-c_{ij}}{\varepsilon}}$, we may even simplify the iterations further:

$$\begin{cases} u_i^{k+1} = \mu_i / (Kv^k)_i \\ v_j^{k+1} = \nu_j / (K^T u^{k+1})_j \\ \gamma^k = \operatorname{diag}(v^k) K \operatorname{diag}(u^k), \end{cases}$$
(5.27)

where diag(x) is the square diagonal matrix with entries x_i . It is also possible to drop the variables u, v and write the iterations purely in term of γ . In practice, this is not advised because of memory requirements: the memory to store u and v is N + M while the memory to store γ is NM. In addition, the use of the variables u and v instead of φ, ψ is not advised, because the iteration (5.27) is less stable numerically than the formula (5.26) for small values of ε . In particular, for (5.26), one may use robust implementation of the LogSumExp function provided in most machine learning frameworks.

The following two properties are very similar to some properties holding for the standard *c*-transform. In the following, we denote $\|\cdot\|_{o,\infty}$ the pseudo-norm of uniform convergence up to addition of a constant:

$$||f||_{o,\infty} = \inf_{a \in \mathbb{R}} ||f+a||_{\infty} = \frac{1}{2} (\sup f - \inf f).$$

This pseudo-norm will be very useful to state convergence results for Sinkhorn's algorithm for solving the regularized optimal transport problem. We first note that the Sinkhorn map is 1-Lipschitz with respect to this norm.

Proposition 31. Let $\psi, \overline{\psi} \in \mathbb{R}^Y$. Then,

(i) for $a \in \mathbb{R}$, $(\psi + a)^{c,\varepsilon} = \psi^{c,\varepsilon} + a$. (ii) $\left\| \psi^{c,\varepsilon} - \overline{\psi}^{c,\varepsilon} \right\|_{\infty,o} \leq \left\| \psi - \overline{\psi} \right\|_{\infty,o}$.

Similar properties hold for the map $\varphi \in \mathbb{R}^X \mapsto \varphi^{\overline{c},\varepsilon}$.

Proof. (i) follows immediately from the definition

(ii) We first show that the map is 1-Lipschitz with respect to the norm of uniform convergence:

$$\begin{split} \psi^{c,\varepsilon}(x) &- \overline{\psi}^{c,\varepsilon}(x) \\ &= \varepsilon \log\left(\langle e^{\frac{\overline{\psi} - c(x,\cdot)}{\varepsilon}} |\nu\rangle \right) - \varepsilon \log\left(\langle e^{\frac{\psi - c(x,\cdot)}{\varepsilon}} |\nu\rangle \right) \\ &= \varepsilon \log\left(\langle e^{\frac{\psi - c(x,\cdot)}{\varepsilon}} e^{\frac{\overline{\psi} - \psi}{\varepsilon}} |\nu\rangle \right) - \varepsilon \log\left(\langle e^{\frac{\psi - c(x,\cdot)}{\varepsilon}} |\nu\rangle \right) \leqslant \left\| \psi - \overline{\psi} \right\|_{\infty} \end{split}$$

Taking the maximum over x leads to $\|\psi^{c,\varepsilon} - \overline{\psi}^{c,\varepsilon}\| \leq \|\psi - \overline{\psi}_{\infty}\|$. The same inequality with $\|\cdot\|_{o,\infty}$ will follow easily using (i) and the definition of the norm $\|\cdot\|_{\infty,o}$ as a minimum.

5.5. Linear convergence of Sinkhorn's algorithm. In order to prove convergence, we need to strengthen the 1-Lipschitz estimation from Proposition 31. This allows to apply Picard's fixed point theorem to get the contraction of the Sinkhorn iteration (5.25). The proof we present in this chapter has been first introduced in course notes of Vialard [41].

Theorem 32 (Convergence of Sinkhorn, [41]). The map S is a contraction for $\|\cdot\|_{o,\infty}$. More precisely,

$$\left\|S(\psi^0) - S(\psi^1)\right\|_{o,\infty} \leqslant \left(1 - e^{-2\frac{\|\varepsilon\|_{o,\infty}}{\varepsilon}}\right) \left\|\psi^0 - \psi^1\right\|_{o,\infty}.$$

In particular, the iterates $(\varphi^{(k)}, \psi^{(k)})$ of Sinkhorn's algorithm (5.24) converge with linear rate to the unique (up to constant) maximizer the regularized dual problem $(\mathbf{KP}^{\varepsilon})$.

Remark 16 (Other convergence proofs). The convergence of Sinkhorn's algorithm is usually proven (e.g. in [39]) using a theorem of Birkhoff [10]. We refer to the recent book by Peyré and Cuturi [32] for this point of view. Other convergence proofs exist, see for instance Berman [7] (in the continuous case), and Altschuler, Weed and Rigollet [1], or Carlier [] and Nutz [] for proofs relying on the strong concavity of $\mathcal{K}^{\varepsilon}$.

Remark 17 (Convergence speed). This theorem shows that the Sinkhorn algorithm converges with linear speed, but the contraction constant has a bad dependency in ε . Denoting $C = \|c\|_{o,\infty}$, to get an error of $\eta > 0$, the number of iterations must satisfy

$$(1 - e^{-2C/\varepsilon})^k \lesssim \eta$$

i.e. $k \gtrsim e^{2C/\varepsilon} \log(1/\eta),$

where the second inequality holds for small values of ε . This bad dependency in ε seems to be a practical obstacle to choosing a very small smoothing parameter. This calls for scaling techniques, as for the auction's algorithm, and was considered by Schmitzer [36, 37].

Remark 18 (Implementation). The numerical implementation of Sinkhorn's algorithm is more complicated than it seems:

- In a naive implementation, the computation of the smoothed c-transforms (5.22) has a cost proportional to Card(X)Card(Y). This can be alleviated for instance when X = Y are grids and when the cost is a ||·||_p norm, using fast convolution techniques (see e.g. [40] or [32, Remark 4.17]), or when the cost is the squared geodesic distance on a Riemannian manifold [16, 40].
- The convergence speed can be slow when the supports of the data X, Y are "far" from each other, and when ε is small. This difficulty is cirvumvented using the ε -scaling techniques mentioned above, often combined with multi-scale (coarse-to-fine) strategies, studied in this context by Benamou, Carlier and Nenna [6] and Schmitzer [36].

OPTIMAL TRANSPORT

• Finally, some numerical difficulties (divisions by zero) can occur when ε is small and the potential ψ is far from the solution.

The book of Cuturi and Peyré present these difficulties in more details and explain how to circumvent them [32]. In addition to the works already cited, we refer to the PhD work of Feydy [14, 19], and especially to the implementation of regularized optimal transport in the library GeomLoss¹.

In order to prove this theorem, we will make use of the following elementary lemma, giving an upper bound on the total variation distance between two Gibbs kernels.

Lemma 33. Let $u_0, u_1 \in \mathcal{C}^0(Y)$ and $\nu \in \mathcal{P}(Y)$. We denote $g_i = e^{u_i}/Z_i\nu$ where $Z_i = \langle e^{u_i} | \nu \rangle$. Then,

$$\forall v \in \mathcal{C}^{0}(Y), \ |\langle v|g_{1} - g_{0} \rangle| \leq 2(1 - e^{-2||u_{0} - u_{1}||_{o,\infty}}) ||v||_{o,\infty}.$$

Proof. Note that by definition the Gibbs kernel g_i does not change if a constant is added to u_i , so that we can assume that

$$\varepsilon := \|u_0 - u_1\|_{o,\infty} = \|u_0 - u_1\|_{\infty}$$

Using the inequality $u_0 - \varepsilon \leq u_1 \leq u_0 + \varepsilon$, one easily shows that

$$e^{u_0-\varepsilon} \leqslant e^{u_1} \leqslant e^{u_0+\varepsilon}$$

Integrating this inequality multiplied by ν , this implies that

$$e^{-\varepsilon}Z_0 \leqslant Z_1 \leqslant e^{\varepsilon}Z_0$$
, i.e. $e^{-\varepsilon}\frac{1}{Z_0} \leqslant \frac{1}{Z_1} \leqslant e^{\varepsilon}\frac{1}{Z_0}$.

Multiplying this last inequality with the first one, we get

$$e^{-2\varepsilon} \frac{e^{u_0}}{Z_0} \leqslant \frac{e^{u_1}}{Z_1} \leqslant e^{2\varepsilon} \frac{e^{u_0}}{Z_0}.$$

Let $v \in \mathcal{C}^0(Y)$ be non-negative. Then,

$$e^{-2\varepsilon}\langle v|g_0\rangle \leqslant \langle v|g_1\rangle \leqslant e^{2\varepsilon}\langle v|g_0\rangle,$$

thus implying

$$|\langle v|g_1 - g_0 \rangle| \leq (1 - e^{-2\varepsilon}) \max(\langle v|g_0 \rangle, \langle v|g_1 \rangle) \leq (1 - e^{-2\varepsilon}) \|v\|_{\infty}.$$

If v is not positive, we apply the previous inequality to $\hat{v} = v - \min_Y v \ge 0$, remarking that $\|\hat{v}\|_{\infty} = 2 \|v\|_{\infty,o}$.

Proof of Theorem 32. Consider $\psi_0, \psi_1 \in \mathcal{C}^0(Y)$. We will first give an upper bound on $\|\psi_1^{c,\varepsilon} - \psi_0^{c,\varepsilon}\|_{a,\infty}$, and to do that we will give an upper bound on

$$A(x,x') = (\psi_1^{c,\varepsilon}(x) - \psi_0^{c,\varepsilon}(x)) - (\psi_1^{c,\varepsilon}(x') - \psi_0^{c,\varepsilon}(x'))$$

which is independent of $x, x' \in X$. For this purpose, we introduce $\psi_t = \psi_0 + tv$ with $v = \psi_1 - \psi_0$, and

$$B(t, x, x') = \psi_t^{c,\varepsilon}(x) - \psi_t^{c,\varepsilon}(x')$$
$$= \varepsilon \log\left(\langle e^{\frac{\psi_t - c(x', \cdot)}{\varepsilon}} | \nu \rangle\right) - \varepsilon \log\left(\langle e^{\frac{\psi_t - c(x', \cdot)}{\varepsilon}} | \nu \rangle\right)$$

¹https://www.kernel-operations.io/geomloss/

Then,

$$\partial_t B(t, x, x') = \langle v | g_{x,t} - g_{x',t} \rangle, \text{ with } g_{x,t} = \frac{e^{\frac{\psi_t - c(x', \cdot)}{\varepsilon}} \nu}{\langle e^{\frac{\psi_t - c(x', \cdot)}{\varepsilon}} | \nu \rangle}.$$

Lemma 33 directly gives us

$$\left|\partial_t B(t, x, x')\right| \leq 2(1 - e^{-2\|c(x', \cdot) - c(x, \cdot)\|_{\infty}}) \|v\|_{\infty, o}$$

We therefore get

$$|A(x,x')| \leq \int_0^1 |\partial_t B(t,x,x')| \leq 2(1 - e^{-2||c||_{\infty,o}}) ||\psi_1 - \psi_0||_{\infty,o}.$$

Taking the supremum over $x, x' \in X$, we obtain

$$\left\|\psi_{1}^{c,\varepsilon}-\psi_{0}^{c,\varepsilon}\right\|_{o,\infty}=\frac{1}{2}\max_{x,x'}\left|A(x,x')\right|\leqslant\left(1-e^{-2\frac{\|\varepsilon\|_{o,\infty}}{\varepsilon}}\right)\left\|\psi_{1}-\psi_{0}\right\|_{o,\infty}.$$

We conclude the proof of the contraction inequality by remarking that the map $\varphi \mapsto \varphi^{\overline{c},\varepsilon}$ is 1-Lipschitz, thanks to Proposition 31.(ii).

6. WASSERSTEIN DISTANCES

6.1. p-Wasserstein spaces over compact metric spaces.

Definition 19 (Wassertein distance). Let (X, d_X) be a compact metric space and $p \ge 1$. The Wasserstein distance between two probability measures $\mu, \nu \in \mathcal{P}(X)$ is defined as

$$W_p(\mu,\nu) = \left(\min_{\gamma \in \Gamma(\mu,\nu)} \langle c_p | \gamma \rangle\right)^{1/p}, \quad c_p(x,y) := d_X(x,y)^p \tag{6.28}$$

Theorem 34 (Kantorovich-Rubinstein). The Wasserstein-1 distances admit the following formulation:

$$W_1(\mu,\nu) = \sup\left\{ \langle f|\mu \rangle - \langle f|\nu \rangle \mid f \in \mathcal{C}^0(X), \ \operatorname{Lip}(f) \leq 1 \right\}.$$
 (6.29)

Proof. Note that for $c = d_X$, $\psi^c(x) = \min_{y \in X} d(x, y) - \psi(y)$ is 1-Lipschitz as a infimum of 1-Lipschitz functions. This implies that the dual problem may be rewritten as

$$\min_{\psi \in \mathcal{C}^0(X) | \operatorname{Lip}(\psi) \leqslant 1} \langle \psi^c | \mu \rangle + \langle \psi | \nu \rangle.$$

If ψ is 1-Lipschitz, then $d(x, y) - \psi(y) \ge -\psi(x)$, so that

$$\psi^c(x) = \inf_y d(x, y) - \psi(y) = -\psi(x)$$

Maximal correlation?

