# The Optimal Transport Problem

## Master Thesis

Oscar Ramirez

# The Optimal Transport Problem

## Master Thesis

by

## Oscar Ramirez

to obtain the degree of Master of Science
in Mathematical Modelling and Engineering,
to be defended publicly on September, 2018.

Project duration:    September, 2016 – September, 2018
Thesis committee:    Prof. Juan Enrique Martinez Legaz,   UAB, supervisor

# Preface

The optimal transport problem was proposed by Gaspard Monge in 1781. Monge was interested in finding a way to transport a fixed amount of sand from one place into another, without losing or gaining mass in the process, such that the transportation cost is optimal. Later, Leonid Kantorovich introduced the proper framework to understand the problem in terms of probability measures, instead of using maps as Monge proposed the first time. The following work is a brief review of the theory developed for this problem, computation of a solution for a discrete optimal transport problem and two results of the developed theory that can be used in applications.

The literature in this topic is really rich and researchers in empirical sciences (such as economy, physics, meteorology, among many others) have found connection between the optimal transport problem and problems related to their respective branches, that at first glance seem disconnected from the transportation problem.

We have divided this text into five chapters. The first one is a reminder about concepts in topology, functional analysis, convex analysis and convergence of probability measures. These concepts are used to find the conditions that assure the existence of a minimizer for the optimal transport problem. If the reader is already familiarized with these concepts please skip directly to the third chapter. The second chapter is a reminder on important results in Linear programming, they are useful to compute a solution for the discrete version of the problem. In case the reader is already familiarized with linear programming, please skip directly to the fourth chapter. The third chapter is the heart of the text, it introduces the problem and we do a brief summary of th theory developed for the problem. The fourth chapter presents the algorithms used to compute a solution for a discrete version of the optimal transport. Finally, in the fifth chapter we discuss how the theory can be used as statistical distance between two probability measures, and an example of how a data assimilation problem can be also understood as an optimal transport problem.

*Oscar Ramirez*
*Barcelona, September 2018*

# Contents

# Notation Table.

| | |
|---|---|
| $\emptyset$ | Empty set |
| $\mathbb{R}$ | Real numbers field. |
| $\overline{\mathbb{R}}$ | $\mathbb{R} \cup \{\pm\infty\}$. That is $[-\infty, \infty]$ |
| $\mathbb{R}_+$ | The set of nonnegative real numbers, that is the interval $[0, \infty)$. |
| $\overline{\mathbb{R}}_+$ | The set of nonnegative extended real numbers, that is the interval $[0, \infty]$ |
| $2^A$ | The Power set of $A$, that is the set of all subsets in $A$. |
| $\mathbb{R}^d$ | The $d$-dimensional Euclidean space. |
| $B \backslash A$ | Given a set $B$ and a subset $A \subset B$, the set $B \backslash A$ is the complement of A in $B$. |
| $A \cup B$ | Union of two sets. |
| $A \cap B$ | Intersection of two sets. |
| $\text{proj}_X$ | Projection of $X$. Given two sets $X$ and $Y$ the projection of $X$ is the function $\text{proj}_X : X \times Y \to X$, defined by $\text{proj}_X(x, y) = x$. |
| $\text{proj}_Y$ | Projection of $Y$. Given two sets $X$ and $Y$ the projection of $Y$ is the function $\text{proj}_Y : X \times Y \to Y$, defined by $\text{proj}_Y(x, y) = y$. |
| $\mathbb{I}_A$ | Indicator function of a set $A \subset X$. If $x \in A$, then $\mathbb{I}_A(x) = 0$. If $x \in X \backslash A$, we have $\mathbb{I}_A(x) = \infty$. |
| $B(x; \epsilon)$ | Open ball with radius $\epsilon$ centered at $x$. |
| $\overline{B}(x; \epsilon)$ | Closed ball with radius $\epsilon$ centered at $x$. |
| $\text{supp}(f)$ | Support of a continuous function $f$. |
| $\text{epi}(f)$ | Epigraph of a function $f$. |
| $\text{graph}(f)$ | Graph of a function $f$. |
| $\partial f$ | Subdifferential of a function $f$. |
| $\text{id}$ | Identity map. That is $\text{id} : X \to X$, defined by $\text{id}(x) = x$. |
| l.s.c. | Lower semicontinuous. |
| $\frac{\delta F}{\delta \rho}$ | First variation of $F : \mathcal{P}(X) \to \mathbb{R}$, that is $\left. \frac{\mathrm{d}}{\mathrm{d}\epsilon} F(\rho + \epsilon \chi) \right|_{\epsilon=0} = \int \frac{\delta F}{\delta \rho} \mathrm{d}\chi$ |
| $\mathrm{D}T(x)$ | Jacobian matrix of a map $T(x)$. |
| $f_{|A}$ | The restriction of a function $f$ to a set $A$. |
| $\mathcal{M}(X)$ | Space of measures on $X$. |
| $\mathcal{M}_+(X)$ | Space of positive measures on $X$. |
| $\mathcal{P}(X)$ | Space of probabilities on $X$. |
| $\mu \ll \nu$ | The measure $\mu$ is absolutely continuous with respect to the measure $\nu$. |
| $\mu \llcorner A$ | A measure $\mu$ is restricted to a set $A$. |
| $\omega_d$ | The Measure of the unite ball in $\mathbb{R}^d$. |
| $T_\# \mu$ | The image measure (or pushforward measure) of $\mu$ through the map $T$. |
| $\text{spt}(\mu)$ | Support of a measure $\mu$. |
| i.i.d. | Independent and identical probability distributions. |
| $\mathcal{L}^p$ | Lebesgue measure on $\mathbb{R}^p$. |
| $\mathcal{H}^p \llcorner A$ | Hausdorff measure of dimension $p$ applied to some set $A \subset \mathbb{R}^d$. |
| $\delta_x$ | The Dirac (used as Kronecker in the discrete setting) mass at point x. |
| $\Pi(\mu, \nu)$ | The set of transport plans from $\mu$ to $\nu$. |
| $W_p$ | Wasserstein distance of order $p$. |
| $\mathbb{W}_p$ | Wasserstein space of order $p$. |
| $\gamma_T$ | The transport plan in $\Pi(\mu, \nu)$ induced by a map $T$. That is $\gamma_T = (\text{id}, T)_\# \mu$ and $T_\# \mu = \nu$. |
| $M(T)$ | Monge cost of a map $T$. |
| $K(\gamma)$ | Kantorovich cost of a plan $\gamma$. |
| $\mu \otimes \nu$ | The product measure of $\mu$ and $\nu$ such that $\mu \otimes \nu(A \times B) = \mu(A)\nu(B)$. |
| $\mathbf{M}^{k \times h}$ | The set of real matrices with $k$ rows and $h$ columns. |

$\mathbf{A}^\top$      Transpose of a matrix $\mathbf{A}$.

$\mathbf{A} \otimes \mathbf{B}$      Kronecker product between two matrices $\mathbf{A}$ and $\mathbf{B}$.

$\mathbf{f} \oplus \mathbf{g}$      The direct sum of two vectors $\mathbf{f}$ and $\mathbf{g}$, a matrix whose entries are given by the sum of the entries of $\mathbf{f}$ and $\mathbf{g}$, that is $(\mathbf{f} \oplus \mathbf{g})_{i,j} = (\mathbf{f})_i + (\mathbf{g})_j$.

$\mathbf{I}_n$      Identity matrix of dimension $n$.

<div align="right">

# 1

</div>

# Preliminaries.

We start this text reminding some definitions and theorems in topology, functional analysis and measure theory, needed to have a suitable framework to discuss the optimal transport problem and its applications. This chapter is meant to work as a brief summary of the definitions, motivations and results necessary to give the formalism of the upcoming chapters, adapted to the notation used in this project. For further the proofs and structured construction of the theory we will present the respective references.

## Topology.

We recall some important notions in topology important to discuss continuity and convergence, for further details about this topic we refer [23], [14].

A set $X$ endowed with a topology $\mathcal{T}$ is a **topological space**. The elements of a topology are called **open** sets. Any set of $X$ that is a complement of a set in $\mathcal{T}$ is called a **closed** set. We call a neighborhood of $x \in X$ to any set in $\mathcal{T}$ that contains $x$.

The **interior** of a set $A$, is defined as the biggest open set contained in $A$. Similarly, the **closure** of a set $A$, is defined as the smallest closed set containing $A$. We use indistinctly the notation $\mathrm{int}\,(A)$ and $A°$ for the interior of a set $A$. In the same way, for the closure we use the notation $\mathrm{clo}\,(A)$ or $\bar{A}$. We remark that a set is open if and only if $A = A°$, and a set is closed if and only if $A = \bar{A}$.

If $A$ is a subset of a topological space $X$, a point $x \in X$ is called a **limit point** of $A$ if every neighborhood of $x$ intersects $A$ in some point different than $x$ itself. A subset of topological space is closed if and only if contains all its limit points.
A subset $D$ of a topological space $X$ is **dense** in $X$ if for any point $x$ in $X$, any neighborhood of $x$ containing at least one point from $D$, different of $x$. Equivalently, $D$ is dense in $X$ if and only if it is identically to its closure in $X$, i.e. $D = \mathrm{clo}\,(D)$. A topological space is called **separable** if it contains a countable, dense subset.

Let $X$ be a topological space endowed with a topology $\mathcal{T}$. If $Z$ is a subset of $X$ the collection $\mathcal{T}_Z = \{Z \cap U : U \in \mathcal{T}\}$ is called the **subspace topology** of $Z$ and $Z$ is called a **subspace** of $X$. That is the topology of $Z$ is composed by all the intersections of the open sets of $X$ with $Z$.

**Definition 1.1** (Continuity). *Let $X$ and $Y$ be topological spaces. A function $f : X \to Y$ is said to be continuous if for each open subset $V$ of $Y$, the set $f^{-1}(V)$ is an open subset of $X$.*

Continuity not only depends upon the function $f$, but also in the topologies specified for the range and the domain of $f$. We say that $f$ is an **homeomorphism** if $f$ is bijective, continuous and $f^{-1}$ is also continuous. The support $\mathrm{supp}\,(f)$ of a continuous real valued function $f$ is the closure of the set $\{x \in X : f(x) \neq 0\}$.

A sequence $(x_n)_{n\in\mathbb{N}}$ in $X$ is a countable indexed set of elements in $X$. We say that a sequence converges to the point $x$ of $X$ if for each neighborhood $U$ of $x$ there is a positive integer $M \in \mathbb{N}$, such that $x_n \in U$, for all $n \geq M$. We use the notation $x_n \to x$, or $\lim_{n\to\infty} x_n = x$, to denote convergence of a sequence to a point $x$.

The notion of sequences is useful for tracking compactness in topological spaces satisfying the *first axiom of countability for topological spaces*[1]. Although for arbitrary topological spaces we can use the notion of nets. A **net** is a generalization of sequences for arbitrary topological spaces. A net in $X$ is a function $f$ from a directed set[2] $\mathcal{J}$ into $X$. If $\alpha \in \mathcal{J}$, we usually denote $f(\alpha) = x_\alpha$. We use the notation $(x_\alpha)_{\alpha\in\mathcal{J}}$ to refer a net. Every sequence is a net but the converse does not hold.

Topologies in which one element is not a closed set, or in which a sequence can converge to more than one point, are not really interesting for practical problems. If such things are allowed the theorems that one can prove are limited. To overcome this situation the mathematician Felix Hausdorff suggested topological spaces where for each pair of points we can find disjoint neighborhoods. We call a **Hausdorff space** a topological space $X$ endowed with a topology $\mathcal{T}$, such that for each pair $x_1$, $x_2$ of distinct points in $X$, there exists a neighborhood $U_1$ of $x_1$, and there exists a neighborhood $U_2$ of $x_2$, such that $U_1$ and $U_2$ are disjoint. In Hausdorff spaces every convergent net converges to at most one point.

A distance $d$ over $X$ is a nonnegative real valued function $d : X \times X \to \mathbb{R}$ satisfying, symmetry, triangle inequality and the property that the distance between any element to itself is zero. That is,

1. $d(x, y) = d(y, x)$, for all $x, y \in X$.

2. Symmetry: $d(x, y) = 0 \iff x = y$, for all $x, y \in X$.

3. Triangle inequality: $d(x, z) \leq d(x, y) + d(y, z)$ for all $x, y, z \in X$.

Given a set $X$ with distance $d$, and a real number $\epsilon > 0$, we call open ball the set,

$$B(x; \epsilon) = \{y | d(x, y) < \epsilon\}$$

The collection of all $\epsilon-$ open balls $B(x; \epsilon)$, for each $x \in X$ and $\epsilon > 0$, is a basis for the topology on $X$, called the metric topology. A set $U$ is open in a metric topology induced by $d$ if and only if for each $y \in U$, there is $\delta > 0$ such that $B(y; \delta) \subset U$.

A topological space $X$ is said to be **metrizable** if there exists a metric $d$ on the set $X$ that induces the topology of $X$. A **metric space** is a metrizable space $X$ together with a specific metric $d$. Every metrizable space satisfies the first axiom of countability.

Given a metric space $X$ the diameter diam of set $A \subset X$ is given by,

$$\text{diam}(A) = \sup\{d(x, y) : \quad x, y \in A\}.$$

We say that $A$ is a bounded subset of a metric space $X$ if there is $M \in \mathbb{R}_+$ such that $\text{diam}(A) < M$.

**Definition 1.2** (Continuity in metric spaces)**.** *The definition 1.1 for metric spaces is equivalent to say: A function $f : X \to Y$ defined from a metrizable space $(X, d_X)$ to a metrizable space $(Y, d_Y)$ is continuous if $\forall \epsilon \in \mathbb{R}$ and $\epsilon > 0$ there is $\delta \in \mathbb{R}$ and $\delta > 0$ such that,*

$$d_X(x, y) < \delta \implies d_Y(f(x), f(y)) < \epsilon$$

---

[1]Topological spaces with a countable topological basis at each of its points, i.e. for any $x \in X$ there is a countable collection $\mathcal{B}$ of neighborhoods of $x$, such that each neighborhood contains at least one of the elements of $\mathcal{B}$

[2]Any set with a partial order $\mathcal{J}$ relation

Let $X$ be topological space and $A \subset X$. If there is a sequence of points of $A$ converging to $x$, then $x \in \operatorname{clo} A$. The converse holds if $X$ is metrizable.

Let $f : X \to Y$ be continuous function between to topological spaces. If $(x_n)_{n \in \mathbb{N}}$ is a sequence converging to $x$, then the sequence $f(x_n)$ converges to $f(x)$. If $X$ satisfies the first axiom of countability we can invert the implication. Therefore for a metrizable space $X$ the converse holds.

A **cover** is a family $\mathcal{B}$ of subsets of $X$ such that the union of all its sets is equal to $X$. If all the sets in the family $\mathcal{A}$ are open we call it an open cover.

In topology we can find different notions of compactness. We start with the usual notion, a topological space $X$ is called **compact** if every open covering contains a finite subcover also covering $X$. A metrizable space $X$ is compact if and only if is sequentially compact. Every compact metric space is separable. The importance of the metrizability consists of the fact that all the topological properties can be characterized by sequences.

Every closed subspace of a compact space is compact. Every compact subspace of a Hausdorff space is closed. The image of a compact space under a continuous map is compact. The product of *finitely many* compact spaces is compact. We call **pre-compact** to any set whose closure is compact.

A metric space $X$ is said to be **totally bounded** if for every $\epsilon > 0$, there is a finite covering of $X$. A metric space $X$ is said to be **complete** if any Cauchy sequence[3] has a limit in $X$. Every compact metric space is complete. A metric space is compact if and only if it is complete and totally bounded.

We call **Polish space** is any topological space that is separable and completely metrizable. There is a subtle difference between complete metric space and completely metrizable space. And the difference lies on the words "*there exists at least a metric...*" in the completely metrizable definition, and "*given a metric*". Complete metrizable is a topological property while completeness is a property of the chosen metric.

Each closed subspace, and each open subspace, of a Polish space is Polish. The product of a finite or infinite sequence of Polish spaces is Polish. Every space in this text is Polish, unless the contrary is stated.

**Theorem 1.1** (Extreme value theorem). *Let $f : X \to \mathbb{R}$ be a continuous real valued function. If $X$ is compact, then there exist points $u$ and $v$ in $X$ such that $f(u) \leq f(x) \leq f(v)$ for all $x \in X$.*

We use the notation $C(X)$ to refer the set of real valued continuous functions $f : X \to \mathbb{R}$ defined on $X$. We denote by $C_b(X)$ the set of bounded continuous functions, that is $f \in C_b(X) \iff f \in C(x)$ and $\sup_{x \in X} |f(x)| < \infty$. We denote by $C_c(X)$ the set of continuous functions with compact support.

**Definition 1.3** (Uniform Continuity). *A function $f$ from the metric space $(X, d_X)$ to the metric space $(Y, d_Y)$ is said to be **uniformly continuous** if for every $\epsilon > 0$, and for every pair of points $x_0$, $x_1$ of $X$ there is a common $\delta > 0$ such that,*

$$d_X(x_0, x_1) < \delta \implies d_Y(f(x_0), f(x_1)) < \epsilon \tag{1.1}$$

If a function $f$ is continuous on a compact set, then $f$ is uniformly continuous.

A space $X$ is said to be **locally compact** at $x$ if there is some compact subspace $K$ of $X$ that contains a neighborhood of $x$. If $X$ is locally compact at each $x$ we call it just locally compact. Equivalently a topological space is locally compact if each of its points has an open neighborhood whose closure is compact. Any compact space is locally compact.

---

[3]Given a metric space $X$, a sequence $(x_n)_{n \in \mathbb{N}}$ is said to be Cauchy, if for every real $\epsilon > 0$ there is $N$ such that for all $m, n > N$ pair $d(x_m, x_n) < \epsilon$

On locally compact spaces $X$. We denote by $C_0(X)$ the set of continuous functions vanishing at infinity, i.e. $f \in C_0(X) \iff f \in C(X)$ and $\forall \epsilon > 0$, $\exists K \subset X$ such that $K$ is compact and $|f(x)| < \epsilon$, for all $x \in X \backslash K$.

We see that for any locally compact topological space $C_0(X) \subset C_b(X) \subset C(X)$. If $X$ is compact we have $C_0(X) = C_b(X) = C(X)$.

A subset of functions $F \subset C(X)$ is said to be **equicontinuous at** $x_0$ if given $\epsilon > 0$, there is a neighborhood $U$ of $x_0$ such that for all $x \in U$ and all $f \in F$, $d(f(x), f(x_0)) < \epsilon$. If $F$ is **equicontinuous** at all $x \in X$ we just call it equicontinuous. The subset $F \subset C(X)$ is **equibounded** if there is a common constant $M$ such that, $|f(x)| \leq M$ for all $f \in F$ and all $x \in X$.

**Theorem 1.2** (Arzelà-Ascoli). *Let $X$ be a compact metric space; The subset $K \subset C(X)$ is pre-compact in the uniform topology (the topology induced by the distance $d(f,g) = \sup_{\alpha \in X} |f(\alpha) - g(\alpha)|$), if and only if it is equicontinuous and equibounded.*

## Functional Analysis

Some notions in functional analysis are needed when we are trying to find optimal values for a given real valued function. We refer for further details of this section [1], [12], [18], [5].

The **epigraph** of a real valued function $f : X \to \overline{\mathbb{R}}$ is the set $\operatorname{epi}(f) = \{(x,t) \in X \times t : f(t) \leq t\}$. In similar way, we define the **graph** $\operatorname{graph}(f)$ of a function $f$ as the set defined by the expression $\operatorname{graph}(f) = \{(x,y) \in X \times Y : f(x) = y\}$.

Given a topological space $X$ we say that **limit inferior** of real valued function $f : X \to \overline{\mathbb{R}}$ as $x$ tends to $x_0$ is defined as

$$\liminf_{x \to x_0} f(x) := \sup_{\substack{x_0 \in A \\ A = \operatorname{int} A}} \left( \inf_{x \in A \backslash \{x_0\}} f(x) \right)$$

**Definition 1.4** (Lower Semicontinuous). *Let $X$ be a topological space and let $f : X \to \overline{\mathbb{R}}$.*

- *The function $f$ is lower semicontinuous (abbreviated as l.s.c.) if the set $\{x \in X : f(x) < t\}$ is closed for every $t \in \mathbb{R}$.*

- *The function $f$ is sequentially lower semicontinuous if the set $\{x \in X : f(x) \leq t\}$ is sequentially closed for every $t \in \mathbb{R}$.*

- *The function $f$ is upper semicontinuous if $-f$ is lower semicontinuous.*

The statement $f$ is lower semicontinuous is equivalent to the epigraph of $f$ is closed. Simultaneously, $f$ is lower semicontinuous if and only if $\forall x_0 \in X$, $f(x_0) \leq \liminf_{x \to x_0}$. In similar way, a function is sequentially lower semicontinuous if and only if its epigraph is sequentially closed. Moreover, $f$ is sequentially lower semicontinuous if and only if for each sequence $(x_n)_{n \in \mathbb{N}}$ converging to $x_0$, $f(x_0) \leq \liminf_n f(x_n)$.

Functions that are lower semicontinuous can be written as the pointwise supremum of a family of lower semicontinuous functions.

**Theorem 1.3.** *Let $X$ be a topological space and let $(f_\alpha)_{\alpha \in J}$ be a net (finite, countable or uncountable) of lower semicontinuous (or sequentially lower semicontinuous) functions defined on $X$ onto the extended real line, $f_\alpha : X \to \overline{\mathbb{R}}$. Then the function $f$ defined by,*

$$f = \sup_{\alpha \in J} f_\alpha$$

*is lower semicontinuous (sequentially lower semicontinuous). In addition if the family is finite, the function $f_-$ defined by,*

$$f_- = \min_{\alpha \in J} f_\alpha$$

*is also lower semicontinuous (sequentially lower semicontinuous).*

**Theorem 1.4** (Weierstrass). *Let $X$ be a topological space, let $K \subset X$ be compact (respectively sequentially compact), and let $f : X \to \overline{\mathbb{R}}$ be a lower semicontinuous (respectively sequentially lower semicontinuous) function. Then there exists $x_0 \in K$ such that,*

$$f(x_0) = \min_{x \in K} f(x)$$

We pay special attention to direct method of calculus of variations and Weierstrass criterion for minimizers.

A **linear space** $V$ over the scalar field $\mathbb{R}$ is a set of points, or vectors, on which are defined operations of vector addition and scalar multiplication with the following properties:

1. The set $V$ is a commutative group with respect to the operation $+$ of vector addition, that is for all $v, w, u \in X$, we have that $v + w = w + v$ and $v + (w + u) = (v + w) + u$, there is a zero vector $0$ such that $v + 0 = v$ for all $v \in V$, and for each $v \in V$ there is a unique vector $-v$ such that $v + (-v) = 0$.


2. For all $v, w \in V$ and $\alpha, \beta \in \mathbb{R}$, we have that $1v = v$, $\alpha(v + w) = \alpha v + \alpha w$, $\alpha(\beta x) = \beta(\alpha v) = (\alpha\beta)v$ and $(\alpha + \beta)v = \alpha v + \beta v$.

**Definition 1.5** (Norms and Seminorms). *A mapping $p : X \to \mathbb{R}$ is called a **seminorm** on a vector space $V$ if it has the following properties:*

- *$p(\lambda v) \le |\alpha| \, p(v) \; \forall v \in V$ and $\forall \alpha \in \mathbb{R}$, (absolutely scalable).*

- *$p(u + v) \le p(u) + p(v), \; \forall u, v \in V$, (subadditive).*

*These conditions imply that $p(x) \ge 0$ for all $x \in X$. If $p$ has the strong condition $p(x) > 0$ for $X \setminus \{0\}$ then $p$ is called a **norm**. If $p$ is a norm we usually use the symbol $\|\cdot\|$, to denote $\|x\| = p(x)$.*

A linear space is a metric space with a distance $d(v, w) = \|v - w\|$. A linear subspace is a subset of a linear space that is also a linear space. A complete vector space is called a **Banach space**.

**Theorem 1.5.** *(Tonelli's Method) Let $X$ be a normed space and let $f : V \to \overline{\mathbb{R}}$ be not identically to $\infty$. The Tonelli's method provides conditions on $V$ and $f$ that ensure the existence of a minimum point for $f$:*

1. *Consider a minimizing sequence $\{x_n\}_{n \in \mathbb{N}} \subset X$, that is a sequence such that*

$$\lim_{n \to \infty} f(v_n) = \inf_{v \in V} f(v)$$

2. *Prove that the sequence $(v_n)_{n \in \mathbb{N}}$ admits a subsequence $\{v_{n_k}\}_{n_k \in \mathbb{N}}$ that converges with respect to some topology $\mathcal{T}$ to some point $v_0 \in V$.*

3. *Establish sequential lower semicontinuity of $f$ with respect $\mathcal{T}$.*

4. *In the view of previous steps, conclude $v_0$ is a minimum of $f$ given that,*

$$\inf_{v \in V} f(v) = \lim_{n \to \infty} f(v_n) = \lim_{n_k \to \infty} f(v_{n_k}) \ge f(v_0) \ge \inf_{v \in V} f(v) \tag{1.2}$$

A map from a vector space $X$ onto a vector space $Y$ is said to be linear if $T(\alpha x + \beta y) = \alpha T(x) + \beta T(y)$ for all $x, y \in X$ and all $\alpha, \beta \in \mathbb{R}$. A linear operator $T$, from $X$ onto $Y$ is continuous if and only if is bounded, i.e. $T(M)$ is bounded for every bounded subset of $M \subset X$. The set $\mathcal{L}(X, Y)$ of linear bounded operators defined on a linear space $(X, \|\cdot\|_X)$ onto the linear space $(Y, \|\cdot\|_Y)$ becomes a linear normed space, with the norm defined by

$$\|T\| := \sup\left\{\|T(x)\|_Y \, ; \|x\|_X \le 1\right\} = \inf\{K : \, \|T(x)_Y \le K \|x\|_X\|, \; \forall x \in X\}.$$

If $Y$ is Banach then $L(X, Y)$ is also Banach. We call the dual of a Banach space $X$ the space $X^* = L(X, \mathbb{R})$ that is also Banach.

Duality.

Two linear spaces $X$ and $Y$, define a **dual system** if a fixed bilinear functional on their product is given $\langle \cdot, \cdot \rangle_{X,Y} : X \times Y \to \overline{\mathbb{R}}$. The dual system is called **separated** if for every $x \in X \backslash \{0\}$ there is $y \in Y$ such that $\langle x, y \rangle_{X,Y} \neq 0$ and for every $y \in Y \backslash \{0\}$ there is $x \in X$ such that $(x, y) \neq 0$.

For each $x$ we can define the application $T_x : Y \to \overline{\mathbb{R}}$ by $T_x(y) = \langle x, y \rangle_{X,Y}$, for all $y \in Y$, we see that $T_x$ is a linear functional on $Y$ and the mapping $x \mapsto T_x$, is linear and injective. Thus the elements of $X$ can be identified with the linear functionals on $Y$. In a similar way, the elements of $Y$ can be considered as linear functionals of $X$, identifying an element $y \in Y$ with $T_y(x) = \langle x, y \rangle_{X,Y}$.

Therefore, each dual system of linear spaces defines a mapping from either of the two linear spaces into the space of linear functionals on the other.

We define the topology $\sigma(X, Y)$ over $X$ in duality with $Y$ (or $Y$-topology of $X$), as the smallest topology which makes the linear functionals $T_y$ continuous. The roles of $X$ and $Y$ are interchangeable here.

A sequence $(x_n)_{n \in \mathbb{N}}$ in $X$ is $\sigma(X, Y)$-convergent to $x_0 \in X$ if and only if the sequence of real values $\left( \langle x_n, y \rangle_{X,Y} \right)_{n \in \mathbb{N}}$ converges to $\langle x_0, y \rangle_{X,Y} \in \mathbb{R}$, for each $y \in Y$. The roles of $X$ and $Y$ are interchangeable here, we have equivalent convention for $\sigma(Y, X)$-convergence.

We can construct a bilinear form $\langle \cdot, \cdot \rangle : X \times X^* \to \mathbb{R}$ defined by $\langle x, x^* \rangle = x^*(x)$, for $x^* \in X^*$ and $x \in X$. The weak topology is the smallest topology for which all the bounded linear functionals are continuous, that is $\sigma(X, X^*)$. We call this topology the weak topology of $X$.

In general, $X$ is not possible to identify a linear operator of $X^*$ with an element from $X$, then $X$ and $X^*$ do not always play a symmetric role in the bilinear form. When we can represent each element of $X^{**}$ with an element of $X$ we call the space $X$ reflexive.

Although, it is not possible to inject $X^{**}$ onto $X$, it is possible to inject $X$ onto $X^{**}$, since the map $x^* \mapsto x^*(x)$ defined over the dual $X^*$ is linear for all $x \in X$. Then we create a topology $\sigma(X^*, X)$ of $X^*$ in duality with $X$, we call this topology weak-star topology of $X^*$.

A bounded linear operator of normed space $X$ is continuous if and only if it is **weakly continuous**, that is continuous on the weak topology.

**Theorem 1.6** (Banach-Alaoglu)**.** *Let $X$ be an arbitrary normed space, and let $X^*$ be its dual. A closed unit ball $B \subset X^*$ is weak-star compact.*

**Corollary 1.1.** *The closed unit ball of the dual of a normed space is weak-star closed.*

If a sequence $(x_n)_{n \in \mathbb{N}} \subset X$ is weak convergent to $x_0 \in X$, then $\|x_0\| \leq \liminf_{n \to \infty} \|x_n\|$, if a sequence $(x_n^*)_{n \in \mathbb{N}} \subset X^*$ is weak-star convergent to $x_0^* \in X^*$, then $\|x_0^*\| \leq \liminf_{n \to \infty} \|x_n^*\|$.

