

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.

**Proposition 2.1.** *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*

$$\bar{\gamma} = \min \left\{ \gamma \in \Gamma(\mu, \nu) : \int_{X \times Y} c(x, y) d\gamma(x, y) \right\}.$$

*Proof.* (to do) □

We will denote the set of probability measures over a space  $X$  by  $\mathcal{P}(X)$ .

**Example 2.2** (Mean and variance in  $\mathbb{R}$ ).

**Definition 2.3.** 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 2.4.** *The definition of  $\mathcal{P}_p(X)$  is independent of the base point  $x_0$*

*Proof.* (to do) □

**Definition 2.5** ( $p$ -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 2.6.**  $W_p$  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 over persistence diagrams

The contents of this thesis are based on [1] and [3].

Along this text, 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** (Chebyshev distance). (To do)  $d_\infty := \max\{|a_x - b_x|, |a_y - b_y|\}$

**Proposition 3.3.** If  $a \in \mathbb{R}_{<}^2$ , then  $d_\infty(a, \Delta) = \inf_{t \in \Delta} d_\infty(a, t) = \frac{a_y - a_x}{2}$ .

*Proof.* (to do) □

**Proposition 3.4.** The upper triangular region of the Euclidean plane with the Chebyshev distance  $(\mathbb{R}_{<}^2, d_\infty)$  is a metric space.

*Proof.* (To do) □

**Definition 3.5** (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  $I'_2 \subseteq I_2$ .

**Definition 3.6.** 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

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

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

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

**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  $(\text{Dgm}_p, \omega_p)$  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  $\omega_p$  is called the *p-Wasserstein distance*.

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

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

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

*Proof.* (To do) □

**Proposition 3.10.** *The space of persistence diagrams with the p-Wasserstein distance  $(\text{Dgm}_p, \omega_p)$  is indeed a metric space.*

*Proof.* (To do) □

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

**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.* 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  $k \in \mathbb{N}$ ,

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

Because of Proposition 3.9, for every  $x \in X$ ,  $\omega_\infty(\eta(x), \emptyset) < \infty$  and  $\eta$  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

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

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  $\eta$  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,” 2000.
- [3] P. Bubenik and A. Wagner, “Embeddings of persistence diagrams into hilbert spaces,” 2020.