Theorem 35 (Properties of W_p). The following properties hold: (i) $W_1 \leq W_p$ for all $p \geq 1$, (ii) W_p is a distance on $\mathcal{P}(X)$, (iii) W_p metrizes weak convergence. *Proof.* (i) The first claim is a consequence of the Jensen's inequality. (ii) To prove the second claim, we note that the stability of optimal transport plans (Theorem 14) implies in particular that the Wasserstein distances W_p^p are weak* continuous with respect to their arguments. To establish the triangle inequality, we let $\mu, \nu, \sigma \in \mathcal{P}(X)$ and we consider empirical measures

$$\mu_N = \frac{1}{N} \sum_{i=1}^N \delta_{x_N^i}, \quad \nu_N = \frac{1}{N} \sum_{i=1}^N \delta_{y_N^i}, \quad \sigma_N = \frac{1}{N} \sum_{i=1}^N \delta_{z_N^i}.$$

converging weak* to μ, ν and σ respectively. Without loss of generality, we can reorder the points so that the optimal transport map between μ_N and ν_N is given by $x_N^i \to y_N^i$, and that the optimal transport map between ν_N and σ_N is $y_N^i \to z_N^i$. Then,

$$W_{p}(\mu_{N},\sigma_{N}) \leq \left(\frac{1}{N}\sum_{1\leq i\leq N} \left\|x_{N}^{i}-z_{N}^{i}\right\|^{p}\right)^{1/p} \\ \leq \left(\frac{1}{N}\sum_{1\leq i\leq N} \left\|x_{N}^{i}-y_{N}^{i}\right\|^{p}\right)^{1/p} + \left(\frac{1}{N}\sum_{1\leq i\leq N} \left\|y_{N}^{i}-z_{N}^{i}\right\|^{p}\right)^{1/p} \\ = W_{p}(\mu_{N},\nu_{N}) + W_{p}(\nu_{N},\sigma_{N})$$

We conclude by taking the limit $N \to +\infty$.

(iii) Since $W_1 \leq W_p$, if a sequence (μ_n) converges to μ with respect to W_p , then it also converges to μ with respect to W_1 . Kantorovich-Rubinstein's formula then implies that for any function $f \in \mathcal{C}^0(X)$ with $\operatorname{Lip}(f) \leq 1$ one has $\lim_{n \to +\infty} \int f d\mu_n = \int f d\mu$, thus proving weak* convergence of (μ_n) towards μ as $n \to +\infty$. Conversely, if μ_n converges weak* to μ , then by the weak* continuity of W_p^p we get

$$\lim_{n \to +\infty} W_p(\mu_n, \mu) = W_1(\mu, \mu) = 0.$$

Theorem 36 (Subdifferential of W_p^p). Let $\mu \in \mathcal{P}(X)$. The function $F = W_p^p(\mu, \cdot)$ is convex and continuous in $\mathcal{P}(X) \times \mathcal{P}(X)$. Its subdifferential is given by

$$\partial F(\nu) = \left\{ \psi \in \mathcal{C}^0(X) \mid \langle \psi^c | \mu \rangle + \langle \psi | \nu \rangle = W_p^p(\mu, \nu) \right\}.$$

In particular, if the dual problem $\max_{\psi} \langle \psi^c | \mu \rangle + \langle \psi | \nu \rangle$ has a unique solution ψ up to an additive constant, then for any measure $\nu' \in \mathcal{P}(X)$ one has

$$\frac{\mathrm{d}}{\mathrm{d}t} F(\nu + t(\nu' - \nu))\big|_{t=0} = \langle \psi | \nu' - \nu \rangle.$$

Proof. Let $(\psi^c, \psi) \in \mathcal{C}^0(X) \times \mathcal{C}^0(X)$ be a maximizer of the dual Kantorovich problem. Then, for all measures $\nu' \in \mathcal{P}(X) \times \mathcal{P}(X)$ one has

$$F(\nu') = W_p^p(\mu, \nu') \ge \langle \psi^c | \mu \rangle + \langle \psi | \nu' \rangle$$

= $\langle \psi^c | \mu \rangle + \langle \psi | \nu \rangle + \langle \psi | \nu' - \nu \rangle$
= $F(\nu) + \langle \psi | \nu' - \nu \rangle$,

thus showing that ψ belong to $\partial F(\nu)$. To prove the converse, we introduce $\tilde{\mathcal{K}}_{\mu}(\psi) = -\int \psi^{c} d\mu$. Then,

$$\tilde{\mathcal{K}}^*_{\mu}(\psi) = \sup_{\nu \in \mathcal{P}(X)} \langle \nu | \psi \rangle + \langle \mu | \psi^c \rangle = \mathbf{W}^p_p(\mu, \nu) = F(\nu).$$

By subdifferential calculus, we have

$$\psi \in \partial F(\nu) \iff \nu \in \partial F^*(\psi) = \partial \tilde{\mathcal{K}}_{\mu}(\psi)$$
$$\iff \psi \in \arg \max (\mathrm{KD}),$$

where the last equivalence comes from Proposition 16.

Remark 19 (Horizontal perturbations in the discrete case). For simplicity, assume that $\mu = \frac{1}{N} \sum_{i} \delta_{x_i}$ and $\nu = \frac{1}{N} \sum_{i} \delta_{y_i}$ and that there exists unique optimal transport maps $S : \mu \to \nu$ and $T : \nu \to \nu$ (which are thus inverse of each other). Let ξ be a smooth and compactly supported vector field. Then, for small values of t, the map (id $+ t\xi$) $\circ S$ is optimal between μ and $\nu_t = (id + t\xi)_{\#}\nu$. Thus,

$$W_p^p(\nu_t,\mu) = \int \|y - (\mathrm{id} + t\xi) \circ T(y)\|^p \,\mathrm{d}\mu(y),$$

directly implying that

$$\begin{aligned} \frac{\mathrm{d}}{\mathrm{d}t} \mathbf{W}_p^p(\nu_t, \mu) &= \int \frac{\mathrm{d}}{\mathrm{d}t} \| y - (\mathrm{id} + t\xi) \circ S(y) \|^p \,\mathrm{d}\mu(y), \\ &= \int p \, \| y - S(y) \|^{p-2} \,\langle \xi \circ S(y) | S(y) - y \rangle \mathrm{d}\mu(y), \\ &= \int p \, \| T(x) - x \|^{p-2} \,\langle \xi(x) | x - T(x) \rangle \mathrm{d}\nu(x) \end{aligned}$$

Letting T be the optimal transport map between μ . More concretely, if we denote

$$\hat{F}: (z_1, \dots, z_N) \in \mathbb{R}^{dN} \mapsto F(\frac{1}{N} \sum_i \delta_{z_i}, \nu),$$

then the previous computation shows that

$$\nabla_{z_i} \hat{F}(x_1, \dots, x_N) = \frac{p}{N} \|T(x_i) - x_i\|^{p-2} (x_i - T(x_i)).$$

6.2. *p*-Wasserstein geodesics on \mathbb{R}^d . In this subsection, we provide a short introduction to the geometry of the Wasserstein space on \mathbb{R}^d . We refer to [2] for a more complete exposition.

Definition 20 (Geodesic). In a metric space (X, d_X) , a curve $\omega : [0, 1] \to X$ is called a constant speed geodesic if

$$\forall s, t \in [0, 1], d_X(\omega_s, \omega_t) \leq |t - s| d_X(\omega_0, \omega_1).$$

A space is called *geodesic* if any pair of points in X is joined by a geodesic.

Remark 20. Let ω be a constant speed geodesic and assume that $s \leq t$. Then, the triangle inequality gives us

$$d_X(\omega_0, \omega_1) \leq d_X(\omega_0, \omega_s) + d_X(\omega_s, \omega_t) + d_X(\omega_t, \omega_1)$$
$$\leq ((1-t) + (s-t) + (1-t))d_X(\omega_0, \omega_1)$$
$$\leq d_X(\omega_0, \omega_1).$$

Thus, all inequalities must in fact be equalities, showing in particular that

 $d_X(\omega_s, \omega_t) = |t - s| d_X(\omega_0, \omega_1).$

Theorem 37 (Geodesics in W_p). Let X be a convex subset of \mathbb{R}^d , let $\mu_0, \mu_1 \in \mathcal{P}(X)$ and let $\gamma \in \Gamma(\mu, \nu)$ be an optimal transport plan for the cost $c_p(x, y) = ||x - y||^p$. Then, the curve $t \in [0, 1] \mapsto \mu_t \in \mathcal{P}(X)$ is a constant speed geodesic between μ_0 and μ_1 , with

 $\mu_t = P_{t\#}\gamma, \quad where \ P_t : (x, y) \mapsto (1 - t)x + ty$

Moreover, all constant speed geodesics between μ_0 and μ_1 are of this form. In particular, if μ_0 or μ_1 are absolutely continuous with respect to the Lebesgue measure, then the geodesic between μ_0 and μ_1 is unique.

Example 6 (Geodesics when a transport map exists). If there exists an optimal transport map T between μ_0 and μ_1 , then the geodesic defined above is $\mu_t = ((1-t)id + tT)_{\#}\mu_0$. In the discrete case, if

$$\mu_0 = \frac{1}{N} \sum_{1 \leqslant i \leqslant N} \delta_{x_0^i} \text{ and } \mu_1 = \frac{1}{N} \sum_{1 \leqslant i \leqslant N} \delta_{x_0^i}$$

are two empirical measures, and if the points are ordered such that

$$W_p^p(\mu_0, \mu_1) = \frac{1}{N} \sum_{1 \le i \le N} \|x_1^i - x_0^i\|^p,$$

a geodesic between μ_0 and μ_1 is given by

$$\mu_t = \frac{1}{N} \sum_{x \in X_0} \delta_{(1-t)x_0^i + tx_1^i}.$$

Thus, μ_t provides an interpolation between the supports of μ_0 and μ_1 .

Remark 21 (Many geodesics). It is quite easy to construct examples of measures μ_0 and μ_1 such that there exists more than one transport map between μ_0 and μ_1 . For instance, take $\mu_0 = \frac{1}{N} \sum_i \delta_{(i/N,0)}$ and $\mu_1 = \frac{1}{N} \sum_i \delta_{(0,i/N)}$. Then, every bijection between the supports of μ_0 and μ_1 is optimal for p = 2, and therefore there exists an countably infinite number of geodesics between μ_0 and μ_1 . In particular, this shows that the space ($\mathcal{P}([0,1]^2, W_2)$ cannot be embedded isometrically into any Banach space.

Proof of Theorem 37. One can observe that $\gamma_{st} = (P_s, P_t)_{\#\gamma}$ has marginals μ_s and μ_t . In particular,

$$W_p^p(\mu_s, \mu_t) \leq \int ||x_s - x_t||^p \, \mathrm{d}\gamma_{st}(x_s, x_t)$$

= $\int ||(1 - s)x + sy - (1 - t)x + sx||^p \, \mathrm{d}\gamma(x, y)$
= $(t - s)^p \int ||x - y||^o \, \mathrm{d}\gamma(x, y) = (t - s)^o \, \mathrm{W}_p^p(\mu_0, \mu_1),$

thus proving that μ_t is a constant speed geodesic.