**Theorem 1.7** (Banach-Alaoglu weak-star sequentially compactness)**.** *Let $X$ be a separable normed vector space, and let $X^*$ be its dual. Then the weak-star topology on a closed ball $B \subset X^*$ is metrizable.*

**Corollary 1.2.** *Let $X$ be a separable normed vector space and let $X^*$ be its topological dual space. Then every bounded sequence $(\phi_n)_{n \in \mathbb{N}} \in X$ has a weak-star convergent subsequence.*

## Convex Analysis.

Convexity plays an important role in optimization. Convex sets and functions have nice properties that can be exploit in the development of theorems in optimization, for further references please check [1], [10].

A subset $E$ of a vector space $V$ is convex if for all $v_1$, $v_2$ in $E$ and all $\alpha \in (0, 1)$ we have that $\alpha v_1 + (1 - \alpha) v_2 = v \in E$. In other words, $E$ contains all line segments between points in $E$. We denote by $[v_1, v_2]$, a line segment with extreme points $v_1, v_2$. A point $x$ in a convex set $C$ is said to be an **extreme point** of $C$ if there are no two distinct points $x_1$ and $x_2$ in $C$ such that $x = \alpha x_1 + (1 - \alpha) x_2$ for some $0 < \alpha < 1$.

A subset of the linear space $X$ is said to be **affine** set if whenever it contains $x_1$ and $x_2$ it also contains $\alpha x_1 + \beta x_2$ for arbitrary $\alpha, \beta \in \mathbb{R}$ satisfying $\alpha + \beta = 1$.

The intersection of many arbitrary convex (affine) sets is again convex (affine) set. A function $f : V \to \overline{\mathbb{R}}$ is convex if and only if epi$(f)$ is a convex set.

A **hyperplane** is a subset $H \subset X$ of a linear space $X$, identified by a pair $(f, k)$ such that $f$ is a real valued linear functional defined on $X$ and $k \in \mathbb{R}$, such that

$$H = \{x \in X : f(x) = k\} \tag{1.3}$$

Let $\alpha \in [0, 1]$, a function $f : V \to \overline{\mathbb{R}}$ defined on a vector space $V$ is said to be **convex** if $f(\alpha v_1 + (1 - \alpha)v_2) \le \alpha f(v_1) + (1 - \alpha)f(v_2)$. We call it **strictly convex** if the inequality holds strictly. We call the function **proper** if it is convex and does not take the value $-\infty$ and it is not equally to $\infty$. We call a function $f$ **concave** if $-f$ is convex.

Note that taking $v_1 \ne v_2$ and $v_2 = 0$, if $f$ is convex for all $v_1$ then $f(\alpha v_1) \le \alpha f(v_1)$ for all $v_1 \in V$; the inequality inverts if $f$ is concave. Every norm is a convex function. If $f$ is concave we have that $f(\alpha v_1 + (1 - \alpha)v_2) \ge \alpha f(v_1) + (1 - \alpha)f(v_2)$, for all $t \in [0, 1]$.

Given any nonempty subset $A$ of $X$, the function $\mathbb{I}_A$ on $X$, defined by

$$\mathbb{I}_A(x) = \begin{cases} 0 & \text{if } x \in A \\ \infty & \text{if } x \ne A \end{cases}$$

is called the indicator function of $A$. The subset $A$ of X is convex if and only if its indicator function $\mathbb{I}_A$ is convex.

A particular class of topological linear spaces with richer properties is the class of **locally convex spaces**; these are topological linear spaces with the property that for every element there exists a base of neighborhoods consisting of convex sets. A locally convex space is Hausdorff if and only if it has a separated family of seminorms.

A convex set is closed if and only if it is weakly closed. The closed unit ball of a normed space is weakly closed.

**Theorem 1.8** (Hahn-Banach theorem). *Let $X$ be a real linear space, let $p$ be a real convex function on $X$ and let $Y$ be a linear subspace of $X$. If a linear functional $T_0$ defined on $Y$ satisfies,*

$$T_0(y) = p(y), \quad \forall y \in Y$$

*then $T_0$ can be extended to a linear functional $T$ defined on all of $X$, satisfying*

$$T(x) \le p(x), \quad \forall x \in X$$

**Theorem 1.9** (Geometric version of Hahn-Banach theorem). *If $A$ is a convex set with a nonempty interior and if $M$ is an affine set which contains no interior point of $A$, then there exists a closed hyperplane which contains $M$ and which again contains no interior point of $A$.*

For a hyperplane $H \subset X$ characterized by a pair $(f, k)$ composed of a linear functional $f$ and a real value $k$, we have two open half spaces, $\{x \in X : f(x) < k\}$, $\{x \in X : f(x) > k\}$. The implications of the above Hanh-Banach theorems, is that a convex set which contains no point of a hyperplane is contained in one of the two open half spaces determined by that hyperplane. Indeed, if $f(x_1) > k$ and $f(x_2) < k$, it exists $\alpha \in [0, 1]$ such that $f(\alpha x_1 + (1 - \alpha)x_2) = k$, hence $x_1$ and $x_2$ cannot be contained in a convex set which is disjoint from the hyperplane $f(x) = k$.

A lower-semicontinuous, proper and convex function $f$ on a reflexive Banach space $X$ takes a minimum value on every bounded, convex and closed subset $M$ of $X$. In other words, $\exists x_0 \in M$ such that

$$f(x_0) = \inf\{f(x) : x \in M\}. \tag{1.4}$$

Moreover, $x_0$ is unique if $f$ is strictly convex. The last equation is also known as a **convex program**.

On the other hand, in relation with concavity we have the modulus of continuity of a continuous function. The following definition presents some examples of characterizing a function according to an intrinsic property.

**Definition 1.6.** *Let $X$ be a metric space with a distance $d_X$ and let $Z \subset X$. A function $f : X \to \mathbb{R}$ onto a metric space $Y$ with a distance $d_Y$ is said to be,*

- ***Lipschitz continuous*** *if*

$$\mathrm{Lip}(f; Z) := \sup \left\{ \frac{d_Y(f(x), f(y))}{d_X(x,y)} : x, y \in Z, x \neq y \right\} < \infty$$

- ***locally Lipschitz continuous*** *if for every compact $K \subset Z$,*

$$\sup \left\{ \frac{d_Y(f(x), f(y))}{d_X(x,y)} : v, w \in K, v \neq w \right\} < \infty$$

- ***Hölder continuous*** *with exponent[4] $0 < \alpha < 1$ if,*

$$\|f\|_{C^{0,\alpha}(E)} := \sup \left\{ \frac{d_Y(f(x), f(y))}{d_X(x,y)} : x, y \in E, x \neq y \right\} < \infty.$$

Given a subset $Z$ of a metric space $X$ and a function $f : Z \to \mathbb{R}$ the **modulus of continuity** $\omega : \mathbb{R}_+ \to \mathbb{R}_+$ defined by $\omega(\delta; f) = \sup\{|f(x) - f(y)| : x, y \in Z, d(x,y) < \delta\}$. Note that $f$ is uniformly continuous if and only if $\omega(\delta; f) \to 0$ as $\delta \downarrow 0$, that is continuous at zero. The modulus of continuity of uniformly continuous function is subadditive. The **concave modulus of continuity** $\overline{\omega} : \mathbb{R}_+ \to \mathbb{R}_+$ is defined as the smallest concave function above $\omega$. If $K = X$ and $f$ is uniformly continuous, then it is always possible to replace $\omega$ with $\overline{\omega}$. For further details of the modulus of continuity we refer [30], [15] and [12].

**Theorem 1.10.** *Let $E$ be a subset of a metric space $X$. Any uniformly continuous bounded function $f : E \to \mathbb{R}$ can be extended to a function $g : X \to \mathbb{R}$ with the same modulus of continuity, supremum and infimum.*

A subset $A \subset X \times X^*$ is called monotone if $\langle x_1 - x_2, y_1 - y_2 \rangle \geq 0$, for any segment $[x_i, y_i] \in A$.

**Definition 1.7** (Convex Conjugate). *Let $f$ a function defined on a linear space $X$, $f : X \to \overline{\mathbb{R}}$. The function $f : X^* \to \overline{\mathbb{R}}$ defined by,*

$$f^*(x^*) = \sup\{\langle x, x^* \rangle - f(x); \ x \in X\}, \quad x^* \in X^*, \tag{1.5}$$

*is called the conjugate function of $f$.*

Any convex, proper and lower-semicontinuous function is bounded from below by an affine function. A lower-semicontinuous convex function is proper if and only its conjugate is proper.

**Theorem 1.11** (Convex Envelope Theorem). *Let $X$ be a locally convex topological space. And let $f : X \to \mathbb{R} \cup \{+\infty\}$ be any functional nonidentically to $+\infty$. Then $f^{**} = f$ if and only if $f$ is convex and lower semicontinuous on $X$.*

A proper function $f$ is convex and lower-semicontinuous on $X$ if and only if it is the supremum of a family of affine continuous functions.

**Definition 1.8** (Subdifferential). *Let $X$ be a real Banach space and let be $X^*$ be its dual with a norm $\|\cdot\|$. Given the proper convex function $f : X \to \mathbb{R} \cup \{\infty\}$, the **subdifferential** of $f$ is the multivalued mapping $\partial f : X \to 2^{X^*}$*

$$\partial f(x) = \{x^* \in X^*; \quad f(x) - f(u) \leq \langle x - u, x^* \rangle, \ \forall u \in X.\}$$

**Theorem 1.12.** *Let $f : X \to \mathbb{R} \cup \{\infty\}$ be a proper convex function. Then, then, $x^* \in \partial f(x) \iff f(x) + f^*(x^*) \leq \langle x, x^* \rangle \iff f(x) + f^*(x^*) = \langle x, x^* \rangle$*

Let $f$ proper convex functional defined on a Banach Space $X$, the minimum points of $f$ are just the solutions to the equation $0 \in \partial f(x)$.

---

[4]Note that we have introduced a norm notation.

## Convergence of Probability measures.

The framework that Kantorovich proposed for the optimal transport problem requires aspects in measure theory and convergence of probability measures. We refer for further details in this section [6],[3], [12] and [21].

Let $X$ be a set, endowed with a $\sigma$-algebra $\mathscr{A}(X)$, we call such a pair a **measurable space**. The Borel $\sigma$-algebra $\mathscr{B}(X)$ of $X$ is the $\sigma$-algebra constructed with all the open sets of the topology of $X$. A **measure** is a function $\mu$ from the $\sigma-$algebra to $\overline{\mathbb{R}}$, satisfying non-negativity, null-empty set ($\mu(\emptyset) = 0$) and countable additivity. A **Radon measure** is a measure on $\mathscr{B}(X)$ of a Hausdorff topological space $X$ that is finite on all compact sets.

**Definition 1.9** (Signed measure). *Let $(X, \mathscr{B}(X))$ be a measurable space. A signed measure is a function $\lambda : \mathscr{B}(X) \to [-\infty, \infty]$, such that,*

1. *$\lambda(\emptyset) = 0$;*

2. *$\lambda$ takes at most one of the two values $\infty$ and $-\infty$ that is, either one of these definitions holds: $\lambda : \mathscr{B}(X) \to (-\infty, \infty]$, or $\lambda : \mathscr{B}(X) \to [-\infty, \infty)$*

3. *For every countable collection $(E_i) \subset \mathscr{B}(X)$ of pairwise disjoint sets we have,*

$$\lambda\left(\bigcup_{n=1}^{\infty} E_n\right) = \sum_{n=1}^{\infty} \lambda(E_n) \tag{1.6}$$

A set function $\lambda : \mathscr{B}(X) \to [-\infty, \infty]$ is a signed measure if and only if it satisfies the second condition of the definition 1.9, it is finitely additive and for every increasing sequence $E_n \subset \mathscr{B}(X)$, $\lambda(\cup_{n=1}^{\infty} E_n) = \lim_{n \to \infty} \lambda(E_n)$.

**Definition 1.10.** *Let $(X, \mathscr{B}(X))$ be a measurable space and let $\gamma$ be a signed measure. A set $E \in \mathscr{B}(X)$ is said to be positive (respectively negative) if $\lambda(F) \geq 0$ (respectively $\lambda(F) \leq 0$) for all $F \subset E$ with $F \in \mathscr{B}(X)$*

If $\lambda(E)$ is bounded, then there exists a positive measurable subset $F \subset E$ with $\lambda(F) > 0$. Moreover, the space can be decomposed as $X = X^+ \cup X^-$, where $X^-$ is negative and $X^+$ is positive. This result is known as the Hahn decomposition theorem. If $(X^+, X_1^+)$ and $(X^-, X_1^-)$ are decompositions of a signed measure $\lambda$, then $\lambda(E) = 0$ for every $E \subset \mathscr{B}(X)$ satisfying

$$E \subset ((X^+ \backslash X_1^+) \cup (X_1^+ \backslash X^+)) \cup ((X^- \backslash X_1^-) \cup (X_1^- \backslash X^-))$$

If $X$ is measurable and $\lambda$ is a signed measure we may uniquely decompose $X$ as $X = X^+ \cup X^-$, where $X^+$ is positive and $X^-$ is negative. For $E \in \mathscr{B}(X)$ define,

$$\lambda^+(E) := \lambda(E \cap X), \quad \lambda^-(E) := -\lambda(E \cap X)$$

Then $\lambda^+$, $\lambda^-$ are measures with at least one of them finite. Moreover $\lambda(E) = \lambda^+(E) - \lambda^-(E)$ and for every $E = \mathscr{B}(X)$

$$\lambda^+(E) = \sup\{\lambda(F) : F \subset E, F \in \mathscr{B}(X), F \text{ is positive}\}$$
$$\lambda^-(E) = -\inf\{\lambda(F) : F \subset E, F \in \mathscr{B}(X), F \text{ is negative}\}$$

The measures are called, respectively, the upper and lower variation of $\lambda$, while the measure

$$\|\lambda\| := \lambda^+ + \lambda^-$$

is called the total variation of $\lambda$. For every $E \in \mathscr{B}(X)$ we have that,

$$\|\lambda\|(E) = \sup\left\{\sum_{n=1}^{\infty} |\lambda(E_n)| : \quad \{E_n\} \subset \mathscr{B}(X) \text{ partition of } E\right\} \tag{1.7}$$

**Theorem 1.13** (Jordan decomposition theorem). *Let $(X, \mathscr{B}(X))$ be a measurable space and let $\lambda$ : $\mathscr{B}(X) \to [-\infty, \infty]$ be a signed measure. Then there exists a unique pair $(\lambda^+, \lambda^-)$ such that $\lambda = \lambda^+ - \lambda^-$.*

A **measurable function** is a function between two measurable spaces such that the preimage of any measurable set is measurable. Let $u$ be a measurable function, $\lambda$ a signed measure, and a measurable set $E \subset \mathscr{B}(X)$, if at least one of the two integrals $\int_E u(x) \mathrm{d}\lambda^+$ or $\int_E u(x) \mathrm{d}\lambda^-$ is finite the Lebesgue integral of $u$ over the measurable set $E$ is defined as

$$\int u \mathrm{d}\lambda := \int_E u \mathrm{d}\lambda^+ - \int_E \mathrm{d}\lambda^-. \tag{1.8}$$

Given two signed measures, $\lambda$ and $\iota$, then $\lambda$ is said to be absolutely continuous with respect to $\iota$, denoted by $\lambda << \iota$, if $\lambda(E) = 0$ whenever $\iota(E) = 0$. If $\lambda << \iota$, there exists a measurable function $u$, such that $\iota(E) = \int_E u \mathrm{d}\lambda$, for any $E$ in $\mathscr{B}(X)$.

A finitely additive signed measure is a function $\lambda : \mathscr{B}(X) \to \mathbb{R}$ such that,

1. $\lambda(\emptyset) = 0$,

2. $\lambda$ is finitely additive, that is ,

$$\lambda(E_1 \cup E_2) = \lambda(E_1) + \lambda(E_2)$$

   for all $E_1, E_2 \in \mathscr{B}(X)$ with $E_1 \cap E_2 = \emptyset$.

3. $\lambda$ is bounded, that is, its total variation norm is finite,

$$\|\lambda\|(X) := \sup \left\{ \sum_{n=1}^{l} |\lambda(E_n)| : \{E_n\} \subset \mathscr{B}(X) \text{finite partition of } X \right\} < \infty. \tag{1.9}$$

The set of finitely additive signed measures is Banach space endowed with a norm $\|\cdot\|(X)$. We denote by $\mathcal{M}(X)$ the set of all finite signed measures on $\mathscr{B}(X)$. The set of all nonnegative, bounded, finitely additive and **regular** measures on $\mathscr{B}(X)$ is denoted by $\mathcal{M}_+(X)$.

**Definition 1.11** (Measure of vectors). *Let $X$ a measurable space. A set function $\lambda = (\lambda_1, ..., \lambda_m)$ : $\mathscr{B}(X) \to \mathbb{R}^m$ is a vectorial measure if each component $\lambda_i$, with $i = 1, ..., m$ is a signed measure. The total variation of $\lambda$ is defined by,*
$$\|\lambda\| := \|\lambda_1\| + ... \|\lambda_m\|.$$

If $X$ is a compact topological space, all the signed measures the regular, finitely additive and bounded measures are countably additive.

**Theorem 1.14** (Riesz representation theorem in $C_b(X)$). *Let $X$ a Polish space. Then every bounded linear functional $L : C_b(X) \to \mathbb{R}$ is represented by a unique regular, finitely additive and bounded measure $\lambda$, in the sense that*

$$L(u) = \int_X u \mathrm{d}\lambda, \quad \forall u \in C_b(X) \tag{1.10}$$

*Moreover, the norm of $L$ coincides with the total variation norm $\|\lambda\|(X)$. Conversely, every functional of the form* (1.10)*, where $\lambda$ is a bounded linear functional on $C_b(X)$*

Let $X$ be a compact Hausdorff space and let $C(X) = C_b(X)$ be the space of all continuous functions, then the dual of $C(X)$ can be identified with $\mathcal{M}(X)$. Let $X$ be a separable metric space, then the weak-star topology on $\mathcal{M}_+(X)$ is metrizable. If $X$ is a compact metric space, any bounded sequence of regular, finitely additive and bounded measures $(\lambda_n)_{n \in \mathbb{N}}$ is sequentially weakly-star compact. In particular if $u \in C_b(X)$ then,

$$\lim_{k \to \infty} \int_X u \mathrm{d}\lambda_{n_k} = \int_X u \mathrm{d}\lambda,$$

for a subsequence $(\lambda_{n_k})_{n_k \in \mathbb{N}}$ and for some $\lambda$ being bounded, regular and finitely additve measures.

**Theorem 1.15** (Luizin). *Let $X$ be a locally compact Hausdorff space, let $\mathcal{A}$ be a $\sigma-$algebra on $X$ containing the Borel $\sigma-$algebra of $X$. Let $\mu$ a regular measure on $(X, \mathcal{A})$ and let $f : X \to \mathbb{R}$ be $\mathcal{A}$-measurable. If $A$ belongs to $\mathcal{A}$ and satisfies $\mu(A) < \infty$ and given a positive number $\epsilon > 0$, then there is a compact $K \subset X$ such that $\mu(X \backslash K) < \epsilon$ and the restriction $f_{|K}$ is continuous in $K$. Moreover, there is a function $g \in C_c(K)$ such that $g(x) = f(x)$ for each $x \in K$.*

It follows that the supremum of a collection of continuous (or lower semicontinuous) functions is lower semicontinuous and that each lower semicontinuous function on a Hausdorff space is Borel measurable. A positive measure $\mu$ is said to be a probability measure if $\mu(X) = 1$. The space of probability measures $\mathcal{P}(X)$ is not closed in $\mathcal{M}(X)$.

A set $\Pi$ of probability measures is precompact if any sequence of probability measures $\mu_n \in \Pi$ contains a subsequence $(\mu_{n_k})_{n_k \in \mathbb{N}}$ which converges weakly-star to a probability measure in $\mathcal{P}(X)$.

We say that a probability measure $\mu$ on $X$ is **tight** if for any $\mu \in \mathcal{P}(X)$ and any positive real number $\epsilon > 0$ there is a compact $K \subset X$ such that $\mu(X \backslash K) \le \epsilon$. We also say that a set of probabilities $\Gamma$ is tight if $\forall \mu \in \Gamma$, $\mu$ is tight.

**Theorem 1.16.** *If $X$ is a complete and separable topological space, then $\mathcal{P}(X)$ is tight. If $X$ is compact any set of probability measures is pre-compact*

**Theorem 1.17** (Prohorov). *Let $X$ be a metric space and let $\{\mu_n\}$ be a sequence of Borel measures, $\mu_n : \mathscr{B}(X) \to \mathbb{R}_+$, such that $\sup \mu_n(X) < \infty$. Assume that the elements of the sequence is tight. Then there exist a subsequence $\{\mu_k\}$ of $\{\mu_n\}$ and a Borel measure $\mu : \mathscr{B}(X) \to \mathbb{R}$, such that $\mu_{n_k}$ converges weakly-star to $\mu$ in the dual of the continuous bounded functions of $X$. In particular if $X$ is a locally compact metric space, $\mu, \mu_n$ are finite Radon measures and $\mu_n$ converges weakly-star to $\mu \in \mathcal{M}(X)$, then $\mu_n$ converges weakly-star in the $C_b(X)^*$ if and only if the sequence is tight.*

## Image Measure.

We state some definitions useful to prove the density of transport plans induced by maps in the set of transportation plans (Lemma 3.5). The proof for the following theorems and lemmas, as well as for Lemma 3.5, are detailed in [29].

**Definition 1.12** (Support of a measure.). *Given a separable metric space $X$, the support of a measure $\mu$ is defined as the smallest closed set on which $\mu$ is concentrated, that is*

$$\mathrm{spt}(\mu) := \bigcap_{\substack{\mu(X \backslash A) = 0 \\ A = \mathrm{clo}\, A}} A \tag{1.11}$$

**Definition 1.13** (Image Measure). *Let $(X, \mathcal{A}_X)$ and $(Y, \mathcal{A}_Y)$ be two measurable spaces. Let $T : X \to Y$ be a measurable map from $X$ to $Y$. Let $\mu$ be a measure $\mu : \mathcal{A}_X \to \overline{\mathbb{R}}_+$, then the image measure (or pushforward measure) $T_\# \mu : \mathcal{A}_Y \to \overline{\mathbb{R}}_+$ is given by,*

$$T_\# \mu(B) = \mu\left(T^{-1}(B)\right), \quad \forall B \in \mathcal{A}_Y.$$

**Lemma 1.1.** *If $\mu$, $\nu$ are two probability measures on the real line $\mathbb{R}$ and $\mu$ is atomless, then there exists at least a map $T$ such that $T_\# \mu = \nu$.*

**Lemma 1.2.** *There exists a Borel map $\sigma_d : \mathbb{R}^d \to \mathbb{R}$ which is injective, its image is a Borel subset of $\mathbb{R}$, and its inverse map is Borel measurable as well.*

**Theorem 1.18.** *If $\mu$ and $\nu$ are two probability measures on $\mathbb{R}^d$ and $\mu$ is atomless, then there exists at least a map $T$ such that $T_\# \mu = \nu$.*

**Theorem 1.19.** *Consider a compact metric space $X$, endowed with a probability measure $\rho \in \mathcal{P}(X)$, a sequence of partitions $G_n$ such that each $G_n = (C_\alpha)_{\alpha \in \mathcal{J}_n}$ is a family of disjoint subsets satisfying $\bigcup_{i \in \mathcal{J}_n} C_{i,n} = X$ for every $n$. Suppose that $size(G_n) := \max_i \left(\mathrm{diam}\left(C_{i,n}\right)\right)$ tends to $0$ as $n \to \infty$ and consider a sequence of probability measures $\rho_n$ on $X$ such that, for every $n$ and $i \in I_p$, we have $\rho_n(C_{i,n})$. Then $\rho_n \rightharpoonup \rho$.*

# 2

# Linear Programming

Linear programming is a well studied branch in mathematics. It deals with optimization of linear functions under linear constraints. The study of linear programming started during the second part of the 1940s, as a technique to solve military oriented problems. We dedicate one chapter to present important results for the finite dimensional part of this branch, since the discrete version of the optimal transport problem can be seen as linear program. The following results are taken from [20].

We can formulate a finite dimensional linear programming problem in its general form as follows:

**Problem 1.** *Given a cost vector $\mathbf{c} \in \mathbb{R}^n$, a linear operator $\mathbf{A} \in \mathbf{M}^{m \times n}$, the problem consists in finding $\mathbf{x} \in \mathbb{R}^n$ such that*

$$\min \qquad \mathbf{c}^\top \mathbf{x} \qquad\qquad (2.1)$$
$$\textit{subject to} \qquad \mathbf{A}\mathbf{x} = \mathbf{b} \qquad\qquad (2.2)$$
$$\mathbf{x} \geq 0 \qquad\qquad (2.3)$$

*We refer to this formulation as the **primal**.*

Where $\mathbf{A}$ is a $m \times n$ matrix, and $\mathbf{b} \in \mathbb{R}^m$ is an m-dimensional column vector. The vector inequality $\mathbf{x} \geq 0$ means that each component is nonnegative. This problem has a solution if $n > m$.

**Definition 2.1** (Basic solutions)**.** *Given the set of $m$ simultaneous linear equations* (2.2) *with $n$ unknowns, let $\mathbf{B}$ be any nonsingular $m \times m$ submatrix made up of columns of $\mathbf{A}$. Then if all $n - m$ components of $\mathbf{x}$ not associated with columns of $\mathbf{B}$ are set equal to zero, the solution to the resulting set of equations is said to be a basic solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$, with respect to the basis $\mathbf{B}$. The components of $\mathbf{x}$ associated with columns of $\mathbf{B}$ are called **basic variables**.*

We assume that the $m$ rows of $\mathbf{A}$ are linearly independent and $m < n$. Under this assumption the problem have at least one basic solution.

**Definition 2.2** (Degenerated basic solutions)**.** *If one or more of the basic variables in a basic solution have value zero, is said to be a **degenerated basic solution**.*

**Definition 2.3** (Feasible solutions.)**.** *A vector $\mathbf{x}$ satisfying the constraints* (2.2) *and* (2.4) *is said to be feasible. A feasible solution that is also basic is said to be a **basic feasible** solution. If the solution is also a degenerated basic solution, it is called a **degenerated basic feasible** solution.*

**Theorem 2.1** (Fundamental theorem of linear programming.)**.** *Given a linear program in the standard form* (2.1)*,* (2.2) *and* (2.3) *where $\mathbf{A}$ is a $m \times n$ matrix of rank $m$,*

- *If there is a feasible solution, there is a basic feasible solution.*

- *If there is an optimal solution, there is an optimal basic feasible solution.*

For a problem having $n$ variables and $m$ constraints there are at most

$$\binom{n}{m} = \frac{n!}{m!\,(n-m)!}$$

basic solutions, the fundamental theorem of linear programming simplifies the problem to a finite number of possibilities. This is a powerful theoretical result, but practical represents an inefficient method to find an optimal solution. This result has an interesting connection to convexity since we are finding the optimal points in the faces of a convex polytope.

**Theorem 2.2.** *Let $\mathbf{A}$ be an $m \times n$ matrix of rank $m$ and $\mathbf{b}$ an $m$-vector. Let $K$ be the convex polytope consisting of all $n$-vectors $\mathbf{x}$ satisfying*

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b} \\ \mathbf{x} &\geq 0 \end{aligned} \tag{2.4}$$

*A vector $\mathbf{x}$ is an extreme point of $K$ if and only if $\mathbf{x}$ is a basic feasible solution of* (2.4).

**Corollary 2.1.** *If the convex set $K$ corresponding to* (2.4) *is nonempty, it has at least one extreme point.*

**Corollary 2.2.** *If there is a finite optimal solution to a linear programming problem, there is a finite optimal solution which is an extreme point of the constraint set.*

**Corollary 2.3.** *The constraint set $K$ corresponding to* (2.4) *possesses at most a finite number of extreme points.*

*Proof.* There is only a finite number of basic solutions generated by selecting $m$ basis vectors and $n$ columns of $\mathbf{A}$. The extreme points of $K$ are a subset of the basic solutions. $\quad\square$

**Corollary 2.4.** *If the convex polytope $K$ corresponding to* (2.4) *is bounded, then $K$ is a convex polyhedron. That is, $K$ consists of points that are convex combinations of a finite number of points.*

## Simplex Method.

The idea of the simplex method is to proceed from one basic feasible solution that belongs to the constraint set of a linear program in standard form to another, in such a way to decrease the value of the objective function continually until a minimum is reached.

Pivoting in a set of simultaneous linear equations is crucial for the development of the algorithm. Remember that the matrix $\mathbf{A}$ has $m$ rows and $n$ columns. Let us write the constraint $\mathbf{Ax} = \mathbf{b}$ as follows,

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n = \mathbf{b}$$

Where $\mathbf{a}_i$ are $m$-dimensional column vectors of the matrix $\mathbf{A}$, for integers $1 \leq i \leq n$. We try to find an expression for $\mathbf{b}$ as a linear combination of the vectors $\mathbf{a}_j$.

If $m < n$ and the vectors $\mathbf{a}_i$ span the space $\mathbb{R}^m$, then the representation of $\mathbf{b}$ using column vectors of $\mathbf{A}$ is not unique but a whole family of different representations. However, $\mathbf{b}$ has a unique representation of $m$ linear independent vectors $\mathbf{a}_j$.

