

# The Optimal Transport Problem

Master Thesis

Oscar Ramirez



Math  
Meets

Y  
Younes Bounhar



# 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

Preface...

*Oscar Ramirez*  
*Barcelona, September 2018*



# Contents

<b>1 Preliminaries.</b>	<b>3</b>
1.1 Definitions and important theorems to remember. . . . .	3
1.1.1 Topology. . . . .	3
1.1.2 Functional Analysis. . . . .	5
1.1.3 Measure Theory . . . . .	7
<b>2 Basics in Convex Analysis.</b>	<b>11</b>
2.1 Convexity in $\mathbb{R}^n$ . . . . .	12
<b>3 Linear Programming</b>	<b>15</b>
3.1 Simplex Method. . . . .	16
3.2 Duality. . . . .	20
3.3 Complementary Slackness. . . . .	21
3.3.1 The Dual simplex Method. . . . .	21
3.3.2 The Primal-Dual Simplex Method. . . . .	21
3.4 Min-Max Theorems. . . . .	21
3.5 Interior Methods. . . . .	21
<b>4 Optimal Transport Theory</b>	<b>23</b>
4.0.1 Existence of a minimizer for Kantorovich's Problem. . . . .	28
4.1 Kantorovich formulation as relaxation . . . . .	29
4.2 Cyclical Monotonicity and Duality. . . . .	31
4.2.1 Duality. . . . .	31
4.3 Properties of Optimal plans. . . . .	34
<b>5 Computation of an Optimal Transport</b>	<b>35</b>
5.1 Linear Programming Formulation. . . . .	35
5.1.1 Simplex Method Algorithm and Duality. . . . .	35
5.1.2 The simplex method is Not polynomial time. . . . .	36
5.1.3 Sinkhorn-Knopp Algorithm. . . . .	36
5.2 Continuous Formulation. . . . .	36
5.2.1 Beckman Problem and Optimal Transport. . . . .	36
5.2.2 Proximal Splitting Algorithms. . . . .	36
<b>6 Applications</b>	<b>37</b>
6.1 Domain Adaptation. . . . .	37
6.2 Isoperimetric Inequality. . . . .	37
6.3 Wasserstein's Distance . . . . .	37
<b>Bibliography</b>	<b>39</b>





# Notation Table.

$\emptyset$	Empty set
$\mathbb{R}$	Real numbers field.
$\overline{\mathbb{R}}$	$\mathbb{R} \cup \{+\infty\}$ . 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 $[0, \infty]$
$\delta_x$	The Dirac mass at point $x$ .
$\mathbb{R}^d$	The $d$ -dimensional Euclidean space.
$\text{id}$	Identity map.
$\mathcal{M}(X)$	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 is absolutely continuous with respect to $\nu$ .
$\mathbb{1}_\Omega$	Indicator function of a set $\Omega$ . If $x \in \Omega$ then $\mathbb{1}_\Omega(x) = 1$ . If $x \in \Omega^c$ , we have $\mathbb{1}_\Omega(x) = 0$ .
$\mu \llcorner A$	A measure $\mu$ restricted to a set $A$ .
$\omega_d$	The Measure of the unite ball in $\mathbb{R}^d$ .
$\wedge$	The min operator, that is $a \wedge b := \min\{a, b\}$ .
$\vee$	The max operator, that is $a \vee b := \max\{a, b\}$ .
$DT(x)$	Jacobian matrix of a map $T(x)$ .
$T_\# \mu$	The image measure (or pushforward measure) of $\mu$ through the map $T$ .
$f _\Omega$	The restriction of a function $f$ to a set $\Omega$ .
$\Pi(\mu, \nu)$	The set of transport plans from $\mu$ to $\nu$ .
$\frac{\delta F}{\delta \rho}$	First variation of $F : \mathcal{P}(X) \rightarrow \mathbb{R}$ , that is $\left. \frac{d}{d\epsilon} F(\rho + \epsilon \chi) \right _{\epsilon=0} = \int \frac{\delta F}{\delta \rho} d\chi$
$W_p$	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.
$M^\top$	Transpose of a matrix $M$ .
i.i.d.	Independent and identical probability distributions.
l.s.c.	Lower semicontinuous.
$\mathcal{L}^p$	Lebesgue measure on $\mathbb{R}^p$
$\mathcal{H} \llcorner A$	Hausdorff measure applied to some set $A \subset \mathbb{R}^d$ .
$B(x, \epsilon)$	Open ball with radius $\epsilon$ centered at $x$ .
$\overline{B}(x, \epsilon)$	Closed ball with radius $\epsilon$ centered at $x$ .



# 1

## Preliminaries.

We start this chapter reminding the basic definitions and theorems in topology and measure theory, since they are needed to have a suitable framework to discuss the optimal transport problem and its applications.

### Definitions and important theorems to remember.

#### Topology.

We start with the definition of topology that is needed to introduce a notion of continuity. We refer to [3] for more details in topology.

**Definition 1.1** (Topology). *A topology on a set  $X$  is a collection  $\mathcal{T}$  of subsets of  $X$  having the following properties*

- *The space  $X$  itself and  $\emptyset$  are in  $\mathcal{T}$ .*
- *The union of the elements of any sub-collection of  $\mathcal{T}$  is in  $\mathcal{T}$ .*
- *The intersection of the elements of any finite sub-collection of  $\mathcal{T}$  is in  $\mathcal{T}$ .*

A pair  $(X, \mathcal{T})$  is called a topological space. The elements of  $\mathcal{T}$  are called open sets. The complements of the sets of  $\mathcal{T}$  are called closed sets. 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  $\text{int}(A)$  and  $A^\circ$  for the interior of a set  $A$ . In the same way, for the closure we use the notation  $\text{clo}(A)$  or  $\bar{A}$ . An equivalent way to define the same ideas is give by the following,

$$\begin{aligned}\text{int}(A) &= \bigcup_{B \text{ is open.}} B \\ \text{clo}(A) &= \bigcap_{B \text{ is closed.}} B\end{aligned}$$

We remark that a set is open if and only if  $A = A^\circ$ , and a set is closed if and only if  $A = \bar{A}$ . We call a neighborhood of  $x$  an element of  $\mathcal{T}$  containing  $x$ .

**Definition 1.2** (Topological Basis.). *Give a set  $X$  endowed with a topology  $\mathcal{T}$ . We call a basis for  $\mathcal{T}$  is a collection  $\mathcal{B}$  of subsets of  $X$  (called basis elements), such that,*

1. *For each  $x \in X$ , there is at least one basis element  $B$  containing  $x$ .*
2. *If  $x$  belongs to the intersection of two basis elements  $B_1$  and  $B_2$ , then there is a basis element  $B_3$  containing  $x$  such that  $B_3 \in B_1 \cap B_2$ .*

**Definition 1.3** (Dense set). *A subset  $D$  of a topological space  $X$  is dense in  $X$  if for any point  $x$  in  $X$ , any neighborhood of  $x$  contains at least one point from  $D$ .*

**Definition 1.4** (Separable space). *A topological space is called separable if it contains a countable, dense subset.*

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. A mathematician Felix Hausdorff suggested to add the following condition:

**Definition 1.5** (Hausdorff space). *A topological space  $X$  is called a Hausdorff space if for each pair  $x_1, x_2$  of distinct points of  $X$ , there exist neighborhoods  $U_1$ , and  $U_2$  of  $x_1$  and  $x_2$ , respectively, that are disjoint.*

**Definition 1.6** (Distance).

**Definition 1.7** (Metric Spaces).

**Definition 1.8** (Completeness). *A metric space  $X$  is called complete if every Cauchy-Sequence of points in  $X$  has a limit that is also in  $X$ .*

**Definition 1.9** (Completely metrizable space).

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.

**Definition 1.10** (Polish space). *We call Polish space to any topological space that is separable and completely metrizable.*

**Definition 1.11** (Sequentially compact). *A subset  $K$  of a metric space  $X$  is said to be compact if from any sequence  $x_n$ , we can extract a converging subsequence  $x_{n_k} \rightarrow x \in K$ .*

**Definition 1.12** (Compactness). *A subset  $K$  of a metric space  $X$  is compact if every open cover of  $K$  has a finite subcover.*

**Theorem 1.1.** *A subset of a metric space is compact if and only if it is sequentially compact.*

**Definition 1.13** (Liminf and Limsup). *Let  $X$  be a Hausdorff space. Let  $\mathcal{V}(x_0)$  be a topological basis of  $X$ , such that all  $V \in \mathcal{V}$  contains  $x_0$ . Let  $f : X \rightarrow \overline{\mathbb{R}}$  a functional valued in  $\overline{\mathbb{R}}$ . We define,*

$$\liminf_{x \rightarrow x_0} f(x) = \sup_{V \in \mathcal{V}(x_0)} \inf_{s \in V} f(s)$$

$$\limsup_{x \rightarrow x_0} f(x) = \inf_{V \in \mathcal{V}(x_0)} \sup_{s \in V} f(s)$$

*The above definitions can be expressed in terms of sequences of real numbers. Let  $(x_n)_{n \in \mathbb{N}}$  be a sequence in  $X$ , the above formulation is equivalent to say.*

$$\liminf_{n \in \mathbb{N}} x_n := \lim_{n \rightarrow \infty} \left( \inf_{m \geq n} x_m \right)$$

*Equivalently for lim sup,*

$$\limsup_{n \in \mathbb{N}} x_n := \lim_{n \rightarrow \infty} \left( \sup_{m \geq n} x_m \right)$$

*Please note that the convergence to some point  $x_0$ ,  $(x_n)_{n \in \mathbb{N}} \rightarrow x_0$  is not required in the last definitions.*

## Functional Analysis.

**Definition 1.14** (Linear Space).

**Definition 1.15** (Banach Space).

**Definition 1.16** (Inner product).

**Definition 1.17** (Hilbert Space).

**Definition 1.18** (Continuity).

For a given be a metric space  $X$ . We denote the set of continuous, real-valued functions  $f : X \rightarrow \mathbb{R}$  by  $C(X)$ .

**Theorem 1.2.** *Let  $K \subset X$  a compact subset of a metric space  $X$ . The space  $C(K)$  is complete.*

A natural norm on spaces of continuous functions is the uniform norm (also called infinity norm), which is defined by,

$$\|f\|_{\infty} = \sup_{x \in X} |f(x)|$$

The norm  $\|f\|_{\infty}$  is finite if and only if  $f$  is bounded. And we use  $C_b(X)$  to refer the space of bounded functions on  $X$ .