Let us now prove that all geodesics are of this form. For every $T \in \mathbb{N}$ and any $i \in \{1, \ldots, T+1\}$, denote $\gamma_{i,i+1}^T$ an optimal transport between μ_{t_i} and $\mu_{t_{i+1}}$, with $t_i = (i-1)/T$. By the gluing lemma recalled below, there exists $\Gamma^T \in \mathcal{P}(X^{T+1})$ whose projection on (X_i, X_{i+1}) agrees with $\gamma_{i,i+1}^T$. Moreover,

$$\left(\int \|x_1 - x_{T+1}\|^2 \,\mathrm{d}\Gamma^T(x_1, \dots, x_{T+1})\right)^{1/2} \leqslant \sum_{j=0}^{T-1} \left(\int \|x_{i+1} - x_i\|^2 \,\mathrm{d}\Gamma^T(x_1, \dots, x_{T+1})\right)^{1/2}$$
$$= \sum_{j=1}^{T+1} \mathrm{W}_2(\mu_{t_i}, \mu_{t_i})$$
$$= \mathrm{W}_2(\mu_0, \mu_1)$$

This implies in particular that $\gamma^T = (\Pi_1, \Pi_{T+1})_{\#} \Gamma^T$, but also that for Γ^T almost every $x = (x_1, \ldots, x_{T+1})$, the points x_1, \ldots, x_{T+1} are aligned, i.e. $x_i = (1 - t_i)x_1 + t_i x_{T+1}$. Thus, we see that $\Gamma^T = (P_0, P_{1/T}, \ldots, P_1)_{\#} \gamma^T$ with $P_t(x, y) = (1 - t)x + ty$. In particular, we have $\mu_t = P_{t\#} \gamma^T$ for all $t \in \{0/T, \ldots, T/T\}$. One can finally check that if γ is a weak*-limit of γ , then for all $t \in [0, 1]$, one has $\mu_t = P_{t\#} \gamma$.

Lemma 38 (Gluing). Let X_1, \ldots, X_N be compact metric spaces, and for any $1 \leq i \leq N-1$ consider a transport plan $\gamma_i \in \Gamma(\mu_i, \mu_{i+1})$. Then, there exists $\gamma \in \mathcal{P}(X_1, \ldots, X_N)$ such that for all $i \in \{1, \ldots, N-1\}$, $\pi_{i,i+1}\gamma = \gamma_i$, where $\pi_{i,i+1}: X_1 \times \cdots \times X_N \to X_i \times X_{i+1}$ is the projection.

Proof. See Lemma 5.3.2 and Remark 5.3.3 in [3].

6.3. Geodesic convexity with respect to W_2 on \mathbb{R}^d .

Definition 21 (Geodesic convexity for sets). A set $S \subseteq \mathcal{P}_2^{\mathrm{ac}}(\mathbb{R}^d)$ is called geodesically convex if for any $\mu_0, \mu_1 \in S$, any W₂-geodesic between μ_0 and μ_1 remains in S.

Example 7 (Geodesically convex subsets of $(\mathcal{P}(X), W_2)$.). Example of geodesically convex subsets of $\mathcal{P}(X)$ include :

(a) the set obtained by translating and shearing a reference measure μ ,

 $\{T_{\#}\mu \mid T(x) = Ax + b, A \text{ symmetric}, A \ge 0\}$

In particular, the set of Gaussians densities is geodesically convex in $\mathcal{P}(\mathbb{R}^d)$. The restriction of the Wasserstein distance on this set can be computed in near closed-form, and called the Bures-Wasserstein metric.

- (b) the set $\mathcal{P}^{\mathrm{ac}}(X)$ of absolutely continuous measures
- (c) the set of probability densities whose density is upper bounded by a constant
- (d) the set of measures of the form $\mu = \frac{1}{N} \sum_{i} \delta_{x_i}$ (where the points x_i are not necessary distinct) is convex under *some* geodesics, namely those induced by bijections (cf Example 6.

Proposition 39. The set $\mathcal{P}^{ac}(X)$ is geodesically convex. More precisely, given $\mu_0 \in \mathcal{P}^{ac}(X)$ and $\mu_1 \in \mathcal{P}(X)$, one has $\mu_t \in \mathcal{P}^{ac}(X)$ for any $t \in [0, 1)$.

Proof. Let $\mu_0 \in \mathcal{P}^{\mathrm{ac}}(X), \mu_1 \in \mathcal{P}(\mathbb{R}^d)$ and $\varphi \in \mathrm{Lip}(X)$ be a convex function so that $\mu_t = ((1-t)\mathrm{id} + t\nabla\varphi)_{\#}\mu_0$ is the unique Wasserstein geodesic between

OPTIMAL TRANSPORT

 $\mu_0 \text{ and } \mu_1. \text{ Define } T_t = (1-t)\mathrm{id} + t\nabla\varphi. \text{ Then, for any } x, y \in \mathrm{spt}(\mu_0),$ $\langle T_t(x) - T_t(y) | x - y \rangle = (1-t) \|x - y\|^2 + t\langle \nabla\varphi(x) - \nabla\varphi(y) | x - y \rangle$ $\geq (1-t) \|x - y\|^2.$

where we used the monotonicity of the gradient of convex functions to get the inequality. In particular, if $x \neq y$ and t < 1, then $T_t(x) \neq T_t(y)$ and the inverse map T_t^{-1} is well-defined. Moreover, the same inequality shows that T_t^{-1} is Lipschitz with constant L = 1/(1-t). In addition, T_t^{-1} transports μ_t to μ_0 , i.e. $\mu_t(B) = \mu_0(T_t^{-1}(B))$ for any Borel set B. Thus, if N is Lebesgue-negligible, $T_t^{-1}(N)$ is also negligible (by the next lemma), so that $\mu_t(N) = \mu_0(T_t^{-1}(N)) = 0$. This implies that $\mu_t \ll \lambda$.

Lemma 40. If N is Lebesgue-negligible, and if S is Lipschitz, then S(N) is Lebesgue-negligible.

Definition 22 (Geodesic convexity for functions). A function $F : \mathcal{P}^{ac}(X)$ to $\mathbb{R} \cup \{+\infty\}$ is geodesically convex if and only if for any $\mu_0, \mu_1 \in \mathcal{P}^{ac}(W)$,

$$F(\mu_t) \leqslant (1-t)F(\mu_0) + tF(\mu_1)$$
 (6.30)

where (μ_t) is the W₂-geodesic. Following McCann, a geodesically convex function is often called displacement convex.

Definition 23 (Internal energy). Let $A : \mathbb{R}^+ \to \mathbb{R} \cup \{+\infty\}$. The *internal* energy associated to A generalizes Boltzmann's functional. It is defined as

$$E_A: \mu \in \mathcal{P}(X) \mapsto \begin{cases} \int_{\Omega} A(\rho(x)) \mathrm{d}x \text{ if } \mu \ll \lambda \text{ and } \rho := \frac{\mathrm{d}\mu}{\mathrm{d}\lambda} \\ +\infty \text{ if not} \end{cases}$$
(6.31)

Theorem 41 (McCann). Let $A : \mathbb{R}_+ \to \mathbb{R}_+ \cup \{+\infty\}$ be such that (i) A(0) = 0 and (ii) $r \mapsto A(r^{-d})r^d$ is convex non-increasing. Then internal energy E_A is displacement convex on $\mathcal{P}(X)$.

We will call conditions (i) and (ii) McCann's conditions. Example of functions A that satisfy such conditions include

- $A(r) = r^q$ for q > 1;
- $A(r) = r \log r;$
- $A(r) = -r^m$ for $m \in [1 1/d, 1)$.

This theorem is a corollary of the more general result below. Indeed, take $\mu_0 = \mu \in \mathcal{P}^{\mathrm{ac}}(X), \, \varphi_0 = \frac{1}{2} \|\cdot\|^2$ and φ_1 a convex function such that $T = \nabla \varphi_1$ is the optimal transport map between μ_0 and μ_1 . Then,

$$\mu_t = ((1-t)\nabla\varphi_0 + t\nabla\varphi_1)_{\#}\mu = ((1-t)\mathrm{id} + tT)_{\#}\mu_0$$

is the unique Wasserstein geodesic between μ_0 and μ_1 .

Theorem 42. Let $\mu \in \mathcal{P}^{ac}(X)$ and let $\varphi_0, \varphi_1 \in \operatorname{Lip}(X)$ be two convex functions such that $\nabla \varphi_i(X) \subseteq X$, and let $\varphi_t = (1-t)\varphi_0 + t\varphi_1$. Assume that $A : \mathbb{R}_+ \to \mathbb{R}_+$ satisfies McCann's conditions. Then

$$t \in [0,1] \mapsto E_A(\nabla \varphi_{t\#} \mu)$$

is convex.

We only prove this theorem when the functions φ_0 and φ_1 are \mathcal{C}^2 and uniformly convex, which implies that the gradients $\nabla \varphi_i$ are diffeomorphisms from X to $\nabla \varphi_i(X)$. The proof in the general case can be found in the article of McCann [26] or in Villani's first book [42].

Lemma 43. Assume that $\mu \in \mathcal{P}^{ac}(X)$ has density ρ and that $\varphi \in \mathcal{C}^2(X)$ is uniformly convex. Then

$$E_A(\nabla \varphi_{\#} \mu) = \int_{\mathbb{R}^d} A\left(\frac{\rho(x)}{\det(\mathbf{D}^2 \varphi(x))}\right) \det(\mathbf{D}^2 \varphi(x)) \mathrm{d}x.$$

Proof. Since T is a diffeomorphism, the measure $T_{\#}\mu$ is absolutely continuous with respect to the Lebesgue measure. We denote σ the density of $T_{\#}\mu$, which satisfies

$$\sigma(T(x))\det(\mathrm{D}T(x)) = \rho(x)$$

Moreover, by the change of variable formula y = T(x) and using det $(DT(x)) = |\det DT(x)|$, which follows from the convexity of T, we get

$$E_A(\nabla \varphi_{\#} \mu) = \int A(\sigma(y)) dy$$

= $\int A(\sigma(T(x))) \det(DT(x)) dx$
= $\int A\left(\frac{\rho(x)}{\det(DT(x))}\right) \det(DT(x)) dx$

Lemma 44. The map $M \mapsto \det(M)^{1/d}$ is concave over the set of symmetric positive d-by-d matrices.

Proof. Recall Hadamard's formula for a symmetric positive matrix M:

$$\det(M) = \min_{e_1,\dots,e_d \text{ orthonormal}} \langle e_1 | M e_1 \rangle \cdots \langle e_d | M e_d \rangle,$$

where the minimum is taken over orthonormal bases. Given a fixed orthonormal basis e_1, \ldots, e_d consider $f(M) = (\langle e_1 | M e_1 \rangle \cdots \langle e_d | M e_d \rangle)^{1/d}$. Then f is concave over the set of matrices M satisfying $\langle e_i | M e_i \rangle \ge 0$ as the composition of the geometric mean $(x \in (\mathbb{R}^+)^d \mapsto (x_1 \cdots x_d)^{1/d})$ with linear functions. Then, $\det(\cdot)^{1/d}$ is concave over the set of symmetric positive matrices, as a minimum of concave functions.

Proof of Theorem 42. If φ_0, φ_1 are \mathcal{C}^2 and uniformly convex, the interpolant $\varphi_t := (1-t)\varphi_0 + t\varphi_1$ is also \mathcal{C}^2 and uniformly convex. Hence, by Lemma 43,

$$E_A(\nabla \varphi_{t\#}\mu) = \int_X B(D(x,t))\rho(x) \mathrm{d}x,$$

where we have set $B(r) = A(r^{-d})r^d$ and $D(x,t) = (\det(\mathbb{D}^2\varphi_t(x))/\rho(x))^{1/d}$. By Lemma 44, for all $x \in X, t \in [0,1] \mapsto D(x,t)$ is concave so that

$$D(x,t) \ge (1-t)D(x,0) + tD(x,1).$$

Hence, since B is non-decreasing and convex,

$$B(D(x,t)) \leq B((1-t)D(x,0) + tD(x,1)) \leq (1-t)B(D(x,0)) + tB(D(x,1)).$$

Integrating this inequality gives the desired convexity result. \Box

Remark 22 (Displacement convexity in the "linear" case). Assume that $T_0(x) = x$ and $T_1(x) = M \cdot x$ where M is a fixed symmetric positive definite matrix, and $\rho \in \mathcal{P}^{\mathrm{ac}}(\Omega)$. Then, the 2-Wasserstein geodesic between $\rho_0 = \rho$ and $\rho_1 = T_{1\#}\rho$ is given by $\rho_t = T_{t\#}$, where $T_t = M_t \cdot x$, with $M_t = (1-t)\mathrm{Id} + tM$. The density of ρ_t satisfies $\rho_t(T_t(x)) \det DT_t(x) = \rho(x)$, where $\det(DT_t(x)) = \det(M_t)$ does not depend on x. Then,

$$E_A(\rho_t) = \int A(\rho_t(y)) dy$$

= $\int A(\rho_t(T_t(x))) \det(M_t) dx$
= $\int A\left(\frac{\rho(x)}{\det(M_t)}\right) \det(M_t) dx$

Since $A(r) = r \ln(r)$, A(r/s)s = rln(r/s) = r(ln(r) - ln(s)). Thus,

$$E_A(\rho_t) = \int A(\rho(x)) dx - \log(\det(M_t))\rho(x)$$

Using that $M \mapsto \log \circ \det(M)$ is concave on the set of symmetric positive definite matrices, we conclude that $t \mapsto E_A(\rho_t)$ is convex.

