

Universidad Autónoma de Madrid

MATHEMATIC ANALYSIS FUNDAMENTALS

OPTIMAL TRANSPORT FOR
TOPOLOGICAL DATA ANALYSIS

End of Course Thesis.
2024-2025.

Author:
Gonzalo Ortega Carpintero

January 2025

1 Introduction

Transport maps were introduced in 1781 by Gaspard Monge to represent the idea of moving earth from one place into another [1][1.1 Historical overview]. In this original formulation of the optimal transport problem, it was enough to consider \mathbb{R}^3 as the ambient space, using the Euclidean distance as the cost function of moving mass between two points.

In the 30's, Leonid Kantorovich reformulated the problem to describe the optimization process of supply and demand distributions of diverse problems. The mass could be divided between different origin and destinations, making it possible to interpret the problem as the way to measure the cost of transforming one probability distribution into another. In this thesis, we will introduce the p -Wasserstein distance as a metric on the probability measures with finite p -moment space. When $p = 1$, the distance will represent the metric introduced in the Kantorovich optimal transport problem, also used and named Earth Mover's distance, used for machine learning algorithms and computer vision problems [2]. When $p = \infty$ it is named the bottleneck distance, and will be the main theme of study of this thesis.

In topological data analysis, diagrams arise to represent the persistence of the homology groups of a data set through time. Those diagrams are named persistence diagrams, and those homology groups, persistence homology groups. We will introduce an analogous p -Wasserstein distance in the space of persistence diagrams and prove that there exists an isometric embedding from a separable metric space into the space of persistence diagrams with the Wasserstein distance.

2 Optimal transport

The main result of optimal transport theory is the solution of Kantorovich's problem for general costs: the existence of an optimal transport plan. Lets start by introducing Monge's and Kantorovich's problems, observing its main key difference. For that, we shall fist define a way to compare probability measures from two different spaces. We will denote the set of probability measures over a space X by $\mathcal{P}(X)$, and the class of Borel-measurable sets by $\mathcal{B}(X)$.

Definition 2.1 (Push-forward measure). Let $T : X \rightarrow Y$ be a Borel map, and $\mu \in \mathcal{P}(X)$. Let $A \in \mathcal{B}$. The *push-forward measure* $T_{\#}\mu \in \mathcal{P}(Y)$ is defined as

$$T_{\#}\mu(A) := \mu(T^{-1}(A)).$$

Now we can introduce transport maps, as functions witch transform one probability measure into an other.

Definition 2.2 (Transport map). Given $\mu \in \mathcal{P}(X)$ and $\nu \in \mathcal{P}(Y)$, a *transport map from μ to ν* is a Borel map $T : X \rightarrow Y$ that satisfies $T_{\#}\mu = \nu$.

Definition 2.3 (Transport plan). Let $\pi_X : (X \times Y) \rightarrow X$ and $\pi_Y : (X \times Y) \rightarrow Y$ such that for every $(x, y) \in (X, Y)$, $\pi_X(x, y) = x$ and $\pi_Y(x, y) = y$. A *transport plan between μ and ν* is a probability measure $\gamma \in \mathcal{P}(X \times Y)$ where

$$(\pi_X)_{\#}\gamma = \mu \text{ and } (\pi_Y)_{\#}\gamma = \nu.$$

The set of all couplings between μ and ν is denoted $\Gamma(\mu, \nu)$.

While the set of transport maps between two given probability measures might be empty, transport plans are a more flexible generalization of them allowing to modulate one measure into the other. In probability theory, transport plans are named *couplings*, and $\Gamma(\mu, \nu)$ is the collection of all probability measures in $X \times Y$ with *marginals* μ and ν [3].

Given this definitions, we can introduce Monge and Kantorovich problems, $C_M(\mu, \nu)$ and $C_K(\mu, \nu)$ respectively, as follows.

Definition 2.4 (Transport problems). Fix $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$ and consider a lower semicontinuous map $c : X \times Y \rightarrow [0, \infty]$. Then

$$\begin{aligned} C_M(\mu, \nu) &:= \inf \left\{ \int_X c(x, T(x)) d\mu(x) : T_{\#}\mu = \nu \right\}, \\ C_K(\mu, \nu) &:= \inf \left\{ \int_{X \times Y} c(x, y) d\gamma(x, y) : \gamma \in \Gamma(\mu, \nu) \right\}. \end{aligned}$$

Next theorem asserts that it actually exists a minimizing transport plan that minimizes Kantorovich problem. This will prove useful to verify that Wasserstein distance exists and it is a well defined metric.