**Definition 1.19.** *Let  $f$  be a real valued function,  $f : X \rightarrow \mathbb{R}$  on a metric space. The **support of a function**,  $\text{supp } f$  is the closure of the set on which  $f$  is nonzero.*

$$\text{supp } f = \text{clo}(\{x \in X : f(x) \neq 0\})$$

We say that  $f$  has compact support if  $\text{supp } f$  is a compact subset of  $X$ , and denote the space of continuous functions on  $X$  with compact support by  $C_c(X)$ .

The space  $C_c(X)$  is a linear subspace of  $C_b(X)$ , but it does not need to be closed.

**Definition 1.20.** *Suppose that  $X$  is a separable and locally compact metric space. We say that a real valued function  $f$  belongs to  $C_0(X)$  if and only if  $f \in C(X)$ , and for every  $\epsilon > 0$ , there exists a compact set  $K \subset X$  such that  $|f| < \epsilon$  on  $X \setminus K$ .*

**Definition 1.21.** *Let  $\mathcal{F}$  be a family of functions from a metric space  $(X, d)$  to a metric space  $(Y, d)$ . The family  $\mathcal{F}$  is equicontinuous if for every  $x \in X$  and  $\epsilon > 0$  there is  $\delta > 0$  such that  $d(x, y) < \delta$  implies  $d(f(x), f(y)) < \epsilon$  for all  $f \in \mathcal{F}$ .*

**Theorem 1.3.** *An equicontinuous family of functions from a compact metric space to a metric space is uniformly equicontinuous.*

**Theorem 1.4** (Ascoli-Arzelà). *Let  $K$  be a compact metric space. A subset  $M$  of the set of continuous functions  $M \subset C(K)$  is compact if and only if it is closed, bounded and equicontinuous. That is any sequence  $(f_n)_{n \in \mathbb{N}}$  in  $M$  admits a subsequence converging to  $f$  in  $M$ .*

**Definition 1.22** (Proper convex function). *Let  $f : X \rightarrow \overline{\mathbb{R}}$ , a function taking values in the extended real number line. We call it proper convex function if  $\exists x \in X$  such that  $f(x) < \infty$ . And  $\forall x \in X$ ,  $f(x) > -\infty$ .*

**Definition 1.23** (Projection of a Cartesian Product.). *Let  $\text{proj}_x : X \times Y \rightarrow X$  be defined by the equation*

$$\text{proj}_x(x, y) = x;$$

*Equivalently, let  $\text{proj}_y : X \times Y \rightarrow Y$  be defined by,*

$$\text{proj}_y(x, y) = y$$

*These maps  $\text{proj}_x$  and  $\text{proj}_y$  are called the projections of  $X \times Y$  onto  $X$  and  $Y$  respectively.*

We can generalize a definition over a general Cartesian product. Given a set  $X$ , we define a  $J$ -tuple of elements of  $X$  to be a function  $\mathbf{x} : J \rightarrow X$ . If  $\alpha$  is an element of  $J$ , we often denote the value of  $\mathbf{x}$  at  $\alpha$  by  $x_\alpha$  rather than  $\mathbf{x}(\alpha)$ ; we call it the  $\alpha$ -th coordinate of  $\mathbf{x}$ . And we often denote the function  $\mathbf{x}$  by the symbol.

$$(x_\alpha)_{\alpha \in J}$$

Let  $\{A_\alpha\}_{\alpha \in J}$  be a set of indexed family of sets; let  $X = \bigcup_{\alpha \in J} A_\alpha$ . The *Cartesian product* of this indexed family, denoted by

$$\prod_{\alpha \in J} A_\alpha$$

is defined to be the set of all  $J$ -tuples  $(x_\alpha)_{\alpha \in J}$  of elements of  $X$  such that  $x_\alpha \in A_\alpha$  for each  $\alpha \in J$ . That is, it is the set of all functions,

$$\mathbf{x} : J \rightarrow \bigcup_{\alpha \in J} A_\alpha$$

such that  $\mathbf{x}(\alpha) \in A_\alpha$  for each  $\alpha \in J$ .

**Definition 1.24** (Lower Semicontinuity). *On a complete metric space  $X$ , a function  $f : X \rightarrow \overline{\mathbb{R}}$  is said to be lower semi-continuous (l.s.c.) if for every sequence  $(x_n)_{n \in \mathbb{N}}$  converging to  $x \in X$ , we have*

$$f(x) \leq \liminf_{n \in \mathbb{N}} f(x_n)$$

We can see from the above definition that any continuous function is lower-semicontinuous. In other words, lower-semicontinuity is a milder requirement than continuity, although it preserves interesting properties that can be exploited in optimization.

**Proposition 1.1.** *Let  $f : X \rightarrow \overline{\mathbb{R}}$  be a convex and lower-semicontinuous function. Assume that there exists  $x_0 \in X$  such that  $f(x_0) = -\infty$ . Then  $f$  is nowhere finite on  $X$ .*

**Theorem 1.5.** *If  $f_\alpha$  is an arbitrary family of lower semi-continuous functions on  $X$ , then  $f = \sup_\alpha f_\alpha$  is also lower-semicontinuous.*

**Definition 1.25** (Lipschitz condition).

**Theorem 1.6.** *Let  $f : X \rightarrow \overline{\mathbb{R}} \setminus -\infty$  be a function bounded from below. Then  $f$  is l.s.c. if and only if there exists a sequence  $(f_n)_{n \in \mathbb{N}}$  of  $k$ -Lipschitz functions such that for every  $x \in X$ ,  $f_n(x)$  converges increasingly to  $f(x)$ .*

**Definition 1.26.** *Let  $F : X \rightarrow \overline{\mathbb{R}}$  be a given functional bounded from below on a metric space  $X$ . Let  $\mathcal{G}$  be the set of lower semicontinuous functions  $G : X \rightarrow \overline{\mathbb{R}}$ , such that  $G \leq F$ . We call a relaxation the supremum of  $\mathcal{G}$ . This functional does exist since the supremum of an arbitrary family of lower semicontinuous functions is also lower semicontinuous. It is possible to have a representation formula as follows:*

$$\bar{F}(x) = \inf \left\{ \liminf_{n \in \mathbb{N}} F(x_n) : x_n \rightarrow x \right\}. \quad (1.1)$$

As consequence of this definition we see that  $F \geq \bar{F}$  implies  $\inf F \geq \inf \bar{F}$ . Let  $l = \inf F$  then  $F \geq l$ . A constant function is lower semicontinuous. Therefore,  $\bar{F} \geq l$  and  $\inf \bar{F} \geq \inf F$ . Implying that the infimum of both  $F$  and its regularization  $\bar{F}$  coincide, i.e.  $\inf \bar{F} = \inf F$ .

**Theorem 1.7** (Maxima and Minima). *Let  $X$  be a compact metric space and  $f : X \rightarrow \mathbb{R}$  is continuous, real-valued function. Then  $f$  is bounded on  $X$  and attains its maximum and minimum. That is, there are  $x, y$  belonging to  $X$  such that,*

$$f(x) = \inf_{z \in X} f(z) \quad \text{and} \quad f(y) = \sup_{z \in X} f(z)$$

Continuity is a strong requirement. Luckily, we can assure the existence of a minimizer of lower-semicontinuous functionals (or maximizer for upper-semicontinuity). The usual procedure to prove existence of a minimizer is making use of Weierstrass' criterion. We take a minimizing sequence and then we prove that the space in which we are trying to find a minimizer element is compact.

**Theorem 1.8** (Weierstrass' criterion for existence of minimizers). *If  $f : X \rightarrow \overline{\mathbb{R}}$  is lower semi-continuous and  $X$  is compact, then there exists  $\hat{x} \in X$  such  $f(\hat{x}) = \min \{f(x) : x \in X\}$ .*

*Proof.* Define  $l := \inf \{f(x) : x \in X\} \in \overline{\mathbb{R}}$ , notice that  $l = +\infty$  only if  $f$  is identically  $+\infty$ , then this case is trivial since any point minimizes  $f$ . By compactness there exists a minimizing sequence  $x_n$ , that is  $f(x_n) \rightarrow l$ . By compactness we can extract a subsequence converging to some  $\hat{x}$  such that  $\hat{x} \in X$ . By lower-semicontinuity of  $f$ , we have that  $f(\hat{x}) \leq \liminf_n f(x_n) = l$ . Since  $l$  is the infimum  $l \leq f(\hat{x})$ . This proves that  $l = f(\hat{x}) \in \mathbb{R}$ .  $\square$

We can apply the above analysis using a notion of upper-semicontinuity and compactness to find the maximum.

**Definition 1.27** (Topological Dual). *If  $X$  is a normed space, the dual space  $X^* = \mathcal{B}(X, \mathbb{R})$ . Consists of all linear and bounded functionals mapping from  $X$  to  $\mathbb{R}$ .*

**Definition 1.28** (Weak compactness in dual spaces). *A sequence  $x_n$  in a Banach space  $X$  is said to be weakly converging to  $x$ , and we write  $x_n \rightharpoonup x$ , if for every  $\xi \in X^*$ . We have  $\langle \xi, x_n \rangle \rightarrow \langle \xi, x \rangle$ . A sequence  $\xi_n \in X^*$  is said to be weakly-\* converging to  $\xi$ , and we write  $\xi_n \rightharpoonup^* \xi$ , if for every  $x \in X$  we have  $\langle \xi_n, x \rangle \rightarrow \langle \xi, x \rangle$ .*

**Theorem 1.9** (Banach-Alaoglu). *If  $X$  is separable and  $\xi_n$  is a bounded sequence in  $X^*$ , then there exists a subsequence  $\xi_{n_k}$  weakly converging to some  $\xi \in X^*$ .*

The Banach-Alaoglu's theorem is a well known result in functional analysis, an equivalent formulation is saying the closed unit ball in  $X^*$  is weak-\* compact.

## Measure Theory

The optimal transport problem theory is based mostly on Measure Theory. We present some abstract objects and theorems needed to develop in proper way the problem. For better understanding on Measure Theory we refer [1].

**Definition 1.29** (Sigma Algebra). *An algebra of sets  $\mathcal{A}$  is a class of subsets of some fixed set  $X$  (called the space) such that,*

- $X$  and  $\emptyset$  belong to  $\mathcal{A}$ .
- If  $A, B \in \mathcal{A}$ , then  $A \cap B \in \mathcal{A}$ ,  $A \cup B \in \mathcal{A}$ ,  $A \setminus B \in \mathcal{A}$ .