Corollary 45 (Brunn-Minkowski's inequality). Let K_0, K_1 be compact subsets of X, and let $K_t = (1-t)K_0 + tK_1$. Then, $\lambda(K_t)^{1/d} \ge (1-t)\lambda(K_0)^{1/d} + t\lambda(K_1)^{1/d}$.

Proof. Assume that $\lambda(K_0), \lambda(K_1) > 0$. Let $\mu_i = \lambda|_{K_i} / \lambda(K_i)$, let μ_t be the geodesic between μ_0 and μ_1 . Then μ_t is absolutely continuous, with density ρ_t , and supported on K_t . The convexity of $A(r) = -r^{1-1/d}$ and Jensen's inequality implies

$$\int_{K_t} A(\rho_t(x)) d\lambda(x) = \lambda(K_t) \int_{K_t} A(\rho_t(x)) \frac{d\lambda(x)}{\lambda(K_t)}$$
$$\leqslant \lambda(K_t) A\left(\int_{K_t} \rho_t(x) \frac{d\lambda(x)}{\lambda(K_t)}\right)$$
$$= \lambda(K_t) A(1/\lambda(K_t)) = -\lambda(K_t)^{1/d}$$

Moreover, for t = 0 and t = 1 we get

$$\int_{K_i} A(\rho_i(x)) \mathrm{d}\lambda(x) = \int_{K_i} \lambda(K_i) A(1/\lambda(K_i)) = -\lambda(K_i)^{1/d}$$

7. QUANTIZATION AND UNIFORM QUANTIZATION OF MEASURES

8. Embedding of the Wasserstein space

We recall that there exists an explicit isometric embedding of $(\mathcal{P}_p(\mathbb{R}), W_p)$ into the space $L^p([0, 1])$, which maps $\mu \in \mathcal{P}(\mathbb{R})$ to its quantile function $T_{\mu} \in L^p([0, 1])$, i.e. the unique non-decreasing map which transports $\lambda_{[0,1]}$ onto ν :

$$\forall \mu, \nu \in \mathcal{P}_p(\mathbb{R}), \quad \mathbf{W}_p(\mu, \nu) = \|T_\mu - T_\nu\|_{\mathbf{L}^p([0,1])}.$$

This embedding is practically very useful, because it allows to simplify many constructions in the 1D Wasserstein space (geodesics, barycenters, etc). However, we already saw (Remark 21) that in $(\mathcal{P}_p(\mathbb{R}^d), W_p)$, there may exist several geodesics between two probability measures, preventing an isometric embedding of this space into a Banach space. A natural question is whether (for instance) there may exist bi-Lipschitz or bi-Hölder embeddings of $(\mathcal{P}_p(\mathbb{R}^d), \mathbb{W}_p)$ into Banach spaces, and we will see that the answer is mostly negative.

8.1. Non-embeddability results.

Definition 24 (Coarse embedding). Let X, Y be metric spaces. A function $f: X \to Y$ is a *coarse embedding* if there exists a non-decreasing functions $\rho_{\pm}: \mathbb{R}_+ \to \mathbb{R}_+$ satisfying

- $\rho_-(d_X(x,y)) \leq d_Y(f(x), f(y)) \leq \rho_+(d_X(x,y))$ $\lim_{t \to +\infty} \rho_-(t) = +\infty.$

In particular, if f is bi-Lipschitz or bi-Hölder (or even uniformly continuous and with uniformly continuous inverse)), the embedding is coarse. However, note that the definition does not imply that f is continuous. The notion of coarse embedding is therefore extremely weak, so that theorem establishing the impossibility of coarse embeddings are usually strong and difficult theorems. Our aim with this sections is not to present a comprehensive review of the literature on embedding metric spaces into Banach spaces, but rather to present some striking impossibility results when the source space is a Wasserstein space over \mathbb{R}^d with d > 1.

The first result is about coarse embeddability into a Hilbert space, and was originally proven by Wagner in the context of persistence diagrams in computational topology [44]. We provide here an (easy) adaptation of the arguments of Wagner to the Wasserstein setting.

Theorem 46 (Wagner). Let p > 2. Then, there is no coarse embedding of the Wasserstein space $(\mathcal{P}_p(\mathbb{R}^d), W_p)$ into a Hilbert space.

The proof of this theorem relies on a characterization of coarse embeddability of a metric space in a Hilbert space through the "uniform" coarse embeddability of finite subsets, due to P. Nowak [29] (see also [18]), and relying on Schoenberg's characterization of *isometric* embeddability into a Hilbert space.

Theorem 47 (Nowak). A metric space (X, d_X) admits a coarse embedding into a Hilbert space if and only if there exists non-decreasing functions ρ_{\pm} : $\mathbb{R}_+ \to \mathbb{R}_+$ satisfying $\lim_{t\to+\infty} \rho_-(t) = +\infty$, and such that for any finite subset $A \subseteq X$, there exists a coarse embedding $f_A: A \to \ell_2$ satisfying

 $\forall x, y \in A, \quad \rho_{-}(\mathbf{d}_X(x, y)) \leq \|f_A(x) - f_A(y)\| \leq \rho_{+}(\mathbf{d}_X(x, y)).$

Nowak's theorem was used for instance to prove that the space $\ell^p = \{a \in$ $\mathbb{R}^{\mathbb{N}} \mid \sum_{i} |a_{i}|^{p} < +\infty$ with p > 2 cannot be coarsely embedded into a Hilbert space [22]. We will use this result to prove that (\mathcal{P}_{p}, W_{p}) also cannot be embedded into a Hilbert space.

Lemma 48. Given any $p \ge 1$, $R \in \mathbb{R}_+$ and $N \in \mathbb{N}$, the map

$$R_{N,N} : ([-R,R]^N, \|\cdot\|_p) \to (\mathcal{P}_p(\mathbb{R}^2), \mathbb{W}_p)$$
$$a = (a_1, \dots, a_N) \mapsto \frac{1}{N} \sum_{i=1}^N \delta_{2N^{1/p}Ri, N^{1/p}a_i}$$

is an isometry.

Φ

Proof. Let $a, b \in [-R, R]^N$ and denote $p_i(a) = (2N^{1/p}Ri, N^{1/p}a_i) \in \mathbb{R}^2$. Then, for all $i \neq j$, one has

$$||p_i(a) - p_j(b)|| \ge 2N^{1/p}R |i - j| \ge 2RN^{1/p}.$$

On the other hand,

$$|p_i(a) - p_i(b)|| = N^{1/p} |a_i - b_i| \leq 2RN^{1/p}$$

Thus, the optimal transport map between $\Phi_{R,N}(a)$ and $\Phi_{R,N}(b)$ simply maps the point $p_i(a)$ to $p_i(b)$, so that

$$W_{p}(\Phi_{R,N}(a), \Phi_{R,N}(b)) = \frac{1}{N} \sum_{i=1}^{N} \|p_{i}(a) - p_{i}(b)\|^{p}$$
$$= \frac{1}{N} \sum_{i=1}^{N} (N^{1/p} |a_{i} - b_{i}|)^{p}$$
$$= \|a - b\|_{p}^{p} \qquad \Box$$

Proof of Theorem 46. Assume that $(\mathcal{P}_p(\mathbb{R}^2), W_p)$ can be embedded into a Hilbert space. Then there exists two functions ρ_{\pm} with $\lim_{t\to+\infty} \rho_{-}(t) = +\infty$ and for any finite subset S of $\mathcal{P}_p(\mathbb{R}^2)$, there exists a map $f_S: S \to \ell_2$ such that

$$\forall \mu, \nu \in S, \quad \rho_{-}(W_{p}(\mu, \nu)) \leq ||f_{S}(\mu) - f_{S}(\nu)|| \leq \rho_{+}(W_{p}(\mu, \nu)).$$
 (8.32)

We now prove using the converse of Theorem 47 that this would imply that ℓ_p can be coarsely embedded into a Hilbert space, which is known to be false [22]. Let A be a finite subset of ℓ_p . Then, there exists N > 0 such that

$$\forall a \in A, \left(\sum_{i=N+1}^{+\infty} |a_i|^p\right)^{1/p} \leq \frac{1}{2}.$$

We let $R = \max_{a \in A} ||a||_{\infty}$ and we consider $\Pi_A : A \to [-R, R]^N$ obtained by keeping only the first N coordinates of a sequence, i.e. $\Pi_A(a) = (a_i)_{1 \leq i \leq N}$. Then,

$$\forall a, b \in A, \quad ||a - b||_p - 1 \leq ||\Pi_A(a) - \Pi_A(b)|| \leq ||a - b||_p.$$

Thus, denoting $S = \Phi_{R,N}(\Pi_A(A))$ and $f_A = f_S \circ \Phi_{R,N} \circ \Pi_A$, we get for all $a, b \in A$, using (8.32) and Lemma 48,

$$\|f_{A}(a) - f_{A}(b)\| \leq \rho_{+}(\mathbb{W}_{p}(\Phi_{R,N} \circ \Pi_{A}(a), \Phi_{R,N} \circ \Pi_{A}(a)))$$

= $\rho_{+}(\|\Pi_{A}(a) - \Pi_{A}(b)\|)$
 $\leq \rho_{+}(\|a - b\|_{p}).$

Similarly, introducing $\tilde{\rho}_{-}(r) = \max(\rho_{-}(r) - 1, 0)$ we get

$$||f_A(a) - f_A(b)|| \ge \tilde{\rho}_-(||a - b||_p).$$

We can use Theorem 47 to conclude that ℓ^p can be coarsely embedded into a Hilbert space, contradicting p > 2.

Below, we report a more difficult negative result from Andoni, Naor and Neiman [4, Theorem 7]. In particular, this result shows that it is not possible construct a bi-Lipschitz or bi-Hölder embedding of $(\mathcal{P}_2(\mathbb{R}^d), W_2), d \ge 3$, into any Hilbert or L^p space.

Theorem 49 (Andoni, Naor, Neiman). For every p > 1, the space $(\mathcal{P}_p(\mathbb{R}^3), \mathbb{W}_p)$ does not admit a coarse embedding into any Banach space of nontrivial type.

8.2. Embedding via slicing. Theorems 46 and 49 show that it is impossible to coarsely embed Wasserstein spaces into a Hilbert space when $d \ge 3$ and p = 2 or $d \ge 2$ and p > 2. In the case d = 1, the map $\mu \mapsto T_{\mu}$ (quantile function) is an $(\mathcal{P}_p(\mathbb{R}), \mathbb{W}_p)$ into the space $L^p([0, 1])$.

An natural idea, initially proposed by Marc Bernot, is to define an easy to compute analogues of the Wasserstein distance in dimension d > 1 using averages of 1D Wasserstein distances. This idea was first exploited in a joint work between Marc Bernot, Julie Delon, Gabriel Peyré and Julien Rabin in the context of texture generation [33]. We also refer to the PhD theses of Nicolas Bonnotte [11] and Kimia Nadjahi [28], which provide the most detailed theoretical study of these distances.

Definition 25 (Sliced Wasserstein). Given a direction θ in the unit sphere \mathcal{S}^{d-1} , we note $P_{\theta}(x) = \langle x | \theta \rangle$ the projection of x on the line spanned by θ . The *p*-sliced Wasserstein distance between two measures $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ are then defined by averaging 1D Wasserstein distances:

$$SW_p^p(\mu,\nu) = \int W_p^p(P_{\theta\#}\mu, P_{\theta\#}\nu) d\sigma(\theta),$$

where σ is the uniform probability measure over the unit sphere \mathcal{S}^{d-1} .

Proposition 50. The sliced Wasserstein distance SW_p enjoys the following properties:

- (i) SW_p is indeed a distance on $\mathcal{P}_p(\mathbb{R}^d)$; in particular, SW_p(μ, ν) = 0 if and only if $\mu = \nu$
- (ii) SW_p is weaker than the Wasserstein distance:

 $\forall \mu, \nu \in \mathcal{P}_p(\mathbb{R}^d), \quad \mathrm{SW}_p(\mu, \nu) \leqslant \mathrm{W}_p(\mu, \nu).$

(iii) the topology induced by SW_p is stronger than the weak* topology (iv) the map

$$\Phi: \mathcal{P}_p(\mathbb{R}^d) \to \mathrm{L}^p(\mathcal{S}^{d-1} \times [0,1])$$
$$\mu \mapsto \left[(\theta, r) \mapsto T_{P_{\theta \neq \mu}}(r) \right]$$

where $T_{P_{\theta\#\mu}}$ is the quantile function of $P_{\theta\#\mu}$, is an isometric embedding of $(\mathcal{P}_p(\mathbb{R}^d), \mathrm{SW}_p)$ into the Banach space $\mathrm{L}^p(\mathcal{S}^{d-1} \times [0, 1])$. We refer to [11, Proposition 5.1.2 and 5.1.3] for a detailed proof of some of these properties.