Theorem 2.5. *Let $c : X \times Y \rightarrow [0, \infty]$ be lower semicontinuous, and let $\mu \in \mathcal{P}(X)$ and $\nu \in \mathcal{P}(Y)$. Then there exists a coupling $\bar{\gamma} \in \Gamma(\mu, \nu)$ that verifies*

$$C_K(\mu, \nu) = \int_{X \times Y} c(x, y) d\bar{\gamma}(x, y).$$

Proof. (to do) □

Example 2.6 (Mean and variance in \mathbb{R}).

Definition 2.7 (Probability measures with finite p -moment). Let (X, d) be a locally compact and separable, metric space. Let $1 \leq p < \infty$. The set of probability measures with finite p -moment is defined As

$$\mathcal{P}_p(X) := \left\{ \sigma \in \mathcal{P}(X) : \int_X d(x, x_0)^p d\mu(x) < \infty \text{ for some } x_0 \in X \right\}.$$

Proposition 2.8. *The definition of $\mathcal{P}_p(X)$ is independent of the base point x_0*

Proof. (to do) □

Definition 2.9 (p -Wasserstein distance). Given $u, v \in \mathcal{P}_p(X)$, the p -Wasserstein distance is defined as

$$W_p(u, v) := \left(\inf_{\gamma \in \Gamma(u, v)} \int_{X \times X} d(x, y)^p d\gamma(x, y) \right)^{\frac{1}{p}}.$$

Proposition 2.10. *W_p is a distance on the space $\mathcal{P}_p(X)$.*

Proof. We will follow the steps made in [1][Theorem 3.1.5]. To prove the triangle inequality, let $\mu_1, \mu_2, \mu_3 \in \mathcal{P}_p(X)$ and

(to do) □

3 Wasserstein distance in persistence diagrams

In last section we have exposed the original optimal transport problem where the objective was to measure distance between probability measures. We will now define a new Wasserstein distance, inspired in the original one, looking forward to measure the distance between persistence diagrams. After making a brief introduction to Algebraic Topology and Topological Data Analysis, we will introduce the required concepts to define our new Wasserstein distance, we will check that it is actually a distance between persistence diagrams, and we will conclude with the main result of this thesis: the existence of an isometric embedding from a separable metric space into the space of persistence diagrams.

We will denote the strict upper triangular region of the Euclidean plane as $\mathbb{R}_{<}^2 := \{(x, y) \in \mathbb{R}^2 : x < y\}$, and the diagonal of the plane as $\Delta := \{(x, y) \in \mathbb{R}^2 : x = y\}$.

Definition 3.1 (Persistence diagram). Let I be a countable set. A *persistence diagram* is a function $D : I \rightarrow \mathbb{R}_{<}^2$.

Definition 3.2 (Partial matching). Let $D_1 : I_1 \rightarrow \mathbb{R}_{<}^2$ and $D_2 : I_2 \rightarrow \mathbb{R}_{<}^2$ be persistence diagrams. A *partial matching* between D_1 and D_2 is the triple (I'_1, I'_2, f) such that $f : I'_1 \rightarrow I'_2$ is a bijection with $I'_1 \subseteq I_1$ and $I'_2 \subseteq I_2$.

Instead of probability measures, now we are actually dealing with countable sets of points in \mathbb{R} . We will make use of the l^p norm at countable spaces to measure the distance between matched pairs and the distance between unmatched pairs and the diagonal Δ . For a more detailed explanation of Lebesgue measures check [4][Definition 3.7]. This norm is named after Pafnuty Chebyshev.

Definition 3.3 (Chebyshev distance). Let $a, b \in \mathbb{R}^2$ with $a = (a_x, a_y)$ and $b = (b_x, b_y)$. The *Chebyshev distance* is defined as

$$d_\infty(a, b) := \|a - b\|_\infty := \max\{|a_x - b_x|, |a_y - b_y|\}.$$

To define our adapted Wasserstein distance we need to check how Chebyshev distance measures distances between points of $\mathbb{R}_{<}^2$ and Δ .

Proposition 3.4. If $a = (a_x, a_y) \in \mathbb{R}_{<}^2$, then $d_\infty(a, \Delta) = \inf_{t \in \Delta} d_\infty(a, t) = \frac{a_y - a_x}{2}$.

Proof. The t which minimizes the distance is the midpoint of a_x and a_y , that is $t = (\frac{a_x + a_y}{2}, \frac{a_x + a_y}{2})$. Then,

$$\left| a_x - \frac{a_x + a_y}{2} \right| = \left| \frac{a_x - a_y}{2} \right| = \left| \frac{a_y - a_x}{2} \right| = \left| a_y - \frac{a_x + a_y}{2} \right|,$$

and as $a_y > a_x$ we have

$$d_\infty(a, t) = \left| \frac{a_y - a_x}{2} \right| = \frac{a_y - a_x}{2}.$$

□

We now verify that the upper triangular region of the Euclidean plane with the Chebyshev distance adapted to measure distances in Δ is a metric space.