*An algebra of sets is a  $\sigma$ -algebra if for any sequence of sets  $A_n \in \mathcal{A}$  we have  $\mathcal{A} \ni \bigcup_{n \in \mathbb{N}} A_n$ .*

**Definition 1.30** (Measure Space). *A pair  $(X, \mathcal{A})$  consisting of a set  $X$  and a  $\sigma$ -algebra  $\mathcal{A}$  of its subsets is called a measurable space.*

**Definition 1.31** (Borel  $\sigma$ -algebra). *The Borel  $\sigma$ -algebra  $\mathcal{B}(\mathbb{R}^n)$  of  $\mathbb{R}^n$  is the  $\sigma$ -algebra generated by all open sets. The sets in a Borel  $\sigma$ -algebra are called Borel sets. For any set  $E \subset \mathbb{R}^n$ , let  $\mathcal{B}(E)$  denote the class of all sets of the form  $E \cap B$ , where  $B \in \mathcal{B}(\mathbb{R}^n)$ .*

Given a topological space, in this text consider the Borel  $\sigma$ -algebra unless stated otherwise.

**Definition 1.32** (Measure). *A real-valued set function  $\mu : \mathcal{A} \rightarrow \overline{\mathbb{R}}$  on a class of sets  $\mathcal{A}$  is called countably additive if*

$$\mu\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} \mu(A_n)$$

*for all pairwise disjoint sets  $A_n$  in  $\mathcal{A}$  such that  $\mathcal{A} \ni \bigcup_{n=1}^{\infty} A_i$ . A countably additive set function defined on an algebra is called a measure.*

Given a measure space  $X$ , we denote  $\mathcal{M}(X)$  and  $\mathcal{M}_+(X)$  to refer the set of finite measures and positive finite measures on  $X$ , respectively.

**Definition 1.33.** A countably additive measure  $\mu$  on a  $\sigma$ -algebra of subsets of a space  $X$  is called a probability measure if  $\mu \geq 0$  and  $\mu(X) = 1$ . A triple  $(X, \mathcal{A}, \mu)$  is called a measure space if  $\mu$  is a nonnegative measure on a  $\sigma$ -algebra  $\mathcal{A}$  of subset of a set  $X$ . If  $\mu$  is a probability measure, then  $(X, \mathcal{A}, \mu)$  is called a probability space.

**Definition 1.34** (Lebesgue Measure).

**Definition 1.35** (Hausdorff Measure).

**Definition 1.36** (Probability). A countably additive measure  $\mu$  on a  $\sigma$ -algebra of subsets of a space  $X$  is called a probability measure if  $\mu \geq 0$  and  $\mu(X) = 1$ .

A triple  $(X, \mathcal{A}, \mu)$  is called a measure space if  $\mu$  is a nonnegative measure on a  $\sigma$ -algebra  $\mathcal{A}$  of subset of a set  $X$ . If  $\mu$  is a probability measure, then  $(X, \mathcal{A}, \mu)$  is called a probability space.

**Definition 1.37** (Product  $\sigma$ -algebra). Let  $(X_1, \mathcal{A}_1, \mu_1)$  and  $(X_2, \mathcal{A}_2, \mu_2)$  be two spaces with finite non-negative measures. On the space  $X_1 \times X_2$  we consider sets of the form  $A_1 \times A_2$ , where  $A_i \in \mathcal{A}_i$ , called measurable rectangles. Let  $\mu_1 \times \mu_2 (A_1 \times A_2) := \mu_1(A_1)\mu_2(A_2)$ .

Let  $\mathcal{A}_1 \otimes \mathcal{A}_2$  denote the  $\sigma$ -algebra generated by all measurable rectangles; this  $\sigma$ -algebra is called the product of the  $\sigma$ -algebras  $\mathcal{A}_1$  and  $\mathcal{A}_2$ .

**Theorem 1.10.** The set function  $\mu_1 \times \mu_2$  is countably additive on the algebra generated by all measurable rectangles and uniquely extends to a countably additive measure, denoted by  $\mu_1 \otimes \mu_2$ .

**Definition 1.38** (Image Measure). Let  $(X, \mathcal{A}_X)$  and  $(Y, \mathcal{A}_Y)$  be two measurable spaces. Let  $T : X \rightarrow Y$  be a measurable map from  $X$  to  $Y$ . Let  $\mu$  be a measure  $\mu : \mathcal{A}_X \rightarrow \overline{\mathbb{R}}_+$ , then the image measure (or pushforward measure)  $T_\# \mu : \mathcal{A}_Y \rightarrow \overline{\mathbb{R}}_+$  is given by,

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

**Theorem 1.11.** Let  $(X, \mathcal{A}_X)$  and  $(Y, \mathcal{A}_Y)$  be two measurable spaces. Let  $\mu$  be a nonnegative measure. A  $\mathcal{A}_Y$ -measurable function  $g$  on  $Y$  is integrable with respect the measure  $\mu \circ f^{-1}$  precisely when the function  $g \circ f$  is integrable with respect to  $\mu$ . In addition we have,

$$\int_Y g(y) \mu \circ f^{-1}(dy) = \int_X g(f(x)) \mu(dx)$$

The space of Borel probability measures on  $X$  is denoted by  $\mathcal{P}(X)$ . The weak topology on  $\mathcal{P}(X)$  is induced by convergence against bounded continuous test functions on  $X$ , that is  $C_b(X)$ .

**Definition 1.39** (Atom and atomless measures). The set  $A \in \mathcal{A}$  is called an atom of the measure  $\mu$  if  $\mu(A) > 0$  and every set  $B \subset A$  from  $\mathcal{A}$  has measure either 0 or  $\mu(A)$ . If there are no atoms, then the measure  $\mu$  is called atomless.

A measure over a set  $\Omega \subset \mathbb{R}$  is atomless if  $\forall x \in \Omega$ , we have  $\mu(\{x\}) = 0$ . The Dirac's measure is not atomless.

**Definition 1.40** (Absolutely continuity and singularity). Let  $\mu$  and  $\nu$  be countably additive measures on a measurable space  $(X, \mathcal{A})$ .

- The measure  $\nu$  is called absolutely continuous with respect to  $\mu$  if  $|\nu|(A) = 0$  for every set  $A$  with  $|\mu|(A) = 0$ . We use the notation  $\nu \ll \mu$ .
- The measure  $\nu$  is called singular with respect to  $\mu$  if there exists a set  $A \in \mathcal{A}$  such that

$$|\mu|(A) = 0 \quad \text{and} \quad |\nu|(X \setminus A) = 0$$

If  $\nu \ll \mu$  and  $\mu \ll \nu$ , then the measures  $\mu$  and  $\nu$  are equivalent. We use the notation  $\mu \sim \nu$  to refer this situation.



The above definition allows us to introduce the Radon-Nikodym theorem that is one of the main results in measure theory.

**Theorem 1.12** (Radon–Nikodym theorem). *Let  $\mu$  and  $\nu$  be two finite measures on a space  $(X, \mathcal{A})$ . The measure  $\nu$  is absolutely continuous with respect to the measure  $\mu$  precisely when there exists a  $\mu$ -integrable function  $f$  such that  $\nu$  is given by*

$$\nu(A) = \int_A f d\mu$$

**Definition 1.41** ( $L^p$  Spaces).

**Theorem 1.13** (Lebesgue dominated convergence theorem). *Suppose that  $\mu$ -integrable functions  $f_n$  converge almost everywhere to a function  $f$ . If there exists a  $\mu$ -integrable function  $\Phi$  such that,*

$$|f_n|(x) \leq \Phi(x), \quad \text{almost everywhere for every } n$$

*then the function is integrable and*

$$\int_X f(x) d\mu(x) = \lim_{n \rightarrow \infty} \int_X f_n(x) d\mu(x)$$

*In addition,*

$$\lim_{n \rightarrow \infty} \int |f(x) - f_n(x)| d\mu(x) = 0$$

**Theorem 1.14** (Monotone Convergence). *Let  $(f_n)_{n \in \mathbb{N}}$  be a sequence of  $\mu$ -integrable functions such that  $f_n(x) \leq f_{n+1}$  almost everywhere for each  $n \in \mathbb{N}$ . Suppose that*

$$\sup_{n \in \mathbb{N}} \int_X f_n(x) d\mu(x) < \infty \quad (1.2)$$

*Then the function  $f(x) = \lim_{n \rightarrow \infty} f_n(x)$  is almost everywhere finite and integrable. In addition the following equality holds,*

$$\int_X f(x) d\mu(x) = \lim_{n \rightarrow \infty} \int_X f_n(x) d\mu(x)$$

**Theorem 1.15** (Fatou's Theorem). *Let  $(f_n)_{n \in \mathbb{N}}$  be a sequence of nonnegative  $\mu$ -integrable functions convergent to a function  $f$  almost everywhere and let*

$$\sup_{n \in \mathbb{N}} \int_X f_n(x) d\mu(x) \leq K < \infty$$

*Then, the function  $f$  is  $\mu$ -integrable and*

$$\int_X f(x) d\mu(x) \leq K$$

*Moreover,*

$$\int_X f(x) d\mu(x) \leq \liminf_{n \rightarrow \infty} \int_X f_n(x) d\mu(x)$$

**Corollary 1.1.** *Let  $(f_n)_{n \in \mathbb{N}}$  be a sequence of nonnegative  $\mu$ -integrable functions such that*

$$\sup_{n \in \mathbb{N}} \int_X f_n(x) d\mu(x) \leq K < \infty$$

*Then the function  $\liminf_{n \rightarrow \infty} f_n$  is  $\mu$ -integrable and one has*

$$\int_X \liminf_{n \rightarrow \infty} f_n(x) d\mu(x) \leq \liminf_{n \rightarrow \infty} \int_X f_n(x) d\mu(x) \leq K \quad (1.3)$$

**Corollary 1.2.** *The dominated convergence theorem and Fatou's theorem remain valid if in place of almost everywhere convergence in their hypotheses we require convergence of  $(f_n)_{n \in \mathbb{N}}$  to  $f$  in measure  $\mu$ .*