Moreover, every vector $\mathbf{a}_j$, with $1 \leq j \leq n$ can be expressed as a linear combination of these basis vectors,

$$\mathbf{a}_j = y_{1,j}\mathbf{a}_1 + y_{2,j}\mathbf{a}_2 + \cdots + y_{m,j}\mathbf{a}_m$$

Without loss of generality we can say that the first $m$ column vectors are linearly independent and therefore they form a basis for $\mathbb{R}^m$. We see that if $\mathbf{a}_j$ is a member of a basis, implies that $y_{j,j} = 1$ and the coefficients $y_{i,j} = 0$ for $i \neq j$. We can use the following tableau to represent the coefficients,

| $\mathbf{a}_1$ | $\mathbf{a}_2$ | ... | $\mathbf{a}_m$ | $\mathbf{a}_{m+1}$ | $\mathbf{a}_{m+2}$ | ... | $\mathbf{a}_n$ | $\mathbf{b}$ |
|---|---|---|---|---|---|---|---|---|
| 1 | 0 | ... | 0 | $y_{1,m+1}$ | $y_{1,m+2}$ | ... | $y_{1,n}$ | $y_{1,0}$ |
| 0 | 1 | ... | 0 | $y_{2,m+1}$ | $y_{2,m+2}$ | ... | $y_{2,n}$ | $y_{2,0}$ |
| $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ |
| 0 | 0 | ... | 1 | $y_{m,m+1}$ | $y_{m,m+2}$ | ... | $y_{m,m}$ | $y_{m,0}$ |

$$\tag{2.5}$$

For simplicity we consider $y_{0,j}$ the representation for $\mathbf{b}$. Consider the process of changing a vector of the basis by another one. Take $\mathbf{a}_k$, with $1 \leq k \leq m$, and we want to substitute it by a vector $\mathbf{a}_l$, with $m + 1 \leq l \leq n$.

Since any vector $\mathbf{a}_j$ can be expressed in terms of the old basis we have,

$$\mathbf{a}_l = y_{kl}\mathbf{a}_k + \sum_{\substack{i=1 \\ i \neq k}}^{m} y_{il}\mathbf{a}_i.$$

From the above we solve for $\mathbf{a}_k$,

$$\mathbf{a}_k = \frac{1}{y_{kl}}\mathbf{a}_l - \sum_{\substack{i=1 \\ i \neq k}}^{m} \frac{y_{il}}{y_{kl}}\mathbf{a}_i.$$

Then we substitute $\mathbf{a}_k$ in the linear combination of the old basis for $\mathbf{a}_j$ by the above equation,

$$\mathbf{a}_j = \frac{y_{kj}}{y_{kl}}\mathbf{a}_l + \sum_{\substack{i=1 \\ i \neq k}}^{m} \left( y_{ij} - \frac{y_{il}}{y_{kl}} \right)\mathbf{a}_i.$$

Therefore, we write a new tableau for the system using the following set of equations,

$$\begin{cases} y'_{k,j} = \frac{y_{k,j}}{y_{k,l}} \\ y'_{i,j} = y_{i,j} - \frac{y_{i,l}}{y_{k,l}} & \text{for } i \neq k \text{ and } 0 \leq i \leq n. \end{cases} \tag{2.6}$$

We can generate a new basic solution from an old one, by the mean of pivoting vectors as explained above . The problem is that the nonnegative constraint can be violated after pivoting operations. Therefore, it is required to control the pair of variables whose roles are going to be interchanged, in order to take in account only basic feasible solutions.

The fundamental theorem of the linear programming shows that it is only necessary to consider basic feasible solutions of the problem. Until this moment we have not considered the possibility of having as result of the pivoting process a degenerated basic feasible solution.

For the sake of simplicity, we assume that every basic feasible solution is non-degenerated. This assumption simplifies the description of the simplex method, however the arguments can be modified to include degenerated basic feasible solutions.

Suppose we have the basic feasible solution $\mathbf{x} = (x_1, x_2, \ldots, x_m, 0, 0, \ldots, 0)$. We are assuming non-degeneracy of the solutions, therefore $x_i > 0$ for $i = 1, \ldots, n$.

Imagine we want to introduce in the basis the vector $\mathbf{a}_l$, with $l > m$. Since the vectors $\mathbf{a}_i$, for $i = 1 \ldots m$, form a basis. Manipulating the basis representation for $\mathbf{b}$ and $\mathbf{a}_l$,

$$\mathbf{b} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_m\mathbf{a}_m + \epsilon\mathbf{a}_l - \epsilon\mathbf{a}_l$$
$$= \left( x_1 - \epsilon y_{1,l} \right)\mathbf{a}_1 + \left( x_2 - \epsilon y_{2,l} \right)\mathbf{a}_2 + \cdots + \left( x_m - \epsilon y_{m,l} \right)\mathbf{a}_m + \epsilon\mathbf{a}_l$$

For simplicity take $\epsilon \geq 0$. Now we have a $m + 1$ representation for $\mathbf{b}$, we see that for $\epsilon = 0$ we have the old basis representation. We are trying to generate a new basic feasible solution, then we set the value of $\epsilon$,

$$\epsilon = \min_{1 \leq i \leq m} \left\{ \frac{x_i}{y_{i,l}} : \quad y_{i,l} > 0 \right\}$$

.

If the minimum is achieved by more than one single index, the new solution is degenerated and any of the vectors with zero component can be regarded as the one leaving the basis.

If all $y_{i,l} \leq 0$ no new basic feasible solution can be obtained. However, we can obtain feasible solutions with arbitrarily large coefficients. That is, the set of feasible solutions is unbounded.

Hence, given a basic feasible solution and arbitrary column vector $\mathbf{a}_l$ of $\mathbf{A}$. We can find either a new basic feasible solution with $\mathbf{a}_l$ as part of its basis and one of the old vectors removed from it, or a set of unbounded feasible solutions.

In summary, the assumption that the coefficients $y_{1,0}, \ldots, y_{m,0}$ are nonnegative, implies that $x_1 = y_{1,0}, x_2 = y_{2,0}, \ldots, x_m = y_{m,0}$ is feasible. We substitute a vector already in the basis by a vector $\mathbf{a}_l$, in such a way that the solution of the new generated coefficients are feasible. We take the smallest ratio to keep the feasibility. In this way we can introduce $\mathbf{a}_l$ as part of the basis creating a new basic feasible solution.

Assume that $\mathbf{A}$ can be written as follows,

$$\mathbf{A} = [\mathbf{B}, \mathbf{D}] \tag{2.7}$$

where $\mathbf{B}$ consists of the first $m$ columns of $\mathbf{A}$ corresponding to the basic variables. These columns are linearly independent and they form a basis for $\mathbb{R}^m$. The matrix $\mathbf{D}$ is a sub-matrix of $\mathbf{A}$ representing the rest of the columns of $\mathbf{A}$.

In order to write the problem in an appropriate way, we write $\mathbf{x}$ and $\mathbf{c}$ as follows,

$$\mathbf{x} = (\mathbf{x_B}, \mathbf{x_D}), \qquad \mathbf{c} = (\mathbf{c_B}, \mathbf{c_D}) \tag{2.8}$$

Where $\mathbf{x_B}$ has $m$ entries and $\mathbf{x_D}$ has $n - m$ entries; in similar way for $\mathbf{c_B}$ and $\mathbf{c_D}$. Then, our primal problem can be written as follows,

$$\begin{aligned} \min \quad & \mathbf{c_B^\top x_B} + \mathbf{c_D^\top x_D} \\ \text{subject to} \quad & \mathbf{Bx_B} + \mathbf{Dx_D} \\ & \mathbf{x_B} \geq 0, \quad \mathbf{x_D} \geq 0 \end{aligned}$$

If $\mathbf{x}$ is a basic feasible solution, the corresponding value is give by

$$z_0 = \mathbf{c_B^T x_B}$$

We can construct nonbasic feasible solution, setting arbitrary values for $\mathbf{x_D} = (x_{m+1}, x_{m+2}, \ldots, x_n)$ and solving for each $x_i$, with $1 \leq i \leq m$,

$$x_i = y_{i,0} - \sum_{j=m+1}^{n} y_{i,j} x_j$$

Let $z$ be a real number given by,

$$z = \mathbf{c^\top x} = z_0 + (c_{m+1} - z_{m+1}) x_{m+1} + (c_{m+2} - z_{m+2}) x_{m+2} + \cdots + (c_n - z_n) x_n. \tag{2.9}$$

where,

$$z_j = y_{1,j} c_1 + y_{2,j} c_2 + \cdots + y_{m,j} c_m, \qquad \text{for } m + 1 \leq j \leq n. \tag{2.10}$$

From this equation, we can determine if there is any advantage in introducing to the basis one of the nonbasic variables.

**Theorem 2.3** (Improvement of basic feasible solution). *Given a non-degenerated basic feasible solution with corresponding objective value $z_0$, suppose that for there is $j$, such that $c_j - z_j < 0$ holds.*

*Then there is a feasible solution with objective value $z < z_0$. If the column $\mathbf{a}_j$ can be substituted for some vector in the original basis to yield a new basic feasible solution, this new solution will have $z < z_0$. If $\mathbf{a}_j$ cannot be substituted to yield a basic feasible solution, then the feasible solutions are unbounded and the objective function can be made arbitrarily small.*

*Proof.* Consider equations (2.9) and (2.10), if $c_j - z_j$ is negative for some $j$, $m + 1 \leq j \leq n$, then changing $x_j$ from zero to a positive value decreases the total cost $z$. Let $(x_1, x_2, \ldots, x_m, 0, \ldots, 0)$ a basic feasible solution with $z_0$ and suppose $c_{m+1} - z_{m+1} < 0$. New feasible solutions can be constructed of

the form $\left(x_1', x_2', \dots, x_{m+1}', 0, 0, \dots, 0\right)$, with $x_m' > 0$, substituting this new solution into equation (2.9) we obtain,

$$z - z_0 = (c_m + 1 - z)x_{m+1}' < 0$$

Hence $z < z_0$ for any such solution. It is clear that we desire to make $x_{m+1}'$ as large as possible. As $x_{m+1}'$ is increased, the other components change their values. Thus $x_{m+1}'$ can be increased until one $x_i' = 0$, for $i \leq m$ in which case we obtain a new basic feasible solution. If no variable $x_i'$ decreases, $x_{m+1}'$ can be increased without bound indicating an unbound solution set and an objective value without lower bound. □

If at any stage $c_j - z_j < 0$ for some $j$, it is possible to make $x_j$ positive and decrease the objective function.

**Theorem 2.4** (Optimality Condition Theorem). *If for some basic feasible solution $c_j - z_j \geq 0$ for all $j$, then that solution is optimal.*

*Proof.* This result comes from equation (2.9), since any other feasible solution must have $x_i \geq 0$ for all $i$, and hence the value $z$ of the objective will satisfy $z - z_0 \geq 0$. □

Since the main role in this method is played by the constants $c_j - z_j$ we refer them as the **relative cost coefficients** and we use the notation $r_j = c_j - z_j$. These coefficients measure the cost of a variable relative to a basis.

We can summarize the simplex algorithm in the following steps:

1. Construct a tableau (2.5) corresponding to a basic feasible solution and compute the relative cost coefficients $r_j$. For this purpose we can use row reduction.

2. If $r_j \geq 0$ for all $j$, then the current basic feasible solution is optimal; stop.

3. Select $l$ such that $r_l < 0$ to determine which nonbasic variable is to become basic.

4. Calculate the ratios $y_{i,0}/y_{i,l}$ for $y_{i,l} > 0$, $i = 1, 2, \dots, m$. If no $y_{i,l} > 0$. then problem is unbounded; stop. Otherwise, select $k$ as the index $i$ corresponding to the minimum ratio.

5. Apply the pivoting procedure to introduce $\mathbf{a}_l$ substituting $\mathbf{a}_k$ in the basis. Return to step 1.

Assuming non-degeneracy it is possible to prove in an easy way the convergence of the algorithm. The process stops only if optimality or unboundedness is discovered. If the algorithm does not find neither optimality or unboundedness, the objective value is strictly decreased. Since there are only a finite number of possible basic feasible solutions, and the basis do not repeat because of the strictly decrease of the objective value. The algorithm must reach a basis satisfying one of the two terminating conditions.

Revised simplex method.
The revised simplex method is a scheme for ordering the computations required by the simplex method, so that unnecessary calculations are avoided. A basic solution has the form $\mathbf{x} = (\mathbf{x_B}, \mathbf{0})$, where $\mathbf{x_B} = \mathbf{B}^{-1}\mathbf{b}$.

For any $\mathbf{x_D}$ the necessary value of $\mathbf{x_B}$ as follows,

$$\mathbf{x_B} = \mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{D}\mathbf{x_D}$$

Therefore, we substitute the above equation in the cost expression,

$$z = \mathbf{c_B^\top}\left(\mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{D}\mathbf{x_D}\right) + \mathbf{c_D^T}\mathbf{x_D}$$
$$= \mathbf{c_B^\top}\mathbf{B}^{-1}\mathbf{b} + \left(\mathbf{c_D^\top} - \mathbf{c_B^\top}\mathbf{B}^{-1}\mathbf{D}\right)\mathbf{x_D}$$

Thus, the vector $\mathbf{r_D} = \mathbf{c_D^\top} - \mathbf{c_B}\mathbf{B}^{-1}\mathbf{D}$ is the relative cost for non-basic variables. The components of this vector are used to determine which vector we will bring into the basis.

In summary,

1. Calculate the current relative cost coefficients $\mathbf{r_D} = \mathbf{c_D^\top} - \mathbf{c_B}\mathbf{B^{-1}D}$. It is more efficient and numerically stable to solve the linear system $\mathbf{v}^\top = \mathbf{c_B^\top}\mathbf{B^{-1}}$, then compute the relative vector $\mathbf{r_D} = \mathbf{c_D^\top} - \mathbf{v}^\top\mathbf{B^{-1}D}$. If $\mathbf{r_D} \geq \mathbf{0}$ then the current solution is optimal, stop.

2. Determine the vector $\mathbf{a}_l$ is to enter the basis by selecting the most negative cost coefficient, and calculate $\mathbf{q} = \mathbf{B^{-1}a}_q$ which gives the vector $\mathbf{a}_q$ in terms of the current basis.

3. If no $y_{i,l} > 0$ then the problem is unbounded; stop. Otherwise calculate the ratios $y_{i,l}/y_{i,l} > 0$ to determine which vector is to leave the basis.

4. Update $\mathbf{B^{-1}}$ and the current solution $\mathbf{B^{-1}}b$. Return to step 1.

## Duality

**Problem 2.** *Given a cost vector $\mathbf{c} \in \mathbb{R}^n$, a linear operator $\mathbf{A} \in M^{m \times n}$ and a column vector. We say that the dual for the primal formulation 1 is given by,*

$$\text{max} \qquad\qquad \boldsymbol{\lambda}^\top \mathbf{b} \qquad\qquad\qquad (2.11)$$
$$\textit{subject to} \qquad\qquad \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top \qquad\qquad\qquad (2.12)$$

**Lemma 2.1** (Weak Duality lemma). *If $\mathbf{x}$ and $\boldsymbol{\lambda}$ are feasible for (2.2) and (2.12), respectively then $\mathbf{c}^\top\mathbf{x} \geq \boldsymbol{\lambda}^\top\mathbf{b}$.*

*Proof.* We see that following inequality holds for equations (2.2), (2.12) and the cone $\mathbf{x} \geq 0$,

$$\boldsymbol{\lambda}^\top\mathbf{b} = \boldsymbol{\lambda}^\top(\mathbf{Ax}) \leq \mathbf{c}^\top\mathbf{x}$$

$\square$

**Corollary 2.5.** *If $\mathbf{x}_0$ and $\boldsymbol{\lambda}_0$ are feasible for the constraints (2.2) and (2.12) respectively and $\mathbf{c}^\top\mathbf{x}_0 = \boldsymbol{\lambda}_0^\top\mathbf{b}$, then $\mathbf{x}_0$ and $\boldsymbol{\lambda}_0$ are optimal for their respective problems.*

This corollary is the result of the Weak Duality lemma. A feasible vector to the primal problem yields an upper bound on the value of the dual problem. In the other hand, a feasible vector to the dual problem yields a lower bound on the value of the primal problem. The values associated with the primal problem are all larger than the values associated with the dual problem. We see that having a feasible pair $\mathbf{x}_0$ and $\boldsymbol{\lambda}_0$ for their respective problems, satisfying the equality means that each problem has reached its optimal value.

**Theorem 2.5** (Duality Theorem). *If the problem (1) has a finite optimal solution then the dual formulation (2) also does. In the same manner, if the dual problem (2) has solution then the primal also does. Moreover, the corresponding values of the objective functions are equal. If either problem has an unbounded objective solution, the other problem has no feasible solution.*

*Proof.* We see from corollary 2.5 that the first condition holds. If the primal is unbounded and $\boldsymbol{\lambda}$ is feasible for the dual we must have, $\boldsymbol{\lambda}\mathbf{b} \leq -M$ for arbitrarily large $M$, leading to a contradiction.

Suppose that the primal problem has a finite optimal solution with value $z_0$. In the space $\mathbb{R}^{m+1}$ define the convex set

$$C = \left\{(r, \mathbf{w}) : r = \alpha z_0 - \mathbf{c}^\top\mathbf{x}, \ \mathbf{w} = \alpha\mathbf{b} - \mathbf{Ax}, \ \mathbf{x} \geq \mathbf{0}, \alpha \geq 0\right\}$$

We see that $C$ is a closed cone convex cone. We need to find a point $(\bar{r}, \bar{\mathbf{w}}) \notin C$, in order to apply the Hahn–Banach separation theorem to prove the existence of a vector $\boldsymbol{\lambda} \in \mathbb{R}^m$ satisfying a condition that allow us to introduce the Weak Duality lemma. Our proposition is the point $(1, \mathbf{0}) \notin C$. We see that, for $\alpha > 0$ and $\mathbf{w} = 0$, $\mathbf{w} = \alpha\mathbf{b} - \mathbf{Ax}_0 = 0$ with $\mathbf{x}_0 \geq 0$, then $\mathbf{x} = \mathbf{x}_0/\alpha$ is feasible for the primal problem.

Hence $r/\alpha = z_0 - \mathbf{c}^\top\mathbf{x} \leq 0$, implying that $r \leq 0$. For $\alpha = 0$, we have $\mathbf{w} = -\mathbf{Ax}_0 = \mathbf{0}$ with $\mathbf{x}_0 \geq \mathbf{0}$ and $\mathbf{c}^\top\mathbf{x}_0 = -1$. If $\mathbf{x}$ is any feasible solution to the primal implies $\mathbf{x} + \beta\mathbf{x}_0$ is feasible for $\beta \geq 0$ and therefore we can obtain an objective value as small as we want, contradicting the fact that the primal has a bounded solution. Therefore, $(1, \mathbf{0}) \notin C$. By the Hahn-Banach's separation theorem we can find

a hyperplane separating $C$ and $(1, \mathbf{0})$. Thus we can find a non zero vector $(s, \boldsymbol{\lambda}) \in \mathbb{R}^{m+1}$ and constant $c$ satisfying

$$s < c = \inf\{sr + \boldsymbol{\lambda}^\mathsf{T}\mathbf{w} : (r, \mathbf{w}) \in C\}.$$

Since $C$ is a cone, it follows that $c \geq 0$. Imagine we have a point $(\tilde{r}, \tilde{\mathbf{w}}) \in C$, such that $s\tilde{r} + \boldsymbol{\lambda}^\mathsf{T}\tilde{\mathbf{w}} < 0$, then for $\beta > 0$ big enough we the point $(\beta\tilde{r}, \beta\tilde{\mathbf{w}})$ can violate the hyperplane inequality. In the other hand, $(0, \mathbf{0}) \in C$, then $c = 0$. Thus, $s < 0$ without loss of generality we can take $s = -1$, resulting for any $(r, \mathbf{w}) \in C$,

$$-r + \boldsymbol{\lambda}^\mathsf{T}\mathbf{w} \geq 0 \tag{2.13}$$

We proved the existence of $\boldsymbol{\lambda} \in \mathbb{R}^m$ holding the above inequality. Using the definition of $C$,

$$\left(\mathbf{c} - \boldsymbol{\lambda}^\mathsf{T}\mathbf{A}\right) - \alpha z_0 + \alpha\boldsymbol{\lambda}^\mathsf{T}\mathbf{b} \geq 0 \tag{2.14}$$

for all $\mathbb{x} \geq \mathbf{0}$, and $\alpha \geq 0$. Setting $\alpha = 0$ we have the inequality $\boldsymbol{\lambda}^\mathsf{T} \leq \mathbf{c}^\mathsf{T}$, which says $\boldsymbol{\lambda}$ is feasible for the dual. Setting $\mathbf{x} = \mathbf{0}$ and $\alpha = 1$ results in $\boldsymbol{\lambda}^\mathsf{T}\mathbf{b} \geq z_0$. Therefore, by means of the Weak Duality lemma we have that $\boldsymbol{\lambda}^\mathsf{T}\mathbf{b} = z_0$ and by corollary 2.5 we have that $\boldsymbol{\lambda}$ is optimal for the dual. □

# Complementary Slackness.

**Theorem 2.6.** *Let $\mathbf{x}$ and $\boldsymbol{\lambda}$ be feasible solutions for the primal and dual programs, respectively. A necessary and sufficient condition that they both be optimal solutions is that for all $i$*

- $x_i > 0 \implies \boldsymbol{\lambda}^\mathsf{T}\mathbf{a}_i = c_i$,

- $x_i = 0 \impliedby \boldsymbol{\lambda}^\mathsf{T}\mathbf{a}_i < c_i$.

*Proof.* If the above conditions hold, then $\left(\boldsymbol{\lambda}^\mathsf{T}\mathbf{A} - \mathbf{c}^\mathsf{T}\right)\mathbf{x} = 0$. By the means of the Weak Duality lemma and corollary 2.5, $\boldsymbol{\lambda}^\mathsf{T}\mathbf{b} = \mathbf{c}^\mathsf{T}\mathbf{x}$ implies two solutions are optimal. Conversely, if the two solutions are optimal, by the means of the Duality Theorem $\boldsymbol{\lambda}^\mathsf{T}\mathbf{b} = \mathbf{c}^\mathsf{T}\mathbf{x}$. Since each component of $\mathbf{x}$ is nonnegative and each component of $\boldsymbol{\lambda}^\mathsf{T}\mathbf{A} - \mathbf{c}^\mathsf{T}$ is nonpositive, and the above conditions must hold. □

### The Dual simplex Method.

For general linear programs the dual simplex method is most frequently used, since it is more efficient to work with the dual tableau instead, making use of the complementary slackness conditions to recover the primal solution.

Given a linear program in its primal form, let $\mathbf{B}$ a basis such that a vector $\boldsymbol{\lambda}$ defined by $\boldsymbol{\lambda}^\mathsf{T} = \mathbf{c}_\mathbf{B}^\mathsf{T}\mathbf{B}^{-1}$ is feasible for the dual. We say that $\mathbf{x}_\mathbf{B} = \mathbf{B}^{-1}\mathbf{b}$ is **dual feasible**. If $\mathbf{x}_\mathbf{B} \geq 0$ then this solution is also primal feasible. The vector $\boldsymbol{\lambda}$ is feasible for the dual, therefore it satisfies $\lambda_j \leq c_j$ for all $j = 1, 2, \ldots, n$. Without loss of generality assume the first columns of $\mathbf{A}$ form a the basis, then

$$\boldsymbol{\lambda}^\mathsf{T}\mathbf{a}_j = c_j, \quad \text{for } j = 1, \ldots, m \tag{2.15}$$

Assuming non-degeneracy there is an inequality,

$$\boldsymbol{\lambda}^\mathsf{T}\mathbf{a}_j < c_j, \quad \text{for } j = m+1, \ldots, n \tag{2.16}$$

We find a new $\bar{\boldsymbol{\lambda}}$ in which the inequality becomes an equality and vice-versa. We need to find $\bar{\boldsymbol{\lambda}}$ such that increases the objective value $\mathbf{b}^\mathsf{T}\boldsymbol{\lambda}$. Let $\boldsymbol{\beta}^i$ the $i$-th row of $\mathbf{B}^{-1}$, note that $\boldsymbol{\beta}^i\mathbf{a}_j = y_{i,j}$. Setting,

$$\bar{\boldsymbol{\lambda}}^\mathsf{T} = \boldsymbol{\lambda}^\mathsf{T} - \epsilon\boldsymbol{\beta}^i \tag{2.17}$$

we have $\bar{\boldsymbol{\lambda}}^\mathsf{T}\mathbf{a}_j = \boldsymbol{\lambda}^\mathsf{T} - \epsilon\boldsymbol{\beta}^i\mathbf{a}_j$.

$$\bar{\boldsymbol{\lambda}}^\mathsf{T}\mathbf{a}_j = c_j \qquad\qquad \text{for } j = 1, 2, \ldots, m \text{ and } i \neq j \tag{2.18}$$

$$\bar{\boldsymbol{\lambda}}^\mathsf{T}\mathbf{a}_j = c_j - \epsilon y_{i,j} \qquad\qquad \text{for } j = m+1, m+2, \ldots, n \tag{2.19}$$

$$\bar{\boldsymbol{\lambda}}^\mathsf{T}\mathbf{a}_i = c_i - \epsilon \tag{2.20}$$

In the other hand,

$$\bar{\boldsymbol{\lambda}}^\top = \boldsymbol{\lambda}^\top \mathbf{b} - \epsilon\, (\mathbf{x_B})_i \tag{2.21}$$

The idea behind the algorithm is:

1. We start with a dual feasible solution $\mathbf{x_B}$, if $\mathbf{x_B} \geq 0$ the solution is optimal. $\mathbf{x_B}$ is not nonnegative we can find $i$ such that the $i$-th component of $\mathbf{x_B}$ is less than zero, i.e. $(\mathbf{x_B})_i < 0$.

2. If all $y_{i,j} \geq 0$, for $j = 1, 2, \ldots, n$, then the dual has no maximum, due to the fact we have feasibility of $\bar{\boldsymbol{\lambda}}$ for any choice of $\epsilon > 0$.
   If $y_{i,j} < 0$ for some $j$, we set

   $$\epsilon_0 = \frac{z_l - c_l}{y_{i,l}} = \min_{j=1,2,\ldots,n} \left\{ \frac{z_j - c_j}{y_{i,j}} : \quad y_{i,j} < 0 \right\}$$

3. Form a new basis $\mathbf{B}$ by pivoting $\mathbf{a}_i$ and $\mathbf{a}_l$. Using this basis determine the new $\mathbf{x_B}$ and return to Step 1.
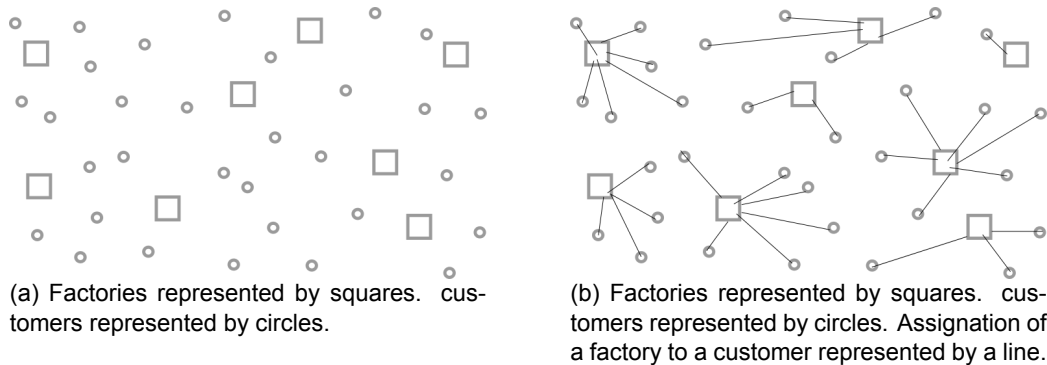
<div style="text-align: right; font-size: 3em;">3</div>

# Optimal Transport Theory

To introduce the optimal transport problem please imagine we are asked by a consortium of factories to design a plan to distribute their products among its many customers in such a way that the transportation costs are minimal.

We can start the approach of this problem considering the customers as members of the set $X$ and the factories as members of a set $Y$. We want to know which factory $y \in Y$ is going to supply a customer $x \in X$, i.e. we represent such assignation of a factory to a customer as map $y = T(x) \in Y$. Therefore, we can estimate the transportation cost $c(x, T(x))$ of supplying a customer $x$ with a factory $y = T(x)$.

We see that our problem is reduced to find an assigning map from the set of customers to the set of factories in such a way that the total cost $C(X, Y) = \sum_{x \in X} c(x, T(x))$ is minimal.

Figure 3.1: Illustration of the problem of Factories supplying customers.



(a) Factories represented by squares. customers represented by circles.

(b) Factories represented by squares. customers represented by circles. Assignation of a factory to a customer represented by a line.

Gaspard Monge was a French mathematician who introduced for the very first time the optimal transport problem as *déblais et remblais* in 1781. Monge was interested in finding a map that distributes an amount of sand or soil extracted from the earth or a mine distributed according to a density $f$, onto a new construction whose density of mass is characterized by a density $g$, in such a way the average displacement is minimal. We see that Monge presented a more continuous flavor of the problem.

We remark that we are not interested in the quantity of mass we are transporting. This information it is not relevant for the problem or has no sense its consideration (for example the factories-customer problem). We are interested in finding a way to assign or distribute elements among two sets. We are interested in applications concerning the transportation of

a finite amount of mass. Therefore, it is reasonable to state our problem in terms of probability measures.

Formally, given two densities of mass $f$ and $g$, Monge was interested in finding a map $T : \mathbb{R}^3 \to \mathbb{R}^3$ pushing the one onto the other,

$$\int_A g(y)\mathrm{dy} = \int_{T^{-1}(A)} f(x)\mathrm{dx}$$

For any Borel subset $A \subset \mathbb{R}^3$. And the transport also should minimize the quantity,
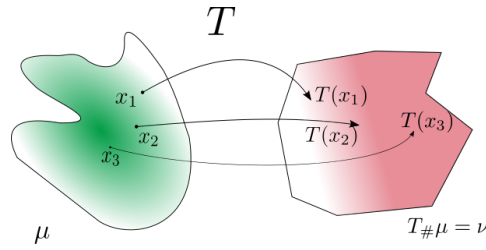
$$\int_{\mathbb{R}^3} |x - T(x)| f(x)\mathrm{dx}$$

Therefore, we need to search for the optimum in the set of measurables maps $T : X \to Y$ such that the condition (3) is translated to,

$$(T_{\#}\mu)(A) = \mu(T^{-1}(A)) \qquad \text{for every measurable set } A \subset X. \tag{3.1}$$

In other words, we need $T_{\#}\mu = \nu$. Notice that given the context for which the problem was formulated, originally it was binded to $\mathbb{R}^3$ or $\mathbb{R}^2$ but we can consider the general case in $\mathbb{R}^d$. In the Euclidean frameworks if we assume $f$, $g$ and $T$ regular enough and $T$ also injective, this equality implies,

$$g(T(x)) \det(\mathrm{D}T(x)) = f(x) \tag{3.2}$$

Figure 3.2: Monge problem. Finding a map.