Proposition 3.5. d_∞ is a distance in $\mathbb{R}_{<}^2$ with the diagonal Δ

Proof. For points $a, b \in \mathbb{R}_{<}^2 \subset \mathbb{R}^2$, d_∞ is a distance as usual Lebesgue norms are well defined. See [4][Chapter 3]. To verify that the metric requirements are fulfilled for $d_\infty(a, \Delta)$, it is enough to consider $t = \frac{a_y - a_x}{2}$ as in Proposition 3.4. □

Definition 3.6 (p -cost). Let $D_1 : I_1 \rightarrow \mathbb{R}_{<}^2$ and $D_2 : I_2 \rightarrow \mathbb{R}_{<}^2$ be persistence diagrams. Let (I'_1, I'_2, f) be a partial matching between them. If $p < \infty$, the p -cost of f is defined as

$$\begin{aligned} \text{cost}_p(f) := & \left(\sum_{i \in I'_1} d_\infty(D_1(i), D_2(f(i)))^p \right. \\ & + \sum_{i \in I_1 \setminus I'_1} d_\infty(D_1(i), \Delta)^p \\ & \left. + \sum_{i \in I_2 \setminus I'_2} d_\infty(D_2(i), \Delta)^p \right)^{\frac{1}{p}}. \end{aligned}$$

For $p = \infty$, the ∞ -cost of f is defined as

$$\text{cost}_\infty(f) := \max \left\{ \sup_{i \in I'_1} d_\infty(D_1(i), D_2(f_i)), \right. \\ \sup_{i \in I_1 \setminus I'_1} d_\infty(D_1(i), \Delta), \\ \left. \sup_{i \in I_2 \setminus I'_2} d_\infty(D_2(i), \Delta) \right\}.$$

Definition 3.7 (p -Wasserstein distance). Let D_1, D_2 be persistence diagrams. Let $1 \leq p \leq \infty$. Define

$$\tilde{\omega}_p(D_1, D_2) = \inf \{ \text{cost}_p(f) : f \text{ is a partial matching between } D_1 \text{ and } D_2 \}.$$

Let \emptyset denote the unique persistence diagram with empty indexing set. Let (Dgm_p, ω_p) be the space of persistence diagrams D that satisfy $\tilde{\omega}_p(D, \emptyset) < \infty$

modulo the equivalence relation $D_1 \sim D_2$ if $\tilde{\omega}_p(D_1, D_2) = 0$. The metric ω_p is called the p -Wasserstein distance.

Definition 3.8 (Bottleneck distance). In the conditions of Definition 3.7, if $p = \infty$, the metric ω_∞ is called the *bottleneck distance*.

Proposition 3.9. *There is only one matching between $D : I \rightarrow \mathbb{R}_{\leq}^2$ and \emptyset . Hence, if $p \leq \infty$,*

$$\tilde{\omega}_p(D, \emptyset) = \left(\sum_{i \in I} d_\infty(D(i), \Delta)^p \right)^{\frac{1}{p}},$$

and, if $p = \infty$,

$$\tilde{\omega}_\infty(D, \emptyset) = \sup_{i \in I} d_\infty(D(i), \Delta)$$

Proof. Let $I' \subseteq D$. If f is a partial matching between D and \emptyset , means that $f(I') = \emptyset$ is a bijection. That is only possible if $I' = \emptyset$ too. Therefore $I \setminus I' = I \setminus \emptyset = I$ and following Definition 3.6 we conclude our proof. \square

Next proposition will prove that, in indeed, the space of persistence diagrams with the p -Wasserstein distance (Dgm_p, ω_p) is a metric space. Its proof is usually omitted in literature, as it based on the simple fact that d_∞ is a distance. We will give, however, an step by step version here.

Proposition 3.10. ω_p is a distance on the space (Dgm_p, ω_p) .

Proof. Let $D_1, D_2, D_3 \in \text{Dgm}_p$, with $1 \leq p \leq \infty$.

First of all, $\omega_p(D_1, D_2) \geq 0$ because $d_\infty \geq 0$. $\omega_p(D_1, D_2) = 0$ if and only if $\tilde{\omega}_p(D_1, D_2) \geq 0$. Thus, because of the equivalence relationship used to define ω_p , it has to be $D_1 \sim D_2$.