**Definition 1.42** (Tightness). *Let  $(X, \mathcal{T})$  a topological space, and let  $\mathcal{A}$  a  $\sigma$ -algebra on  $X$  that contains the topology  $\mathcal{T}$ . Let  $M$  be a collection of measures defined on  $\mathcal{A}$ . The collection  $M$  is called tight if for every  $\epsilon > 0$  there is a compact subset  $K_\epsilon$  of  $X$  such that, for all measures  $\mu \in M$  we have,*

$$|\mu|(X \setminus K_\epsilon) < \epsilon$$

**Definition 1.43.** *A sequence  $\mu_n$  probability measures over  $X$  is said to be tight if for every  $\epsilon > 0$ , there exists a compact subset  $K \subset X$  such that  $\mu_n(X \setminus K) < \epsilon$  for every  $n$ .*

**Theorem 1.16** (Prokhorov). *Suppose that  $\mu_n$  is a tight sequence of probability measures over a Polish space  $X$ . Then there exists  $\mu \in \mathcal{P}(X)$  and a subsequence  $\mu_{n_k}$  such that  $\mu_{n_k} \rightarrow \mu$ , in duality with  $C_b(X)$ . Conversely, every sequence  $\mu_{n_k} \rightarrow \mu$  is tight.*

**Theorem 1.17** (Prokhorov). *If  $X$  is a Polish space, then a set  $P \subset \mathcal{P}(X)$  is precompact for the weak topology if and only if it is tight.*

**Definition 1.44.** *Let  $(X, \mathcal{A}, \mu)$  be a probability space. Then every Borel-measurable mapping  $\mathcal{X} : X \rightarrow \mathbb{R}$  with for all  $B \in \mathcal{B}(\mathbb{R})$  is a random variable, denoted by  $\mathcal{X} : (X, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$*

**Definition 1.45** (Duality between  $C_0$  and  $\mathcal{M}$ ).

**Explanation about notions of convergence with bounded functionals and vanishing in infinity functions.** If  $X$  is compact we have  $C_0(X) = C_b(X) = C(X)$  if  $X$  and both notions of convergence coincide.

**Theorem 1.18** (Rademacher). *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a Lipschitz continuous function. Then the set of points where  $f$  is not differentiable is negligible for the Lebesgue measure.*

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

**Theorem 1.19.** *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 1.20.** *Consider on a compact metric space  $X$ , endowed with a probability  $\rho \in \mathcal{P}(X)$ , a sequence of partitions  $G_n$ , each  $G_n$  being a family of disjoint subsets,  $\bigcup_{i \in I_n} C_{i,n} = X$  for every  $n$ . Suppose that  $\text{size}(G_n) := \max_i (\text{diam}(C_{i,n}))$  tends to 0 as  $n \rightarrow \infty$  and consider a sequence of probability measures  $\rho_n$  on  $X$  such that, for every  $n$  and  $i \in I_n$ , we have  $\rho_n(C_{i,n})$ . Then  $\rho_n \rightarrow \rho$ .*

# 2

## Basics in Convex Analysis.

**Definition 2.1** (Variety).

**Definition 2.2** (Convexity). *Let  $X$  a linear space.*

**Definition 2.3** (Graph and Epigraph).

**Definition 2.4** (Infimal Convolution).

**Theorem 2.1** (Hahn-Banach separation theorem).

**Theorem 2.2** (Geometrical version of Hahn-Banach Theorem). *Let  $M$  be a vector subspace of the topological vector space  $X$ . Suppose  $K$  is a non-empty convex open subset of  $X$  with  $K \cap M = \emptyset$ . Then there is a closed hyperplane  $N \in X$  containing  $M$  with  $K \cap N = \emptyset$ .*

**Definition 2.5** (Cone).

**Definition 2.6** (Extreme Point). *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  $0 < \alpha < 1$ .*

**Definition 2.7** (Convex conjugate function). *Let  $X$  be a Banach space, let  $f : X \rightarrow \overline{\mathbb{R}}$  be a functional over  $X$ . We call the convex conjugate to the function  $f^* : X^* \rightarrow \overline{\mathbb{R}}$ , defined as*

$$f^*(x^*) = \sup_{x \in X} \{ \langle x^*, x \rangle - f(x) \}$$

**Proposition 2.1.** *The convex conjugate  $f^* : X^* \rightarrow \overline{\mathbb{R}}$  of a function  $f : X \rightarrow \overline{\mathbb{R}}$  is convex.*

*Proof.* Let  $x^*, y^*$  elements of the dual space  $X^*$ , and  $t \in [0, 1]$ ,

$$\begin{aligned} f^*(tx^* + (1-t)y^*) &= \sup_{x \in X} \{ \langle tx^* + (1-t)y^*, x \rangle - f(x) \} \\ &= \sup_{x \in X} \{ \langle tx^* + (1-t)y^*, x \rangle - tf(x) - (1-t)f(x) \} \\ &= \sup_{x \in X} \{ t \langle x^*, x \rangle + (1-t) \langle y^*, x \rangle - tf(x) - (1-t)f(x) \} \\ &\leq \sup_{x, y \in X} \{ t \langle x^*, x \rangle + (1-t) \langle y^*, y \rangle - tf(x) - (1-t)f(y) \} \\ &= t \sup_{x \in X} \{ \langle x^*, x \rangle - f(x) \} + (1-t) \sup_{y \in X} \{ \langle y^*, y \rangle - f(y) \} \\ &= tf^*(x^*) + (1-t)f^*(y^*). \end{aligned}$$

Therefore  $f^*$  is convex regardless the convexity of  $f$ . □

**Theorem 2.3.** *A function  $f : \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$  is convex and lower-semicontinuous if and only if  $f^{**} = f$ .*

**Lemma 2.1** (Convex envelope theorem). *Let  $X$  be a reflexive Banach Space. Then the convex conjugate function  $f^*$  is the maximum convex functional below  $f$  (also called convex envelope), i.e. if  $g$  is convex functional and  $g(x) \leq f(x)$ ,  $\forall x \in X$ . Then,  $f^{**}(x) \leq f(x)$ , and  $g(x) \leq f^{**}(u)$ ,  $\forall x \in U$ . In particular  $f^{**} = f$  if and only if  $f$  is convex.*

**Definition 2.8** (Subdifferential). *Given a proper convex function  $f : X \rightarrow (-\infty, \infty]$ , the subdifferential of such a function is the mapping  $\partial f : X \rightarrow X^*$  defined by,*

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

**Theorem 2.4.** *The epigraph of a convex and lower semicontinuous function is a closed convex set in  $\mathbb{R}^d \times \mathbb{R}$ , and can be written as the intersection of the half-spaces which contain it.*

**Here we write the proof for the identity for the projection onto an affine set**

**Definition 2.9** (Projection onto a Set).

**Theorem 2.5.**

An important concept in convex programming is duality.

**Definition 2.10** (Duality).

## Convexity in $\mathbb{R}^n$

Since many results are developed in finite dimensions we state some results related to convexity in  $\mathbb{R}^n$ .

**Theorem 2.6.** *Let  $C$  be a convex set and let  $y$  be a point exterior to the closure of  $C$ . Then there is a vector  $\mathbf{a}$  such that  $\mathbf{a}^\top \mathbf{y} < \inf_{\mathbf{x} \in C} \mathbf{a}^\top \mathbf{x}$ .*

**Definition 2.11.** *A hyperplane in  $\mathbb{R}^n$  is an  $(n - 1)$ -dimensional linear variety.*

**Proposition 2.2.** *Let  $\mathbf{a}$  be a nonzero  $n$ -dimensional column vector, and let  $c$  be a real number. The set*

$$H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} = c\}$$

*is a hyperplane in  $\mathbb{R}^n$ .*

**Proposition 2.3.** *Let  $H$  be a hyperplane in  $\mathbb{R}^n$ . Then there is a nonzero  $n$ -dimensional vector and a constant  $c$  such that,*

$$H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}^\top \mathbf{x} = c\}$$

**Definition 2.12.** *Let  $\mathbf{a}$  be a nonzero vector in  $\mathbb{R}^n$  and let  $c$  be a real number. Corresponding to the hyperplane  $H = \{\mathbf{x} : \mathbf{a}^\top \mathbf{x} = c\}$  are the positive and negative closed half spaces*

$$H_+ = \{\mathbf{x} : \mathbf{a}^\top \mathbf{x} \geq c\}$$

$$H_- = \{\mathbf{x} : \mathbf{a}^\top \mathbf{x} \leq c\}$$

*and the positive and negative open half spaces*

$$\overset{\circ}{H}_+ = \{\mathbf{x} : \mathbf{a}^\top \mathbf{x} > c\}$$

$$\overset{\circ}{H}_- = \{\mathbf{x} : \mathbf{a}^\top \mathbf{x} < c\}$$

Half spaces are convex sets and the union of  $H_+$  and  $H_-$  is the whole space.

**Definition 2.13.** *A set which can be expressed as the intersection of a finite number of closed half spaces is said to be a convex polytope.*

Convex polytopes are the sets obtained as the family of solutions to a set of linear inequalities of the form,

$$\begin{aligned} \mathbf{a}_1^\top \mathbf{x} &\leq b_1 \\ \mathbf{a}_2^\top \mathbf{x} &\leq b_2 \\ &\vdots \\ \mathbf{a}_m^\top \mathbf{x} &\leq b_m, \end{aligned}$$

since each individual inequality defines a half space and the solution family is the intersection of these half space.

**Definition 2.14.** A nonempty bounded polytope is called a **polyhedron**.

**Theorem 2.7.** Let  $C$  be convex set and let  $\mathbf{y}$  be a point exterior to the closure of  $C$ . Then there is a vector  $\mathbf{a}$  such that  $\mathbf{a}^\top \mathbf{y} < \inf_{\mathbf{x} \in C} \mathbf{a}^\top \mathbf{x}$ .

**Theorem 2.8.** Let  $C$  be a convex set and let  $\mathbf{y}$  be a boundary point of  $C$ . Then there is a hyperplane containing  $\mathbf{y}$  and containing  $C$  in one of its closed half spaces.

**Definition 2.15.** A hyperplane containing a convex set  $C$  in one of its closed half spaces and containing a boundary point of  $C$  is said to be a **supporting hyperplane** of  $C$ .

**Theorem 2.9.** Let  $B$  and  $C$  be convex sets with no common relative interior points. Then there is a hyperplane separating  $B$  and  $D$ . In particular, there is a nonzero vector  $\mathbf{a}$  such that  $\sup_{\mathbf{b} \in B} \mathbf{a}^\top \mathbf{b} \leq \inf_{\mathbf{c} \in C} \mathbf{a}^\top \mathbf{c}$ .