The equation (3.2) is nonlinear in $T$ making the analysis of the Monge's Problem really difficult. Moreover, the constraint makes this problem hard to handle since it is not close even under weak convergence.

To appreciate this fact, consider $\mu = \mathscr{L}^1 \llcorner [0,1]$ and the hat functions $h_k$ defined as follow,

$$h_k(x) = \begin{cases} 2kx & x \in \left[0, \frac{1}{2k}\right] \\ 2 - 2kx & x \in \left(\frac{1}{2k}, \frac{1}{k}\right] \\ 0 & \text{otherwise} \end{cases}$$

Then take the sequence $f_n : [0,1] \to [0,1]$,

$$f_n(x) = \sum_{i=0}^{n-1} h_n\left(x - \frac{i}{n}\right) \tag{3.3}$$

We see that the sequence satisfies $f_{n\#}\mu = \mu$. It is easy to check that $\mu\left(f_n^{-1}(A)\right) = \mathscr{L}^1(A)$ for every open set $A \in [0,1]$. In the other hand, this is sequence of oscillating function converges weakly to its mean value $f_n \rightharpoonup f = \frac{1}{2}$ by means of the Riemmann-Lebesgue theorem[1], which makes $f_{\#}\mu \neq \mathscr{L}^1 \llcorner [0,1]$.

---

[1] Let $u \in L^p_{\mathrm{loc}}(\mathbb{R}^N)$, be $Q-$ periodic, i.e. $u(x + e_i) = u(x)$ for any canonical basis $e_i$ of $\mathbb{R}^N$, with $1 \leq i \leq N$. For every $\epsilon > 0$ and any $x \in \mathbb{R}^N$, we set $u_\epsilon(x) = u(x/\epsilon)$. Then $u_\epsilon \rightharpoonup \bar{u}$ in $L(E)$ for every bounded measurable set $E \subset \mathbb{R}^N$ where $\bar{u} = const = \int_{Q(0,1)} u(y)\mathrm{dy}$. We can find the proof of this theorem in [12] in the section of *Weak convergence for $L^p$ spaces*.
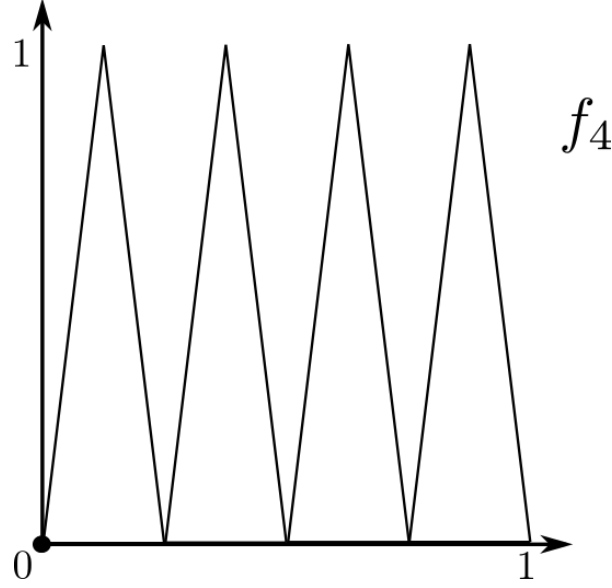
Figure 3.3: $f_n$ constructed using hat functions. The pictures shows the case $n = 4$.

**Problem 3.** *Given two probability measures $\mu \in \mathcal{P}(X)$ and $\nu \in \mathcal{P}(Y)$ and a cost function $c : X \times Y \to \{0, +\infty\}$, the Monge's problem consists in finding a map $T : X \to Y$*

$$\inf\left\{ M(T) := \int_X c(x, T(x))\mathrm{d}\mu(x) : \quad T_\#\mu = \nu \right\} \tag{MP}$$

Monge analyzed geometric properties of the solution to this problem. Although, the question for the existence of an optimal map stayed open until a Russian mathematician named Leonid Vitaliyevich Kantorovich introduced in the paper [16] a suitable framework to study its optimality conditions and prove the existence of a minimizer.

When we formulate our factories-customer problem through finding an assignation map, we are excluding the situations in which one customer can be supplied by two or more factories, or in the case of the Monge's problem we are ignoring the possibility of splitting a unit of mass into small pieces that can be assigned simultaneously to different places.

The idea behind Kantorovich's formulation is to consider the transportation maps from one space to another as transportation plans, that is joint probability measures with their marginals given by the initial and final configurations.

Instead of assigning an element of $Y$ to each element of the set $X$, we can see the problem from a different perspective and assign a weight to the importance of the point $(x, y) \in X \times Y$. We would like to know how much of our total material is going to be distributed from $x$ to $y$, in such a way to be consistent with the information we have about the initial and final configuration of the material. That is, we would like to know the optimal way to concentrate mass to the points $(x, y)$ in such a way we are not creating neither destroying mass.

Designing the transportation strategy using the above procedure is called a transport plan. In terms of probability theory, we are constructing a joint probability measure for $X \times Y$ with marginals given by the measures $\mu \in \mathcal{P}(X)$ and $\nu \in \mathcal{P}(Y)$.

Please note that in contrast to a map, we can always assign to a point $x \in Y$ as many points in $Y$ as we want, just considering the constraints given by the densities $\mu$ and $\nu$. We introduce the following notation to give the necessary formalism to this approach.

**Definition 3.1** (Coupling). *Let $\mu$ and $\nu$ be probability measures over the measurable spaces $(X, \mathcal{A}_X)$ and $(Y, \mathcal{A}_Y)$. Finding a coupling between $\mu$ and $\nu$ means to construct a probability measure $\gamma$ on the space $X \times Y$ (precisely on the product $\sigma$-algebra $\mathcal{A}_X \otimes \mathcal{A}_Y$), $\mu$ and $\nu$ are admitted as marginals on $X$ and $Y$ respectively. That is $\gamma \geq 0$, $\gamma(X \times Y) = 1$, $\mathrm{proj}_{x_\#}\gamma = \mu$ and $\mathrm{proj}_{y_\#}\gamma = \nu$.*

The above definition is equivalently to say that coupling two measures means to find a probability measure $\gamma$, such that for all measurable sets $A \subset X$ and $B \subset Y$, one has $\gamma[A \times Y] = \mu[A]$, $\gamma[A \times X] = \nu[B]$.

Moreover, for all integrable (nonnegative measurable) functions $\phi, \psi$ on $X$ and $Y$,

$$\int_{X \times Y} (\phi(x) + \psi(y)) \, \mathrm{d}\gamma(x, y) = \int_X \phi \mathrm{d}\mu + \int_Y \psi \mathrm{d}\nu$$

Since definition 3.1 is given for measures on probabilistic spaces, we can rephrase it in terms of stochastic variables. Let $(X, \mu)$ and $(Y, \nu)$ be two probability spaces. Coupling $\mu$ and $\nu$ means constructing two random variables $\mathcal{X}$ and $\mathcal{Y}$ on some probability space, such that $\mathrm{law}(\mathcal{X}) = \mu$, $\mathrm{law}(\mathcal{Y}) = \nu$. The couple $(\mathcal{X}, \mathcal{Y})$ is called a coupling of $(\mu, \nu)$.

We use the notation $\Pi(\mu, \nu)$ to refer the **set of couplings** of $\mu$ and $\nu$. That is,

$$\Pi(\mu, \nu) = \left\{ \gamma \in \mathcal{P}(X \times Y) : \left(\mathrm{proj}_x\right)_\# \gamma = \mu \text{ and } \left(\mathrm{proj}_y\right)_\# \gamma = \nu \right\} \tag{3.4}$$

**Lemma 3.1** (Existence of a coupling). *Let $\mu$ and $\nu$ be probability measures over the measurable spaces $(X, \mathcal{A}_X)$ and $(Y, \mathcal{A}_Y)$. Then, there exists $\exists \gamma \in \mathcal{P}(X \times Y)$, such that $\gamma \in \Pi(\mu, \nu)$.*

*Proof.* Take $\gamma = \mu \otimes \nu$.                                                                    $\square$

Notice that this approach to solve the problem is more general, since we can always create a transportation plan given a transportation map, i.e.

$$(\mathrm{id}, T)_\# \mu = \gamma \in \mathcal{P}(X \times Y)$$

If $T$ is a transportation map it is easy to check that indeed $\left(\mathrm{proj}_x\right)_\# \gamma = \mu$ and $\left(\mathrm{proj}_y\right)_\# \gamma = \nu$. This inspires a definition for a coupling between two measures generated by a transport map.

**Definition 3.2** (Deterministic Coupling). *Let $(X, \mu)$ and $(Y, \nu)$ be two probabilistic spaces. If there exists a measurable map $T : X \to Y$ such that $T_\# \mu = \nu$. We call the measure $(\mathrm{id}, T)_\# \mu = \gamma \in \mathcal{P}(X \times Y)$ a deterministic coupling of $\mu$ and $\nu$.*

For the sake of simplicity, we refer as $\gamma_T$ a transportation plan generated from a transportation map $T$.

In terms of stochastic variables, a coupling $(\mathcal{X}, \mathcal{Y})$ is said to be deterministic if there exists a measurable function $T : X \to Y$ such that $\mathcal{Y} = T(\mathcal{X})$. Equivalently, $(\mathcal{X}, \mathcal{Y})$ is a deterministic coupling of $\mu$ and $\nu$, if its law $\gamma = \mathrm{law}((\mathcal{X}, \mathcal{Y}))$ is concentrated on the graph of a measurable map $T : X \to Y$. Other way to rephrase it is saying that $\mu = \mathrm{law}(\mathcal{X})$, $\mathcal{Y} = T(\mathcal{X})$, where $T$ is a change of variables from $\mu$ to $\nu$, for all $\nu$-integrable (nonnegative measurable) function $\phi$,

$$\int_Y \phi(y) \mathrm{d}\nu(y) = \int_X \phi(T(x)) \mathrm{d}\mu(x).$$

The increasing rearrangement on $\mathbb{R}$ is an example of a coupling between two probability measures over one dimensional euclidean space. Let $\mu$, $\nu$ be two probability measures on $\mathbb{R}$. Define their cumulative distribution functions by,

$$F(x) = \int_{-\infty}^x \mathrm{d}\mu, \qquad G(y) = \int_{-\infty}^y \mathrm{d}\nu$$

Cumulative distributions not always are invertible, since they are not always strictly increasing. Although we can define their pseudo-inverses as follow,

$$F^{-1}(t) = \inf\{x \in \mathbb{R}; F(x) > t\}, \qquad (3.5)$$

$$G^{-1}(t) = \inf\{y \in \mathbb{R}; G(y) > t\}. \qquad (3.6)$$

Then, we set the map $T$ as $T = G^{-1} \circ F$. If $\mu$ is atomless then $T_\# \mu = \nu$.

The increasing rearrangement coupling is useful to construct the *Knothe-Rosenblatt coupling* between two Stochastic variables $\mathbb{R}^n$. Let $\mu$ and $\nu$ be two probability measures on $\mathbb{R}^n$, such that $\mu$ is absolutely continuous with respect to Lebesgue measure. This coupling is constructed in the following way:

1. Take the marginal of the first projection on the first variable; this gives probability measures $\mu_1(\mathrm{dx}_1)$, $\nu_1(\mathrm{dy}_1)$ on $\mathbb{R}$, with $\mu_1$ being atomless. Then define $y_1 = T_1(x_1)$ by the composition of the pseudo-inverse functions of the increasing rearrangement, with $F$ and $G$ considered as they are in (3.5) and (3.6) respectively.

2. Now take the marginal on the first two variables and disintegrate it with respect to the first variable. This gives probability measures $\mu_2(\mathrm{dx}_1\mathrm{dx}_2) = \mu_1(\mathrm{dx}_1)\mu_2(\mathrm{dx}_2|x_1)$, $\nu_2(\mathrm{dy}_1\mathrm{dy}_2) = \nu_1(\mathrm{dy}_1)\nu_2(\mathrm{dy}_2|y_1)$. For each given $y_1 \in \mathbb{R}$, we set $y_1 = T_1(x_1)$, and then we define $y_2 = T_2(x_2; x_1)$ under the increasing rearrangement formula of $\mu(\mathrm{dx}_2|x_1)$ into $\nu(\mathrm{dy}_2|y_1)$.

3. We repeat the construction, adding one variable after another. For example, after the assignation $x_1 \to y_1$ has been determined, the conditional probability of $x_2$ is seen as a one-dimensional probability on a small slice of width $\mathrm{dx}_1$, and it can be transported to the conditional probability of $y_2$ seen as one dimensional probability of a slice of width $\mathrm{dy}_1$. After $n$ constructions, this procedure maps $\mathcal{Y} = T(\mathcal{X})$.

The *Knothe-Rosenblatt coupling* has the property that its Jacobian matrix for the change of variable $T$ is upper triangular with positive entries on the diagonal.

**Lemma 3.2** (Gluing lemma). *If $\mathcal{Z}$ is a function of $\mathcal{Y}$ and $\mathcal{Y}$ is a function of $\mathcal{X}$, then $\mathcal{Z}$ is a function of $\mathcal{X}$. Let $(X_i, \mu_i)$, $i = 1, 2, 3$, be Polish probability spaces. If $(\mathcal{X}_1, \mathcal{X}_2)$ is a coupling of $(\mu_1, \mu_2)$ and $(\mathcal{Y}_2, \mathcal{Y}_3)$ is a coupling of $(\mu_2, \mu_3)$, then it is possible to construct a triple of random variables $(\mathcal{Z}_1, \mathcal{Z}_2, \mathcal{Z}_3)$ such that $(\mathcal{Z}_1, \mathcal{Z}_2)$ has the same law as $(\mathcal{X}_1, \mathcal{X}_2)$ and $(\mathcal{Z}_2, \mathcal{Z}_3)$ has the same law as $(\mathcal{Y}_2, \mathcal{Y}_3)$.*

To understand the above lemma in terms of transportation plans, consider a plan $\gamma_{1,2,3}$ for the Cartesian product $X_1 \times X_2 \times X_3$, with marginals given by $\mu_1$, $\mu_2$ and $\mu_3$ respectively. Let $\gamma_{1,2}$ be a coupling between $\mu_1$ and $\mu_2$, and let $\gamma_{2,3}$ be a coupling between $\mu_2$ and $\mu_3$. Informally, we can disintegrate $\gamma_{1,2}$ and $\gamma_{2,3}$ as follows,

$$\gamma_{1,2}(\mathrm{dx}_1\mathrm{dx}_2) = \gamma_{1,2}(\mathrm{dx}_1|x_2)\mu_2(\mathrm{dx}_2)$$

$$\gamma_{2,3}(\mathrm{dx}_2\mathrm{dx}_3) = \gamma_{2,3}(\mathrm{dx}_3|x_2)\mu_2(\mathrm{dx}_2),$$

and reconstruct $\gamma_{1,2,3}$ as follows,

$$\gamma_{1,2,3}(\mathrm{dx}_1\mathrm{dx}_2\mathrm{dx}_3) = \gamma_{1,2}(dx_1|x_2)\mu(x_2)\gamma_{2,3}(\mathrm{dx}_2\mathrm{dx}_3)$$

We use indistinctly the term coupling and transportation plan. Now we introduce the Kantorovich's problem.

**Problem 4.** *Given $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$, and $c : X \times Y \to [0, +\infty]$, we consider the problem*

$$\Theta(\mu, \nu) = \inf\left\{K(\gamma) := \int_{X \times Y} c\,\mathrm{d}\gamma : \quad \gamma \in \Pi(\mu, \nu)\right\} \qquad \text{(KP)}$$

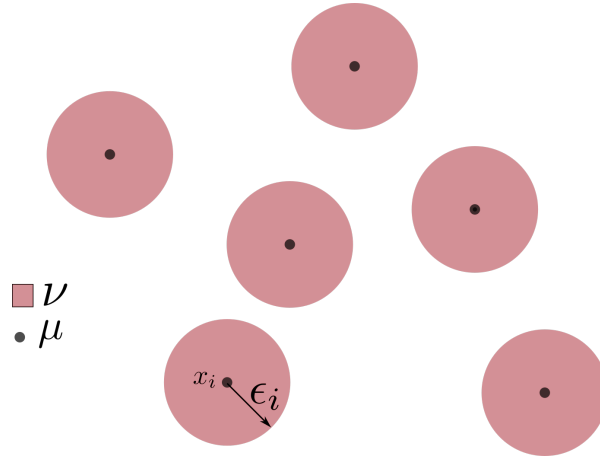*where $\Pi(\mu, \nu)$ is the set of transport plans.*

It is a fact for The Kantorovich's formulation that it is always possible to find a transport plan, to see this fact it is enough to take $\gamma = \mu \otimes \nu$, as shown in lemma 3.1. On the contrary, it is not always possible to find transportation maps (deterministic couplings).

Consider a measure $\mu$ on $\mathbb{R}^d$, concentrated on $N$ different atoms $x_i \in \mathbb{R}^d$,

$$\mu = \frac{1}{N} \sum_{i=0}^{N-1} \delta_{x_i}$$

Where $\delta_{x_i}$ is the Dirac mass at point $x_i$. Consider $N$ open balls on $\mathbb{R}^d$ centered at $x_i$ with radius $\epsilon_i > 0$, such that they disjoint pairwise. Let $D = \cup_{i=0}^{N-1} B(x_i; \epsilon_i)$ be the union of these balls. Let $\nu$ be a the Hausdorff measure of over $D \subset \mathbb{R}^d$. That is $\nu = \mathcal{H} \llcorner D$. We see that it is impossible to couple $\mu$ and $\nu$ deterministically; since there is no map $T$, such that $T_\# \mu = \nu$.

Figure 3.4: Transportation maps. There is no deterministic coupling for $\mu$ and $\nu$, but there is a transportation plan.



## Existence of a minimizer for Kantorovich's Problem.

Properties of transportation plans.

**Lemma 3.3.** *Let $X$ be a metric space. If $f : X \to \overline{\mathbb{R}}$ is a lower semi-continuous function, bounded from below, then the functional $J : \mathcal{M}_+(X) \to \overline{\mathbb{R}}$ defined on the space of finite positive measures on $X$, given by*

$$J(\mu) = \int f \, d\mu$$

*is lower semi-continuous for the weak convergence of measures.*

*Proof.* Let $(f_n)_{n \in \mathbb{N}}$ be a sequence of continuous and bounded functions, converging increasingly to $f$. Consider the functionals $J_n : \mathcal{M}_+(X) \to \overline{\mathbb{R}}$, defined as

$$J_n(\mu) = \int f_n \, d\mu$$

Every $J_n$ is continuous for the weak convergence. We set $J(\mu) = \int f \, d\mu$. We see that $J_n(\mu) \le J(\mu)$ for any $\mu$. Since our functions are bounded, and $f$ is bounded from below, and our measures are finite, we can make use of monotone convergence theorem, $J_n(\mu) \to J(\mu)$, having as result $J(\mu) = \sup_n J_n(\mu)$. Since we have that $J(\mu)$ is the supremum of continuous functions we can assure that that $J$ is lower semicontinuous. $\qquad \square$

**Theorem 3.1** (Lower-semicontinuity of the cost function). *Let $X$ and $Y$ two Polish spaces, and $c : X \times Y \to \mathbb{R} \cup \{+\infty\}$ is a real valued lower semicontinuous function bounded from below. Then the*

functional $K : \mathcal{P}(X \times Y) \to \mathbb{R} \cup \{+\infty\}$,

$$K(\gamma) := \int_{X \times Y} c \, \mathrm{d}\gamma, \tag{3.7}$$

is lower semicontinuous.

*Proof.* This is a consequence of lemma 3.3 setting $f = c$ over $X \times Y$. $\square$

*The beauty of Kantorovich's formulation lies on the fact that the set of transport plans is compact under weak convergence making it a suitable framework where we can use the Weierstrass' criterion to show the existence of a minimizer.*

**Theorem 3.2.** *Let $X$ and $Y$ be compact metric spaces, $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$ and a cost function $c : X \times Y \to \mathbb{R}$ a continuous function. Then* (KP) *admits a solution.*

*Proof.* To prove the existence we make use of the Weierstrass' criterion for existence of minimizers. Therefore, we need to prove that $K(\gamma)$ is at least lower semicontinuous and compactness of the space $\Pi(\mu, \nu)$ under some topology.

We choose as a notion of convergence the weak convergence of probability measures in duality with $C_b(X \times Y)$. This immediately implies continuity for $K(\gamma)$ by definition, since $c$ is already in $C(X \times Y)$ defined in a compact space $X \times Y$.

Now take a sequence $(\gamma_n)_{n \in \mathbb{N}} \in \Pi(\mu, \nu)$. Since they are probability measures for all $n$ they are bounded in the dual of $C(X \times Y)$. Weak-star compactness in the dual of compact metric spaces guarantees the existence of a convergent subsequence $\gamma_{n_k} \rightharpoonup \gamma$. Let us fix $\phi \in C(X)$ and using $\int \phi(x) \mathrm{d}\gamma_{n_k} = \int \phi \mathrm{d}\mu$ and taking the limit we have $\int_{X \times Y} \phi(x) \mathrm{d}\gamma = \int_X \phi \mathrm{d}\mu$. In this way we prove that $\gamma_\#\left(\mathrm{proj}_x\right) = \mu$. We can repeat this argument for $\nu$, fixing $\psi \in C(Y)$ and taking the limit of $\int_{X \times Y} \psi(y) \mathrm{d}\gamma_{n_k} = \int \psi \mathrm{d}\nu$, implies $\int_{X \times Y} \psi(y) \mathrm{d}\gamma = \int \psi \mathrm{d}\nu$.

This proves that $\gamma_\#\left(\mathrm{proj}_y\right) = \nu$. Hence, the limit $\gamma \in \Pi(\mu, \nu)$ showing that the set of couplings of $\mu$ and $\nu$ is sequentially compact. $\square$

*Continuity for the cost function and compactness of the metric spaces can be demanding requirements. However we can substitute them by milder conditions for the existence of a minimizer.*

**Theorem 3.3.** *Let $X$ and $Y$ be compact metric spaces, $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$, and $c : X \times Y \to \overline{\mathbb{R}}$ be lower semi-continuous and bounded from below. Then Kantorovich's problem admits a solution.*

*Proof.* By theorem 3.1, the functional $K(\gamma) = \int c \mathrm{d}\gamma$ is lower semicontinuous. We apply again Weierstrass criterion proving existence of a minimizer. $\square$

**Lemma 3.4.** *The set of couplings $\Pi(\mu, \nu)$ between two probability measures $\mu$ and $\nu$ defined over two Polish spaces $X$ and $Y$ respectively, is tight.*

*Proof.* Fix $\epsilon > 0$ and find two compact sets $K_x \subset X$ and $K_y \subset Y$ such that $\mu(X \backslash K_x) < \epsilon$, and $\nu(Y \backslash K_y) < \epsilon$. Then the set $K_X \times K_Y$ is compact in $X \times Y$ and, for any $\gamma_n \in \Pi(\mu, \nu)$, we have,

$$\gamma_n\left((X \times Y) \backslash (K_X \times K_Y)\right) \le \gamma_n\left((X \backslash K_X) \times Y\right) + \gamma_n\left(X \times (Y \backslash K_y)\right)$$
$$= \mu(X \backslash K_X) + \nu(Y \backslash K_Y)$$
$$= 2\epsilon$$

Given the arbitrary way to choose $\epsilon$, $\Pi(\mu, \nu)$ is tight. $\square$

**Theorem 3.4.** *Let $X$ and $Y$ be Polish spaces, and $c : X \times Y \to \overline{\mathbb{R}}_+$, a real valued lower semi-continuous cost function on the space $X \times Y$. Then the Kantorovich's problem* (KP) *admits a solution.*

*Proof.* By Prokhorov and the tightness of the set of transport plans between two probability measures, we have compactness. Theorem (3.1) shows lower semicontinuity of the cost function. Applying Weierstrass criterion we obtain the existence of an optimizer. $\square$

### *Properties of Optimal plans*

**Theorem 3.5** (Convexity of optimal plans). *The set of solutions $\bar{\gamma} \in \Pi(\mu, \nu)$ for the Kantorovich's problem is a convex set.*

*Proof.* We see immediately that if $\gamma_1$ and $\gamma_2$ solve the Kantorovich's problem, for any $t \in [0, 1]$, the plan $\gamma = t\gamma_1 + (1 - t)\gamma_2$, also solves the problem.  □

An interesting property of an optimal coupling between two measures, is that optimality remains after restricting the plan to a non zero measure subset of $X \times Y$.

**Theorem 3.6** (Optimality is inherited by restriction). *Let $(X, \mu)$ and $(Y, \nu)$ be two Polish spaces, $a \in L^1(\mu)$, $b \in L^1(\nu)$, let $c : X \times Y \to \mathbb{R} \cup \{+\infty\}$ be a measurable cost function such that $c(x, y) \geq a(x) + b(y)$ for all $x$, $y$; and let $\Theta(\mu, \nu)$ be the optimal transport cost from $\mu$ to $\nu$. Assume that $\Theta(\mu, \nu) < \infty$ and let $\gamma \in \Pi(\mu, \nu)$ be an optimal transport plan. Let $\tilde{\gamma}$ be a nonnegative measure on $X \times Y$, such that $\tilde{\gamma} \leq \gamma$ and $\tilde{\gamma}(X \times Y) > 0$. Then the joint probability measure,*

$$\hat{\gamma} = \frac{\tilde{\gamma}}{\tilde{\gamma}(X \times Y)}$$

*is an optimal plan between its marginals $\hat{\mu} = (\text{proj}_X)_\# \hat{\gamma}$ and $\hat{\nu} = (\text{proj}_Y)_\# \hat{\gamma}$.*

*Proof.* We proceed by contradiction; take a transportation plan $\hat{\gamma}$ such that for a given cost function $c$, it is not optimal. Since $\hat{\gamma}$ is not optimal we can find another plan $\bar{\gamma}$ such that $(\text{proj}_X)_\# \bar{\gamma} = (\text{proj}_X)_\# \hat{\gamma} = \hat{\mu}$ and $(\text{proj}_Y)_\# \bar{\gamma} = (\text{proj}_Y)_\# \hat{\gamma} = \hat{\nu}$ and $K(\bar{\gamma}) < K(\hat{\gamma})$.

Let $\alpha = \tilde{\gamma}(X \times Y) > 0$ be the measure of the space under $\tilde{\gamma}$ which is greater than zero by definition. Let $\gamma' = (\gamma - \tilde{\gamma}) + \alpha\bar{\gamma}$ be a measure over $X \times Y$. We see that $\gamma' > 0$ by construction since $\gamma$

$$\begin{aligned}
\gamma' &= (\gamma - \tilde{\gamma}) + \alpha\bar{\gamma} \\
&= \gamma - \frac{\tilde{\gamma}(X \times Y)}{\tilde{\gamma}(X \times Y)}\tilde{\gamma} + \alpha\bar{\gamma} \\
&= \gamma - \tilde{\gamma}(X \times Y)\hat{\gamma} + \tilde{\gamma}(X \times Y)\bar{\gamma} \\
&= \gamma + \alpha\left(\bar{\gamma} - \hat{\gamma}\right) \\
&< \gamma
\end{aligned}$$

$\bar{\gamma}$ and $\tilde{\gamma}$ share the same marginals, therefore $(\text{proj}_X)_\# \gamma' = (\text{proj}_X)_\# \gamma$ and $(\text{proj}_Y)_\# \gamma' = (\text{proj}_Y)_\# \gamma$ giving as result that $\gamma' \in \Pi(\mu, \nu)$, contradicting the fact that $\gamma$ is optimal.  □

The last theorem tells us that transferring part of the initial mass to part of the final using the optimal plan designed for the total mass is also optimal.

**Corollary 3.1.** *Under the framework of the last theorem, if $\gamma$ is the unique optimal transference plan between $\mu$ and $\nu$, then also $\hat{\gamma}$ is the unique optimal transference plan between $\hat{\mu}$ and $\hat{\nu}$.*

*Proof.* Let $\gamma$ be the unique optimal coupling between $\mu$ and $\nu$. Let $\bar{\gamma}$ be any optimal transfer plan coupling $\hat{\mu}$ and $\hat{\nu}$. Define $\gamma'$ as we did in the last proof $\gamma' = (\gamma - \tilde{\gamma}) + \alpha\bar{\gamma}$, with $\alpha = \tilde{\gamma}(X \times Y)$. This implies that $K(\gamma') = K(\gamma)$. Since the coupling for $\mu$ and $\nu$ is unique $\gamma' = \gamma$, which implies $\tilde{\gamma} = \alpha\gamma$ then $\tilde{\gamma} = \hat{\gamma}$  □

## Kantorovich's formulation as relaxation of Monge's formulation.

There are situations in which is possible to find a deterministic coupling between two measures, but not an optimal one for a cost function $c : X \times Y \to \overline{\mathbb{R}}$. A common example, popular in the literature, is the following: consider as cost function $c : \mathbb{R}^2 \times \mathbb{R}^2 \to \mathbb{R}$, the Euclidean distance $c(x, y) = |x - y|$, the measure $\mu = \mathcal{H} \llcorner D$ as the Hausdorff measure for the segment $D = \left\{(0, t)^\top \in \mathbb{R}^2 : \text{ for } t \in [0, 1]\right\}$.