To check symmetry, note that every partial matching f is bijective, therefore f^{-1} is a partial matching. But, for all $i \in I'_1$, exists $j \in I'_2$ such that $f(i) = j$ and

$$d_\infty(D_1(i), D_2(f(i))) = d_\infty(D_2(f(i)), D_1(i)) = d_\infty(D_2(j), D_1(f^{-1}(j))).$$

Then, $\text{cost}_p(f) = \text{cost}_p(f^{-1})$ and we have

$$\begin{aligned} \omega_p(D_1, D_2) &= \inf \{ \text{cost}_p(f) : f \text{ is a partial matching between } D_1 \text{ and } D_2 \} \\ &= \inf \{ \text{cost}_p(f^{-1}) : f^{-1} \text{ is a partial matching between } D_2 \text{ and } D_1 \} \\ &= \omega_p(D_2, D_1). \end{aligned}$$

Finally, let's prove the triangle inequality. If $f : I'_1 \rightarrow I'_2$ is a partial matching between D_1 and D_2 and $g : I'_2 \rightarrow I'_3$ is a partial matching between D_2 and D_3 , $g \circ f : I'_1 \rightarrow I'_3$ is a partial matching between D_1 and D_3 as both f and g are bijective. Computing the cost of the matchings for $p < \infty$, we notice that

$$\begin{aligned} & \sum_{i \in I'_1} d_\infty(D_1(i), D_2(f(i))) + \sum_{i \in I_1 \setminus I'_1} d_\infty(D_1(i), \Delta) + \sum_{i \in I_2 \setminus I'_2} d_\infty(D_2(i), \Delta) \\ & + \sum_{i \in I'_2} d_\infty(D_2(i), D_3(g(i))) + \sum_{i \in I_2 \setminus I'_2} d_\infty(D_2(i), \Delta) + \sum_{i \in I_3 \setminus I'_3} d_\infty(D_3(i), \Delta) \\ & \geq \sum_{i \in I'_1} d_\infty(D_1(i), D_3(g \circ f(i))) + \sum_{i \in I_1 \setminus I'_1} d_\infty(D_1(i), \Delta) + \sum_{i \in I_3 \setminus I'_3} d_\infty(D_3(i), \Delta) \end{aligned}$$

as $d_\infty(D_1(i), D_2(f(i))) + d_\infty(D_2(f(i)), D_3(g(f(i)))) \geq d_\infty(D_1(i), D_3(g \circ f(i)))$ using the triangle inequality of d_∞ . Therefore, for all partial matchings f and g as described, we have $\text{cost}_p(f) + \text{cost}_p(g) \geq \text{cost}_p(g \circ f)$. Using the same reasoning, for $p = \infty$ we also obtain $\text{cost}_\infty(f) + \text{cost}_\infty(g) \geq \text{cost}_\infty(g \circ f)$. Hence, we have verified that

$$\omega_p(D_1, D_2) + \omega_p(D_2, D_3) \geq \omega_p(D_1, D_3).$$

□

Definition 3.11 (Isometric embedding). Let $(X, d_X), (Y, d_Y)$ be metric spaces. An *isometric embedding* $\eta : (X, d_X) \rightarrow (Y, d_Y)$ is a mapping that satisfies

$$d_X(x_1, x_2) = d_Y(\eta(x_1), \eta(x_2))$$

for all $x_1, x_2 \in X$.

Definition 3.12 (Ball in persistence diagrams). Let $1 \leq p \leq \infty$. Let $D_0 \in \text{Dgm}_p$. The *ball* at the space of persistence diagrams is defined as $B_p(D_0, r) := \{D \in \text{Dgm}_p : w_p(D, D_0) < r\}$.

Theorem 3.13 (Isometric embedding of metric spaces into persistence diagrams). *Let (X, d) be a separable, bounded metric space. Then there exists an isometric embedding to the space of persistence diagrams $\eta : (X, d) \rightarrow (\text{Dgm}_\infty, \omega_\infty)$ such that $\eta(X) \subseteq B(\emptyset, \frac{3c}{c}) \setminus B(\emptyset, c)$.*

Proof. We will follow the procedure followed in [5][Theorem 19]. As (X, d) is bounded, we can let $c > \sup\{d(x, y) : x, y \in X\}$. As (X, d) is separable, we can

take $\{x_k\}_{k=1}^\infty$, a countable, dense subset of (X, d) . Consider

$$\begin{aligned}\eta : (X, d) &\rightarrow (\text{Dgm}_\infty, \omega_\infty) \\ x &\mapsto \{(2c(k-1), 2ck + d(x, x_k))\}_{k=1}^\infty\end{aligned}$$