**Theorem 2.10.** Let  $C$  be a convex set,  $H$  a supporting hyperplane of  $C$ , and  $T$  the intersection of  $H$  and  $C$ . Every extreme point of  $T$  is an extreme point of  $C$ .

**Theorem 2.11.** A closed bounded convex set in  $\mathbb{R}^n$  is equal to the closed convex hull of its extreme point.

**Theorem 2.12.** A convex polyhedron can be described either as a bounded intersection of a finite number of closed half spaces, or as the convex hull of a finite number of points.



# 3

## Linear Programming

Linear programming is a well studied branch of the mathematics that studies the optimization of linear functions under linear constraints. The study of linear programming started during the second part of the 1940s, as a technique military oriented problems.

We can formulate the 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

$$\min \quad \mathbf{c}^T \mathbf{x} \quad (3.1)$$

$$\text{subject to} \quad \mathbf{Ax} = \mathbf{b} \quad (3.2)$$

$$\mathbf{x} \geq 0 \quad (3.3)$$

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 3.1** (Basic solutions). Given the set of  $m$  simultaneous linear equations (3.2) 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  $\mathbf{Ax} = \mathbf{b}$ , 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 3.2** (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 3.3** (Feasible solutions.). A vector  $\mathbf{x}$  satisfying the constraints (3.2) and (3.4) 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 3.1** (Fundamental theorem of linear programming.). Given a linear program in the standard form (3.1), (3.2) and (3.3) 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.

Since 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 3.2.** 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

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b} \\ \mathbf{x} &\geq 0 \end{aligned} \quad (3.4)$$

A vector  $\mathbf{x}$  is an extreme point of  $K$  if and only if  $\mathbf{x}$  is a basic feasible solution of (3.4).

**Corollary 3.1.** If the convex set  $K$  corresponding to (3.4) is nonempty, it has at least one extreme point.

**Corollary 3.2.** 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 3.3.** The constraint set  $K$  corresponding to (3.4) 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.  $\square$

**Corollary 3.4.** If the convex polytope  $K$  corresponding to (3.4) 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 belong to the constraint set of a problem in standard form to another, in such a way as 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  $A$  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,

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

(3.5)

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,

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

From which 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,

$$\begin{cases} y'_{k,j} = \frac{y_{k,j}}{y_{k,l}} \\ y'_{i,j} = y_{i,j} - \frac{y_{i,l}}{y_{k,l}} \end{cases} \quad \text{for } i \neq k \text{ and } 0 \leq i \leq n \quad (3.6)$$

Pivoting vectors as explained above we can generate a new basic solution from an old one. The problem is that the nonnegative constrain 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 nondegenerated. 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  $\mathbf{x} = (x_1, x_2, \dots, x_m, 0, 0, \dots, 0)$ . We are assuming nondegeneracy of the solutions, therefore  $x_i > 0$  for  $i = 1, \dots, m$ .

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

$$\begin{aligned} \mathbf{b} &= x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \dots + x_m\mathbf{a}_m + \epsilon\mathbf{a}_l - \epsilon\mathbf{a}_l \\ &= (x_1 - \epsilon y_{1,l})\mathbf{a}_1 + (x_2 - \epsilon y_{2,l})\mathbf{a}_2 + \dots + (x_m - \epsilon y_{m,l})\mathbf{a}_m + \epsilon\mathbf{a}_l \end{aligned}$$

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  $\epsilon$ ,

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

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, under the assumption that the coefficients  $y_{1,0}, \dots, y_{m,0}$  are nonnegative, implying that  $x_1 = y_{1,0}, x_2 = y_{2,0}, \dots, x_m = y_{m,0}$  is feasible. We substitute a vector already in the basis by a vector  $\mathbf{a}_l$ , in such a way that the solution 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,

$$\mathbf{A} = [\mathbf{B}, \mathbf{D}] \quad (3.7)$$

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,

$$\mathbf{x} = (\mathbf{x}_B, \mathbf{x}_D), \quad \mathbf{c} = (\mathbf{c}_B, \mathbf{c}_D) \quad (3.8)$$

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,

$$\begin{aligned} \min & \quad \mathbf{c}_B^T \mathbf{x}_B + \mathbf{c}_D^T \mathbf{x}_D \\ \text{subject to} & \quad \mathbf{B} \mathbf{x}_B + \mathbf{D} \mathbf{x}_D \\ & \quad \mathbf{x}_B \geq 0, \quad \mathbf{x}_D \geq 0 \end{aligned}$$

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

$$z_0 = \mathbf{c}_B^T \mathbf{x}_B$$

We can construct nonbasic feasible solution, setting arbitrary values for  $\mathbf{x}_D = (x_{m+1}, x_{m+2}, \dots, x_n)$  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,

$$z = \mathbf{c}^T \mathbf{x} = z_0 + (c_{m+1} - z_{m+1}) x_{m+1} + (c_{m+2} - z_{m+2}) x_{m+2} + \dots + (c_n - z_n) x_n. \quad (3.9)$$

where,

$$z_j = y_{1,j} c_1 + y_{2,j} c_2 + \dots + y_{m,j} c_m, \quad \text{for } m+1 \leq j \leq n. \quad (3.10)$$

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

**Theorem 3.3** (Improvement of basic feasible solution). *Given a nondegenerated 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 (3.9) and (3.10), 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  $(x_1, x_2, \dots, x_m, 0, \dots, 0)$  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  $(x'_1, x'_2, \dots, x'_{m+1}, 0, 0, \dots, 0)$ , with  $x'_m > 0$ , substituting this new solution into equation (3.9) we obtain,

$$z - z_0 = (c_{m+1} - z_{m+1}) x'_{m+1} < 0$$

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.  $\square$

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 3.4** (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 (3.9), since any other feasible solution must have  $x_i \geq 0$  for all  $i$ , and hence the value  $z$  of the objective will satisfy  $z - z_0 \geq 0$ .  $\square$

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 (3.5) corresponding to a basic feasible solution. The relative cost coefficients  $r_j$  can be found by 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.

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,

$$\begin{aligned} z &= \mathbf{c}_B^T (\mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{D}\mathbf{x}_D) + \mathbf{c}_D^T \mathbf{x}_D \\ &= \mathbf{c}_B^T \mathbf{B}^{-1}\mathbf{b} + (\mathbf{c}_D^T - \mathbf{c}_B^T \mathbf{B}^{-1}\mathbf{D}) \mathbf{x}_D \end{aligned}$$

Thus, the vector  $\mathbf{r}_D = \mathbf{c}_D^T - \mathbf{c}_B^T \mathbf{B}^{-1}\mathbf{D}$  is the relative cost for non-basic variables. The components of this vector are used to determine which vector bring into the basis.

1. Calculate the current relative cost coefficients  $\mathbf{r}_D = \mathbf{c}_D^T - \mathbf{c}_B^T \mathbf{B}^{-1}\mathbf{D}$ . It is more efficient and numerically stable to solve the linear system  $\mathbf{v}^T = \mathbf{c}_B^T \mathbf{B}^{-1}$ , then compute the relative vector  $\mathbf{r}_D = \mathbf{c}_D^T - \mathbf{v}^T \mathbf{B}^{-1}\mathbf{D}$ . 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  $\mathbf{q} = \mathbf{B}^{-1}\mathbf{a}_q$  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  $\mathbf{B}^{-1}\mathbf{b}$ . 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,

$$\max \quad \boldsymbol{\lambda}^\top \mathbf{b} \quad (3.11)$$

$$\text{subject to} \quad \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c} \quad (3.12)$$

**Lemma 3.1** (Weak Duality lemma). *If  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are feasible for (3.2) and (3.12), respectively then  $\mathbf{c}^\top \mathbf{x} \geq \boldsymbol{\lambda}^\top \mathbf{b}$ .*

*Proof.* We see that following inequality holds for equations (3.2), (3.12) and the cone  $\mathbf{x} \geq 0$ ,

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

□

**Corollary 3.5.** *If  $\mathbf{x}_0$  and  $\boldsymbol{\lambda}_0$  are feasible for the (3.2) and (3.12) 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 3.5** (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 3.5 that the first condition holds. If the primal is unbounded and  $\boldsymbol{\lambda}$  is feasible for the dual we must have,  $\boldsymbol{\lambda}^\top \mathbf{b} \leq -M$  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

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

We see that  $C$  is a closed cone convex cone. We need to find a point  $(\tilde{r}, \tilde{\mathbf{w}}) \notin C$ , 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  $\mathbf{w} = \mathbf{0}$ ,  $\mathbf{w} = \alpha \mathbf{b} - \mathbf{Ax}_0 = \mathbf{0}$  with  $\mathbf{x}_0 \geq \mathbf{0}$ , 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

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

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  $s\tilde{r} + \boldsymbol{\lambda}^\top \tilde{\mathbf{w}} < 0$ , 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  $(r, \mathbf{w}) \in C$ ,

$$-r + \boldsymbol{\lambda}^\top \mathbf{w} \geq 0 \quad (3.13)$$

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

$$(\mathbf{c} - \boldsymbol{\lambda}^\top \mathbf{A}) - \alpha z_0 + \alpha \boldsymbol{\lambda}^\top \mathbf{b} \geq 0 \quad (3.14)$$

for all  $\mathbf{x} \geq \mathbf{0}$ , and  $\alpha \geq 0$ . Setting  $\alpha = 0$  we have the inequality  $\boldsymbol{\lambda}^\top \leq \mathbf{c}^\top$ , which says  $\boldsymbol{\lambda}$  is feasible for the dual. Setting  $\mathbf{x} = \mathbf{0}$  and  $\alpha = 1$  results in  $\boldsymbol{\lambda}^\top \mathbf{b} \geq z_0$ . Therefore, by means of the Weak Duality lemma we have that  $\boldsymbol{\lambda}^\top \mathbf{b} = z_0$  and by corollary 3.5 we have that  $\boldsymbol{\lambda}$  is optimal for the dual. □

## Complementary Slackness.

**Theorem 3.6.** *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  $i$ .*

- $x_i > 0 \Rightarrow \boldsymbol{\lambda}^\top \mathbf{a}_i = c_i$
- $x_i = 0 \Leftarrow \boldsymbol{\lambda}^\top \mathbf{a}_i < c_i$

*Proof.* If the above conditions hold, then  $(\boldsymbol{\lambda}^\top \mathbf{A} - \mathbf{c}^\top) \mathbf{x} = 0$ . By the means of the Weak Duality lemma and corollary 3.5,  $\boldsymbol{\lambda}^\top \mathbf{b} = \mathbf{c}^\top \mathbf{x}$  implies two solutions are optimal. Conversely, if the two solutions are optimal, by the means of the Duality Theorem  $\boldsymbol{\lambda}^\top \mathbf{b} = \mathbf{c}^\top \mathbf{x}$ . Since each component of  $\mathbf{x}$  is nonnegative and each component of  $\boldsymbol{\lambda}^\top \mathbf{A} - \mathbf{c}^\top$  is nonpositive, and the above conditions must hold.  $\square$

## The Dual simplex Method.

For general linear programs the dual simplex method is most frequently used,

## The Primal-Dual Simplex Method.

This method begins with a feasible solution to the dual problem that is improved at each step by optimizing an associated restricted primal problem. As the method progresses it can be regarded as striving to achieve the complementary slackness conditions for optimality.