Let $D_1$ and $D_2$ be the segments given by,

$$D_1 = \left\{(-1, t)^\top \in \mathbb{R}^2 : \text{ for } t \in [0, 1]\right\}$$

$$D_2 = \left\{(+1, t)^\top \in \mathbb{R}^2 : \text{ for } t \in [0, 1]\right\}$$

And we set the measure $\nu$ as follows,

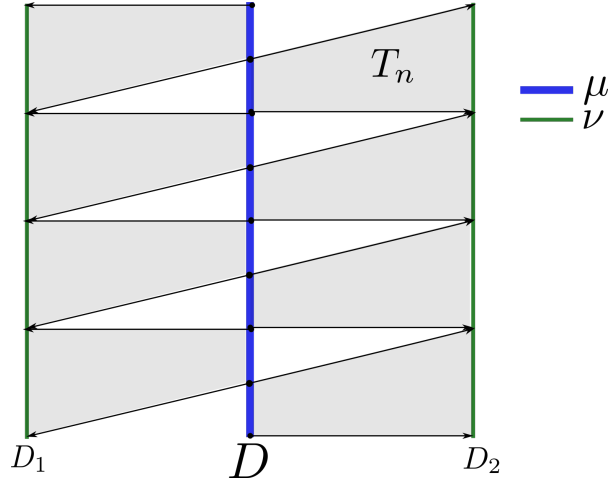$$\nu = \frac{\mathscr{H} \llcorner D_1 + \mathscr{H} \llcorner D_2}{2}$$



Figure 3.5: There is a deterministic coupling for $\mu$ and $\nu$, but no optimal one. The map $T_n$ shown in this picture with $n = 4$.

There are many ways to construct a transportation map for this situation. Consider the maps $T_n$ constructed splitting the segment $D$ into $2n$ equal parts and the segments $D_1$ and $D_2$ in $n$ equal parts. We label the parts of the segment $D$ with the integer numbers from 0 to $2n - 1$. Then the map $T_n$ assign the parts of $D$ labeled with even numbers to the right hand side segment $D_2$ and the parts labeled with odd numbers to the left right side segment $D_1$.

Formally, let $k = 0, \ldots, 2n - 1$ be an integer used to label the equal parts of $D$,

$$T_n\left(\begin{pmatrix} 0 \\ t \end{pmatrix} \in D\right) = \begin{cases} \begin{pmatrix} 1 \\ 2t - \frac{k}{2n} \end{pmatrix} & k \text{ even and } t \in \left[\frac{k}{2n}, \frac{k+1}{2n}\right), \\ \begin{pmatrix} -1 \\ 2t - \frac{k+1}{2n} \end{pmatrix} & k \text{ odd and } t \in \left(\frac{k}{2n}, \frac{k+1}{2n}\right]. \end{cases}$$

We can find an upper boundary for the total cost $C(T_n)$,

$$M(T_n) = \int_D \left| x - T_n(x) \right| \mathrm{d}\mu(x)$$

$$= 2n \int_0^{\frac{1}{2n}} \sqrt{1 + 4t^2} \mathrm{d}t$$

$$\leq 2n \left( \int_0^{\frac{1}{2n}} 1 + 4t^2 \mathrm{d}t \right)^{1/2} \left( \int_0^{\frac{1}{2n}} \mathrm{d}t \right)^{1/2}$$

$$= \sqrt{1 + \frac{1}{3n^3}}$$

$$\leq 1 + \frac{1}{n}$$

Let $T_n$ be an assignation map. We see that we can always find a cheaper map $T_{n+1}$ for any $n \in \mathbb{N}$. This sequence of transportation maps converges weakly to the plan $\gamma_{T_n} \rightharpoonup \gamma_T = \frac{\gamma_{T^+}}{2} + \frac{\gamma_{T^-}}{2}$. Where $T^+$ and $T^-$ are given by:

$$T^+(x) = x + \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$T^-(x) = x - \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

The idea is that the mass of each point $x \in D$ is split in two and equally distributed among $D_1$ and $D_2$ assigning one half of the mass respectively. Note that this distribution is an optimal plan for the cost function $c(x, y) = |x - y|$. Because of the triangle inequality, sending the mass from $x \in D$ to any other point of $D_1$ and $D_2$ different than those assigned by the maps $T^{\pm}$, implies a higher cost.

From the last example we see that a sequence of deterministic couplings converges to a transportation plan that is a solution for Kantorovich's problem (KP), but clearly it is not for Monge's problem (MP). We also gave one example where (MP) has no solution. Assume for a moment that Monge's situation where indeed does exist a solution for Monge's problem, then the following question arises: Is there any situation where Monge's problem and Kantorovich's problem have the same solution?

**Lemma 3.5.** *On a compact subset $\Omega \subset \mathbb{R}^d$, the set of plans $\gamma_T$ induced by a transport is dense in the set of plans $\Pi(\mu, \nu)$ whenever $\mu$ is atomless.*

**Theorem 3.7.** *On a compact subset $\Omega \subset \mathbb{R}^d$, $K(\gamma)$ is the relaxation of $J(\gamma)$. In particular, $\inf J = \min K$, and hence Monge and Kantorovich problems have the same infimum.*

*Proof.* Since $K$ is continuous, then it is lower semicontinuous, and since we have $K \leq J$, then $K$ is necessarily smaller than the relaxation of $J$. We only need to prove that, for each $\gamma$, we can find a sequence of transports $T_n$ such that $\gamma_{T_n} \to \gamma$ and $J(\gamma_{T_n}) \to K(\gamma)$, so that the infimum in the sequential characterization of the relaxed functional will be smaller than $K$, thus providing the equality.

Actually, since for $\gamma = \gamma_{T_n}$ be two functionals $K$ and $J$ coincide, and since $K$ is continuous we only need to produce a sequence $T_n$ such that $\gamma_{T_n} \to \gamma$. It is possible to do the last step due to the density of transport plans generated by a map $\gamma_{T_n}$ in the set of transport plans $\Pi(\mu, \nu)$. $\square$

To understand the relation between both formulations we come back to the factories-customers example. The consortium instead of having the policy of assigning a factory to each costumer, they prefer a relaxation of the problem. Now, they consider their products as mass distributed across the city in different places given by the production rate of each factory that they need to redistribute to a given configuration, that is the costumers location with their respective demand.

Approaching the problem in this way gives us the flexibility to supply each customer's demand with the production of many factories.

## Cyclical Monotonicity and Duality.

Imagine that the consortium changed its policy and it has decided not to be responsible any longer for the transportation of the goods, letting the customers to solve this problem by themselves (assume that the consortium has the monopoly of the goods and the customers have no choice but to adhere to this policy). An entrepreneur feeling that he can ship the goods more efficiently than the consortium did, he intend to buy the goods at the factories and sell them at the customers' stores. Then, he must negotiate with the consortium the prices $-\phi(x)$ that he is able to pay at each factory for the goods, and the selling prices $\psi(y)$ at each customers' store. In order to succeed, he needs to be competitive and should do it better than the consortium did. Therefore, he must be able to cover with the difference of the sale prices, the transportation costs and they should be less than the consortium's costs $\psi(y) + \phi(x) \leq c(x, y)$. He is subject to this constraint and he should negotiate with the consortium and the customers the prices $\phi(x)$ and $\psi(y)$ in order to obtain the maximum profit.

## Duality

We see that for any $\gamma \in \mathcal{M}_+(X \times Y)$ we have,

$$\sup_{\phi,\psi}\left(\int_X \phi \, d\mu + \int_Y \psi \, d\nu - \int_{X\times Y}(\phi(x)+\psi(y))\,d\gamma\right) = \mathbb{I}_{\Pi(\mu,\nu)}(\gamma) = \begin{cases} 0 & \text{if } \gamma \in \Pi(\mu,\nu) \\ +\infty & \text{otherwise.} \end{cases} \tag{3.8}$$

where the supremum is taken among all bounded and continuous functions $\phi$, $\psi$.

Note that the result of this problem is 0 if $\gamma$ satisfies the constraint of being a probability measure over $X \times Y$ with given marginals $\mu$ and $\nu$, and we obtain $\infty$ if $\gamma$ does not satisfy this constraint.

Therefore, we can rewrite the Kantorovich's transport problem as an unconstrained minimization problem,

$$K(\gamma) = \min_\gamma \left(\int_{X\times Y} c \, d\gamma + \sup_{\phi,\psi}\left(\int_X \phi \, d\mu + \int_Y \psi \, d\nu - \int_{X\times Y}(\phi(x)+\psi(y))\,d\gamma\right)\right) \tag{3.9}$$

$$= \sup_{\phi,\psi}\left(\int_X \phi \, d\mu + \int_Y \psi \, d\nu\right) + \min\left(\int_{X\times Y} c(x,y)\,d\gamma - \sup_{\phi,\psi}\left(\int_X \phi(x)\,d\mu + \int_Y \psi(y)\,d\nu\right)\right) \tag{3.10}$$

Note that,

$$\inf_\gamma \int_{X\times Y}(c(x,y)-(\phi(x)+\psi(y)))\,d\gamma = \begin{cases} 0 & \text{if } \phi(x)+\psi(y) \leq c(x,y), \quad \forall(x,y) \in X \times Y \\ -\infty & \text{otherwise.} \end{cases} \tag{3.11}$$

In the other hand, equation (3.10) it is not really useful if we are not able to exchange the min and sup; for a moment suppose that the conditions that allow to exchange them do exist, we can rewrite the equation (3.10) as follows,

$$K(\gamma) = \sup_{\phi,\psi}\left(\int_X \phi \, d\mu + \int_Y \psi \, d\nu + \inf_\gamma \int_{X\times Y}(c(x,y)-(\phi(x)+\psi(y)))\,d\gamma\right) \tag{3.12}$$

If it exists $(x,y) \in X \times Y$ such that $\phi(x)+\psi(y) > c$, we can find measures $\gamma$ concentrated on the set where the strict inequality holds and mass tending to infinity, sending the value of the integral to $-\infty$. The above equation motivates the dual formulation of the problem,

**Problem 5.** *Given $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$, and the cost function $c : X \times Y \to \mathbb{R}_+$ we refer as the dual formulation of the transport problem* (DP),

$$\Delta = \sup\left\{\int_X \phi \, d\mu + \int_Y \psi \, d\nu : \phi \in C_b(X), \, \psi \in C_b(Y), \, \phi(x)+\psi(y) \leq c(x,y), \, \forall(x,y) \in X \times Y\right\} \quad \text{(DP)}$$

Notice that for any $\gamma \in \Pi(\mu,\nu)$,

$$\int_X \phi \, d\mu + \int_Y \psi \, d\nu = \int_{X\times Y} \phi(x)+\psi(y)\,d\gamma \leq \int_{X\times Y} c(x,y)\,d\gamma \tag{3.13}$$

Then we see that the objective value of the dual problem is less or equal than the primal Kantorovich's problem, as long as the pair $(\phi,\psi)$ is admissible.

Since the set of admissible maps is not compact we cannot assure the existence of a maximizer. To find the necessary conditions to assure existence of the minimizer or a duality equality we need to characterize functions by means of $c-$concavity.

**Definition 3.3.** *Let $X$, $Y$ be two sets, let $c : X \times Y \to \mathbb{R}$, a real valued cost function bounded from below. Given a function $\zeta : X \to \mathbb{R} \cup \{-\infty\}$, we define its c-concave transform of $\zeta$, the function $\zeta^c : Y \to \overline{\mathbb{R}}$ by*

$$\zeta^c(y) = \inf_{x\in X}(c(x,y)-\zeta(x)) \tag{3.14}$$

*In similar way, we can define the $\bar{c}-$transform of $\xi : Y \to \mathbb{R} \cup \{-\infty\}$ by*

$$\xi^{\bar{c}}(x) = \inf_{y\in Y}(c(x,y)-\xi(y)) \tag{3.15}$$

*A function $\psi$ defined on $Y$ is said to be $\bar{c}-$concave if it is not identically to $-\infty$ and there exists $\zeta$ defined on $X$, such that $\psi = \zeta^c$. Similarly, a function $\phi$ defined on $X$ is said to be $c-$concave if it is not identically to $-\infty$ and there exists $\xi$ defined on $Y$ such that $\phi = \xi^{\bar{c}}$.*

We use a bar to denote the difference between the transformation respect to the first and the second parameter of the cost function. This notation becomes trivial if we deal with symmetric cost functions. For a given cost function $c : X \times Y \to \mathbb{R}$. We denote by $c - conc(X)$ the set of $c-$concave functions defined on $X$. In similar way, we denote by $\bar{c} - conc(Y)$ the set of $\bar{c}-$concave functions defined on $Y$.

The following theorem encompasses the importance of using $c$-characterization for functions, it represents a generalization for the convex envelop theorem.

**Theorem 3.8.** *Suppose that $c$ is real valued. For any $\phi : X \to \mathbb{R} \cup \{-\infty\}$, we have $\phi^{c\bar{c}} \geq \phi$. We have the equality $\phi^{c\bar{c}} = \phi$ if and only if $\phi$ is $c-$concave. Moreover, $\phi^{c\bar{c}}$ is the smallest $c-$concave function larger than $\phi$.*

*Proof.* Let $\phi^{c\bar{c}}$ be the $c$-transform of the $\bar{c}$-transform of $\phi$,

$$\phi^{c\bar{c}}(x) = \inf_{y \in Y} \left( c(x,y) - \phi^c(y) \right)$$

$$= \inf_{y \in Y} \left( c(x,y) - \inf_{\tilde{x} \in X} \left( c(\tilde{x},y) - \phi(\tilde{x}) \right) \right).$$

Note that $\forall (x,y) \in X \times Y$,

$$\inf_{\tilde{x} \in X} \left( c(\tilde{x},y) - \phi(\tilde{x}) \right) \leq c(x,y) - \phi(x).$$

Therefore,

$$\phi^{c\bar{c}}(x) \geq \inf_{y \in Y} \left( c(x,y) - c(x,y) + \phi(x) \right) = \phi(x)$$

As we did for $\phi$, we can repeat the above arguments for any $\xi : Y \to \mathbb{R} \cup \{-\infty\}$, having as result $\xi^{\bar{c}c} \geq \xi$. If $\phi$ is $c-$concave, there exists $\xi$ allowing to write $\phi$ as $\phi = \xi^{\bar{c}}$. Therefore, using the fact that $\xi^{\bar{c}c} \leq \xi$,

$$\phi^{c\bar{c}}(x) = \inf_{y \in Y} \left( c(x,y) - \phi^c(y) \right)$$

$$\leq \inf_{y \in Y} \left( c(x,y) - \xi \right)$$

$$= \xi^{\bar{c}} = \phi$$

Proving in this way that if $\phi \in c - conc(X)$, then $\phi^{c\bar{c}} = \phi$.

To prove the implication in the other direction, take any c-concave function $\overline{\phi} = \chi^{\bar{c}}$ larger than $\phi$. Taking the $c-$concave transform of $\overline{\phi}$ and the assumption $\overline{\phi} \geq \phi$ imply,

$$\chi^{\bar{c}c}(y) = \inf_{x \in X} \left( c(x,y) - \chi^{\bar{c}}(x) \right)$$

$$= \inf_{x \in X} \left( c(x,y) - \overline{\phi}(x) \right)$$

$$\leq \inf_{x \in X} \left( c(x,y) - \phi(x) \right)$$

$$= \phi^c.$$

Therefore $\overline{\phi} = \chi^{\bar{c}c} \leq \phi^c$. Since $\chi^{\bar{c}c} \geq \chi$, we have that $\chi \leq \phi^c$. Taking the $\bar{c}$-transform and using the last inequality,

$$\phi^{c\bar{c}}(x) = \inf_{y \in Y} \left( c(x,y) - \phi^c \right)$$

$$\leq \inf_{y \in Y} \left( c(x,y) - \chi \right)$$

$$\leq \chi^{\bar{c}} = \overline{\phi}$$

Hence $\phi^{c\bar{c}}$ is smaller than any $c-$concave function $\overline{\phi}$, larger than $\phi$. This proves that if the equality $\phi^{c\bar{c}} = \phi$ holds, then $\phi$ is a $c-$concave function.                                                    $\square$

An interesting property of $c-$concave transforms is that taking $c-$, $\bar{c}-$ and $c-$ concave transforms consecutively, is equivalent to applying just one a $c-$concave transform. Consider any $\phi : X \to \mathbb{R} \cup \{-\infty\}$, and take the $c\bar{c}c-$transform as follows,

$$\phi^{c\bar{c}c}(y) = \inf_{x \in X} \left( c(x,y) - \phi^{c\bar{c}}(x) \right)$$

$$= \inf_{x \in X} \left( c(x,y) - \inf_{\tilde{y} \in Y} \left( c(x,\tilde{y}) - \phi^c(\tilde{y}) \right) \right)$$

$$= \inf_{x \in X} \left( c(x,y) - \inf_{\tilde{y} \in Y} \left( c(x,\tilde{y}) - \inf_{\tilde{x} \in X} \left( c(\tilde{x},\tilde{y}) - \phi(\tilde{x}) \right) \right) \right)$$

$$= \inf_{x \in X} \left( c(x,y) - \inf_{\tilde{y} \in Y} \inf_{\tilde{x} \in X} \left( c(x,\tilde{y}) - c(\tilde{x},\tilde{y}) + \phi(\tilde{x}) \right) \right)$$

$$= \inf_{x \in X} \inf_{\tilde{y} \in Y} \inf_{\tilde{x} \in X} \left( c(x,y) - c(x,\tilde{y}) + c(\tilde{x},\tilde{y}) - \phi(\tilde{x}) \right)$$

For any $y \in Y$, we are taking the infimum among all $\tilde{y} \in Y$. Since $y \in Y$ we take $\tilde{y} = y$,

$$\phi^{c\bar{c}c}(y) \leq \inf_{\tilde{x} \in X} \left( c(\tilde{x},y) - \phi(\tilde{x}) \right) = \phi^c(y)$$

Having as result that for any $\phi$, $\phi^{c\bar{c}c} = \phi^c$. We can use this fact to prove the theorem 3.8, as we can find it in [34] proceeding through $c-$convexity-$c-$concavity.

This result makes use only of their definition to be proven, since $c-$concave functions are exactly defined as $c-$transform of something. Its convex version 1.11 needs the Hahn Banach theorem due to the fact that convex functions are not defined via sup of affine functions, but via the convexity inequality.

**Lemma 3.6** (Improvement through c-transforms). *Let $X$ and $Y$ be two Polish spaces. Let $c : X \times Y \to \mathbb{R}$ be a real valued cost function bounded from below. Let $\phi : X \to \mathbb{R}$ and $\psi : Y \to \mathbb{R}$ be two bounded real valued functions such that $\forall (x,y) \in X \times Y$ we have $\phi(x) + \psi(y) \leq c(x,y)$. Then,*

1. *$\phi^c(y) + \phi(x) \leq c(x,y)$ and $\psi^{\bar{c}}(x) + \psi(y) \leq c(x,y)$.*

2. *$\phi(x) + \psi(y) \leq \phi^c(y) + \phi(x)$ and $\phi(x) + \psi(y) \leq \psi^{\bar{c}}(x) + \psi(y)$.*

*Proof.*     1. We see that

$$c(x,y) = \phi(x) + c(x,y) - \phi(x)$$
$$\geq \phi(x) + \inf_{z \in X} \left( c(z,y) - \phi(z) \right) = \phi(x) + \phi^c(y)$$

And the same for $c(x,y) = \psi(y) + \psi^{\bar{c}}(y)$.

2. Note that $\forall x \in X$,

$$\psi(y) \leq c(x,y) - \phi(x) \implies \psi(y) \leq \inf_{x \in X} \left( c(x,y) - \phi(x) \right) = \phi^c(y).$$

Similarly $\forall y \in Y$,

$$\phi(x) \leq c(x,y) - \psi(y) \implies \phi(x) \leq \inf_{y \in Y} \left( c(x,y) - \psi(y) \right) = \psi^{\bar{c}}(x).$$

$\square$

Every lower semicontinuous function on a Polish space is always measurable. The last lemma allows to substitute a pair $(\phi, \psi)$ by a pair $(\phi, \phi^c)$ in order to increase the objective value of the problem (DP). Following this procedure we can apply again the same result substituting $(\phi, \phi^c)$ by $\left( \phi^{c\bar{c}}, \phi^c \right)$.

The functions $\phi$ and $\psi$ are known in the literature related with optimal transport as Kantorovich's potentials. Recall the factories-customers problem, with an entrepreneur in charge

of the transportation. A function $c$ represents the cost for an agent to move a unit of mass from a factory $y \in Y$ to a customer $x \in X$. It is natural that our entrepreneur will take the merchandise to supply each customer from the closest factory. However, a more natural situation is the possibility that each customer $x \in X$ negotiates the price $\phi(x)$ that is able to pay for each unit of product. The entrepreneur will try to obtain the maximum profit possible. For this purpose, at each factory he needs to know which customers supply from the factory, such that the profit $\zeta(y)$ is maximum, i.e. $\zeta(y) = \sup_{x \in X} \phi(x) - c(x, y) = -\inf_{x \in X} c(x, y) - \phi(y) = -\phi^c(x)$, this is exactly the negative of the $c-$concave transform of the price $\phi(x)$ that each customer is able to pay. We can imagine a similar situation but this time we do not hire the entrepreneur, we just set the prices $\psi(y)$ at each factory $y \in Y$ and we wait for the customer to come. It is natural that each customer will try to minimize the purchase costs $\psi_-(y)$ plus transportation costs, that is[2] $\inf_{y \in Y} c(x, y) + \psi_-(y)$. From the perspective of the factory $\psi(y) = -\psi_-(y)$, leading to the $\bar{c}-$concave transform of the prices we are setting at each factory.

**Lemma 3.7.** *Let $X$ and $Y$ be two compact metric spaces. And let $c : X \times Y \to \mathbb{R}$ be a uniform continuous function with continuity modulus $\omega$. Then $\phi^c(x)$ and $\psi^{\bar{c}}$ are uniformly continuous with continuity modulus $\omega$, for any $\phi \in C(X)$ and $\psi \in C(Y)$.*

*Proof.* Given $\phi$, we set $g_x(y) = c(x, y) - \phi(x)$. The function $c(x, y)$ is uniformly continuous, we know that the modulus of continuity for $c$ is subadditive and nondecreasing, then

$$\begin{aligned}
\left| g_x(y) - g_x(y') \right| &= \left| c(x, y) - g(x) - c(x, y') + g(x) \right| \\
&= \left| c(x, y) - c(x, y') \right| \\
&< \omega(d((x, y), (x, y'))) \\
&\leq \omega(d(x, x) + d(y, y')) \\
&\leq \omega(d(y, y'))
\end{aligned}$$

Then, $g_x$ share the same modulus of continuity, for all $x$. By definition of $c-$concave transform $\phi^c(y) = \inf_{x \in X} g_x(y)$. If $g_x(y) \leq g_x(y') + \omega(d(y, y'))$ we take the infimum in both sides and we get $\phi^c(y) \leq \phi^c(y') + \omega(d(y, y'))$. Repeating the same argument exchanging $y$ and $y'$ in the inequality, we have $\phi^c(y') \leq \phi^c(y) + \omega(d(y, y'))$. Therefore $\phi^c$ shares the same modulus of continuity than $c$ implying that $\phi^c$ is uniformly continuous. The proof follows the same structure for the $\bar{c}-$concave transform.  □

**Theorem 3.9.** *Let $A \subset X$ and $B \subset Y$ be two compact subsets of two Polish Spaces $X$ and $Y$. Let $\mu$ and $\nu$ be probability measures defined on $A$ and $B$ respectively. And let $c : A \times B \to \mathbb{R}$ be a continuous and finite cost function. Then the problem,*

$$\sup \left\{ \int_A \phi \, d\mu + \int_B \psi \, d\nu : \quad \phi \in C(A), \ \psi \in C(B), \ \phi(a) + \psi(b) \leq c(a, b), \ \forall (a, b) \in A \times B \right\}$$

*admits a feasible solution $(\phi, \psi)$.*

*Proof.* The cost function is defined over the compact set $A \times B$, then is bounded and uniformly continuous. We can take a maximizing sequence $(\phi_n, \psi_n)_{n \in \mathbb{N}}$, by means of $c-$transforms we can improve it. The new sequence is uniformly continuous with same modulus of continuity than $c$. We still denote the new improved sequence by $\phi_n$ and $\psi_n$. Each pair $\phi_n, \psi_n$ is continuous defined on an compact set, implying that both functions are bounded. We can subtract from $\phi_n$ its minimum and add it to $\psi_n$, this does not change the value of the functional. Therefore without any loss of generality we can substitute the sequence by a new one such that $\min \phi_n = 0$ and $\psi_n$ is increased by the respective minimum of its partner. Then We have that $\max \phi_n \leq \omega(X)$.

Since $\psi_n$ is the $c-$concave transform, we have for all $y \in Y$, $\psi_n(y) \inf_{x \in X}(c(x, y) - \phi_n(x))$, therefore

$$\min_{(x,y) \in X \times Y} (c(x, y)) - \omega(X) \leq \psi_n(y) \leq \max c, \tag{3.16}$$

Proving that all the members of the sequence are bounded by the same constants. By the means of the Arzelà-Ascoli theorem the sequence $(\phi_n, \psi_n)$ converges uniformly to a pair of continuous functions

---

[2]This is known as a $\bar{c}-$convex transform of $\psi$.

in $(\phi, \psi)$. Then by uniform convergence,

$$\int_A \phi_n \mathrm{d}\mu + \int_B \psi_n \mathrm{d}\nu \to \int_A \phi \mathrm{d}\mu + \int_B \psi \mathrm{d}\nu \tag{3.17}$$

Uniform convergence implies pointwise convergence therefore, $\phi_n(x) + \psi_n(y) \leq c(x,y)$ implies that $\phi(x) + \phi(y) \leq c(x,y)$. Finding in this way an admissible pair for the problem.    □

Now we are able to introduce $c-$monotonicity, we will see the importance of this construction in order to formalize the relation between the dual problem and primal problem of Kantorovich's optimal transport problem formulation.

We start considering again our factories-costumers transportation problem. This time assume that the consortium is in charge of the transportation and it has already a plan to transport the products from factories to stores.

Let $c(x, y)$ be the transportation cost of sending a unit of product from a factory $y$ to a customer $x$. Currently, the plan is expensive and we would like to make it cheaper. For this purpose, we start redesigning the plan. We take a customer $x_1$ whose demand is in part supplied by a factory $y_1$. We take one unit of product from the demand supplied by $y_1$ and now we send it from a factory $y_2$. We send a product from $y_2$ to $x_1$ in such a way is cheaper than doing it from $y_1$ to $x_1$, after the above rearrangement we see that we earned $c(x_1, y_2) - c(x_1, y_1)$.

The supply of each factory is finite, then we have just taken one product from some store that is being supply by $y_2$, and we have sent it to $x_1$, leaving one product idle in $y_1$ and some store with a missing unit of product, without loss of generality let us call it $x_2$. Since we know that each customer needs to satisfy its demand, we take one product from a factory $y_3$ and we send it to the store $x_2$. Again we have one product missing in some store $x_3$, we continue this process redirecting one unit of product from a factory $y_{i+1}$ to a store $x_i$ until we finally have no choice to send the idle product in factory $y_1$ to a store $x_N$.

Therefore, at the end if we have,

$$c(x_1, y_1) + c(x_2, y_2) + \cdots + c(x_N, y_N) > c(x_1, y_2) + c(x_2, y_3) + \cdots + c(x_N, y_1),$$

we have made an improvement to the transportation cost. Note that if we are not able to find any rearrangement that decreases the cost we have already an optimal plan. This example allow us to introduce the next definition.

**Definition 3.4.** *Let $X, Y$ be arbitrary sets, and $c : X \times Y \to (-\infty, \infty]$ be a cost function. A subset $\Gamma \subset X \times Y$ is said to be c-cyclically monotone if $\forall N \in \mathbb{N}$ and any family of points $(x_1, y_1), (x_2, y_2), \ldots (x_N, y_N)$ of $\Gamma$, the following inequality holds,*

$$\sum_{i=1}^{N} c(x_i, y_i) \leq \sum_{i=1}^{N} c(x_i, y_{i+1}),$$

*considering $N + 1 = 1$.*

Since any permutation $\sigma$ over the set $\{1, \ldots, N\}$ can be written as a product of disjoint cycles, we have that this property satisfies,

$$\sum_{i=1}^{N} c(x_i, y_i) \leq \sum_{i=1}^{N} c(x_i, y_{\sigma(i)}) \tag{3.18}$$

**Definition 3.5** (*$c-$superdifferential*). *Let $c : X \times Y \to \mathbb{R} \cup \{+\infty\}$ be a cost function. Let $\xi : Y \to \mathbb{R} \cup \{-\infty\}$ be a real valued function defined on $Y$. Using $c-$concavity characterization for functions we call the $\bar{c}-$superdifferential of a function $\xi$ the c-cyclically monotone set,*

$$\partial^{\bar{c}} \xi = \left\{ (x, y) \subset X \times Y; \quad \xi(y) + \xi^{\bar{c}}(x) = c(x, y) \right\} \tag{3.19}$$

Note that for any $\bar{c}-$concave function we can find $\phi$ such that $\xi = \phi^c$, then we see that

$$\partial^{\bar{c}}\phi^c = \partial^{\bar{c}}\xi = \left\{(x,y) \subset X \times Y; \quad \phi^c(y) + \phi^{c\bar{c}}(x) = c(x,y)\right\}$$

We would like to have similar characterization for functions defined on the variable parameter of the cost function. We call $c-$superdifferential of a function $\phi : X \to \mathbb{R} \cup \{-\infty\}$ the set,

$$\partial^c\phi = \{(x,y) \in X \times Y; \quad \phi(x) + \phi^c(y) = c(x,y)\}$$

If $\phi$ is a $c-$concave function, then $\phi = \phi^{c\bar{c}}$, and taking the $\bar{c}$-superdifferential of $\phi^c$ we obtain,

$$\begin{aligned}
\partial^{\bar{c}}\phi^c &= \left\{(x,y) \subset X \times Y; \quad \phi^c(y) + \phi^{c\bar{c}}(x) = c(x,y)\right\} \\
&= \{(x,y) \subset X \times Y; \quad \phi^c(y) + \phi(x) = c(x,y)\} \\
&= \partial^c\phi.
\end{aligned}$$

We recall Rockafellar's theorem the subdifferentials of convex functions on $\mathbb{R}^n$ are characterized in terms of a cyclical monotonicity property.