For any $x \in X$ and $k \in \mathbb{N}$,

$$\begin{aligned}d_\infty((2c(k-1), 2ck + d(x, x_k)), \Delta) &= \frac{2ck + d(x, x_k) - 2c(k-1)}{2} \\ &= c + \frac{d(x, x_k)}{2} \\ &< c + \frac{c}{2} = \frac{3c}{2}.\end{aligned}$$

Because of Proposition 3.9, for every $x \in X$, $\tilde{\omega}_\infty(\eta(x), \emptyset) < \infty$ and η is well defined. Note that

$$\omega_\infty(\eta(x), \emptyset) = \sup_{1 \leq k < \infty} d_\infty((2c(k-1), 2ck + d(x, x_k)), \Delta),$$

so $\eta(x) \in B(\emptyset, \frac{3c}{c}) \setminus B(\emptyset, c)$.

Let $\eta(x)$ and $\eta(y)$ two equivalence classes of $(\text{Dgm}_\infty, \omega_\infty)$. Choose the representative diagrams $D_x : \mathbb{N} \rightarrow \mathbb{R}_<^2$ and $D_y : \mathbb{N} \rightarrow \mathbb{R}_<^2$ and consider the partial matching $(\mathbb{N}, \mathbb{N}, \text{id}_\mathbb{N})$. With it, for every $k \in \mathbb{N}$, $(2c(k-1), 2ck + d(x, x_k))$ is matched with $(2c(k-1), 2ck + d(y, x_k))$. The Chebyshev distance between those points is

$$\begin{aligned}d_\infty(D_x(k), D_y(k)) &= \max \{|2c(k-1) - 2c(k-1)|, \\ &\quad |2ck + d(x, x_k) - b_y - (2ck + d(y, x_k))|\} \\ &= \max\{0, |d(x, x_k) - d(y, x_k)|\} \\ &= |d(x, x_k) - d(y, x_k)|.\end{aligned}$$

Hence, because of the triangle inequality, the cost of this partial matching is

$$\text{cost}_\infty(\text{id}_\mathbb{N}) = \sup_k |d(x, x_k) - d(y, x_k)| \leq d(x, y).$$

Since $\{x_k\}_{k=1}^\infty$ is dense, for every $\epsilon > 0$, there exist a $k \in \mathbb{N}$ such that $d(x, x_k) \leq \epsilon$, so

$$\begin{aligned}|d(x, x_k) - d(y, x_k)| &\geq d(y, x_k) - d(x, x_k) \\ &= d(y, x_k) + d(x, x_k) - d(x, x_k) - d(x, x_k) \\ &\geq d(x, y) - 2d(x, x_k) \\ &> d(x, y) - 2\epsilon.\end{aligned}$$

Therefore, $\sup_k |d(x, x_k) - d(y, x_k)| \geq d(x, y)$ and

$$\text{cost}_\infty(\text{id}_\mathbb{N}) = \sup_k |d(x, x_k) - d(y, x_k)| = d(x, y).$$

Suppose $I, J \subseteq \mathbb{N}$ and (I, J, f) is a different partial matching between D_x and D_y . Then there exist a $k \in \mathbb{N}$ such that either $k \notin I$ or $k \in I$ and $f(k) = k \neq k$. If $k \notin I$, then

$$\text{cost}_\infty(f) \geq d_\infty((2c(k-1), 2ck + d(x, x_k)), \Delta) \geq c.$$

If $k \in I$ and $f(k) = k \neq k$, then

$$\text{cost}_\infty(f) \geq \|(2c(k-1), 2ck + d(x, x_k)) - (2c(k'-1), 2ck' + d(x, x_{k'}))\|_\infty \geq 2\epsilon.$$

Hence, $\text{cost}_\infty(f) \geq c > d(x, y)$ and $d(x, y) = \omega_\infty(\eta(x), \eta(y))$, proving that η is an isometric embedding of a metric space into the space of persistence diagrams. \square

References

- [1] A. Figalli and F. Glaudo, *An invitation to Optimal Transport, Wasserstein distances, and gradient flows.* EMS Press, 2020.
- [2] Y. Rubner, C. Tomasi, and L. J. Guibas, “The Earth Mover’s Distance as a metric for image retrieval,” *International Journal of Computer Vision*, 2000.
- [3] C. R. Givens and R. M. Shortt, “A class of Wasserstein metrics for probability distributions,” *Michigan Technological University*, 1984.
- [4] W. Rudin, *Real and Complex Analysis.* McGraw-Hill, 1987.
- [5] P. Bubenik and A. Wagner, “Embeddings of persistence diagrams into Hilbert spaces,” *Journal of Applied and Computational Topology*, 2020.