The primal-dual method was developed for solving a linear program arising in network flow problems, and it continues to be the most efficient procedure for these problems.

Consider the linear program (3.1), (3.2) and (3.3). Also, consider its dual stated by the equations (3.11) and (3.12).

Given a feasible solution  $\boldsymbol{\lambda}$  to the dual, define the subset  $P$  of  $1, 2, \dots, n$  by  $i \in P$  if  $\boldsymbol{\lambda} \mathbf{a}_i = c_i$  where  $\mathbf{a}_i$  is the  $i$ -th column of  $\mathbf{A}$ . Thus, since  $\boldsymbol{\lambda}$  is dual feasible.

## Min-Max Theorems

$$\max_{\mathbf{y} \in Y} \left( \min_{\mathbf{x} \in X} (\mathbf{x}^\top \mathbf{A} \mathbf{y} + \mathbf{B} \mathbf{x} + \mathbf{C} \mathbf{y}) \right) = \min_{\mathbf{x} \in X} \left( \max_{\mathbf{y} \in Y} (\mathbf{x}^\top \mathbf{A} \mathbf{y} + \mathbf{B} \mathbf{x} + \mathbf{C} \mathbf{y}) \right) \quad (3.15)$$

## Interior Methods.



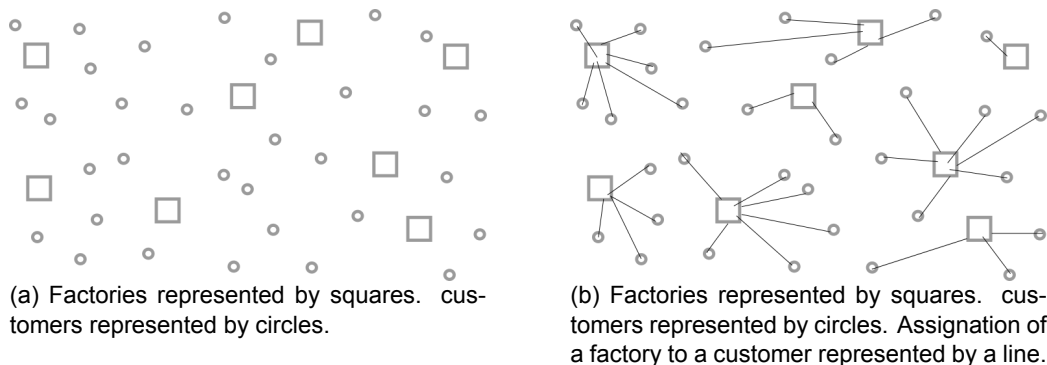
# Optimal Transport Theory

To introduce the optimal transport problem please imagine we are asked by a consortium of factories to design a plan for distributing 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 4.1: Illustration of the problem of Factories supplying customers.



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  $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  pushing the one onto the other,

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

For any Borel subset  $A \subset \mathbb{R}^3$ . And the transport also should minimize the quantity,

$$\int_{\mathbb{R}^3} |x - T(x)| f(x)dx$$

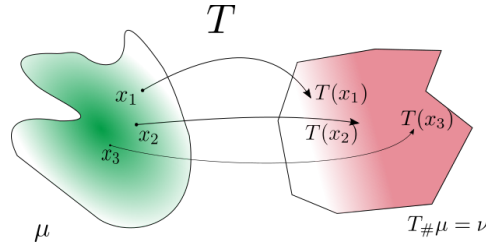
Therefore, we need to search for the optimum in the set of measurable maps  $T : X \rightarrow Y$  such that the condition (4) is translated to,

$$(T_{\#}\mu)(A) = \mu(T^{-1}(A)) \quad \text{for every measurable set } A \subset X. \quad (4.1)$$

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,

$$g(T(x)) \det(DT(x)) = f(x) \quad (4.2)$$

Figure 4.2: Monge problem. Finding a map.



The equation (4.2) is nonlinear in  $T$  making difficult the analysis of the Monge's Problem. Moreover, the constrain makes this problem hard to handle since it is not close even under weak convergence.

To appreciate this fact, consider  $\mu = \mathcal{L}^1 \llcorner [0, 1]$  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  $f_n : [0, 1] \rightarrow [0, 1]$ ,

$$f_n(x) = \sum_{i=0}^{n-1} h_n\left(x - \frac{i}{n}\right) \quad (4.3)$$

We see that the sequence satisfies  $f_{n\#}\mu = \mu$ . It is easy to check that  $\mu(f_n^{-1}(A)) = \mathcal{L}^1(A)$  for every open set  $A \subset [0, 1]$ . In the other hand, the sequence converges weakly to  $f_n \rightharpoonup f = \frac{1}{2}$ , which obviously makes  $f_{\#}\mu \neq \mathcal{L}^1 \llcorner [0, 1]$ .



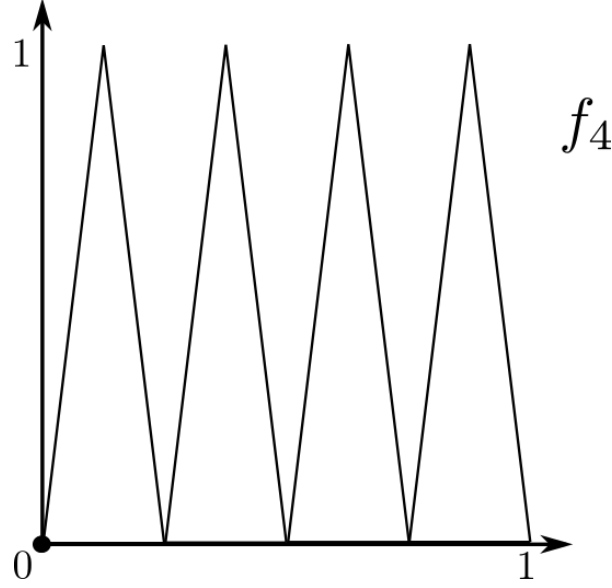


Figure 4.3:  $f_n$  constructed using hat functions. The picture shows the case  $n = 4$ .

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

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

Monge analyzed geometric properties of the solution to this problem. Although, the question of the existence of an optimal map stayed open until a Russian mathematician named Leonid Vitaliyevich Kantorovich introduced in the paper [2] a suitable framework to study its optimality conditions and prove existence of a minimizer.

When we formulate our factories-customer problem through finding an assignment 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 distributed from  $x$  to  $y$ , in such a way to be consistent with information we have the initial and final material configuration. 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  $x \in X$  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 4.1** (Coupling). *Let  $\mu$  and  $\nu$  be probability measures of a probability space  $(X, \mathcal{A}_X)$  and  $(Y, \mathcal{A}_Y)$ . Finding a coupling between  $\mu$  and  $\nu$  means to construct a measure  $\gamma$  on the space  $X \times Y$  (precisely on the product  $\sigma$ -algebra  $\mathcal{A}_X \otimes \mathcal{A}_Y$ ) such that  $\mu$  and  $\nu$  are admitted as marginals on  $X$  and  $Y$  respectively. That is  $\text{proj}_{x\#} \gamma = \mu$  and  $\text{proj}_{y\#} \gamma = \nu$ .*

The above definition is equivalent 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  $Y$ ,

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

Since definition 4.1 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  $\text{law}(\mathcal{X}) = \mu$ ,  $\text{law}(\mathcal{Y}) = \nu$ . The couple  $(\mathcal{X}, \mathcal{Y})$  is called a coupling of  $(\mu, \nu)$ .

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.

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

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

**Definition 4.2** (Deterministic Coupling). *Let  $(X, \mu)$  and  $(Y, \nu)$  be two probabilistic spaces. If there exists a measurable map  $T : X \rightarrow Y$  such that  $T_\# \mu = \nu$ . We call the measure  $(\text{id}, T)_\# \mu = \gamma \in \mathcal{P}(X \times Y)$  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  $T : X \rightarrow Y$  such that  $\mathcal{Y} = T(\mathcal{X})$ . Equivalently,  $(\mathcal{X}, \mathcal{Y})$  is a deterministic coupling of  $\mu$  and  $\nu$ , if its law  $\gamma = \text{law}((\mathcal{X}, \mathcal{Y}))$  is concentrated on the graph of a measurable map  $T : X \rightarrow Y$ . Other way to rephrase it is saying that  $\mu = \text{law}(\mathcal{X})$ ,  $\mathcal{Y} = T(\mathcal{X})$ , where  $T$  is a change of variables from  $\mu$  to  $\nu$ , for all  $\nu$ -integrable (nonnegative measurable) function  $\phi$ ,

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

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

$$\Pi(\mu, \nu) = \left\{ \gamma \in \mathcal{P}(X \times Y) : (\text{proj}_x)_\# \gamma = \mu \text{ and } (\text{proj}_y)_\# \gamma = \nu \right\} \quad (4.4)$$

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,

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

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

$$F^{-1}(t) = \inf\{x \in \mathbb{R}; F(x) > t\}, \quad (4.5)$$

$$G^{-1}(t) = \inf\{y \in \mathbb{R}; G(y) > t\}. \quad (4.6)$$

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  $\mu_1(dx_1)$ ,  $\nu_1(dy_1)$  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 (4.5) and (4.6) respectively.
2. Now take the marginal on the first two variables and disintegrate it with respect to the first variable. This gives probability measures  $\mu_2(dx_1dx_2) = \mu_1(dx_1)\mu_2(dx_2|x_1)$ ,  $\nu_2(dy_1dy_2) = \nu_1(dy_1)\nu_2(dy_2|y_1)$ . 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  $\mu(dx_2|x_1)$  into  $\nu(dy_2|y_1)$ .
3. We repeat the construction, adding one variable after another. For example, after the assignation  $x_1 \rightarrow y_1$  has been determined, the conditional probability of  $x_2$  is seen as a one-dimensional probability on a small slice of width  $dx_1$ , and it can be transported to the conditional probability of  $y_2$  seen as one dimensional probability of a slice of width  $dy_1$ . After  $n$  constructions, this procedure maps  $\mathcal{Y} = T(\mathcal{X})$ .