Proof. (i) We note that for any $\theta \in \mathcal{S}^{d-1}$ and $k \in \mathbb{R}$,

$$\langle \mu | e^{i \langle k\theta | \cdot \rangle} \rangle = \int_{\mathbb{R}^d} e^{ikP_{\theta}(x)} \mathrm{d}x = \int_{\mathbb{R}} e^{ikt} \mathrm{d}P_{\theta \#} \mu(t) = \langle P_{\theta \#} \mu | e^{ik \cdot} \rangle.$$

where we used the change of variable $t = P_{\theta}(x)$. If $SW_p(\mu, \nu) = 0$, then for all θ , $P_{\theta \#}\mu = P_{\theta \#}\nu$, so that by the above computation the Fourier transform of μ, ν agree. This implies that $\mu = \nu$.

(ii) The upper bound of SW_p in terms of W_p is obtained by remarking that if $\gamma \in \Gamma(\mu, \nu)$, then $(P_{\theta}, P_{\theta})_{\#}\gamma$ belongs to $\Gamma(P_{\theta\#}\mu, P_{\theta\#}\nu)$. (iii) This is a consequence of the Cramér-Wold theorem and the fact that W_p topologizes weak convergence on $\mathcal{P}_p(\mathbb{R})$.

From the last point of the previous proposition, and form the Theorem 49, we deduce that the sliced Wasserstein distance SW_p cannot be coarsely equivalent to the Wasserstein distance W_p on $\mathcal{P}_p(\mathbb{R}^d)$ when $d \ge 3$ and p > 1. However, Bonotte [11, Theorem 5.1.5] was still able to prove that SW_p and W_p are bi-Hölder equivalent to each other when the probability measures are on a fixed compact set.

Theorem 51 (Bonotte). Let K be a compact subset of \mathbb{R}^d with diam $(K) \leq R$. Then for all $p \geq 1$, there exists constants $C_{d,p,R} > 0$ such that

$$\forall \mu, \nu \in \mathcal{P}(K), \quad \mathrm{SW}_p(\mu, \nu) \leqslant \mathrm{W}_p(\mu, \nu) \leqslant C_{d, p, R} \, \mathrm{SW}_p(\mu, \nu)^{\overline{p(d+1)}}$$

Remark 23 (Non-convex image). The image of $\mathcal{P}_p(\mathbb{R}^d)$ under the embedding Φ described above is not necessarily a convex subset of $L^p(\mathcal{S}^{d-1} \times [0,1])$ (this is discussed in detail in [33]). This complicates many tasks that one could wish to accomplish using this embedding, e.g. defining a notion of barycenter between measures. In order to define the barycenter between μ_1, \ldots, μ_N with weights $\alpha_1, \ldots, \alpha_N > 0$, one could start by taking a weighted average of the images $\Phi(\mu_i)$,

$$\frac{1}{\sum_i \alpha_i} \sum_{i=1}^N \alpha_i \Phi(\mu_i).$$

However, this average might not belong to the range of Φ , so that we may not be able to take its inverse image.

8.3. "Linearization" of the quadratic Wasserstein distance. We fix a supported probability density $\rho \in \mathcal{P}^{\mathrm{ac}}(\mathbb{R}^d)$, supported over a compact convex set X and bounded from above and below positive constants. Given $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we call *Brenier map* $T_{\rho \to \mu}$ the quadratic optimal transport map between ρ and μ . In practice, since ρ is fixed, we will often denote T_{μ} the Brenier map.

Remark 24 (Relation to quantile function). Note that in dimension 1, if ρ is the Lebesgue measure on [0, 1], the map $T_{\rho \to \mu}$ coincides with the quantile function. In fact, maps of the form $T_{\rho \to \mu}$ have been suggested as a analogue of the quantile function in [15].

Proposition 52. The mapping $\mu \mapsto T_{\mu}$ enjoys the following properties:

- (i) $\mu \mapsto T_{\mu}$ is injective;
- (ii) $\mu \mapsto T_{\mu}$ is reverse-Lipschitz:

$$\forall \mu, \nu \in \mathcal{P}(Y), \ W_2(\mu, \nu) \leq \|T_\mu - T_\nu\|_{L^2(\rho, \mathbb{R}^d)}.$$
 (8.33)

- (iii) $\mu \mapsto T_{\mu}$ is continuous.
- (iv) the image of $\mu \mapsto T_{\mu}$ is a convex subset of $L^{2}(\rho, \mathbb{R}^{d})$.

We note that the arguments used to prove the general continuity result (iii) are non-quantitative.

Proof. (i) The injectivity comes from the fact that $\mu = T_{\mu \#} \rho$. (ii) If we denote $\gamma := (T_{\mu}, T_{\nu})_{\#} \rho$, then $\gamma \in \Pi(\mu, \nu)$. The change of variable formula gives

$$\begin{aligned} \mathbf{W}_{2}^{2}(\mu,\nu) &\leqslant \int_{Y \times Y} \|y - y'\|_{2}^{2} d\gamma(y,y') \\ &= \int_{X} \|T_{\mu}(x) - T_{\nu}(x)\|_{2}^{2} \rho(x) \mathrm{d}x = \|T_{\mu} - T_{\nu}\|_{\mathrm{L}^{2}(\rho)}^{2}. \end{aligned}$$

(iii) If a sequence of probability measures $(\mu_n)_n$ converges to some μ in $(\mathcal{P}_2(\mathbb{R}^d), \mathbb{W}_2)$, then T_{μ_n} converges to T_{μ} in $L^2(\rho, \mathbb{R}^d)$. This continuity property is for instance implied by Corollary 5.23 in [43], together with the dominated convergence theorem.

These properties of the map $\mu \mapsto T_{\mu}$ motivated its use to embed the metric space $(\mathcal{P}_2(\mathbb{R}^d), \mathbb{W}_2)$ into the Hilbert space $L^2(\rho, \mathbb{R}^d)$. This approach is often referred to as the *Linearized Optimal Transport* (LOT) [45] framework and has shown great results in applications to image processing [45, 24, 5, 13, 31, 23].

Remark 25 (Convex image). A practical benefit of the embedding is to enable the use of the classical Hilbertian statistical toolbox on families of probability measures while keeping some features of the Wasserstein geometry. A particularly nice feature of the embedding $\mu \mapsto T_{\mu}$ is that its image in $L^2(\rho, \mathbb{R}^d)$ is convex, i.e. barycenters of optimal transport maps are optimal transport maps, and that the inverse image of the embedding is very easy to compute.

Remark 26 (Relation to generalized geodesics). Working with this embedding is equivalent to replacing the Wasserstein distance by the distance

$$W_{2,\rho}(\mu,\nu) = \|T_{\mu} - T_{\nu}\|_{L^{2}(\rho,\mathbb{R}^{d})}.$$

We note that the geodesic curves with respect to the distance $W_{2,\rho}$ are called the *generalized geodesics* in the book of Ambrosio, Gigli, Savaré [3].

Remark 27 ($\mu \mapsto T_{\mu}$ as a Riemannian logarithm). The choice of the Brenier map between a reference measure ρ and a measure μ as an embedding of μ may also be motivated by the Riemannian interpretation of the Wasserstein geometry [30, 3]. In this interpretation, the tangent space to $\mathcal{P}_2(\mathbb{R}^d)$ at ρ is included in $L^2(\rho, \mathbb{R}^d)$. The Brenier map minus the identity, T_{μ} – id, can be regarded as the vector in the tangent space at ρ which supports the Wasserstein geodesic from ρ to μ . In the Riemannian language again, the map $\mu \mapsto T_{\mu}$ – id would be called a *logarithm*, i.e. the inverse of the Riemannian

OPTIMAL TRANSPORT

exponential map: it sends a probability measure μ in the (curved) manifold $\mathcal{P}_2(\mathbb{R}^d)$ to a vector T_{μ} – id belonging to the linear space $L^2(\rho, \mathbb{R}^d)$. This establishes a connection between the linearized optimal transport framework idea and similar strategies used to extend statistical inference notions such as principal component analysis to manifold-valued data.

It is quite natural to expect that the embedding $\mu \mapsto T_{\mu}$ retains some of the geometry of the underlying space, or equivalently that the metric $W_{2,\rho}$ is comparable, in some coarse sense, to the Wasserstein distance (however, not on the whole space $\mathcal{P}_2(\mathbb{R}^d)$, by Theorem 49!). A first question in this direction is to estimate the Hölder exponent of the map $\mu \mapsto T_{\mu}$ when restricted to a suitable family of probability measures.

We first note that $\mu \mapsto T_{\mu}$ cannot be better than $\frac{1}{2}$ -Hölder (another example of this fact can be found in [21]).

Lemma 53. Let ρ be uniform on the unit disc $X \subseteq \mathbb{R}^2$. Then, there is a curve $\theta \in [0, 2\pi] \to \mu_{\theta} \in \mathcal{P}(X)$ and C > 0 such that

$$||T_{\mu_{\theta}} - T_{\mu_{0}}||_{L^{2}(\rho)} \ge C \operatorname{W}_{2}(\mu_{\theta}, \mu_{0})^{1/2}$$

Proof. Given $\theta \in \mathbb{R}$, we denote $x_{\theta} = (\cos \theta, \sin(\theta))$ and $\mu_{\theta} = \frac{1}{2}(\delta_{x_{\theta}} + \delta_{-x_{\theta}})$. Then, the optimal transport map between ρ and μ_{θ} is given by

$$T_{\mu_{\theta}}(x) = \begin{cases} x_{\theta} & \text{if } \langle x | x_{\theta} \rangle \ge 0\\ -x_{\theta} & \text{if not.} \end{cases}$$
(8.34)

One can easily check that for θ one has $W_2(\mu_0, \mu_\theta) \leq |\theta|$. For $\theta > 0$ we set

$$D_{\theta} = \{ x \in \mathbb{R}^2 \mid \langle x | x_0 \rangle \ge 0 \text{ and } \langle x | x_{\theta} \rangle \le 0 \}.$$
(8.35)

Then, on D_{θ} , $T_{\mu_{\theta}} \equiv x_{-\theta}$ and $T_{\mu_0} \equiv x_0$, giving

$$\|T_{\mu_{\theta}} - T_{\mu_{0}}\|_{\mathrm{L}^{2}(\rho)}^{2} \ge \int_{D_{\theta}} \|x_{-\theta} - x_{0}\|^{2} \,\mathrm{d}x = |D_{\theta}| \,\|x_{-\theta} - x_{0}\|^{2} \,.$$
(8.36)

Moreover, if $|\theta| \leq \frac{\pi}{2}$ one has $||x_{-\theta} - x_0||^2 \ge 2$. This gives

$$\|T_{\mu_{\theta}} - T_{\mu_{0}}\|_{\mathrm{L}^{2}(\rho)}^{2} \geq 2 |D_{\theta}| \geq \frac{|\theta|}{\pi}.$$

9. STABILITY OF QUADRATIC OPTIMAL TRANSPORT MAPS

We are interested in establishing quantitative continuity estimates for the map $\mu \in \mathcal{P}(Y) \mapsto T_{\mu}$, where T_{μ} is the optimal transport map between a reference probability density ρ on \mathbb{R}^d and μ , and where Y is a (fixed) compact subset of \mathbb{R}^d . We will rely on the following assumptions and notations:

Definition 26. We fix a supported probability density $\rho \in \mathcal{P}^{\mathrm{ac}}(\mathbb{R}^d)$, supported over a compact convex set X and bounded from above and below positive constants. Given $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, we call

- Brenier map T_{μ} the optimal transport map between ρ and μ ;
- Brenier potential the unique lower semi-continuous convex function $\varphi_{\rho \to \mu} \in \dot{\mathbf{L}}^2(\rho)$ such that $T_{\mu} = \nabla \varphi_{\mu}$ and $\int_X \varphi_{\mu} d\rho = 0$; • dual potential the convex conjugate of φ_{μ} , denoted $\psi_{\mu} = \varphi_{\mu}^*$.