**Theorem 3.10** (Rockafellar). *Let $\Gamma$ be a cyclically monotone set. In order that there exists a closed proper convex function $f$ on $\mathbb{R}^n$ such that $\Gamma \subset \partial f(x)$ for every $x$, it is necessary and sufficient that $\Gamma$ be cyclically monotone.*

The theorem 3.10 is a well known result in convex analysis. It basically states that every cyclically monotone set is contained in the graph of the subdifferential of a convex function. Note that a $\bar{c}-$concave function $\xi$ has the property $\xi = \xi^{\bar{c}c}$ (theorem 3.8),

We can provide an extension of Rockafellar's theorem in terms of $c-$concave functions. We can say that every $c-$cyclically monotone set is contained in the graph of the $c-$superdifferential of a $c-$concave function.

**Theorem 3.11.** *If $\Gamma$ is a not empty, c-cyclically monotone set in $X \times Y$ and $c : X \times Y \to \mathbb{R}$, then there is a $c-$concave function $\phi : X \to \mathbb{R} \cup \{-\infty\}$ and not everywhere $-\infty$ such that,*

$$\Gamma \subset \partial^c\phi = \{(x,y) \in X \times Y : \quad \phi(x) + \phi^c(y) = c(x,y)\} \tag{3.20}$$

**Theorem 3.12.** *Let $X$ and $Y$ be two Polish spaces. If $\gamma$ is an optimal transport plan for the cost $c : X \times Y \to \mathbb{R}$ and $c$ is continuous then $\mathrm{spt}(\gamma)$ is a c-cyclically monotone set.*

*Proof.* We proceed by contradiction. Suppose that $\mathrm{spt}(\gamma)$ is not $c-$cyclically monotone. Then there are a natural $k \in \mathbb{N}$, a permutation $\sigma$ and a set of pairs $\{(x_1,y_1),(x_2,y_2),...,(x_k,y_k)\} \subset \mathrm{spt}(\gamma)$ such that,

$$\sum_{i=1}^{k} c(x_i,y_i) > \sum_{i=1}^{k} c(x_i,y_{\sigma(i)})$$

Take $\epsilon \in \mathbb{R}$ satisfying

$$0 < k\epsilon < \sum_{i}^{k} c(x_i,y_i) - c(x_i,y_{\sigma(i)}).$$

Note that $\epsilon$ satisfying the above condition allows to write the inequality as follows,

$$\sum_{i=1}^{k} c(x_i,y_{\sigma(i)}) < \sum_{i=1}^{k}\left(c(x_i,y_{\sigma(i)}) + \frac{\epsilon}{2}\right) < \sum_{i=}^{k}\left(c(x_i,y_i) - \frac{\epsilon}{2}\right) < \sum_{i}^{k} c(x_i,y_i)$$

Since $c$ is continuous, there exists $r$ such that for any $i = 1,...,k$ and any $B(x_i;r) \times B(y_i;r)$ we have $c(x_i,y_i) - \epsilon < c(x,y)$. Similarly, we have $c(x,y) < c(x_i,y_{\sigma(i)}) + \epsilon$, $\forall (x,y) \in B(x_i;r) \times B(y_{\sigma(i)};r)$.

Now, consider the neighborhood $V_i = B(x_i; r) \times B(y_{\sigma(i)}; r)$. Given that $(x_i, y_i) \in \mathrm{spt}(\gamma)$, we have for all $i = 1, \ldots, k$ that $\gamma(V_i) > 0$. We set $\gamma_i = \frac{\gamma \llcorner V_i}{\gamma(V_i)}$, and $\mu_i = (\mathrm{proj}_X)_\# \gamma_i$ and $\nu_i = (\mathrm{proj}_Y)_\# \gamma_i$. Set $0 < \epsilon_0 < \frac{1}{k} \min_i \gamma(V_i)$.

Lemma 3.1 allows to construct for every $i$ an arbitrary coupling $\tilde{\gamma} \in \Pi(\mu_i, \nu_{\sigma(i)})$.

Set $\hat{\gamma} := \gamma - \epsilon_0 \sum_{i=1}^k \gamma_i + \epsilon_0 \sum_{i=1}^k \tilde{\gamma}_i$. Given that $\hat{\gamma}$ is a probability measure, we have that $\hat{\gamma} > 0$. Note that,

$$\epsilon_0 \gamma_i = \epsilon_0 \left( \frac{\gamma \llcorner V_i}{\gamma(V_i)} \right) < \frac{1}{k} \left( \min_i \gamma(V_i) \right) \left( \frac{\gamma \llcorner V_i}{\gamma(V_i)} \right) \le \frac{1}{k} (\gamma \llcorner V_i) \le \frac{\gamma}{k}$$

Then $\gamma - \sum_{i=1}^k \epsilon_0 \gamma_i > 0$, implying that $\hat{\gamma} > 0$. We see that $\hat{\gamma} = \Pi(\mu, \nu)$,

$$(\mathrm{proj}_X)_\# \hat{\gamma} = \mu - \epsilon_0 \sum_{i=1}^k \mu_i + \epsilon_0 \sum_{i=1}^k \mu_i = \mu$$

$$(\mathrm{proj}_Y)_\# \hat{\gamma} = \nu - \epsilon_0 \sum_{i=1}^k \nu_i + \epsilon_0 \sum_{i=1}^k \nu_{\sigma(i)} = \nu$$

Note that $\gamma_i$ is concentrated on $V_i = B(x_i; r) \times B(y_i; r)$, and $\tilde{\gamma}_i$ is concentrated on $B(x_i; r) \times B(y_{\sigma i}; r)$. And both are probability measures having total mass one, then we check have,

$$\int c \mathrm{d}\gamma - \int c \mathrm{d}\hat{\gamma} = \int c \mathrm{d}\gamma - \int c \mathrm{d}\gamma + \epsilon_0 \sum_{i=1}^k \int c \mathrm{d}\gamma_i - \epsilon_0 \sum_{i=1}^k \int c \mathrm{d}\tilde{\gamma}_i$$

$$= \epsilon_0 \sum_{i=1}^k \int c \mathrm{d}\gamma_i - \epsilon_0 \sum_{i=1}^k \int c \mathrm{d}\tilde{\gamma}_i$$

$$\ge \epsilon_0 \sum_{i=1}^k \int \left( c(x_i, y_i) - \frac{\epsilon}{2} \right) \mathrm{d}\gamma_i - \epsilon_0 \sum_{i=1}^k \int \left( c(x_i, y_{\sigma(i)}) + \frac{\epsilon}{2} \right) \mathrm{d}\tilde{\gamma}_i$$

$$= \epsilon_0 \left( \sum_{i=1}^k c(x_i, y_i) - \sum_{i=1}^k c(x_i, y_{\sigma(i)}) - k\epsilon \right) > 0$$

Therefore, $\hat{\gamma} < \gamma$ contradicting the assumption that $\gamma$ is optimal. Then $\mathrm{spt}(\gamma)$ must be $c-$cyclically monotone. $\qquad \square$

**Theorem 3.13.** *Suppose that $X$ and $Y$ are Polish spaces and suppose that $c : X \times Y \to \mathbb{R}$ is uniformly continuous and bounded. Then the problem* (DP) *admits a solution* $(\phi, \psi)$*. Moreover* $\psi = \phi^c$ *and the objective value of* (DP) *is equal to the objective value of the Kantorovich's problem* (KP)*.*

*Proof.* First consider the minimization problem (KP). Since $c$ is uniformly continuous, it is continuous, then (KP) admits a solution $\gamma$, and $\mathrm{spt}(\gamma)$ is a $c-$cyclically monotone set. Any $c-$cyclically monotone set is contained in the $c-$superdifferential of a $c-$concave function $\phi$. The uniform continuity of $c$ and $\phi$ being a $c-$concave function implies that $\phi$ and $\phi^c$ are continuous.

Now we check for boundedness of $\phi$ and $\phi^c$. Since $c$ is bounded we can take $(x_0, y_0)$ from $\mathrm{spt}(\gamma)$, such that $\phi(x_0) < \infty$ and $\phi^c(y_0) < \infty$. So,

$$\phi^c(y) = \inf_{x \in X} (c(x, y) - \phi(x)) \le \|c\|_\infty - \phi(x_0) \quad \forall y \in Y$$

$$\phi(x) = \phi^{c\bar{c}}(x) = \inf_{y \in Y} (c(x, y) - \phi^c(y)) \le \|c\|_\infty - \phi^c(y_0) \quad \forall x \in X$$

Proving that $\phi$ and $\phi^c$ are bounded form above. We get a lower bound, for both functions, injecting the above inequalities back into the definitions and $c$ bounded. So, $\forall x \in X$

$$\phi(x) = \inf_{y \in Y} (c(x,y) - \phi^c(y))$$

$$\geq \inf_{y \in Y} \left( c(x,y) - \left( \|c\|_\infty - \phi(x_0) \right) \right)$$

$$\geq - \|c\|_\infty + \phi(x_0) + \inf_{x \in X} \inf_{y \in Y} c(x,y)$$

$$\geq -2 \|c\|_\infty + \phi(x_0)$$

and all $y \in Y$,

$$\phi^c(y) = \inf_{x \in X} (c(x,y) - \phi(x))$$

$$\geq \inf_{x \in X} \left( c(x,y) - \left( \|c\|_\infty - \phi^c(y_0) \right) \right)$$

$$\geq - \|c\|_\infty + \phi^c(y_0) + \inf_{y \in Y} \inf_{x \in X} c(x,y)$$

$$\geq -2 \|c\|_\infty + \phi^c(y_0).$$

Having upper and lower bounds for $\phi$ and $\phi^c$, we can integrate $\phi$ and $\phi^c$ with respect to $\mu$ and $\nu$ respectively,

$$\int_X \phi \mathrm{d}\mu + \int_y \phi^c \mathrm{d}\nu = \int_{X \times Y} (\phi + \phi^c) \, \mathrm{d}\gamma = \int_{X \times Y} c(x,y) \mathrm{d}\gamma$$

This equality holds because of $\gamma$ being optimal is concentrated on a $c-$cyclically monotone set satisfying $\phi(x) + \phi^c(y) = c(x,y)$.

By definition of $\Delta$ in the dual problem (DP),

$$\sup (\mathrm{DP}) = \Delta \geq \int_X \phi \mathrm{d}\mu + \int_Y \phi^c \mathrm{d}\nu = \int_{X \times Y} c \mathrm{d}\gamma = \min (\mathrm{KP}) \tag{3.21}$$

We have that $\sup (\mathrm{DP}) \leq \min (\mathrm{KP})$ and we have an admissible optimal pair $(\phi, \phi^c)$, hence the desired equality for both problems $\max (\mathrm{DP}) = \min (\mathrm{KP})$ holds. $\qquad \square$

We have proved that the optimal plan is concentrated in a $c-$cyclically monotone set, impressively we can prove that a given $\gamma$ concentrated in a $c-$cyclically monotone set is optimal for its marginals.

**Theorem 3.14.** *Let $X$ and $Y$ two Polish Spaces, a cost function $c : X \times Y \to \mathbb{R}$ is given and it is uniformly continuous and bounded. Let $\gamma \in \mathcal{P}(X \times Y)$ a probability measure over $X \times Y$ such that $\mathrm{spt}(\gamma)$ is $c-$cyclically monotone. Then $\gamma$ is an optimal coupling for the measures $\mu = (\mathrm{proj}_X)_\# \gamma$ and $\nu = (\mathrm{proj}_Y)_\# \gamma$.*

*Proof.* We can find a $c-$concave function $\phi$ such that $\mathrm{spt}(\gamma)$ is contained in $\partial^c \phi$. Both $\phi$ and $\phi^c$ are $c-$concave and $\bar{c}-$concave, respectively. By continuity of $c$ we obtain that $\phi$ and $\phi^c$ are continuous. Note that they are also bounded. Therefore, we have satisfied the conditions to use the duality result,

$$\min (\mathrm{KP}) \leq \int_{X \times Y} c \mathrm{d}\gamma = \int_{X \times Y} \phi(x) + \phi^c(y) \mathrm{d}\gamma = \int_X \phi(x) \mathrm{d}\mu + \int_Y \phi^c(y) \mathrm{d}\nu \leq \max (\mathrm{DP}) = \min (\mathrm{KP})$$

Proving in this way that $\gamma$ is optimal. $\qquad \square$

If $c$ is lower semicontinuous we cannot assure the existence of an optimal pair, but we can assure duality.

**Lemma 3.8** (Stability). *Let $X$ and $Y$ be Polish spaces. Let $\mu$ and $\nu$ be two probability measures defined on $X$ and $Y$ respectively. Let $(c_n)_{n \in \mathbb{N}}$ be a sequence of lower-semicontinuous bounded from below functions converging uniformly increasingly to a cost function $c$ also bounded from below. Let $\gamma_n$ be a solution for the Kantorovich's problem for $c_n$,*

$$\gamma_n = \operatorname*{arg\,min}_{\gamma \in \Pi(\mu,\nu)} K(\gamma_n) = \operatorname*{arg\,min}_{\gamma \in \Pi(\mu,\nu)} \left( \int_{X \times Y} c_n \mathrm{d}\gamma \right)$$

*and let $\gamma$ be a solution for the Kantorovich's problem with $c$ as cost function. Then,*

$$\lim_{n \to \infty} K(\gamma_n) = K(\gamma).$$

**Theorem 3.15** (Duality and $c$ l.s.c.)**.** *Let $X$ and $Y$ be Polish spaces. If $c : X \times Y \to \mathbb{R} \cup \{+\infty\}$ is bounded from below and lower semicontinuous, then the equality $\sup (\text{DP}) = \min (\text{KP})$ holds.*

## Duality by convex analysis.

We can prove duality between the two problems introducing a perturbation (3.22), and then using convex analysis to prove a relation between the dual and primal transport problems. We can find an equivalent proof in [4] for cost functions satisfying the Lipschitz condition. As it is presented in [29] we do it for uniformly continuous cost functions, since we already have proved the existence of a minimizer.

**Definition 3.6.** *Suppose that $X$ and $Y$ are compact metric spaces and $c : X \times Y \to \mathbb{R}$ is uniformly continuous. For every $p \in C(X \times Y)$, let $H_\gamma : C(X \times Y) \to \overline{\mathbb{R}}$ be a perturbation of the problem,*

$$H_\gamma(p) = -\max \left\{ \int_X \phi \mathrm{d}\mu + \int_Y \psi \mathrm{d}\nu : \quad \phi(x) + \psi(y) \leq c(x, y) - p(x, y) \right\} \tag{3.22}$$

**Lemma 3.9.** *Under the setting of definition 3.6, $H_\gamma$ is convex.*

*Proof.* Take $p_0, p_1 \in C(X \times Y)$, with their optimal potentials $(\phi_0, \psi_0)$ and $(\phi_1, \psi_1)$. For $\alpha \in [0, 1]$, define $p_\alpha = \alpha p_0 + (1 - \alpha) p_1$, $\phi_\alpha = \alpha \phi_0 + (1 - \alpha) \phi_1$ and $\psi_\alpha = \alpha \psi_0 + (1 - \alpha) \psi_1$. The pair $(\phi_\alpha, \psi_\alpha)$ is admissible for $-H(p_\alpha)$. We have,

$$H(p_\alpha) \leq - \left( \int_X \phi_t \mathrm{d}\mu + \int_Y \psi_t \mathrm{d}\nu \right) = \alpha H(p_0) + (1 - \alpha) H(p_1).$$

**Lemma 3.10.** *Under the setting of definition 3.6, $H_\gamma$ is lower semicontinuous for the uniform convergence on the compact set $X \times Y$.*

For lower semicontinuity, take a sequence $(p_n)_{n \in \mathbb{N}}$ converging uniformly to $p$, $p_n \to p$. Extract a subsequence $(p_{n_k})_{n_k \in \mathbb{N}}$ approaching to limit inferior of $H(p_n)$. By the means of uniform convergence and the converse implication of Arzelà-Ascoli theorem, the sequence $(p_{n_k})_{n_k \in \mathbb{N}}$ is equicontinuous and bounded. Thus, its corresponding optimal potentials $(\phi_{n_k}, \psi_{n_k})$ are also equicontinuous and bounded, with $\phi_{n_k} \to \phi$ and $\psi_{n_k} \to \psi$ uniformly. The sequence satisfies $\phi_{n_k}(x) + \psi_{n_k} \leq c(x, y) - p_{n_k}(x, y)$, therefore the limit satisfies the same inequality, making it a feasible solution. Hence,

$$H(p) \leq - \left( \int_X \phi \mathrm{d}\mu + \int_Y \psi \mathrm{d}\nu \right) = \lim_{n \to \infty} H(p_{n_k}) = \liminf_{n \in \mathbb{N}} H(p_n). \tag{3.23}$$

$\square$

Taking the Legendre transform of $H^* : \mathcal{M}(X \times Y) \to \mathbb{R} \cup \{+\infty\}$. For every $\gamma \in \mathcal{M}(X \times Y)$, we have

$$\begin{aligned}
H^*(\gamma) &= \sup_{p \in C(X \times Y)} \left\{ \int_{X \times Y} p \mathrm{d}\gamma - H(p) \right\} \\
&= \sup_{p \in C(X \times Y)} \left\{ \int_{X \times Y} p \mathrm{d}\gamma + \max \left\{ \int_X \phi \mathrm{d}\mu + \int_Y \psi \mathrm{d}\nu : \quad \phi(x) + \psi(y) \leq c(x, y) - p(x, y) \right\} \right\} \\
&= \sup_{p \in C(X \times Y)} \left\{ \int_{X \times Y} p \mathrm{d}\gamma + \sup \left\{ \int_X \phi \mathrm{d}\mu + \int_Y \psi \mathrm{d}\nu : \quad \phi(x) + \psi(y) \leq c(x, y) - p(x, y) \right\} \right\}
\end{aligned}$$

Note that, if $\gamma \neq \mathcal{M}_+(X \times Y)$, then there is $p \leq 0$ such that $\int p_0 \mathrm{d}\gamma$, and one can take $\phi = 0$, $\psi = 0$, $p = c + n p_0$, and, for $n \to \infty$, we get $H^*(\gamma) = +\infty$. In the other hand, $\mathcal{M}_+(X \times Y)$, we should choose the largest possible $p$, that is $p(x, y) = c(x, y) - \phi(x) - \psi(y)$. This yields,

$$H^*(\gamma) = \sup_{\psi, \phi} \int_{X \times Y} c(x, y) \mathrm{d}\gamma. \tag{3.24}$$

Therefore,

$$H^*(\gamma) = \sup_{\phi,\psi} \int_{X \times Y} c(x,y)\mathrm{d}\gamma + \int_\phi \phi \mathrm{d}\mu - \int_\phi \phi(x)\mathrm{d}\gamma + \int_Y \phi \mathrm{d}\nu - \int_{X \times Y} \psi(y)\mathrm{d}\gamma,$$

this is exactly the same equation (3.12), then we can rewrite the constraints as follows,

$$H^*(\gamma) = \begin{cases} K(\gamma) & \text{if } \gamma \in \Pi(\mu,\nu) \\ +\infty & \text{otherwise} \end{cases} \tag{3.25}$$

We have that,

$$\max(\mathrm{DP}) = -H(0) = -H^{**}(0), \tag{3.26}$$

since $H$ is convex and lower semicontinuous. Moreover,

$$H^{**}(0) = \sup(\langle 0,\gamma \rangle - H^*(\gamma) = -\inf(H^*)) = -\min(\mathrm{KP}). \tag{3.27}$$

In this way we have duality between the two problems (KP) and (DP) if $c$ is uniformly continuous defined on the $X \times Y$ Cartesian product of compact subsets of two Polish spaces. Note that we did not rely on the concept of $c-$cyclically monotonicity for proving the equality between the two problems.

### Brenier's theorem

The following result is one of the most important results in the optimal transport theory. It brings consistence to many applications and allows to formulate many problems in terms of transportation maps.

**Theorem 3.16** (Brenier's theorem)**.** *Let $\mu$, $\nu$ probability measures on $\mathbb{R}^d$. Assume that $\mu$ is absolutely continuous with respect the Lebesgue measure, $\mu << \mathscr{L}^d$. Therefore,*

1. *There exists a unique transport plan.*

2. *The optimal transport plan is induced by a map.*

3. *The map is the gradient of a convex function.*

**Theorem 3.17.** *Let $\mu$ and $\nu$ probability measures on a compact domain $\Omega \subset \mathbb{R}^d$, there exists an optimal plan $\gamma$ for the cost $c(x,y) = h(|x-y|)$, with $h$ a strictly convex. The plan is unique and it has the form $\mu_\#(\mathrm{id},T)$, provided $\mu$ is absolutely continuous and the boundary $\partial\Omega$ is negligible. Moreover, there exists a function $\phi$, such that $T$ and $\phi$ are related by,*

$$T(x) = x - (\nabla h)^{-1}(\nabla \phi(x)) \tag{3.28}$$

**Theorem 3.18** (Distance between two Gaussians)**.** *The optimal transport plan for a cost function $c(x,y) = |x-y|^2$ between two Gaussians in $\mathbb{R}^d$ is a deterministic coupling, given by an affine map. That is, $\mu$ is a probability measure for the space $X = \mathbb{R}^n$, defined by,*

$$\mu(A) = \frac{1}{\sqrt{2\pi \det(\boldsymbol{\Sigma}_\mu)}} \int_A e^{(\mathbf{x}-\mathbf{m}_\mu)^\top \boldsymbol{\Sigma}_\mu^{-1}(\mathbf{x}-\mathbf{m}_\mu)} \mathrm{d}\mathbf{x}, \quad \forall A \in X,$$

*where $\boldsymbol{\Sigma}_\mu$ is its covariance matrix and $\mathbf{m}_\mu$ is its first moment. Let $\nu$ be a probability measure for the space $Y = \mathbb{R}^m$ defined by,*

$$\nu(B) = \frac{1}{\sqrt{2\pi \det(\boldsymbol{\Sigma}_\nu)}} \int_B e^{(\mathbf{y}-\mathbf{m}_\nu)^\top \boldsymbol{\Sigma}_\nu^{-1}(\mathbf{y}-\mathbf{m}_\nu)} \mathrm{d}\mathbf{y}, \quad \forall B \in Y$$
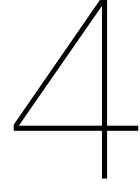
*where $\boldsymbol{\Sigma}_\mu$ is its covariance matrix and $\mathbf{m}_\mu$ is its first moment.*

*Then the plan $\gamma$ that solves the Kantorovich's problem is given by a translation map, that is $\gamma = \mu_\#(\mathrm{id},T)$, where $T$ is defined by,*

$$T(x) = \mathbf{m}_\nu + \mathbf{A}(\mathbf{x} - \mathbf{m}_\mu), \quad \mathbf{x} \in \mathbb{R}^n, \tag{3.29}$$

*and $\mathbf{A}$ is given by,*

$$\mathbf{A} = \boldsymbol{\Sigma}_\mu^{-\frac{1}{2}} \left( \boldsymbol{\Sigma}_\mu^{\frac{1}{2}} \boldsymbol{\Sigma}_\nu \boldsymbol{\Sigma}_\mu^{\frac{1}{2}} \right)^{\frac{1}{2}} \boldsymbol{\Sigma}_\alpha^{-\frac{1}{2}}. \tag{3.30}$$

<div align="right">

# 4

</div>

# Computation and Applications of an Optimal Transport

The approximation of an optimal transport is a computationally expensive problem. In this chapter we present the discrete formulation of the Kantorovich's problem. We can see that the discrete Kantorovich's problem is a subcase of a general linear program.

## Discrete Formulation and Linear Programming.

We denote by $\mathbb{1}_n = (1, 1, \dots, 1)^\top \in \mathbb{R}^n$ the vector composed of $n$ 1's. The vector composed of $n$ 0's is denoted by $\mathbb{0}_n = (0, 0, \dots, 0)^\top \in \mathbb{R}^n$. We call a simplex $\Sigma_n$ the convex set of vectors $\boldsymbol{p} = (p_1, \dots, p_n)^\top \in \mathbb{R}^n$ such that,

$$\Sigma_n = \left\{ \mathbf{p} \in \mathbb{R}^n : \sum_{i=1}^n p_i = 1 \text{ and } \mathbf{p} \geq 0 \right\}. \tag{4.1}$$

Let $X = \{x_1, x_2, \dots, x_n\}$ be a finite set of $n$ elements. Let $\mu$ a probability measure defined over $X$. Since $X$ is a finite set, we can express the probability measure $\mu$ as follows,

$$\mu = \sum_{i=1}^n a_i \delta_{x_i}, \tag{4.2}$$

where $\mathbf{a} = (a_1, a_2, \dots, a_n)^\top \in \Sigma_n$. In this section we represent the Kronocker's delta $\delta_{x_i} : X \to \{0, 1\}$ defined over a finite set $X$ by the equation,

$$\forall x_i \in X \quad \delta_{x_i}(x) = \begin{cases} 1 & x = x_i \\ 0 & x \neq x_i \end{cases}.$$

Hence, we have that $\mu(x_i) = a_i$ for all $x_i \in X$. Using the same notation, let $Y = \{y_1, y_2, \dots, y_m\}$ be a finite set of $m$ elements. Let $\nu$ be a probability measure defined over $Y$ defined by,

$$\nu = \sum_{i=1}^m b_i \delta_{y_i}, \tag{4.3}$$

where $\mathbf{b} = \{b_1, b_2, \dots, b_m\}^\top \in \Sigma_m$, and $\delta_{y_i} : Y \to \{0, 1\}$ is a Kronecker's delta defined on $Y$. We see that $\mu$ and $\nu$ are identified by a vector $\mathbf{a} \in \Sigma_n$ and $\mathbf{b} \in \Sigma_m$ respectively.

We can construct a coupling between $\mu$ and $\nu$, using matrix notation. Let $\mathbf{M}^{n \times m} \ni \boldsymbol{\gamma}$ be a matrix of $n \times m$ real entries, i.e. $(\boldsymbol{\gamma})_{i,j}$ with its marginals given by $\mathbf{a}$ and $\mathbf{b}$,

$$\boldsymbol{\gamma} \mathbb{1}_m = \sum_{j=1}^m \gamma_{i,j} = a_i \quad \text{and} \quad \boldsymbol{\gamma}^\top \mathbb{1}_n = \sum_{i=1}^n \gamma_{i,j} = b_j \tag{4.4}$$

From the equation 4.4, and the fact that $\mathbf{a} \in \Sigma_n$, we can see that the sum of all the entries of $\boldsymbol{\gamma}$ is also 1,

$$\sum_{i=1}^{n}\sum_{j=1}^{m} \gamma_{i,j} = \sum_{i=1}^{n} a_i = \sum_{j=1}^{m}\sum_{i=1}^{n} \gamma_{i,j} = \sum_{j=1}^{m} b_j = 1 \tag{4.5}$$

Therefore we can see $\boldsymbol{\gamma}$ as a joint probability measure over the space $X \times Y$ with marginals given by $\mu$ and $\nu$.

Let $c : X \times Y \to \mathbb{R}$ a cost function. The $X$ and $Y$ are finite sets with $n$ and $m$ elements respectively, therefore we can use a matrix $\mathbf{C} \in \mathbf{M}^{n \times m}$ to represent the cost function.

$$(\mathbf{C})_{i,j} = c_{i,j} = c(x_i, y_j). \tag{4.6}$$

We see that the dimensions of $\mathbf{C}$ are equal to the dimensions of $\boldsymbol{\gamma}$. Then given a transportation plan between $\mu$ and $\nu$ the total cost is given by,

$$\langle \mathbf{C}, \boldsymbol{\gamma} \rangle = \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}} c_{i,j}\gamma_{i,j} = \mathrm{tr}\left(\mathbf{C}^{\mathsf{T}}\boldsymbol{\gamma}\right), \tag{4.7}$$

where $\langle \cdot, \cdot \rangle$ in this setting is the Frobenius inner product [1] of the two matrices.

## Linear Programming Notation.

As mentioned before, the discrete formulation is a specific case of linear programming in finite dimensions. Given the discrete analysis presented above the transportation problem is reduced to,

$$\Theta(\mu, \nu) = \min_{\boldsymbol{\gamma} \in \Pi(\mu,\nu)} \langle \mathbf{C}, \boldsymbol{\gamma} \rangle \tag{4.8}$$

The vectorization $\mathrm{vec}(\mathbf{R})$ of a matrix $\mathbf{R} \in \mathbf{M}^{n \times m}$ is a linear transformation which converts the matrix into a column vector $\mathbf{m} \in \mathbb{R}^{mn}$. If $\mathbf{z} \in \mathbf{M}^{1 \times n}$ is one row matrix, we see that its vectorization $\mathrm{vec}(\mathbf{z}) = \mathbf{z}^{\mathsf{T}} \in \mathbb{R}^n$ is a column vector. Given two matrices $\mathbf{R} \in \mathbf{M}^{n \times k}$ and $\mathbf{P} \in \mathbf{M}^{k \times m}$ with real entries, this operation has the properties:

$$\mathrm{vec}(\mathbf{RP}) = (\mathbf{I}_m \otimes \mathbf{R})\mathrm{vec}(\mathbf{P}) = (\mathbf{P}^{\mathsf{T}} \otimes \mathbf{I}_n)\mathrm{vec}(\mathbf{R}) \tag{4.9}$$

Let $\mathbf{R}$, $\mathbf{P}$ and $\mathbf{Q}$ be matrices such that we can compute the product $\mathbf{RPQ}$ and let $\mathbf{S}$ be the result of the product. We can write the equation $\mathbf{RPQ} = \mathbf{S}$ as,

$$\left(\mathbf{Q}^{\mathsf{T}} \otimes \mathbf{R}\right)\mathrm{vec}(\mathbf{P}) = \mathrm{vec}(\mathbf{S}) \tag{4.10}$$

And the Frobenius inner product between two matrices $\mathbf{R}$ and $\mathbf{P}$ in the space of real matrices $\mathbf{M}^{n \times m}$ can be written as,

$$\mathrm{tr}(\mathbf{R}^{\mathsf{T}}\mathbf{P}) = \mathrm{vec}(\mathbf{R})^{\mathsf{T}}\mathrm{vec}(\mathbf{P}) \tag{4.11}$$

Applying the vectorization property for the Frobenius product (4.8) we have that,

$$\langle \mathbf{C}, \boldsymbol{\gamma} \rangle = \mathbf{c}^{\mathsf{T}}\mathbf{p}.$$

The $nm$-dimensional vectors $\mathbf{c} = \mathrm{vec}(\mathbf{C})$ and $\mathbf{p} = \mathrm{vec}(\boldsymbol{\gamma})$ are the vectorization of $\mathbf{C}$ and $\boldsymbol{\gamma}$ equal to the stacked columns contained in the cost matrix and transportation plan respectively.