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

**Lemma 4.1** (Gluing lemma). *If  $Z$  is a function of  $\mathcal{Y}$  and  $\mathcal{Y}$  is a function of  $\mathcal{X}$ , then  $Z$  is a function of  $\mathcal{X}$ . Let  $(X_i, \mu_i)$ ,  $i = 1, 2, 3$ , be Polish probability spaces. If  $(X_1, X_2)$  is a coupling of  $(\mu_1, \mu_2)$  and  $(Y_2, Y_3)$  is a coupling of  $(\mu_2, \mu_3)$ , then it is possible to construct a triple of random variables  $(Z_1, Z_2, Z_3)$  such that  $(Z_1, Z_2)$  has the same law as  $(X_1, X_2)$  and  $(Z_2, Z_3)$  has the same law as  $(Y_2, Y_3)$ .*

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

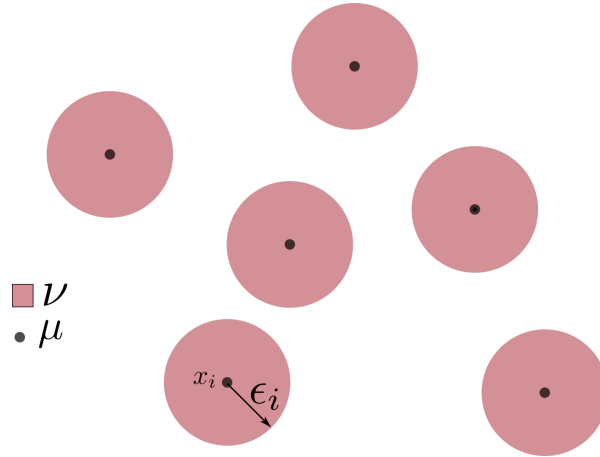
$$\inf \left\{ K(\gamma) := \int_{X \times Y} c d\gamma : \quad \gamma \in \Pi(\mu, \nu) \right\} \quad (\text{KP})$$

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  $\gamma = \mu \otimes \nu$ . Such a thing it is not always possible with transportation maps (deterministic couplings). For example, 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  $D = \bigcup_{i=0}^{N-1} B(x_i, \epsilon_i)$  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 4.4: 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.

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 to use the Weierstrass' criterion to show the existence of a minimizer.

**Theorem 4.1.** *Let  $X$  and  $Y$  be compact metric spaces,  $\mu \in \mathcal{P}(X)$ ,  $\nu \in \mathcal{P}(Y)$  and a cost function  $c : X \times Y \rightarrow \mathbb{R}$  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  $C(X \times Y)$ .

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-\* compactness in dual spaces guarantees the existence of a convergent subsequence  $\gamma_{n_k} \rightarrow \gamma$ . Let us fix  $\phi \in C(X)$  and using  $\int \phi(x) d\gamma_{n_k} = \int \phi d\mu$  and taking the limit we have  $\int_{X \times Y} \phi(x) d\gamma = \int_X \phi d\mu$ . And in this way we prove that  $\gamma_{\#}(\text{proj}_X) = \mu$ .

We can repeat this argument for  $\nu$ , fixing  $\psi \in C(Y)$  and taking the limit of  $\int_{X \times Y} \psi(y) d\gamma_{n_k} = \int \psi d\nu$ , implies  $\int_{X \times Y} \psi(y) d\gamma = \int \psi d\nu$ . This proves that  $\gamma_{\#}(\text{proj}_Y) = \nu$ . 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.

**Lemma 4.2.** *Let  $X$  be a metric space. If  $f : X \rightarrow \overline{\mathbb{R}}$  is a lower semi-continuous function, bounded from below, then the functional  $J : \mathcal{M}_+(X) \rightarrow \overline{\mathbb{R}}$  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  $J_n : \mathcal{M}_+(X) \rightarrow \overline{\mathbb{R}}$ , 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  $J_n(\mu) \leq J(\mu)$  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,  $J_n(\mu) \rightarrow J(\mu)$ , 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  $J$  is lower semicontinuous.  $\square$

**Theorem 4.2.** *Let  $X$  and  $Y$  be compact metric spaces,  $\mu \in \mathcal{P}(X)$ ,  $\nu \in \mathcal{P}(Y)$ , and  $c : X \times Y \rightarrow \overline{\mathbb{R}}$  be lower semi-continuous and bounded from below. Then Kantorovich's problem admits a solution.*

*Proof.* We make use of lemma 4.2, setting  $f = c$  on the space  $X \times Y$ . We apply again Weierstrass criterion proving existence of a minimizer.  $\square$

**Theorem 4.3.** *Let  $X$  and  $Y$  be Polish spaces, and  $c : X \times Y \rightarrow \overline{\mathbb{R}}_+$ , a real valued lower semi-continuous cost function on the space  $X \times Y$ . Then (KP) admits a solution.*

**Lemma 4.3.** *Let  $X$  and  $Y$  be Polish spaces,  $\mu \in \mathcal{P}(X)$ ,  $\nu \in \mathcal{P}(Y)$  and  $c : X \times Y \rightarrow [0, +\infty]$  lower semicontinuous. Then the Kantorovich's problem admits a solution.*

*Proof.* Fix  $\epsilon > 0$  and find two compact sets  $K_X \subset X$  and  $K_Y \subset Y$  such that  $\mu(X \setminus K_X) < \epsilon$ , and  $\nu(Y \setminus K_Y) < \epsilon$ . 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,

$$\gamma_n((X \times Y) \setminus (K_X \times K_Y)) \leq \gamma_n((X \setminus K_X) \times Y) + \gamma_n(X \times (Y \setminus K_Y)) \quad (4.7)$$

$$= \mu(X \setminus K_X) + \nu(Y \setminus K_Y) \quad (4.8)$$

$$= 2\epsilon \quad (4.9)$$

Given the arbitrary way to choose  $\epsilon$ , this shows tightness of all sequences in  $\Pi(\mu, \nu)$  and hence compactness.  $\square$

## Kantorovich formulation as relaxation

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

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

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

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

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

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

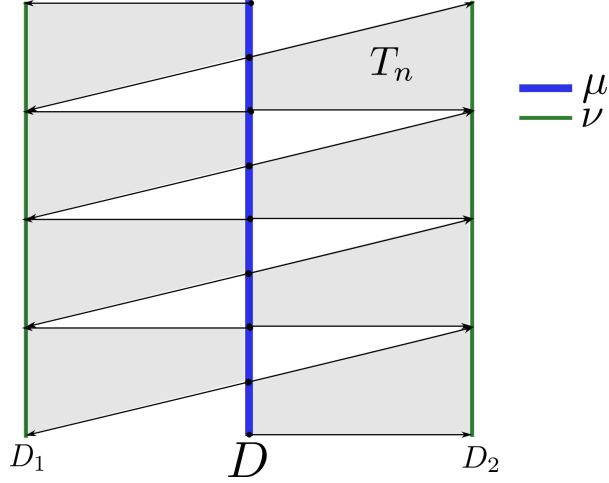


Figure 4.5: There is a deterministic coupling for  $\mu$  and  $\nu$ , but no optimal one. The map  $T_n$  shown in this picture with  $n = 4$ .

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  $k = 0, \dots, 2n - 1$  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  $\mathcal{C}(T_n)$ ,

$$\begin{aligned} \mathcal{C}(T_n) &= \int_D |x - T_n(x)| d\mu(x) \\ &= 2n \int_0^{\frac{1}{2n}} \sqrt{1 + 4t^2} dt \\ &\leq 2n \left( \int_0^{\frac{1}{2n}} 1 + 4t^2 dt \right)^{1/2} \left( \int_0^{\frac{1}{2n}} dt \right)^{1/2} \\ &= \sqrt{1 + \frac{1}{3n^3}} \\ &\leq 1 + \frac{1}{n} \end{aligned}$$

Let  $\gamma_{T_n}$  be deterministic coupling generated by  $T_n$ . We see that we can find always find a cheaper plan  $\gamma_{T_{n+1}}$  for any  $n \in \mathbb{N}$ . This sequence of transportation plans 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:

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

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 4.4.** *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.*

*Proof.* Fix  $n$  and consider any partition of  $\Omega$  into sets  $K_{i,n}$  of diameter smaller than  $1/(2n)$  (for instances, small cubes). The sets  $C_{i,j,n} := K_{i,n} \times K_{j,n}$  make a partition of  $\Omega \times \Omega$  with size smaller than  $1/n$ .

Take any measure  $\gamma \in \Pi(\mu, \nu)$ . □

**Theorem 4.4.** *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  $\gamma_{T_n} \rightarrow \gamma$  and  $J(\gamma_{T_n}) \rightarrow K(\gamma)$ , 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  $\gamma_{T_n} \rightarrow \gamma$ . The last step is possible because of the density of transport plans generated by a map  $\gamma_{T^n}$  in the set of transport plans  $\Pi(\mu, \nu)$ . □

## Cyclical Monotonicity and Duality.

### Duality

Imagine instead 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 selling 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 need to be competitive and should do it better than the consortium did. Therefore, he must be able that 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.

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

$$\sup_{\phi, \psi} \left( \int_X \phi d\mu + \int_Y \psi d\nu - \int_{X \times Y} (\phi(x) + \psi(y)) d\gamma \right) = \begin{cases} 0 & \text{if } \gamma \in \Pi(\mu, \nu) \\ +\infty & \text{otherwise.} \end{cases} \quad (4.10)$$

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 constrain 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 constrain.