Our main tool to prove these continuity estimates will be Kantorovich duality with respect to the cost $c(x, y) = -\langle x | y \rangle$. More precisely we will rely on the semi-dual problem

$$\min_{\psi \in \mathcal{C}^0(Y)} \mathcal{K}(\psi) + \langle \psi | \mu \rangle.$$
(9.37)

where \mathcal{K} is the Kantorovich functional

$$\mathcal{K}(\psi) = \int \psi^* \mathrm{d}\rho.$$

We have already established (see Proposition 17) that $\nabla \mathcal{K}(\psi) = -\nabla \psi_{\#}^* \rho$, so that the optimality condition for the problem (9.37) is

$$\nabla \psi_{\mu \#}^* \rho = \mu.$$

If ψ_{μ} satisfies this condition, then $\varphi_{\mu} = \psi_{\mu}^{*}$ is the Brenier potential and $T_{\mu} = \nabla \varphi_{\mu}$ is the Brenier map. Adding a constant from φ_{μ} if necessary, we may assume that $\langle \varphi_{\mu} | \rho \rangle = 0$; the same constant is then substracted from ψ_{μ} and $(\varphi_{\mu}, \psi_{\mu}, T_{\mu})$ are then uniquely defined.

9.1. Stability near a regular configuration. We state a first positive result, which is a slight variant of a known stability result due to Ambrosio and reported in [21].

Theorem 54. Let
$$\mu, \nu \in \mathcal{P}(Y)$$
 and assume that T_{μ} is K-Lipschitz. Then,
 $\|T_{\mu} - T_{\nu}\|_{L^{2}(\rho)} \leq 2\sqrt{M_{X}K} \operatorname{W}_{1}(\mu, \nu)^{1/2},$ (9.38)
where $M_{X} = \max_{x \in X} \|x\|.$

We deduce this theorem from the following elementary lemma, which can be regarded as a strong concavity estimate for the Kantorovich functional \mathcal{K} , as it can be rephrased as

$$\langle \psi_{\nu} - \psi_{\mu} | \mathcal{K}(\psi_{\nu}) - \mathcal{K}(\psi_{\mu}) \rangle \ge \frac{1}{2K} \left\| T_{\mu} - T_{\nu} \right\|_{\mathrm{L}^{2}(\rho)}^{2}$$

Lemma 55. Under the assumptions of Theorem 54,

$$\|T_{\mu} - T_{\nu}\|_{\mathrm{L}^{2}(\rho)}^{2} \leqslant 2K \int_{Y} (\psi_{\nu} - \psi_{\mu}) \mathrm{d}(\mu - \nu)$$
(9.39)

Proof. From convex analysis, we know that the map $T_{\mu} = \nabla \varphi_{\mu}$ is K-Lipschitz if and only if the dual ψ_{μ} is $\frac{1}{K}$ -strongly convex. We denote $A = \int_{Y} \psi_{\nu} d(\mu - \nu)$ and $B = \int_{Y} \psi_{\mu} d(\nu - \mu)$. Using that $(\nabla \varphi_{\mu})_{\#} \rho = \mu$ and $(\nabla \varphi_{\nu})_{\#} \rho = \nu$, we get

$$A = \int_{X} (\psi_{\nu}(\nabla\varphi_{\mu}) - \psi_{\nu}(\nabla\varphi_{\nu})) d\rho$$
$$= \int_{X} (\psi_{\nu}(\nabla\psi_{\mu}^{*}) - \psi_{\nu}(\nabla\psi_{\nu}^{*})) d\rho \qquad (9.40)$$

We now use the inequality $\psi_{\nu}(y) - \psi_{\nu}(z) \ge \langle y - z | v \rangle$, which holds for all v in the subdifferential $\partial \psi_{\nu}(z)$. The convex functions ψ_{ν}, ψ_{μ} are differentiable ρ -almost everywhere. Taking $z = \nabla \psi_{\nu}^{*}(x)$ and $y = \nabla \psi_{\mu}^{*}(x)$, and using $x \in \partial \psi_{\nu}(z)$, we obtain

$$A \ge \int_X \langle \mathrm{id}, \nabla \psi^*_\mu - \nabla \psi^*_\nu \rangle \mathrm{d}\rho \tag{9.41}$$

Using the strong convexity of ψ_{μ} , we get a similar lower bound on B, with an extra quadratic term

$$B = \int_{X} (\psi_{\mu}(\nabla\psi_{\nu}^{*}) - \psi_{\mu}(\nabla\psi_{\mu}^{*})) d\rho$$

$$\geqslant \int_{X} \left(\langle \operatorname{id}, \nabla\psi_{\nu}^{*} - \nabla\psi_{\mu}^{*} \rangle + \frac{1}{2K} \|\nabla\psi_{\nu}^{*} - \nabla\psi_{\mu}^{*}\|_{2}^{2} \right) d\rho.$$
(9.42)

Summing up the lower bounds on A and B, we get:

$$\int_{Y} (\psi_{\nu} - \psi_{\mu}) d(\mu - \nu) \geq \frac{1}{2K} \int_{X} \|\nabla \psi_{\nu}^{*} - \nabla \psi_{\mu}^{*}\|_{2}^{2} d\rho$$
$$= \frac{1}{2K} \|T_{\nu} - T_{\mu}\|_{L^{2}(\rho)}^{2}.$$

Proof of Theorem 54. Being defined as the convex conjugate of $\varphi_{\nu} : X \to \mathbb{R}$, we know that ψ_{ν} is Lipschitz with constant $\leq M_X$. Combining this with Lemma 55,

$$\|T_{\mu} - T_{\nu}\|_{L^{2}(\rho)}^{2} \leqslant 2K \int_{Y} (\psi_{\nu} - \psi_{\mu}) d(\mu - \nu)$$
$$\leqslant 2K \max_{\operatorname{Lip}(f) \leqslant M_{X}} \int_{Y} f d(\mu - \nu)$$
$$= 2KM_{X} \max_{\operatorname{Lip}(f) \leqslant 1} \int_{Y} f d(\mu - \nu)$$
$$= 2KM_{X} \operatorname{W}_{1}(\mu, \nu), \qquad (9.43)$$

where we used Kantorovich-Rubinstein's theorem to get the last equality. \Box

9.2. Stability of potentials for entropy-regularized quadratic optimal transport. In this section, we fix a reference probability measure σ in $\mathcal{P}(Y)$, with support equal to Y. Given $\varepsilon > 0$, and $\psi \in \mathcal{C}^0(Y)$, we define the ε -regularized convex conjugate as

$$\psi^{*,\varepsilon}(x) = \varepsilon \log\left(\int_Y e^{\frac{\langle x|y\rangle - \psi(y)}{\varepsilon} \mathrm{d}\sigma(y)}\right),$$

and the ε -regularized Kantorovich as

$$\mathcal{K}_{\varepsilon}(\psi) = \int_{X} \psi^{*,\varepsilon} \mathrm{d}\rho.$$

Lemma 56 (Convergence as $\varepsilon \to 0$). For any $\psi \in \mathcal{C}^0(Y)$,

- $\psi^{*,\varepsilon}$ converges pointwise to ψ as $\varepsilon \to 0$;
- $\lim_{\varepsilon} \mathcal{K}_{\varepsilon}(\psi) = \mathcal{K}(\psi)$;

We now look at the gradient of $\mathcal{K}_{\varepsilon}$. To each potential $\psi \in \mathcal{C}^0(Y)$ and any point x in the source domain X, we will associate a probability density $\hat{\mu}_{\varepsilon}^x[\psi]$ (with respect to σ) and the corresponding probability measure $\mu_{\varepsilon}^x[\psi]$ on Y:

$$\hat{\mu}_{\varepsilon}^{x}[\psi] = \frac{e^{\frac{\langle x|y\rangle - \psi(y)}{\varepsilon}}}{\int_{Y} e^{\frac{\langle x|z\rangle - \psi(z)}{\varepsilon}} \mathrm{d}\sigma(z)}$$
$$\mu_{\varepsilon}^{x}[\psi] = \hat{\mu}_{\varepsilon}^{x}[\psi] \mathrm{d}\sigma$$

We also define

$$\mu_{\varepsilon}[\psi] = \int_{X} \mu_{\varepsilon}^{x}[\psi] \mathrm{d}\rho(x).$$

Remark 28 (Limit as $\varepsilon \to 0$). Note that if $\nabla \psi^*$ is differentiable at x, then the maximum of $y \mapsto \langle x | y \rangle - \psi(y)$ is attained at the point $\nabla \psi^*(x)$. If this is the case, then

$$\lim_{\varepsilon \to 0} \mu_{\varepsilon}^{x}[\psi] = \delta_{\nabla \psi^{*}(x)}.$$

Thus, at least formally, $\mu_{\varepsilon}[\psi]$ is the analogue of

$$\nabla \psi_{\#}^* \rho = \int \delta_{\nabla \psi^*(x)} \mathrm{d}\rho(x).$$

Lemma 57 (Gradient of $\mathcal{K}_{\varepsilon}$). The smoothed Kantorovich functional $\mathcal{K}_{\varepsilon}$ is differentiable at any $\psi \in \mathcal{C}^0(Y)$ with

$$\nabla \mathcal{K}_{\varepsilon}(\psi) = -\mu_{\varepsilon}[\psi],$$

i.e. for all $v \in \mathcal{C}^0(Y)$,

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{K}_{\varepsilon}(\psi+tv) = -\langle \mu_{\varepsilon}[\psi]|v\rangle = -\int_{X} \langle \hat{\mu}_{\varepsilon}^{x}[\psi]|v\rangle \mathrm{d}\rho(x).$$

Moreover,

$$\lim_{\varepsilon} \nabla \mathcal{K}_{\varepsilon}(\psi) = \nabla \mathcal{K}(\psi)$$

Proof. Let $\psi_t = \psi + tv$. Then, by definition of $\mu_{\varepsilon}^x[\psi]$ we have

$$\frac{\mathrm{d}}{\mathrm{d}t}\psi_{t}^{*,\varepsilon}(x) = \varepsilon \frac{\mathrm{d}}{\mathrm{d}t}\log\left(\int_{Y} e^{\frac{\langle x|y\rangle - \psi_{t}(y)}{\varepsilon}}\mathrm{d}\sigma(y)\right) \\
= \varepsilon \frac{\int_{Y} \frac{\mathrm{d}}{\mathrm{d}t} e^{\frac{\langle x|y\rangle - \psi_{t}(y)}{\varepsilon}}\mathrm{d}\sigma(y)}{\int_{Y} e^{\frac{\langle x|y\rangle - \psi_{t}(y)}{\varepsilon}}\mathrm{d}\sigma(y)} \tag{9.44}$$

$$= -\langle v|\mu_{\varepsilon}^{x}[\psi]\rangle$$

We conclude by differentiating under the integral defining $\mathcal{K}_{\varepsilon}(\psi)$ that

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}t} \mathcal{K}_{\varepsilon}(\psi + tv) &= \varepsilon \int_{X} \frac{\mathrm{d}}{\mathrm{d}t} \psi_{t}^{*,\varepsilon}(x) \mathrm{d}\rho(x). \\ &= -\int_{X} \langle v | \mu_{\varepsilon}^{x}[\psi] \rangle \mathrm{d}\rho(x). \\ &= - \langle v | \mu_{\varepsilon} \rangle. \end{split}$$

We consider the set of continuous functions with bounded oscillation, where $osc(\psi) = \sup \psi - \inf \psi$:

$$\mathcal{C}^0_R(Y) = \{ \psi \in \mathcal{C}^0(Y) \mid \operatorname{osc}(\psi) \leqslant R \}.$$

Theorem 58. The functional $\mathcal{K}_{\varepsilon}$ is strongly convex on $\mathcal{C}^{0}_{R}(Y)$. More precisely, if $\psi_{0}, \psi_{1} \in \mathcal{C}^{0}_{R}(Y)$, then $\langle \nabla \mathcal{K}_{\varepsilon}(\psi_{1}) - \nabla \mathcal{K}_{\varepsilon}(\psi_{0}) | \psi_{1} - \psi_{0} \rangle \geq c \operatorname{Var}_{\sigma}(\psi_{1} - \psi_{0}).$

where
$$c_{\varepsilon} = \frac{1}{\varepsilon} e^{-\frac{1}{\varepsilon} \frac{1}{\varepsilon} - \frac{1}{\varepsilon}}$$
 and $R_Z = \max_{z \in Z} \|z\|$.

Remark 29 (Stability of dual potentials). In particular, if ψ_0, ψ_1 are Lipschitz with constant K, then

$$\operatorname{Var}_{\sigma}(\psi_{1} - \psi_{0}) \leqslant \frac{1}{c_{\varepsilon}} \langle \mu_{\varepsilon}[\psi_{1}] - \mu_{\varepsilon}[\psi_{0}] | \psi_{1} - \psi_{0} \rangle$$
$$\leqslant \frac{2}{c_{\varepsilon}} K \operatorname{W}_{1}(\mu_{\varepsilon}[\psi_{1}], \mu_{\varepsilon}[\psi_{0}]),$$

where we used Kantorovich-Rubinstein's theorem (Theorem 34) to get the second inequality. Note that as $\varepsilon \to 0$, c_{ε} tends to infinity, so that this inequality does not translate into a stability inequality for the unregularized case $\varepsilon = 0$. In a similar spirit (but using different techniques, involving the so-called Hilbert metric), Giulia Luise et al [25, Theorem C.4] proves an estimate of the form

$$\operatorname{osc}(\psi_1 - \psi_0) \leqslant c_{\varepsilon} \|\mu_{\varepsilon}[\psi_0] - \mu_{\varepsilon}[\psi_1]\|_{\mathrm{TV}},$$

with $\lim_{\varepsilon \to 0} c_{\varepsilon} = +\infty$. Note that the dependence is *Lipschitz* in their case, which is important for some applications.