Consider the constraint for $\mathbf{b}$ take the transpose and right multiply both sides by the identity,

$$\mathbb{1}_n^{\mathsf{T}}\boldsymbol{\gamma}\mathbf{I}_m = \mathbf{b}^{\mathsf{T}}\mathbf{I}_m = \mathbf{b}^{\mathsf{T}}$$

---

[1] For matrices with real entries the product is defined as $\langle A, B \rangle = \mathrm{tr}(A^{\mathsf{T}}B)$.

Using (4.9) in the constraints (4.4), and the vectorization of a product of three matrices (4.10) in the last equation we obtain,

$$\boldsymbol{\gamma} \mathbb{1}_m = \left(\mathbb{1}_m^{\mathsf{T}} \otimes \mathbf{I}_n\right) \text{vec}(\boldsymbol{\gamma}) = \left(\mathbb{1}_m^{\mathsf{T}} \otimes \mathbf{I}_n\right) \mathbf{p} = \mathbf{a} \tag{4.12}$$

$$\text{vec}(\mathbf{b}^{\mathsf{T}}) = \left(\mathbf{I}_m \otimes \mathbb{1}_n^{\mathsf{T}}\right) \text{vec}(\boldsymbol{\gamma}) = \left(\mathbf{I}_m \otimes \mathbb{1}_n^{\mathsf{T}}\right) \mathbf{p} = \mathbf{b} \tag{4.13}$$

Therefore, we can write the equation 4.8 as a linear program,

$$\min_{\mathbf{p} \in \mathbb{R}^{mn}} \quad \mathbf{c}^{\mathsf{T}} \mathbf{p} \tag{4.14}$$

$$\text{subject to} \quad \mathbf{A}\mathbf{p} = \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \tag{4.15}$$

$$\mathbf{p} \geq 0 \qquad , \tag{4.16}$$

where the linear operator of the constraint is given by a $(m+n) \times mn$ matrix,

$$\mathbf{A} = \begin{pmatrix} \mathbb{1}_m^{\mathsf{T}} \otimes \mathbf{I}_n \\ \mathbf{I}_m \otimes \mathbb{1}_n^{\mathsf{T}} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_n & \cdots & \cdots & \mathbf{I}_n \\ \mathbb{1}_n^{\mathsf{T}} & \mathbb{0}_n^{\mathsf{T}} & \cdots & \mathbb{0}_n^{\mathsf{T}} \\ \mathbb{0}_n^{\mathsf{T}} & \mathbb{1}_n^{\mathsf{T}} & \cdots & \mathbb{0}_n^{\mathsf{T}} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{0}_n^{\mathsf{T}} & \vdots & \ddots & \vdots \end{pmatrix} \tag{4.17}$$

Then the dual formulation of the problem is given by,

$$\max_{\mathbf{h} \in \mathbb{R}^{m+n}} \mathbf{h}^{\mathsf{T}} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \tag{4.18}$$

$$\text{subject to } \mathbf{h}^{\mathsf{T}} \mathbf{A} \leq \mathbf{c}^{\mathsf{T}}. \tag{4.19}$$

## C-transform improvement and Duality.

We know that the primal and the dual problem reach the share the same optimal value. Without loss of generality we can write $\mathbf{h} = (\mathbf{f}, \mathbf{g})^{\mathsf{T}}$ as two vectors $\mathbf{f} \in \mathbb{R}^n$ and $\mathbf{g} \in \mathbb{R}^m$. Note matrix product $\mathbf{d}$, between $(\mathbf{f} \ \mathbf{g}) = (f_1, f_2, \ldots f_n, g_1, g_2 \ldots, g_m)$ and the matrix $\mathbf{A}$ can be written column by column using $\mathbf{A}^k$ to denote the $(k+1)$-th column of $\mathbf{A}$ for $0 \leq k < nm$ as follows,

$$\mathbf{d}^{k+1} = (\mathbf{f} \ \mathbf{g})\mathbf{A}^k = f_{i+1} + g_{j+1}, \tag{4.20}$$

where the indexes $i$ and $j$ are related by,

$$\begin{array}{ll} i & = k \bmod n \\ jn & = k - i \end{array} \tag{4.21}$$

We remember that each natural $0 \leq k \leq mn - 1$, is represented uniquely by $k = nj + i$ with $0 \leq j \leq m - 1$ and $0 \leq i \leq n - 1$, hence the product $\mathbf{h}^{\mathsf{T}} \mathbf{A}$ is exactly the direct sum $\mathbf{f} \oplus \mathbf{g}$.

In this way, the dual formulation of the Kantorovich's problem for the discrete is given by,

$$\max_{(\mathbf{f}, \mathbf{g}) \in \mathbb{R}^{n+m}} \langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{g}, \mathbf{b} \rangle \tag{4.22}$$

$$\text{subject to } \mathbf{f} \oplus \mathbf{g} \leq \mathbf{c}^{\mathsf{T}} \tag{4.23}$$

This result reminds us the dual formulation for the continuous case. We see that the complexity of our problem is reduced from $nm$ to $(m+n)$.

It is interesting analyzing the $\mathbf{C}-$ and $\overline{\mathbf{C}}-$ transforms of our variables $\mathbf{f}$ and $\mathbf{g}$. In equivalent way that we did before for the general formulation we define the transformations as follows,

$$\left(\mathbf{f}^{\mathbf{C}}\right)_j = \min_{1 \leq i \leq n} \left((\mathbf{C})_{i,j} - (\mathbf{f})_i\right) \tag{4.24}$$

$$\left(\mathbf{g}^{\overline{\mathbf{C}}}\right)_i = \min_{1 \leq j \leq m} \left((\mathbf{C})_{i,j} - (\mathbf{g})_j\right). \tag{4.25}$$

As we did before we can check the improvement of a solution through $\mathbf{C}-$ and $\overline{\mathbf{C}}-$transforms. We check that for each $\left(\mathbf{f}^{\mathbf{C}}\right)_j$ the constraint is satisfied,

$$0 = -(\mathbf{f}^{\mathbf{C}})_j + \min_{1 \leq i \leq n}\left((\mathbf{C})_{i,j} - (\mathbf{f})_i\right)$$

$$0 = \min_{1 \leq i \leq n}\left((\mathbf{C})_{i,j} - (\mathbf{f}^{\mathbf{C}})_j - (\mathbf{f})_i\right)$$

$$0 \leq (\mathbf{C})_{i,j} - (\mathbf{f}^{\mathbf{C}})_j - (\mathbf{f})_i$$

For a given pair $(\mathbf{f}, \mathbf{g})$ satisfying the constraint we have for all $1 \leq i \leq n$ and $1 \leq j \leq m$,

$$\begin{aligned}
(\mathbf{g})_j &\leq \mathbf{C}_{i,j} - (\mathbf{f})_i \\
&\leq \min_{1 \leq i \leq n}\left(\mathbf{C}_{i,j} - (\mathbf{f})_i\right) \\
&\leq \left(\mathbf{f}^{\mathbf{C}}\right)_j
\end{aligned}$$

Since $\mathbf{a} \in \Sigma_n$ and $\mathbf{b} \in \Sigma_m$ are positive,

$$\langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{g}, \mathbf{b} \rangle \leq \langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{f}^{\mathbf{C}}, \mathbf{b} \rangle. \tag{4.26}$$

Using similar arguments we can prove that the $\overline{\mathbf{C}}-$transform of $\mathbf{g}$ also satisfies the constraint and improves the objective value,

$$\langle \mathbf{f}, \mathbf{a} \rangle + \langle \mathbf{g}, \mathbf{b} \rangle \leq \left\langle \mathbf{g}^{\overline{\mathbf{C}}}, \mathbf{a} \right\rangle + \langle \mathbf{g}, \mathbf{b} \rangle \tag{4.27}$$

We see that $\mathbf{f} \leq \mathbf{x}$ implies $\mathbf{x}^{\mathbf{C}} \leq \mathbf{f}^{\mathbf{C}}$. And exactly as we did for the continuous case we have $\mathbf{f} \leq \mathbf{f}^{\mathbf{C}\overline{\mathbf{C}}}$ and $\mathbf{g} \leq \mathbf{f}^{\overline{\mathbf{C}}\mathbf{C}}$. And the process stabilizes in $\mathbf{f} = \mathbf{f}^{\mathbf{C}\overline{\mathbf{C}}\mathbf{C}}$. The complementary slackness conditions for the discrete transportation problem can be read as follows,

**Proposition 4.1.** *Let $\gamma^\star$ and $(\mathbf{f}^\star, \mathbf{g}^\star)$ be optimal solutions for the primal and dual problems, respectively. Then, $\left(\gamma_{i,j}\right)\left((\mathbf{C})_{i,j} - \mathbf{f}_i - \mathbf{g}_j\right) = 0$ holds. That is,*

   *1. If $(\gamma^\star)_{i,j} > 0$ then necessarily $(\mathbf{f}^\star)_i + (\mathbf{g}^\star)_j = (\mathbf{C})_{i,j}$*

   *2. If $(\mathbf{f}^\star)_i + (\mathbf{g}^\star)_j < (\mathbf{C})_{i,j}$ then necessarily $(\gamma^\star)_{i,j} = 0$*

## The simplex method is not polynomial time.

The linear program for the optimal transport problem can be solved by any of the algorithms presented before in this text. Although neither of the methods above developed are polynomial-time. Victor Klee and George Minty exhibited a class of linear programs requiring an exponential number of iterations when solved by the conventional simplex method. As it is proved in [17] the following linear program under the pivoting rule of the least reduced cost, requires $2^n - 1$ iterations before optimality is reached and, hence, they are exponential.

$$\min_{\mathbf{x} \in \mathbb{R}^{2n}} \quad \sum_{j=1}^{n} -\epsilon^{n-1} x_j$$

$$\text{subject to} \quad \begin{cases}
x_1 + x_{n+1} &= 1, \\
2\left(\epsilon x_1 + x_2 + x_{n+2}\right) &= 1 \\
2\left(\epsilon^2 x_1 + \epsilon x_2 + x_3 + x_{n+3}\right) &= 1 \\
\quad \vdots \quad \vdots \\
2\sum_{j=1}^{n} \epsilon^{n-j} x_j + x_{2n} &= 1
\end{cases}$$

## Network Simplex Method.

In the other hand, we have abused of the terminology calling simplex method to the criterion of choosing the basis according to the value of the least reduced cost, but indeed the simplex

method is a set of different ways to choose a basis in order to find an optimal solution. The nature of the transportation problem allows to implement different versions of the simplex algorithm that exploit **the mass conservation property**, also known as *balanced transportation*, to develop more efficient algorithms.

$$
\begin{array}{cccc|c}
\gamma_{1,1} & \gamma_{2,1} & \cdots & \gamma_{1,m} & a_1 \\
\gamma_{2,1} & \gamma_{2,2} & \cdots & \gamma_{2,m} & a_2 \\
\vdots & \ddots & \ddots & \vdots & \vdots \\
\gamma_{n,1} & \gamma_{n,2} & \cdots & \gamma_{n,m} & a_n \\
\hline
b_1 & b_2 & \cdots & b_m & 1
\end{array}
$$

The most common way to choose the basis in for optimal transport problems is the so called **Network simplex method**, which is polynomial time. The main idea behind this method is that a transportation plan can be seen as the graph of a network connecting sources and sinks, with a flow through the network satisfying that each source $x_i$ is flowing out $a_i$ units of mass, and sink $y_j$ is flowing into $b_j$ units of mass. Each edge connecting a source with a sink is associated with a cost $(\mathbf{C})_{i,j}$. The Network simplex methods needs to be initialized with a feasible solution, this step can be done in a simple and efficient way using the North-west Corner Rule. The rule starts assigning the highest possible value to $(\boldsymbol{\gamma})_{1,1} = \min\{a_1, b_1\}$. At each step the entry $(\boldsymbol{\gamma})_{i,j}$ is chosen in order to saturate either the row at $i$, the column at $j$ or both. The rule finish until we reach the position $(\boldsymbol{\gamma})_{n,m}$. Briefly, the North west rule can be summarized as follows:

- Initialize $i = 1$, $j = 1$. And set $r \leftarrow a_1$ and $c \leftarrow b_1$.

- While $i \leq n$, $j \leq m$,

  1. Set $t \leftarrow \min\{c, r\}$.

  2. Set $(\boldsymbol{\gamma})_{i,j} \leftarrow t$

  3. Set $r \leftarrow r - t$ and $c \leftarrow c - t$.

  4. If $r = 0$ then increment $i \leftarrow \min\{m, i+1\}$, and update $r \leftarrow a_i$.

  5. If $c = 0$ then increment $j \leftarrow \min\{m, j+1\}$, and update $c \leftarrow b_j$.

The algorithm constructs the graph of an extreme point matrix $\boldsymbol{\gamma}$. That is the transportation plan at each stage is a basic feasible solution of the constraints. It can be proved that the graph of a matrix $\boldsymbol{\gamma}$ that is a basic feasible solution has no cycles, then the graph is a tree or a disjoint union of trees (called forest), and it has no more than $m + n - 1$ non-zero entries.

The Network simplex method relies on the fact that complementary slackness conditions, $\gamma_{i,j} > 0 \implies (\mathbf{f})_i + (\mathbf{g})_j = (\mathbf{C})_{i,j}$, allows to find a pair $(\mathbf{f}, \mathbf{g})$ (not necessarily dual feasible) for each feasible matrix $\boldsymbol{\gamma}$ and by corollary 2.5 if the pair is dual feasible then we have reached optimality. If the pair $(\mathbf{f}, \mathbf{g})$ is not feasible then it is modified to get closer to feasibility.

We can compute a complementary pair $(\mathbf{f}, \mathbf{g})$ starting from one some the edges satisfying $(\gamma)_{i'_1, j'_1} > 0$, then we set $(\mathbf{f})_{i'_1} = 0$ and then we compute $(\mathbf{g})_{j'_1} = (\mathbf{C})_{i'_1, j'_1} - \mathbf{f}_{i'_1}$. Imagine that does exist an edge of the graph $\boldsymbol{\gamma}$ such that the corresponding entry is strictly greater than zero and it is connecting a source $x_{i'_2}$ to the sink $y_{j'_1}$, then we set $(\mathbf{f})_{i'_2} = (\mathbf{C})_{i'_2, j'_1} - (\mathbf{g})_{j'_1}$, we repeat this process until we cover the whole tree. Since the graph is composed of a disjoint union of trees (If the graph is a tree and assuming the non-degeneracy condition we will cover all the nodes), we repeat for all the disconnected trees covering in this way all the elements of $\mathbf{f}$ and $\mathbf{g}$ (again assuming the non-degeneracy condition).

(a) Graph of the network connecting sources and sinks.

(b) Feasible transportation plan.

(c) Basic Feasible transportation plan.



$(\mathbf{C})_{i,j}$        $(\mathbf{C})_{i,j}$        $(\mathbf{C})_{i,j}$

The network simplex algorithm maintains a basic feasible solution at each stage. Basically, given a basic feasible solution $\boldsymbol{\gamma}$ we compute a complementary pair $(\mathbf{f},\mathbf{g})$ and we check for feasibility. If the complementary pair is not feasible we create a new edge in the graph from one of the pairs $(x_i, y_j)$ violating $\mathbf{C}_{i,j} < (\mathbf{f})_i + (\mathbf{g})_j$ to the edges to the graph. Two situations arise:

- The graph keeps its property of being a forest, in this case we only recompute a complementary pair $(\mathbf{f},\mathbf{g})$ for the new graph. This situation does not represent any change in $\boldsymbol{\gamma}$, it is just the result of the arbitrary way we have found a pair $(\mathbf{f},\mathbf{g})$.

- The graph has a cycle, in this case we need to remove a pair from the graph of $\boldsymbol{\gamma}$, to ensure that we have a forest and modify $\boldsymbol{\gamma}$ in order to be consistent with the new graph. Let $i'_1 \to j'_1 \to i'_2 \to j'_2 \to i'_3 \to \dots \to i'_K \to j'_K \to i'_1$ the directed path of the cycle we have created. Increase the flow of all the edges $(i', j')$ and decrease it by the same amount in the edges $(j', i')$. Let $\tilde{\boldsymbol{\gamma}}$ the new matrix resulting from this operation, that is,

$$
\begin{aligned}
(\tilde{\boldsymbol{\gamma}})_{i'_k, j'_k} &= (\boldsymbol{\gamma})_{i'_k, j'_k} + \epsilon & \forall k = 1, \dots, K \\
(\tilde{\boldsymbol{\gamma}})_{i'_{k+1}, j'_k} &= (\boldsymbol{\gamma})_{i'_{k+1}, j'_k} - \epsilon & \forall k = 1, \dots, K-1 \\
(\tilde{\boldsymbol{\gamma}})_{i'_1, j'_k} &= (\boldsymbol{\gamma})_{i'_1, j'_K} - \epsilon &
\end{aligned}
$$

Take $\epsilon = \min\{\gamma_{i_2,j_1}, \dots, \gamma_{i_K,j_{K-1}}, \gamma_{1,K}\}$, in order to keep the constraint $\tilde{\boldsymbol{\gamma}} \geq 0$. Having in this way a tree, since we have removed an edge from a cycle. Finally, we update the transport plan $\boldsymbol{\gamma} \leftarrow \tilde{\boldsymbol{\gamma}}$.

Briefly, we initialize the algorithm with a feasible solution $\boldsymbol{\gamma}$ using the North-West rule, resulting into a feasible solution, (not necessarily basic feasible) and we compute a complementary pair $(\mathbf{f},\mathbf{g})$ for $\boldsymbol{\gamma}$. We check for feasibility of the computed pair. If the pair is feasible we have an optimal solution. If the pair $(\mathbf{f},\mathbf{g})$ is not feasible we search for cycles and we make them disappear using the above explained procedure. We compute again the complementary pair $(\mathbf{f},\mathbf{g}$ and we add to the graph of $\boldsymbol{\gamma}$ the edge corresponding to the indexes that violate the complementary condition, and we continue with the above explained procedure.

We recommend [8], [25] and [32] for a detailed description of the algorithm (for example how to deal with degenerated solutions) and the proofs for the convergence and polynomial order of the Network simplex method.

## Sinkhorn-Knopp Algorithm.

The *Shannon-Neumann anecdote* is a famous conversation between these two great mathematicians Claude Shannon[2] and John Neumann[3], that occurred in the time period fall 1940 to spring 1941 in New Jersey in which von Neumann suggested Shannon to use the term entropy in one of his papers. In April of 1961, Myron Tribus visits Shannon at his office at MIT and questions him about the reason behind his "entropy" name adoption.

---

[2] Claude Elwood Shannon (April 30, 1916 – February 24, 2001). American mathematician and electrical engineer known as the *Father of the information theory*.

[3] John von Neumann (December 28, 1903 – February 8, 1957). Hungarian-American mathematician and chemical engineering known for applications in many fields covering from computer sciences to physics and mathematics.

Tribus recounts the Shannon interview [33] as such:

> I thought of calling it "information". But the word was overly used, so I decided to call it "uncertainty". When I discussed it with John von Neumann, he had a better idea: (...) "You should call it entropy, for two reasons. In first place your uncertainty has been used in statistical mechanics under that name, so it already has a name. In second place, and more important, no one knows what entropy really is, so in a debate you will always have the advantage."

Informally, we can say that the term introduced by Shannon as **Entropy** quantifies how much information there is in a random signal, language. In relation to our problem, in [8] is discussed the following: "In practice actual traffic patterns in a network do not agree with those predicted by the solution of the optimal transport problem. Indeed, the former are more diffuse than the latter, which tend to rely on a few routes as a result of the sparsity of optimal couplings".

The maximum entropy of a set of events is reached when the probability of occurrence is uniform, hence we can interpret the entropy also as a measure of the chaos inside of a system.

In general for any vector representing a probability measure the Shannon's Entropy $H(\mathbf{p})$ : $\mathbb{R}^n \to \overline{\mathbb{R}}$ can be expressed as,

$$H(\mathbf{p}) = \begin{cases} -\infty & \exists i \in \mathbb{N},\ 1 \le i \le n \text{ and } \mathbf{p}_i = 0 \\ -\sum_{i=1}^n p_i\left(\log(p_i) - 1\right) & \mathbf{p} > 0 \end{cases} \tag{4.28}$$

We can consider $\mathbf{p} = \mathrm{vec}(\boldsymbol{\gamma})$ as the vectorization of a transportation plan matrix. The entropy of the coupling is equivalent to then entropy of its vectorization, that is $H(\boldsymbol{\gamma}) = H(\mathbf{p})$. Note that given a matrix with positive entries the Hessian has the form $\delta^2 H(\mathbf{p}) = -\mathrm{diag}\left(\frac{1}{\gamma_{i,j}}\right)$, which implies that is strongly concave in $\mathbf{p}$.

The idea is to introduce a strictly convex regularization to the transportation cost in order to have a strictly convex program and hence a unique solution for the problem.

$$\Theta(\mu, \nu; \epsilon) = \min_{\boldsymbol{\gamma} \in \Pi(\mu,\nu)} \langle \mathbf{C}, \boldsymbol{\gamma} \rangle - \epsilon H(\boldsymbol{\gamma}), \tag{4.29}$$

where $\epsilon > 0$, and $\Pi(\mu, \nu)$ is the set of couplings between two probability measures $\mu$ and $\nu$ on finite dimensional spaces, identified by vectors $\mathbf{a}$ and $\mathbf{b}$ respectively.

This regularization produces a blurred prediction of the transportation plan given marginals and transportation costs.

**Theorem 4.1.** *Let $\tilde{\gamma}^\epsilon$ the solution for* (4.29) *converges to the solution $\gamma^\star$ that has the maximum entropy among the solutions $\boldsymbol{\gamma}$ of the problem* (4.8) *as $\epsilon \to 0$.*

*Proof.* Consider a sequence $(\epsilon_\alpha)_{\alpha \in \mathbb{N}} \to 0$ converging to zero, such that $\epsilon_\alpha > 0$ for any $\alpha \in \mathbb{N}$. Let $\boldsymbol{\gamma}_\alpha$ the solution of the $\epsilon-$convex program for $\epsilon = \epsilon_\alpha$. The set of couplings $\Pi(\mu, \nu)$ is a closed convex subset of the simplex $\Sigma_{nm}$, therefore we can extract a convergent subsequence (for the sake of simplicity we use the same index of the sequence), such $\boldsymbol{\gamma}_\alpha \to \boldsymbol{\gamma}^\star$ as $\alpha \to \infty$. We are in a closed set therefore $\boldsymbol{\gamma}^\star \in \Pi(\mu, \nu)$.

In the other hand consider any $\boldsymbol{\gamma}$, such that $\Theta(\mu, \nu) = \langle \mathbf{C}, \boldsymbol{\gamma} \rangle$. Since $\boldsymbol{\gamma}$ is optimal for (4.8),

$$0 \le \langle \mathbf{C}, \boldsymbol{\gamma}_\alpha \rangle - \langle \mathbf{C}, \boldsymbol{\gamma} \rangle, \tag{4.30}$$

and $\boldsymbol{\gamma}_\alpha$ is optimal for (4.29), then

$$0 \le (\langle \mathbf{C}, \boldsymbol{\gamma} \rangle - \epsilon_\alpha H(\boldsymbol{\gamma})) - (\langle \mathbf{C}, \boldsymbol{\gamma}_\alpha \rangle - \epsilon_\alpha H(\boldsymbol{\gamma}_\alpha))$$
$$0 \le (\langle \mathbf{C}, \boldsymbol{\gamma} \rangle - \langle \mathbf{C}, \boldsymbol{\gamma}_\alpha \rangle) + \epsilon_\alpha \left(H(\boldsymbol{\gamma}_\alpha) - H(\boldsymbol{\gamma})\right). \tag{4.31}$$

Multiplying both sides of the inequality (4.30) by two and adding it to the inequality (4.31) we have the following inequality,

$$0 \leq \langle \mathbf{C}, \boldsymbol{\gamma}_\alpha \rangle - \langle \mathbf{C}, \boldsymbol{\gamma} \rangle \leq \epsilon_\alpha \left( H(\boldsymbol{\gamma}_\alpha) - H(\boldsymbol{\gamma}) \right). \tag{4.32}$$

We have that $H$ and $\langle \mathbf{C}, \cdot \rangle$ are continuous functions, therefore taking the limit $\alpha \to \infty$ and by the above inequality we obtain $\langle \mathbf{C}, \boldsymbol{\gamma} \rangle = \langle \mathbf{C}, \boldsymbol{\gamma}^\star \rangle$, implying that $\boldsymbol{\gamma}^\star$ is optimal and feasible. Moreover, dividing by $\epsilon_\alpha$ we obtain that $H(\boldsymbol{\gamma}) \leq H(\boldsymbol{\gamma}^\star)$ proving that in the limit we have an optimal feasible solution for (4.8) with higher entropy that anyone else that is also optimal.                                                                                       □

Sikhorn Algorithm

There are many methods that we can use to solve a convex program, for example one of the different versions of the Newton's methods or Backward propagation algorithm. Although recently the so called **Sinkhorn method** has taken the attention of many researches since it has a high level of parallelism making it GPU[4] friendly.

Consider the Lagrangian for the strictly convex program (4.29),

$$\Lambda(\boldsymbol{\gamma}; \mathbf{f}, \mathbf{g}) = \langle \mathbf{C}, \boldsymbol{\gamma} \rangle - \epsilon H(\boldsymbol{\gamma}) - \langle \mathbf{f}, \boldsymbol{\gamma} \mathbb{1}_m - \mathbf{a} \rangle - \langle \mathbf{g}, \boldsymbol{\gamma}^\top \mathbb{1}_n - \mathbf{b} \rangle \tag{4.33}$$

where $\mathbf{f}$ and $\mathbf{g}$ are the Lagrange multipliers of the constraints (4.4), the necessary condition for optimality should be satisfied,

$$\frac{\delta \Lambda(\boldsymbol{\gamma}; \mathbf{f}, \mathbf{g})}{\delta (\boldsymbol{\gamma})_{i,j}} = (\mathbf{C})_{i,j} + \epsilon \log(\boldsymbol{\gamma}_{i,j}) - (\mathbf{f})_i - (\mathbf{g})_j = 0, \tag{4.34}$$

which results for each entry of $\boldsymbol{\gamma}$,

$$(\boldsymbol{\gamma})_{i,j} = \exp \left( \frac{(\mathbf{f})_i + (\mathbf{g})_j - (\mathbf{C})_{i,j}}{\epsilon} \right) = e^{(\mathbf{f})_i/\epsilon} e^{-(\mathbf{C})_{i,j}/\epsilon} e^{(\mathbf{g})_j/\epsilon} \tag{4.35}$$

Let $\mathbf{u} \in \mathbb{R}^n$ and $\mathbf{v} \in \mathbb{R}^m$ be two real vectors, and $\mathbf{K} \in \mathbf{M}^{nm}$ a $n \times m$ matrix, for $\forall 1 \leq i \leq n$ and $\forall 1 \leq j \leq m$ we set,

$$(\mathbf{u})_i = \exp\left( (\mathbf{f})_i \right) - \exp\left( \epsilon \right) \tag{4.36}$$

$$(\mathbf{v})_j = \exp\left( (\mathbf{g})_j \right) - \exp\left( \epsilon \right) \tag{4.37}$$

$$(\mathbf{K})_{i,j} = -\exp\left( (\mathbf{C})_{i,j} \right) + \exp\left( \epsilon \right) \tag{4.38}$$

Hence we can write $\boldsymbol{\gamma}$ in a closed form as follows,

$$(\boldsymbol{\gamma}) = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \tag{4.39}$$

and the constraints,

$$\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \mathbb{1}_m = \mathbf{a} \tag{4.40}$$

$$\mathbf{v} \odot \left( \mathbf{K}^\top \mathbf{u} \right) = \text{diag}(\mathbf{v}) \mathbf{K}^\top \text{diag}(\mathbf{u}) \mathbb{1}_n = \mathbf{b}, \tag{4.41}$$

where $\odot$ is the product entry by entry multiplication between two vectors. Also note that $\text{diag}(\mathbf{u}) \mathbb{1}_n = \mathbf{u}$ and $\text{diag}(\mathbf{v}) \mathbb{1}_m = \mathbf{v}$. The equations (4.40) and (4.41) are known as the *matrix scaling problem.*.