Given $\psi \in \mathcal{C}^0(Y)$ and a direction $v \in \mathcal{C}^0(Y)$, we define

$$\langle \mathbf{D}^{2}\mathcal{K}_{\varepsilon}(\psi)v|v\rangle = \lim_{t\to 0} \langle \nabla\mathcal{K}_{\varepsilon}(\psi+tv) - \nabla\mathcal{K}_{\varepsilon}(\psi)|v\rangle.$$

Lemma 59. $\langle D^2 \mathcal{K}_{\varepsilon}(\psi) v | v \rangle = \frac{1}{\varepsilon} \int_X \operatorname{Var}_{\mu_{\varepsilon}^x[\psi]}(v) d\rho(x).$

Proof. Let $\psi_t = \psi + tv$. By integration under the integral, we have

$$\frac{\mathrm{d}}{\mathrm{d}t} \langle \nabla \mathcal{K}_{\varepsilon}(\psi_t) | v \rangle = -\int_{X \times Y} \langle v | \frac{\mathrm{d}}{\mathrm{d}t} \mu_{\varepsilon}^x[\psi] \rangle \mathrm{d}\rho(x).$$

Let us compute the derivative of the density $\hat{\mu}_{\varepsilon}^{x}[\psi]$:

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}t} \ \hat{\mu}_{\varepsilon}^{x}[\psi_{t}](y)|_{t=0} &= \frac{\mathrm{d}}{\mathrm{d}t} \frac{e^{\frac{\langle x|y\rangle - \psi_{t}(y)}{\varepsilon}}}{\int_{Y} e^{\frac{\langle x|z\rangle - \psi_{t}(z)}{\varepsilon}} \mathrm{d}\sigma(z)} \\ &= \frac{-\frac{v(y)}{\varepsilon} e^{\frac{\langle x|y\rangle - \psi_{t}(y)}{\varepsilon}} \int_{Y} e^{\frac{\langle x|z\rangle - \psi_{t}(z)}{\varepsilon}} \mathrm{d}\sigma(z) + \frac{1}{\varepsilon} e^{\frac{\langle x|y\rangle - \psi_{t}(y)}{\varepsilon}} \int_{Y} e^{\frac{\langle x|z\rangle - \psi_{t}(z)}{\varepsilon}} v(z) \mathrm{d}\sigma(z)}{(\int_{Y} e^{\frac{\langle x|z\rangle - \psi_{t}(z)}{\varepsilon}} \mathrm{d}\sigma(z))^{2}} \\ &= -\frac{v(y)}{\varepsilon} \hat{\mu}_{\varepsilon}^{x}[\psi](y) + \frac{1}{\varepsilon} \hat{\mu}_{\varepsilon}^{x}[\psi](y) \langle \mu_{\varepsilon}^{x}|\psi \rangle \end{split}$$

Thus,

$$\begin{split} \langle \mathbf{D}^{2} \mathcal{K}_{\varepsilon}(\psi) v | v \rangle &= \frac{1}{\varepsilon} \int_{X} \langle v^{2} | \mu_{\varepsilon}^{x}[\psi] \rangle - (\langle v | \mu_{\varepsilon}^{x}[\psi] \rangle)^{2} \mathrm{d}\rho(x) \\ &= \frac{1}{\varepsilon} \int_{X} \mathrm{Var}_{\mu_{\varepsilon}^{x}[\psi]}(v) \mathrm{d}\rho(x). \end{split}$$

Proof of Theorem 58. Let $v = \psi_1 - \psi_0$ and $\psi_t = \psi_0 + tv$. By Taylor's formula, we have

$$\begin{split} \langle \nabla \mathcal{K}_{\varepsilon}(\psi_{1}) - \nabla \mathcal{K}_{\varepsilon}(\psi_{0}) | \psi_{1} - \psi_{0} \rangle &= \langle \nabla \mathcal{K}_{\varepsilon}(\psi_{1}) - \nabla \mathcal{K}_{\varepsilon}(\psi_{0}) | v \rangle \\ &= \int_{0}^{1} \langle \mathrm{D}^{2} \mathcal{K}_{\varepsilon}(\psi_{t}) v | v \rangle \mathrm{d}t \\ &\geqslant \frac{1}{\varepsilon} \int_{0}^{1} \int_{X} \mathrm{Var}_{\mu_{\varepsilon}^{x}[\psi_{t}]}(v) \mathrm{d}\rho(x) \mathrm{d}t \end{split}$$

Recall that $\mu_{\varepsilon}^{x}[\psi_{t}]$ has density $\hat{\mu}_{\varepsilon}^{x}[\psi_{t}]$ with respect to σ :

$$\hat{\mu}_{\varepsilon}^{x}[\psi_{t}](y) = \frac{e^{\frac{\langle x|y\rangle - \psi_{t}(y)}{\varepsilon}}}{\int_{Y} e^{\frac{\langle x|z\rangle - \psi_{t}(z)}{\varepsilon}} \mathrm{d}\sigma(z)} \geqslant \frac{e^{\frac{\langle x|y\rangle - \sup\psi_{t}}{\varepsilon}}}{\int_{Y} e^{\frac{\langle x|z\rangle - \inf\psi_{t}}{\varepsilon}} \mathrm{d}\sigma(z)} \geqslant \frac{e^{\frac{-R_{X}R_{Y} - \sup\psi_{t}}{\varepsilon}}}{e^{\frac{R_{X}R_{Y} - \inf\psi_{t}}{\varepsilon}}},$$

where $R_X = \max_{x \in X} \|x\|$, $R_Y = \max_{y \in Y} \|y\|$ with respect to σ , so that with $c = e^{-\frac{2R_XR_Y+R}{\varepsilon}}$ we get

$$\hat{\mu}^x_{\varepsilon}[\psi_t](y) \geqslant c\sigma.$$

From this, we deduce that $\operatorname{Var}_{\hat{\gamma}_{\varepsilon}^{x}}(v) \ge c \operatorname{Var}_{\sigma}(v)$, which allows to conclude.

10. HÖLDER STABILITY OF DUAL POTENTIALS

In this section, we show how to prove Hölder estimates for the dual potentials. The proof is taken from [17], and relies on strong convexity estimates for $\mathcal{K}_{\varepsilon}$, with a constant that does not degrade as $\varepsilon \to 0$. The main idea is to deduce strong convexity of $\mathcal{K}_{\varepsilon}$ from mere convexity of

$$I_{\varepsilon}(\psi) = \log\left(\int_X e^{-\psi^{*,\varepsilon}} \mathrm{d}x\right).$$

Given $\psi \in \mathcal{C}^0(Y)$, we will denote $\rho_{\varepsilon}[\psi]$ the Gibbs measure of $\psi^{*,\varepsilon}$, i.e.

$$\rho_{\varepsilon}[\psi](x) = \frac{e^{-\psi^{*,\varepsilon}(x)}}{\int_X e^{-\psi^{*,\varepsilon}(z)} \mathrm{d}z}.$$

Proposition 60. I_{ε} is concave and

$$\nabla I_{\varepsilon}(\psi) = \int_{X} \mu_{\varepsilon}^{x}[\psi] \rho_{\varepsilon}[\psi](x) \mathrm{d}x,$$
$$- \left(\frac{1}{\varepsilon} + 1\right) \int_{\varepsilon} \operatorname{Var}_{\varepsilon} \psi_{\varepsilon}(v) \rho_{\varepsilon}[\psi_{\varepsilon}](x) \mathrm{d}x,$$

$$\langle \mathbf{D}^2 I_{\varepsilon}(\psi) v | v \rangle = -\left(\frac{1}{\varepsilon} + 1\right) \int_X \operatorname{Var}_{\mu_{\varepsilon}^x[\psi_t]}(v) \rho_{\varepsilon}[\psi_t](x) \mathrm{d}x + \operatorname{Var}_{\mu_{\varepsilon}^I[\psi_t]}(v)$$
where $\mu_{\varepsilon}^I[\psi] = \int \mu_{\varepsilon}^x[\psi] \rho_{\varepsilon}[\psi](x) \mathrm{d}x.$

Proof. Let $\psi_t = \psi + tv$. Then,

$$\begin{aligned} \frac{\mathrm{d}}{\mathrm{d}t} e^{I_{\varepsilon}(\psi_t)} &= \frac{\mathrm{d}}{\mathrm{d}t} \int_X e^{-\psi_t^{*,\varepsilon}} \mathrm{d}x \\ &= \int_X \frac{\mathrm{d}}{\mathrm{d}t} e^{-\psi_t^{*,\varepsilon}} \mathrm{d}x \\ &= -\int_X \left(\frac{\mathrm{d}}{\mathrm{d}t} \psi_t^{*,\varepsilon}(x)\right) e^{-\psi_t^{*,\varepsilon}} \mathrm{d}x \\ &= -\int_X \left(\frac{\mathrm{d}}{\mathrm{d}t} \psi_t^{*,\varepsilon}(x)\right) e^{-\psi_t^{*,\varepsilon}} \mathrm{d}x \\ &= \int_X \langle v | \mu_{\varepsilon}^x [\psi_t] \rangle e^{-\psi_t^{*,\varepsilon}} \mathrm{d}x \end{aligned}$$

Thus,

$$\nabla I_{\varepsilon}(\psi_t) = \frac{1}{I_{\varepsilon}(\psi_t)} \int_X \mu_{\varepsilon}^x[\psi_t] e^{-\psi_t^{*,\varepsilon}} \mathrm{d}x = \int_X \mu_{\varepsilon}^x[\psi_t] \rho_{\varepsilon}[\psi_t] \mathrm{d}x$$

,

We now compute the second derivative:

$$\begin{split} \langle \mathbf{D}^{2}I_{\varepsilon}(\psi_{t})v|v\rangle &= \frac{\mathrm{d}}{\mathrm{d}t} \langle \nabla I_{\varepsilon}(\psi_{t})|v\rangle \\ &= \frac{\mathrm{d}}{\mathrm{d}t} \int_{X} \int_{Y} v(y)\hat{\mu}_{\varepsilon}^{x}[\psi_{t}](y)\rho_{\varepsilon}[\psi_{t}](x)\mathrm{d}\sigma(y)\mathrm{d}x \\ &= \int_{X} \int_{Y} v(y) \left[\left(\frac{\mathrm{d}}{\mathrm{d}t}\hat{\mu}_{\varepsilon}^{x}[\psi_{t}](y)\right)\rho_{\varepsilon}[\psi_{t}](x) + \hat{\mu}_{\varepsilon}^{x}[\psi_{t}](y) \left(\frac{\mathrm{d}}{\mathrm{d}t}\rho_{\varepsilon}[\psi_{t}](x)\right) \right] \mathrm{d}\sigma(y)\mathrm{d}x \\ &= \int_{X} \int_{Y} v(y) \left(\frac{\mathrm{d}}{\mathrm{d}t}\hat{\mu}_{\varepsilon}^{x}[\psi_{t}](y)\right)\rho_{\varepsilon}[\psi_{t}](x)\mathrm{d}\sigma(y)\mathrm{d}x \\ &+ \int_{X} \int_{Y} v(y)\hat{\mu}_{\varepsilon}^{x}[\psi_{t}](y) \left(\frac{\mathrm{d}}{\mathrm{d}t}\rho_{\varepsilon}[\psi_{t}](x)\right)\mathrm{d}\sigma(y)\mathrm{d}x \end{split}$$

Following the computations already performed for $\mathcal{K}_{\varepsilon}$, we can see that the first term of the sum is equal to

$$-\frac{1}{\varepsilon}\int_X \operatorname{Var}_{\mu_{\varepsilon}^x[\psi_t]}(v) \mathrm{d}\rho_{\varepsilon}[\psi].$$

We turn to the second term. Let us first differentiate $\rho_{\varepsilon}[\psi_t](x)$ with respect to t:

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}t}\rho_{\varepsilon}[\psi_{t}](x) &= \frac{\mathrm{d}}{\mathrm{d}t} \frac{e^{-\psi_{t}^{*,\varepsilon}(x)}}{\int_{X} e^{-\psi_{t}^{*,\varepsilon}(z)}\mathrm{d}z} \\ &= \frac{\frac{\mathrm{d}}{\mathrm{d}t} e^{-\psi_{t}^{*,\varepsilon}(x)}}{\int_{X} e^{-\psi_{t}^{*,\varepsilon}(z)}\mathrm{d}z} - \frac{e^{-\psi_{t}^{*,\varepsilon}(x)}\frac{\mathrm{d}}{\mathrm{d}t}\int_{X} e^{-\psi_{t}^{*,\varepsilon}(z)}\mathrm{d}z}{(\int_{X} e^{-\psi_{t}^{*,\varepsilon}(z)}\mathrm{d}z)^{2}} \\ &= -\frac{\left(\frac{\mathrm{d}}{\mathrm{d}t}\psi_{t}^{*,\varepsilon}(x)\right)e^{-\psi_{t}^{*,\varepsilon}(x)}}{\int_{X} e^{-\psi_{t}^{*,\varepsilon}(z)}\mathrm{d}z} + \frac{e^{-\psi_{t}^{*,\varepsilon}(x)}\int_{X}\left(\frac{\mathrm{d}}{\mathrm{d}t}\psi_{t}^{*,\varepsilon}(z)\right)e^{-\psi_{t}^{*,\varepsilon}(z)}\mathrm{d}z}{(\int_{X} e^{-\psi^{*,\varepsilon}(z)}\mathrm{d}z)^{2}} \\ &= \langle \mu_{\varepsilon}^{x}[\psi_{t}]|v\rangle\rho_{\varepsilon}[\psi_{t}](x) - \rho_{\varepsilon}[\psi_{t}](x)\int_{X}\langle \mu_{\varepsilon}^{z}[\psi_{t}]|v\rangle\rho_{\varepsilon}[\psi_{t}](z)\mathrm{d}z} \\ &= \langle \mu_{\varepsilon}^{x}[\psi_{t}]|v\rangle\rho_{\varepsilon}[\psi_{t}](x) - \rho_{\varepsilon}[\psi_{t}](x)\langle \mu_{\varepsilon}^{I}[\psi_{t}]|v\rangle, \end{split}$$

where we used (9.44). Then,

$$\begin{split} &\int_X \int_Y v(y)\hat{\mu}_{\varepsilon}^x[\psi_t](y) \left(\frac{\mathrm{d}}{\mathrm{d}t}\rho_{\varepsilon}[\psi_t](x)\right) \mathrm{d}\sigma(y)\mathrm{d}x \\ &= \int_X \langle v|\mu_{\varepsilon}^x[\psi_t] \rangle \left(\langle \mu_{\varepsilon}^x[\psi_t]|v \rangle \rho_{\varepsilon}[\psi_t](x) - \rho_{\varepsilon}[\psi_t](x) \langle \mu_{\varepsilon}^I[\psi_t]|v \rangle \right) \right) \mathrm{d}x \\ &= \int \left(\langle v|\mu_{\varepsilon}^x[\psi_t] \rangle \right)^2 \rho_{\varepsilon}[\psi_t](x)\mathrm{d}x - \left(\langle \mu_{\varepsilon}^I[\psi_t]|v \rangle \right)^2 \\ &= \operatorname{Var}_{\rho_{\varepsilon}[\psi_t]}(x \mapsto \langle v|\mu_{\varepsilon}^x[\psi_t] \rangle) \\ &= \operatorname{Var}_{\mu_{\varepsilon}^I[\psi_t]}(v) - \int_X \operatorname{Var}_{\mu_{\varepsilon}^x[\psi_t]}(v) \rho_{\varepsilon}[\psi_t](x)\mathrm{d}x, \end{split}$$

where we used a variance decomposition formula to get the last line.

Concavity comes from the Prékopa-Leindler inequality

Let $c_t = \langle v | \mu_{\varepsilon}[\psi_t] \rangle$ and $c = \int_0^1 c_t dt$. Then,

$$\begin{aligned} \operatorname{Var}_{\rho}(\psi_{1}^{*,\varepsilon} - \psi_{0}^{*,\varepsilon}) &\leqslant \int_{X} (\psi_{1}^{*,\varepsilon}(x) - \psi_{0}^{*,\varepsilon}(x) - c)^{2} \mathrm{d}\rho(x) \\ &\leqslant \int_{X} \left| \int_{0}^{1} \frac{\mathrm{d}}{\mathrm{d}t} \psi_{t}^{*,\varepsilon}(x) - c_{t} \right|^{2} \mathrm{d}\rho(x) \\ &= \int_{0}^{1} \int_{X} |\langle \mu_{\varepsilon}^{x}[\psi_{t}] - c_{t}|v\rangle|^{2} \mathrm{d}\rho(x) \\ &= \int_{0}^{1} \int_{X} \operatorname{Var}_{\mu_{\varepsilon}[\psi_{t}]}(v) \mathrm{d}t \end{aligned}$$

References

- Jason Altschuler, Jonathan Weed, and Philippe Rigollet, Near-linear time approximation algorithms for optimal transport via sinkhorn iteration, Advances in Neural Information Processing Systems, 2017, pp. 1964–1974.
- 2. Luigi Ambrosio and Nicola Gigli, A user's guide to optimal transport, Modelling and optimisation of flows on networks, Springer, 2013, pp. 1–155.
- 3. Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré, Gradient flows: in metric spaces and in the space of probability measures, Springer Science & Business Media, 2008.
- Alexandr Andoni, Assaf Naor, and Ofer Neiman, Snowflake universality of wasserstein spaces, Annales Scientifiques de l'Ecole Normale Superieure 51 (2018), no. 3, 657–700.
- 5. Saurav Basu, Soheil Kolouri, and Gustavo K. Rohde, Detecting and visualizing cell phenotype differences from microscopy images using transport-based morphometry, Proceedings of the National Academy of Sciences 111 (2014), no. 9, 3448-3453.
- 6. Jean-David Benamou, Guillaume Carlier, and Luca Nenna, A numerical method to solve multi-marginal optimal transport problems with coulomb cost, Splitting Methods in Communication, Imaging, Science, and Engineering, Springer, 2016, pp. 577-601.
- Robert J. Berman, The Sinkhorn algorithm, parabolic optimal transport and geometric Monge-Ampère equations, arXiv preprint arXiv:1712.03082, 2017.
- D.P. Bertsekas, A new algorithm for the assignment problem, Mathematical Programming 21 (1981), no. 1, 152-171.
- 9. D.P. Bertsekas and J. Eckstein, Dual coordinate step methods for linear network flow problems, Mathematical Programming 42 (1988), no. 1, 203-243.
- Garrett Birkhoff, Tres observaciones sobre el algebra lineal, Univ. Nac. Tucuman, Ser. A 5 (1946), 147–154.
- 11. Nicolas Bonnotte, Unidimensional and evolution methods for optimal transportation, Ph.D. thesis, Paris 11, 2013.
- 12. Yann Brenier, Polar factorization and monotone rearrangement of vector-valued functions, Communications on pure and applied mathematics 44 (1991), no. 4, 375-417.
- Tianji Cai, Junyi Cheng, Nathaniel Craig, and Katy Craig, Linearized optimal transport for collider events, Phys. Rev. D 102 (2020), 116019.
- 14. Benjamin Charlier, Jean Feydy, Joan Alexis Glaunes, and Alain Trouvé, An efficient kernel product for automatic differentiation libraries, with applications to measure transport, Working version, 2017.
- Victor Chernozhukov, Alfred Galichon, Marc Hallin, Marc Henry, et al., Mongekantorovich depth, quantiles, ranks and signs, The Annals of Statistics 45 (2017), no. 1, 223-256.
- Keenan Crane, Clarisse Weischedel, and Max Wardetzky, Geodesics in heat: A new approach to computing distance based on heat flow, ACM Transactions on Graphics (TOG) 32 (2013), no. 5, 152.
- Alex Delalande, Nearly tight convergence bounds for semi-discrete entropic optimal transport, International Conference on Artificial Intelligence and Statistics, PMLR, 2022, pp. 1619-1642.

OPTIMAL TRANSPORT

- AN Dranishnikov, G Gong, V Lafforgue, and G Yu, Uniform embeddings into hilbert space and a question of gromov, Canadian Mathematical Bulletin 45 (2002), no. 1, 60-70.
- Jean Feydy, Pierre Roussillon, Alain Trouvé, and Pietro Gori, Fast and scalable optimal transport for brain tractograms, International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2019, pp. 636-644.
- Wilfrid Gangbo and Robert J McCann, The geometry of optimal transportation, Acta Mathematica 177 (1996), no. 2, 113-161.
- Nicola Gigli, On hölder continuity-in-time of the optimal transport map towards measures along a curve, Proceedings of the Edinburgh Mathematical Society 54 (2011), no. 2, 401-409.
- 22. William B Johnson and N Lovasoa Randrianarivony, $\ell^p(p > 2)$ does not coarsely embed into a hilbert space, Proceedings of the American Mathematical Society (2006), 1045–1050.
- Soheil Kolouri and Gustavo K. Rohde, Transport-based single frame super resolution of very low resolution face images, 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015, pp. 4876-4884.
- Soheil Kolouri, Akif B. Tosun, John A. Ozolek, and Gustavo K. Rohde, A continuous linear optimal transport approach for pattern analysis in image datasets, Pattern Recognition 51 (2016), 453-462.
- 25. Giulia Luise, Saverio Salzo, Massimiliano Pontil, and Carlo Ciliberto, *Sinkhorn barycenters with free support via frank-wolfe algorithm*, Advances in neural information processing systems **32** (2019).
- Robert J McCann, A convexity principle for interacting gases, Advances in mathematics 128 (1997), no. 1, 153-179.
- Quentin Merigot and Boris Thibert, Optimal transport: discretization and algorithms, Handbook of Numerical Analysis, vol. 22, Elsevier, 2021, pp. 133-212.
- 28. Kimia Nadjahi, Sliced-wasserstein distance for large-scale machine learning: theory, methodology and extensions, Ph.D. thesis, Institut Polytechnique de Paris, 2021.
- Piotr Nowak, Coarse embeddings of metric spaces into banach spaces, Proceedings of the American Mathematical Society 133 (2005), no. 9, 2589-2596.
- 30. Felix Otto, The geometry of dissipative evolution equations: the porous medium equation, Communications in Partial Differential Equations 26 (2001), 101-174.
- 31. S. Park and M. Thorpe, Representing and learning high dimensional data with the optimal transport map from a probabilistic viewpoint, 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 7864-7872.
- 32. Gabriel Peyré and Marco Cuturi, *Computational optimal transport*, Foundations and Trends® in Machine Learning **11** (2019), no. 5-6, 355-607.
- 33. Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernot, Wasserstein barycenter and its application to texture mixing, International Conference on Scale Space and Variational Methods in Computer Vision, Springer, 2011, pp. 435-446.
- 34. Filippo Santambrogio, Optimal transport for applied mathematicians, Springer, 2015.
- Giuseppe Savaré and Giacomo E Sodini, A simple relaxation approach to duality for optimal transport problems in completely regular spaces, Journal of Convex Analysis 29 (2022), no. 1, 1-12.
- Bernhard Schmitzer, A sparse multiscale algorithm for dense optimal transport, Journal of Mathematical Imaging and Vision 56 (2016), no. 2, 238-259.
- 37. _____, Stabilized sparse scaling algorithms for entropy regularized transport problems, SIAM Journal on Scientific Computing **41** (2019), no. 3, A1443–A1481.
- Richard Sinkhorn, A relationship between arbitrary positive matrices and doubly stochastic matrices, The annals of mathematical statistics 35 (1964), no. 2, 876–879.
- Richard Sinkhorn and Paul Knopp, Concerning nonnegative matrices and doubly stochastic matrices, Pacific Journal of Mathematics 21 (1967), no. 2, 343-348.
- 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 (TOG) **34** (2015), no. 4, 66.

OPTIMAL TRANSPORT

- 41. François-Xavier Vialard, An elementary introduction to entropic regularization and proximal methods for numerical optimal transport, Lecture, May 2019.
- 42. Cédric Villani, *Topics in optimal transportation*, no. 58, American Mathematical Soc., 2003.
- 43. _____, Optimal transport: old and new, vol. 338, Springer Science & Business Media, 2008.
- 44. Alexander Wagner, Nonembeddability of persistence diagrams with p > 2 wasserstein metric, Proceedings of the American Mathematical Society **149** (2021), no. 6, 2673-2677.
- 45. Wei Wang, Dejan Slepčev, Saurav Basu, John A. Ozolek, and Gustavo K. Rohde, A linear optimal transportation framework for quantifying and visualizing variations in sets of images, Int. J. Comput. Vision 101 (2013), no. 2, 254–269.