The above equations are the entry-wise product between two vectors, therefore if $\mathbf{u}$ and $\mathbf{v}$ are solutions of these equations they should satisfy,

$$\mathbf{u} = \mathbf{a} \oslash (\mathbf{K}\mathbf{v}) \tag{4.42}$$

$$\mathbf{v} = \mathbf{b} \oslash (\mathbf{K}^\top \mathbf{u}), \tag{4.43}$$

---

[4]Graphics Processing Unit.

where $\oslash$ is the division entry-wise. Sinkhorn algorithm exploits the form of the equations in order to find a solution for this condition. The algorithm is really simple, at step $k$ we compute:

$$\mathbf{u}^{k+1} := \mathbf{a} \oslash (\mathbf{K}\mathbf{v}^{k}) \tag{4.44}$$

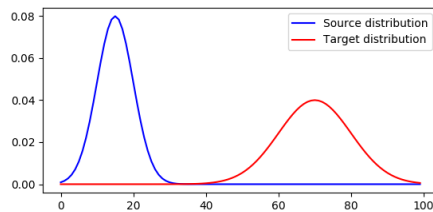$$\mathbf{v}^{k+1} := \mathbf{b} \oslash (\mathbf{K}\mathbf{u}^{k+1}), \tag{4.45}$$

starting from $\mathbf{v}^0 = \mathbb{1}_m$, until reach convergence under some norm. Although the algorithm seems simple, the proof for its convergence it is not a trivial task, and it makes use of the Hilbert's Projective metric and the fact that a positive matrix is a strict contraction on the cone of positive vectors, for a detailed proof and details of the order of convergence we refer [13]. Moreover, the simplicity of the algorithm allows to make the computations in parallel since we have the division entry-wise of two vectors.

## Examples.

The following examples were generated using the library of Optimal Transport [28] developed in Python by Rémi Flamary and Nicolas Courty. The numerical examples takes as cost function $c(x,y) = \frac{1}{2}|x-y|^2$. The source is given by a normal distribution $\mu \sim \mathcal{N}(15,5)$, with mean $m_\mu = 15$ and standard deviation $\sigma_\mu = 5$, and the target is given by a normal distribution $\nu \sim \mathcal{N}(70,10)$, with mean $m_\nu = 70$ and $\sigma_\nu = 10$. Note that the figure 4.3a resembles a linear map, and the figure is just a blurred version of this map.

Figure 4.2: Numerical examples, $c(x,y) = \frac{1}{2}|x-y|^2$.

(a) Two Normal distributions.

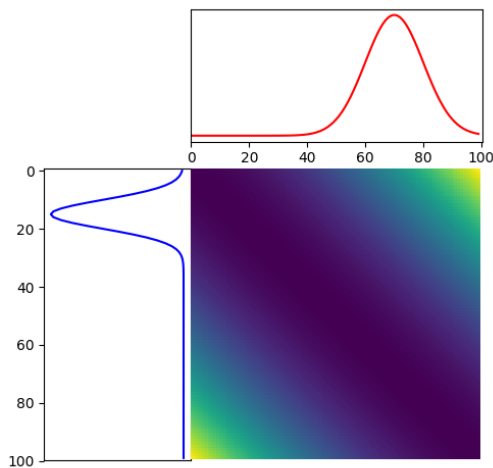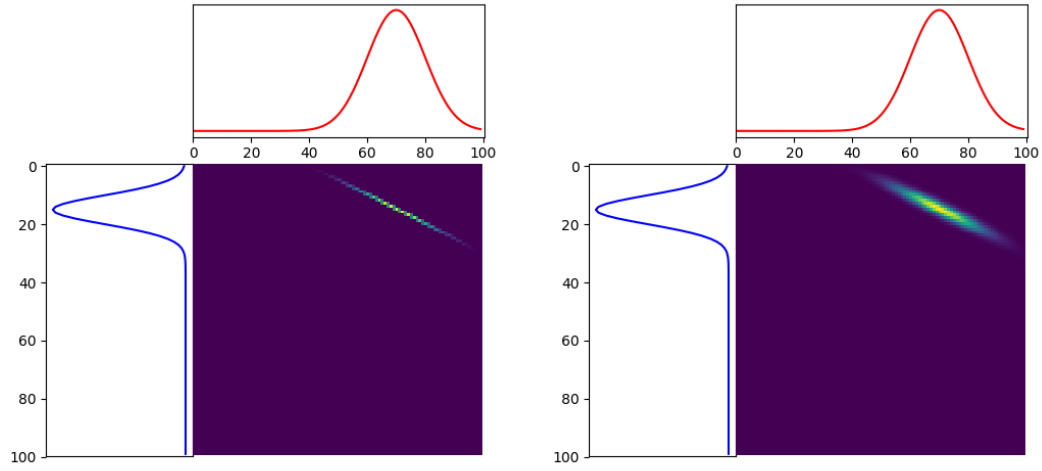(b) Cost matrix, the darker the color the smaller the distance between two points.

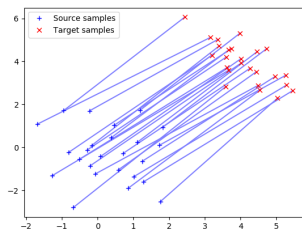Figure 4.3: Optimal couplings for two Gaussians.

(a) Optimal transport plan, found using Network (b) Optimal regularized transport, found using
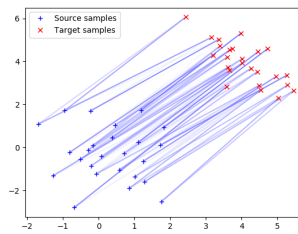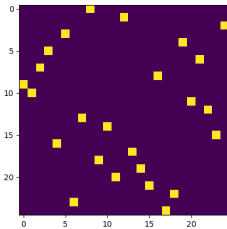simplex algorithm.                                Sinkhorn Algorithm.



We present in figure 4.4 the transportation problem for two data sets $X$ and $Y$ in $\mathbb{R}^2$. Using also the squared euclidean norm for $\mathbb{R}^2$ as cost function, we compute a solution for the linear program ($\epsilon = 0$) and the solutions for the regularized problems with $\epsilon = 1 \times 10^{-3}$ and $\epsilon = 1 \times 10^{-1}$.
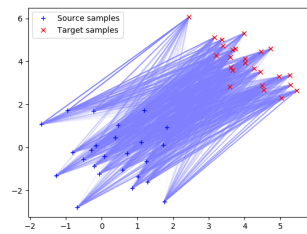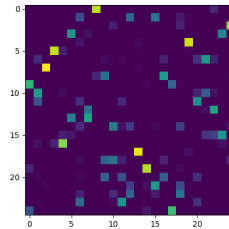
Figure 4.4: Transportation for 2D data sets

(a) Solved using Network Simplex (b) Solved using Sinkhorn method, (c) Solved using Sinkhorn method,
method, $\epsilon = 0$                          $\epsilon = 1 \times 10^{-3}$                          $\epsilon = 1 \times 10^{-1}$



(d) Solution in matrix form,          (e) Solution in matrix form,          (f) Solution in matrix form,
$\epsilon = 0$                              $\epsilon = 1 \times 10^{-3}$                              $\epsilon = 1 \times 10^{-1}$

# 5

# Applications

The applications of the optimal transport are many and we can find them in many discipline of sciences. Optimal couplings are used for color transferring and image segmentation, [26]. In [9] is presented a way to generate high-quality blue noise point distributions of arbitrary density functions through optimal transports. In [27], Wasserstein's Barycenters are used as different way to "average" images. In the field of machine learning, the domain adaptation in unsupervised learning can be done with transportation plans [7]. Yann Brenier, and Jean-David Benamou presented a way to compute a dynamical version of Monge-Kantorovich mass transfer problem using fluid mechanics [2]. In the field of mathematics optimal transports are studied in Riemannian manifolds [34], [11]. In physics, there are many problems where optimal couplings have found a place, for example there exists a relation with Maxwellian distributions and the Monge-Kantorovich problem [31]. In [19] is presented a connection between the Kantorovich's problem and Schrödinger's problem, just to mention some examples among many others.

We present the Wasserstein's distances as an application in statistics since they can be used as a Statistical Distance between two probability distributions. As a second example we present an application of the Brenier's theorem for data assimilation of a dynamic having a linear model and a sensor.

## Wasserstein's Distances as Statistical Distance.

Optimal couplings can be used to metrize the space of probability measures. Since we can induce a distance using couplings, they are useful to use as a measure of the difference between two probability distributions, since we use already a distance of a Polish space to define it, they inherit some of the "geometrical" properties of the space where they are defined.

**Definition 5.1.** *Let $(X, d)$ be a Polish space, and let $p \in [1, \infty)$. For any two probability measures $\mu$, $\nu$ on $X$, the Wasserstein distance $W_p : \mathcal{P}(X) \times \mathcal{P}(X) \to \mathbb{R}$ of order $p$ between $\mu$ and $\nu$ is defined by,*

$$W_p(\mu, \nu) = \left( \inf_{\gamma \in \Pi(\mu, \nu)} \int_{X \times X} d(x, y)^p \mathrm{d}\gamma(x, y) \right)^{\frac{1}{p}} \tag{5.1}$$

*We call a Wasserstein space of order $p$ to the set,*

$$\mathbb{W}_p(X) = \left\{ \mu \in \mathcal{P}(X) : \quad \int_X d(x_0, x)\mu(\mathrm{dx}) < +\infty \right\} \tag{5.2}$$

*where $x_0 \in X$ is arbitrary. This spaces does not depend on the choice of the point $x_0$.*

In words, the Wasserstein space is the space of probability measures with finite moment of order $p$. The case $p = 1$ is also known as the Kantorovich distance. Since $\mu$ and $\nu$ are

defined on the same space, and $d(x, y) = d(y, x)$, we see that $W_p(\mu, \nu) = W_p(\nu, \mu)$. Assume $W_p(\mu, \nu) = 0$; then there exists a transport plan which is entirely concentrated on the diagonal $x = y$ in $X \times X$, therefore $\nu = \mathrm{id}_{\#}\mu$. It is possible to prove the triangle inequality, that is given $\mu_1, \mu_2$ and $\mu_3$, then $W_p(\mu_1, \mu_3) \leq W_p(\mu_1, \mu_2) + W_p(\mu_2, \mu_3)$. We proceed as shown in [34] making use of the Gluing lemma and Minkowski inequality,

$$W_p(\mu_1, \mu_3) = \left( \int_{X \times X} d(x, z)^p \gamma_{1,3}(\mathrm{dxdz}) \right)^{\frac{1}{p}} \tag{5.3}$$

$$\leq \left( \int_{X \times X} (d(x, y) + d(y, z))^p \, \gamma_{1,2,3}(\mathrm{dxdydz}) \right)^{\frac{1}{p}} \tag{5.4}$$

$$\leq \left( \int_{X \times X} d(x, y)^p \gamma_{1,2}(\mathrm{dx|dy})\mu_2(\mathrm{dy}) \right)^{\frac{1}{p}} + \left( \int_{X \times X} d(y, z)^p \gamma_{2,3}(\mathrm{dy|dz})\mu_3(\mathrm{dz}) \right)^{\frac{1}{p}} \tag{5.5}$$

$$\leq \left( \int_{X \times X} d(x, y)^p \gamma_{1,2}(\mathrm{dxdy}) \right)^{\frac{1}{p}} + \left( \int_{X \times X} d(x, z)^p \gamma_{2,3}(\mathrm{dx|dz}) \right)^{\frac{1}{p}} \tag{5.6}$$

$$= W_p(\mu_1, \mu_2) + W_p(\mu_2, \mu_3) \tag{5.7}$$

We can find a more structured proof for the triangle inequality in [29] in the section of Wasserstein Spaces. Note that the inequality,

$$d(x, y) \leq 2^{p-1} \left( d(x, x_0)^p + d(x_0, y)^p \right) \tag{5.8}$$

shows that $W_p$ is finite on $\mathbb{W}_p$, since $W_p(\mu, \nu) \leq 2^{p-1} \left( \int_X d(x, x_0)^p \mathrm{d}\mu(x) + \int_X d(x, x_0)^p \mathrm{d}\nu(x) \right) < \infty$.

In the following we just state some interesting topological results in $\mathbb{W}_p$ spaces that have implications in several applications, the reader can find the respective proofs in [34],

**Definition 5.2** (Weak Convergence in $\mathbb{W}_p$)**.** *Let $(X, d)$ be a Polish space, and $p \in [1, \infty)$. Let $(\mu_k)_{k \in \mathbb{N}}$ be a sequence of probability measures in $\mathbb{W}_p(X)$ and let $\mu$ be another element of $\mathbb{W}_p(X)$. Then $(\mu_k)_{k \in \mathbb{N}}$ is weakly convergent in $\mathbb{W}_p(X)$ if any one of the following properties is satisfied for any $x_0 \in X$:*

- *$\mu_k \to \mu$ and $\int_X d(x_0, x)^p \mathrm{d}\mu_k(x) \to \int_X d(x_0, x)^p \mathrm{d}\mu(x)$;*

- *$\mu_k \to \mu$ and $\limsup_{k \to \infty} \int_X d(x_0, x)^p \mathrm{d}\mu_k(x) \leq \int_X d(x_0, x)^p \mathrm{d}\mu(x)$*

- *$\mu_k \to \mu$ and $\lim_{R \to \infty} \limsup_{k \to \infty} \int_{d(x_0, x) \geq R} d(x_0, x)^p \mathrm{d}\mu_k = 0$*

- *For all continuous functions $\phi$ with $|\phi| \leq C(1 + d(x_0, x)^p)$, $C \in \mathbb{R}$, one has*

$$\int \phi(x) \mathrm{d}\mu_k(x) \to \int \phi(x) \mathrm{d}\mu(x) \tag{5.9}$$

**Theorem 5.1.** *Let $(X, d)$ be a Polish space, and $p \in [1, \infty)$; then the Wasserstein distance $W_p$ metrizes the weak convergence in $\mathbb{W}_p$. That is, if $(\mu_k)_{k \in \mathbb{N}}$ is a sequence of measures in $\mathbb{W}_p$ and $\mu$ is another measure in $\mathcal{P}(X)$, then, the statement $\mu_k$ converges weakly in $\mathbb{W}_(X)$ to $\mu$ and $W_{\mu_k, \mu} \to 0$ are equivalent.*

If $(X, d)$ is a Polish space, and $p \in [1, \infty)$, the $W_p$ is continuous on $\mathcal{P}(X)$. Explicitly if $\mu_k$ (respect to $\nu_k$) converges to $\mu$ weakly (respect to $\nu$ ) in $\mathbb{W}_p$ as $k \to \infty$ then $W_p(\mu_k, \nu_k) \to W_p(\mu, \nu)$. The Wasserstein distance is lower semicontinuous on $\mathcal{P}(X)$. If $\tilde{d}$ is a bounded distance inducing the same topology[1] as $d$, then the convergence in Wasserstein sense for the distance $\tilde{d}$ is equivalent to the usual weak convergence of probability measures in $\mathcal{P}(X)$. Any Cauchy Sequence in $\mathbb{W}_p$ is tight.

The definition of Wasserstein distances makes them convenient to use in problems related with optimal transport, specially many problems coming from partial differential equations.

To conclude this section we introduce the Wasserstein Barycenter problem.

---

[1]For example $\tilde{d} = \frac{d}{1+d}$. Check [23] for the proof that $\tilde{d}$ is bounded and induces the same topology.

**Problem 6.** *Let $X$ a Polish space, and let $\mu_1, \mu_2, \dots, \mu_n$ a set of probability measures defined on $X$. Find $\tilde{\mu} \in \mathbb{W}_p$, such that*

$$\sum_{i=1}^{n} W_p(\tilde{\mu}, \mu_i) \leq \sum_{i=1}^{n} W_p(\mu, \mu_i), \quad \forall \mu \in \mathbb{W}_p(X). \tag{5.10}$$

*That is find the measure $\mu$ that is closer in average to all the given measures in the sense of $\mathbb{W}_p$ metric.*

## Track of a Dynamic (Kalman Filter).

An interesting application of optimal transport is the design of filters for data taken from noisy environments. For the sake of simplicity consider an dynamical system,

$$\mathbf{x}_{n+1} = \mathbf{A}\mathbf{x}_n + \mathbf{B}u_n \tag{5.11}$$

where $\mathbf{x}_n \in \mathbb{R}^n$ is a time series vector constrained to the dynamic given by the constant matrices $\mathbf{A} \in \mathbf{M}^{n \times n}$, $\mathbf{B} \in \mathbf{M}^{n \times m}$ and a control variable $\mathbf{u}_n \in \mathbb{R}^n$.

The above equation is just an idealization of the reality, usually the models ignore small perturbations in order to have a simple and useful model that describes the reality without getting too far from it, as Einsteins says: "Everything should be as simple as possible but not simpler". In practice, many models satisfy this condition, therefore we can assume that the original state $\mathbf{y}_n$ is not the ideal state but behaves really similar. It is reasonable to consider $\mathbf{y}$ as the ideal state plus some noise[2]. For the sake of simplicity, we find convenient to consider it as a random variable normally distributed $\mathbf{y}_n \sim \mathcal{N}(\mathbf{x}_n, \Sigma_x)$ with expected value $\mathbf{x}_n$ and a covariance matrix $\Sigma_x$.

Imagine we are trying to read the data from the original state $\mathbf{y}_n$ using a sensor. In practice, the environment is noisy and sensors are not perfect. This conditions get reflected as noise in the acquisition of data. We can consider the data obtained from the sensor as a random variable normally distributed $\mathbf{z}_n$ with expected value $\mathbf{y}_n$, without loss of generality this is equivalent to say $\mathbf{z}_n = \mathbf{y}_n + \boldsymbol{\sigma}$, where $\boldsymbol{\sigma}$ is a random variable normally distributed with zero mean and covariance matrix $\Sigma_z$.

We remark that are able to read the sensor, so we would like to find a joint probability measure that given a reading of the sensor allows to predict the next state of the dynamic as close as possible.

We need to choose a measure for the difference between the prediction and reality. For simplicity we choose a standard deviation. This is nothing less that an optimal transport formulation of the problem, with the marginals given by normal distributions of our two random variables and a cost function $c(x, y) = |x - y|^2$. Moreover, from the Brenier's theorem we know that this joint probability measure is a deterministic coupling given by a map $T$. Hence we reduce it to an optimal map problem,

$$\min_{T} \mathbb{E}\left[\left\|\mathbf{y}_{n+1} - T(\mathbf{z}_n)\right\|^2\right] \tag{5.12}$$

$$\text{subject to } \nu = T_{\#}\mu \tag{5.13}$$

where $\mu$ and $\nu$ are the normal distributions corresponding to the random variables $\mathbf{z}_n$ and $\mathbf{y}_{n+1}$ respectively. Both variables are normally distributed, and the cost function is the squared euclidean distance, therefore the $T$ is an affine map given by,

$$T(\mathbf{z}_n) = \mathbf{x}_n + \Sigma \mathbf{A}^\top \Sigma_z^{-1} (\mathbf{z}_n - \mathbf{A}\mathbf{x}_n) \tag{5.14}$$

$$\text{where, } \Sigma = \left(\mathbf{A}^\top \Sigma_z^{-1} \mathbf{A} + \Sigma_x^{-1}\right)^{-1}$$

The above equation is the map generating a deterministic coupling between a normally distributed random variable under a linear map and an arbitrary random variable that is also normally distributed. Please check [22] and [24] for a detailed derivation of the equations.

---

[2]Please note that this is an assumption, in many applications this models in good way the phenomenon. Although there is a huge set of problems that this approach result flawed or insufficient.

Please note that computing this map is not possible, since we need the ideal state and this is something that we are not able to obtain, even if our sensor is perfect. Although, we can do an approximation. We proceed as Kalman did in the development of the filter bearing his name.

We start with a guess of the initial state, $\mathbf{r}_0$ and we read our first sample from the sensor $\mathbf{z}_0$, we use our initial guess to do an approximation of the ideal state, that is $\mathbf{x}_0 \approx \mathbf{A}\mathbf{r}_0 + \mathbf{B}\mathbf{u}_0 = \tilde{\mathbf{x}}_0$. And we save $\mathbf{r}_1 = T(z_0)$ since we consider that this is the best approximation we can get from the dynamic under standard deviation measure. We continue in this way, taking a sample from the sensor and using the data to get good approximations $\tilde{\mathbf{x}}_{n+1} = \mathbf{A}\mathbf{r}_n + \mathbf{B}\mathbf{u}_n$ of the reality,

$$T(\mathbf{z}_n) = \tilde{\mathbf{x}}_n + \Sigma \mathbf{A}^\top \Sigma_z^{-1} (\mathbf{z}_n - \mathbf{A}\tilde{\mathbf{x}}_n) \tag{5.15}$$

$$\mathbf{r}_{n+1} = T(\mathbf{z}_n) \tag{5.16}$$

If our assumptions are consistent with the reality $\mathbf{r}_n$ is the best approximation we can do sampling $\mathbf{z}_{n-1}$ from the sensor at step $n-1$, and our filtered data should not be far from the reality.

As example consider the circuit shown in figure 5.1a whose dynamic is given by,

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1/(RC) \\ -R/L & -R/L \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ R/L \end{bmatrix} v_{in} = A_c x + B v_{in} \tag{5.17}$$
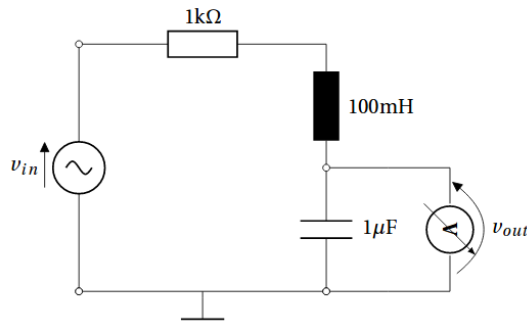
In this hypothetical situation consider that the sensor is able to take samples at a frequency of 10kHz, that is a sampling period $T =$ 1e-4 seconds and the source $v_{in} = \sin(2\pi f t)$, where $f = 100$Hz. The values for the electrical components are $R = 1$kΩ for the resistance, $L = 100$m$H$ for the inductance, and $C = $ 1e-6F for the capacitance. The discrete version of this dynamic is given by,

$$\mathbf{x}_{n+1} = \begin{bmatrix} 0,9056542 & 0,0009066 \\ -0,9065617 & -0,0009075 \end{bmatrix} \mathbf{x}_n + \begin{bmatrix} 0,0943458 \\ 0,9065617 \end{bmatrix} u_n \tag{5.18}$$
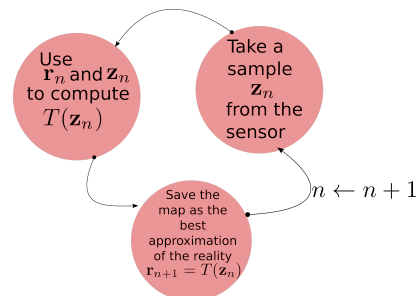
The above equation was obtained using the equations,

$$\mathbf{A} = e^{A_c T}, \quad \mathbf{B} = A_c^{-1}(\mathbf{A} - \mathbf{I})B, \tag{5.19}$$

and $u_n$ is just the sample of the input $v_i n$ at time $nT$ corresponding to the $n$ sample, i.e. $u_n = \sin(2\pi f n T)$.
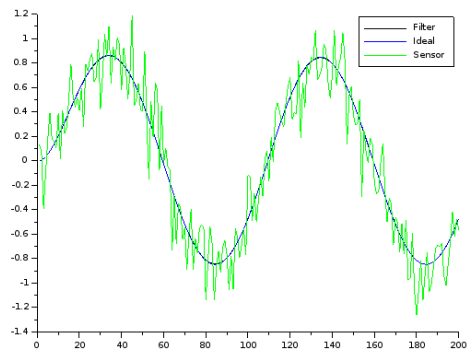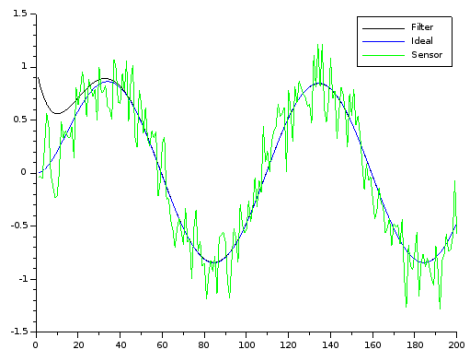


(a) RLC Circuit.

(b) Procedure to track a dynamic.

The figure 5.2 shows the results obtained for the described dynamical system. Note that we have two simulations, one for an initial guess equal to the ideal state and one for an initial guess totally different to the ideal state. We can appreciate that even if we started from a bad guess for the initial state, the coupling works properly.

Figure 5.2: Filter with Optimal Transport coupling. $\boldsymbol{\Sigma}_z = 0.2$, $\boldsymbol{\Sigma}_\chi = 0.01$.

(a) Initial guess $\mathbf{r}_0 = \mathbf{x}_0$                                  (b) Initial guess $\mathbf{r}_0 \neq \mathbf{x}_0$

# Bibliography

[1] Viorel Barbu and Teodor Precupanu. *Convexity and Optimization in Banach Spaces*. Springer Monographs in Mathematics. Springer, 4 edition, 2012.

[2] Jean-David Benamou and Yann Brenier. A computational fluid mechanics solution to the monge-kantorovich mass transfer problem. *Numerische Mathematik*, 84:375–393, 2000.

[3] Patrick Billingsley. *Convergence of Probability Measures*. Probability and Statistics. John Wiley and Sons., 2 edition, 1999.

[4] Guy Bouchitté and Giuseppe Buttazzo. Characterization of optimal shapes and masses through monge-kantorovich equation. *J. Eur. Math. Soc.*, 3:139–168, 2001.

[5] Francis Clarke. *Functional Analysis, Calculus of Variations and Optimal Control*, volume 264 of *Graduate Texts in Mathematics*. Springer, 2013.

[6] Donald L. Cohn. *Measure Theory*. Birkhäuser, 2 edition, 2006.

[7] N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy. Optimal transport for domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(9):1853–1865, Sept 2017. ISSN 0162-8828. doi: 10.1109/TPAMI.2016.2615921.

[8] Gabriel Peyré; Marco Cuturi. Computational optimal transport, 2018. URL https://optimaltransport.github.io/.

[9] F. de Goes, K. Breeden, V. Ostromoukhov, and M. Desbrun. Blue noise through optimal transport. *ACM Trans. Graph. (SIGGRAPH Asia)*, 31, 2012.

[10] Ivar Ekeland and Roger Témam. *Convex Analysis and Variational Problems*, volume 28 of *Classics in Applied Mathematics*. S.I.A.M., 1999.

[11] Alessio Figalli and Cédric Villani. Nonlinear pde's and applications. *Nonlinear PDE's and applications*, 2028:171–217, 2011.

[12] Irene Fonseca and Giovanni Leoni. *Modern methods in the calculus of variations: Lp spaces*. Springer, 2007. ISBN 978-0-387-69006-3.

[13] Joel Franklin and Jens Lorenz. On the scaling of multidimensional matrices. *LINEAR ALGEBRA AND ITS APPLICATIONS*, 114-115:717–735, 1989.

[14] Morris W. Hirsch. *Differential Topology*. springer, 1976. ISBN 3-540-90148-5.

[15] A. Hulanicki, P. Wojtaszczyk, and W. Żelazko, editors. *Selected Papers of Antoni Zygmund*, volume 3 of *Mathematics and Its Applications (East European Series)*. KLUWER ACADEMIC PUBLISHERS, 1989.

[16] L. Kantorovich. On the translocation of masses. *Dokl. Acad. Nauk. USSR*, 37(7-8):227–229, 1942.

[17] Dimitrios Zissopoulos Konstatinos Paparrizos, Nikolaos Samaras. Linear programming: Klee–minty examples. In Pardalos P. Floudas C., editor, *Encyclopedia of Optimization*. Springer, Boston, 2 edition, 2008.

[18] Erwin Kreyzig. *NTRODUCTORY FUNCTIONAL ANALYSIS WITH APPLICATIONS*. Wiley's classics Library. JOHN WILEY & SONS, 1978.

[19] Christian Léonard. From the schrödinger problem to the monge–kantorovich problem. *Functional Analysis*, 262:1879–1920, 2012.

[20] David G. Luenberger and Yinyu Ye. *Linear and Nonlinear Programming*. Springer, 2007.

[21] Milan Merkle. Topics in weak convergence of probability measures. *Zb. radova Mat. Inst. Beograd*, 9(17):235–274, 2000.

[22] El Moselhy, Tarek A., and Youssef M. Marzouk. Bayesian inference with optimal maps. *ournal of Computational Physics 231*, 23:7815–7850, 2012.

[23] Munkres and J.R. *Introduction to Topology*. Series in Topology. Prentice-Hall, 2nd. edition, 2000. ISBN 9780131784499.

[24] Dean S. Oliver. Minimization for conditional simulation: Relationship to optimal transport. *Journal of Computational Physics*, 265:1–15, 2014.

[25] James B. Orlin. A polynomial time primal network simplex algorithm for minimum cost flows. *Mathematical Programming*, 78:109–129, 1997.

[26] Nicolas Papadakis. *Optimal Transport for Image Processing*. PhD thesis, Signal and Image Processing. Université de Bordeaux, l'Institut de Mathématiques de Bordeaux (UMR 5251), 12 2015.

[27] Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernot. Wasserstein barycenter and its application to texture mixing. *Scale Space and Variational Methods in Computer Vision*, 6667:435–446, 2012.

[28] Flamary Rémi and Nicolas Courty. Pot python optimal transport library, 2017. URL https://github.com/rflamary/POT.

[29] Filippo Santambrogio. *Optimal Transport for Applied Mathematicians.*, volume 87 of *Progress in Nonlinear Differential Equations and Their Applications*. Springer, 2015.

[30] Karl-Georg Steffens. *The History of Approximation Theory. From Euler to Bernstein*. Birkhäuser, 2006.

[31] Hitoshi Tanaka. Probabilistic treatment of the boltzmann equation of maxwellian molecules. *Wahrscheinlichkeitstheorie verw Gebiete*, 46:67–105, 1978.

[32] Robert E. Tarjan. Dynamic trees as search trees via euler tours, applied to the network simplex algorithm. *Mathematical Programming*, 78:169–177, 1997.

[33] Myron Tribus. Information theory and thermodynamics. *Heat Transfer, Thermodynamics and Education: Boelter Anniversary*, page 354, 1964.

[34] Cédric Villani. *Optimal Transport: Old and New*, volume 338 of Grundlehren der mathematischen Wissenschaften. Springer Science and Business Media, 2008.