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

$$K(\gamma) = \min_{\gamma} \left( \int_{X \times Y} c d\gamma + \int_X \phi d\mu + \int_Y \psi d\nu - \int_{X \times Y} (\phi(x) + \psi(y)) d\gamma \right) \quad (4.11)$$

$$= \sup_{\phi, \psi} \left( \int_X \phi d\mu + \int_Y \psi d\nu \right) + \min_{\gamma} \left( \int_{X \times Y} c(x, y) d\gamma - \sup_{\phi, \psi} \left( \int_X \phi(x) d\mu + \int_Y \psi(y) d\nu \right) \right) \quad (4.12)$$

Note that,

$$\inf_{\gamma} \int_{X \times Y} (c(x, y) - (\phi(x) + \psi(y))) d\gamma = \begin{cases} 0 & \text{if } \phi(x) + \psi(y) \leq c(x, y), \quad \forall (x, y) \in X \times Y \\ -\infty & \text{otherwise.} \end{cases} \quad (4.13)$$

In the other hand, equation (4.12) 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 (4.12) as follows,

$$K(\gamma) = \sup_{\phi, \psi} \left( \int_X \phi d\mu + \int_Y \psi d\nu + \inf_{\gamma} \int_{X \times Y} (c(x, y) - (\phi(x) + \psi(y))) d\gamma \right) \quad (4.14)$$

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  $c : X \times Y \rightarrow \mathbb{R}_+$  we refer as the dual formulation of the transport problem (DP),

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

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

$$\int_X \phi d\mu + \int_Y \psi d\nu = \int_{X \times Y} \phi(x) + \psi(y) d\gamma \leq \int_{X \times Y} c(x, y) d\gamma \quad (4.15)$$

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 tackle this problem we need to introduce the concept

**Definition 4.3.** Let  $X, Y$  be two sets, let  $c : X \times Y \rightarrow \mathbb{R} \cup \{-\infty\}$ , a cost function over from  $X$  to  $Y$ . Given a function  $\omega : X \rightarrow \overline{\mathbb{R}}$ , we define its  $c$ -concave transform of  $\omega$ , the function  $\omega^c : Y \rightarrow \overline{\mathbb{R}}$  by

$$\omega^c(y) = \inf_{x \in X} (c(x, y) - \omega(x)) \quad (4.16)$$

In similar way, we can define the  $\bar{c}$ -transform of  $v : Y \rightarrow \overline{\mathbb{R}}$  by

$$v^{\bar{c}}(x) = \inf_{y \in Y} (c(x, y) - v(y)) \quad (4.17)$$

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

Using the above definitions we call the superdifferential of a function  $\omega$  the  $c$ -cyclically monotone set,

$$\partial^{\bar{c}} v = \{(x, y) \subset X \times Y; \quad v(y) + v^{\bar{c}}(x) = c(x, y)\} \quad (4.18)$$



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  $c : X \times Y \rightarrow \mathbb{R}$ . We denote by  $c - \text{conc}(X)$  the set of  $c$ -concave functions with respect to  $X$ . In similar way, we denote by  $\bar{c} - \text{conc}(Y)$  the set of  $\bar{c}$ -concave functions with respect to  $Y$ .

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

**Theorem 4.5.** *Suppose that  $c$  is real valued. For any  $\phi : X \rightarrow \mathbb{R} \cup \{-\infty\}$ , 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  $\phi$ ,

$$\begin{aligned}\phi^{c\bar{c}}(x) &= \inf_{y \in Y} (c(x, y) - \phi^{\bar{c}}(y)) \\ &= \inf_{y \in Y} \left( c(x, y) - \inf_{\tilde{x} \in X} (c(\tilde{x}, y) - \phi(\tilde{x})) \right).\end{aligned}$$

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

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

Therefore,

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

We can just repeat the above arguments for some  $v : Y \rightarrow \mathbb{R} \cup \{-\infty\}$ , having as result  $\square$

Consider a similar situation to the factories-customers example, but in this new hypothetical situation the consortium has already a fixed transportation plan. They know that the costs are high and they want to make them cheaper.

**Definition 4.4.** *Let  $X, Y$  be arbitrary sets, and  $c : X \times Y \rightarrow (-\infty, \infty]$  be a cost function. A subset  $\Gamma \subset X \times Y$  is said to be  $c$ -cyclically monotone if, for any  $N \in \mathbb{N}$ , and any family of points  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$  of  $\Gamma$ , the inequality*

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

considering  $N + 1 = 1$ .

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

$$\sum_{i=1}^N c(x_i, y_i) \leq \sum_{i=1}^N c(x_i, y_{\sigma(i)}) \quad (4.19)$$

The subdifferentials of convex functions on  $\mathbb{R}^n$  are characterized in terms of a monotonicity property. **Restate Rockafellar**

**Theorem 4.6** (Rockafellar). *Let  $\Gamma$  be a  $c$ -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 4.6 is a well known theorem in convex analysis. It basically states that every cyclically monotone set is contained in the graph of the subdifferential of a convex function.

**Definition 4.5** (Support of transport plan.). *Given a separable metric space  $X$ , the support of a measure  $\gamma$  is defined as the smallest closed set on which  $\gamma$  is concentrated,*

$$\text{spt}(\gamma) := \bigcap_{\substack{A \\ \gamma(X \setminus A) = 0}} A \quad (4.20)$$

**Theorem 4.7.** *If  $\gamma$  is an optimal transport plan for the cost  $c$  and  $c$  is continuous then  $\text{spt}(\gamma)$  is a  $c$ -cyclical monotone set.*

**Theorem 4.8.** *Suppose that  $X$  and  $Y$  are Polish spaces and that  $c : X \times Y \rightarrow \mathbb{R}$  is uniformly continuous and bounded. Then the problem admits a solution  $(\phi, \phi^c)$  and their objective values coincide.*

**Theorem 4.9.** *Suppose that  $X$  and  $Y$  are compact and  $c$  is continuous. Then there exists a solution  $(\phi, \psi)$  to the dual problem (DP) and  $\phi$  is a  $c$ -concave function over  $X$  and  $\psi = \phi^c$ .*

$$\max(DP) = \max_{\phi \text{ is } c\text{-conc}(X)} \left( \int_X \phi d\mu + \int_Y \phi^c d\nu \right) \quad (4.21)$$

**Theorem 4.10.** *Let  $\mu, \nu$  be probabilities over  $\mathbb{R}^d$  and  $c(x, y) = \frac{1}{2} |x - y|^2$ . Suppose  $\int |x|^2 dx < \infty$  and  $\int |y|^2 dy < \infty$ . Consider*

**Definition 4.6** (Perturbation.). *Suppose that  $X$  and  $Y$  are compact metric spaces and  $c : X \times Y \rightarrow \mathbb{R}$  is continuous. For every  $p \in C(X \times Y)$ , let  $H : C(X \times Y) \rightarrow \overline{\mathbb{R}}$  be a perturbation of the problem ,*

$$H(p) = - \max \left\{ \int_X \phi d\mu + \int_Y \psi d\nu : \phi(x) + \psi(y) \leq c(x, y) - p(x, y) \right\} \quad (4.22)$$

## Properties of Optimal plans.

**Theorem 4.11** (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.  $\square$

**Lemma 4.5.** *Suppose that*

**Theorem 4.12.** *The optimal transport between two gaussians for cost function  $c(x, y) = |x - y|^2$  is given by a translation map.*

*Proof.*  $\square$

# Computation of an Optimal Transport

The approximation of an optimal transport is a challenging problem, computationally speaking. We have found a rich literature on it, and many recent advances in this topic have arisen in the very last years.

## Linear Programming Formulation.

Let  $X$  and  $Y$  be two finite sets having  $n$  and  $m$  elements respectively. Let  $\mu$  a probability measure defined over  $X$ ,

$$\mu = \sum_{i=1}^n a_i \delta_{x_i}, \quad (5.1)$$

where  $X = \{x_1, x_2, \dots, x_n\}$  and  $0 \leq \mathbf{a} = \{a_1, a_2, \dots, a_n\}$ , and  $\sum_{i=1}^n a_i = 1$ .

Let  $\nu$  a probability measure defined over  $Y$ ,

$$\nu = \sum_{i=1}^m b_i \delta_{y_i}, \quad (5.2)$$

where  $Y = \{y_1, y_2, \dots, y_m\}$  and  $0 \leq \mathbf{b} = \{b_1, b_2, \dots, b_m\}$ , and  $\sum_{i=1}^m b_i = 1$ . Let  $(\boldsymbol{\gamma})_{i,j} = \gamma_{i,j}$  be a joint probability distribution with marginals given by  $\mu$  and  $\nu$ . That is,

$$P(x_i|Y) = \sum_{j=1}^m \gamma_{i,j} = a_i \quad (5.3)$$

$$P(y_j|X) = \sum_{i=1}^n \gamma_{i,j} = b_j \quad (5.4)$$

Let  $c : X \times Y \rightarrow \mathbb{R}$  a cost function. Since we  $X$  and  $Y$  are finite dimensional, we find convenient to use the following notation,

$$(\mathbf{C})_{i,j} = c_{i,j} = c(x_i, y_j) \quad (5.5)$$

Given the matrix nature of the formulation for discrete measures we can compute the total cost by the Frobenius inner product of the two matrices, since  $\mathbf{C}$  and  $\boldsymbol{\gamma}$  have the same dimensions.

$$\langle \mathbf{C}, \boldsymbol{\gamma} \rangle = \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}} c_{i,j} \gamma_{i,j} = \text{tr}(\mathbf{C}^\top \boldsymbol{\gamma}) \quad (5.6)$$

## Simplex Method Algorithm and Duality.

We can solve the problem using the simplex method.

$$\arg \min_{\boldsymbol{\gamma} \in \Pi(\mu, \nu)} \langle \mathbf{C}, \boldsymbol{\gamma} \rangle \quad (5.7)$$

### The simplex method is Not polynomial time.

Consider the following example

$$\max \sum_{j=1}^n 10^{n-j} x_j \quad (5.8)$$

$$\text{subject to } 2 \sum_{j=1}^{i-1} 10^{i-j} x_j + x_i \leq \quad (5.9)$$

### Sinkhorn-Knopp Algorithm.

Consider the problem adding a regularization.

$$\arg \min_{\gamma \in \Pi(\mu, \nu)} \langle \mathbf{C}, \gamma \rangle + H(\gamma) \quad (5.10)$$

Where  $H(\gamma)$  is the Shannon's Entropy defined for a matrix,

$$H(\gamma) = - \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}} \gamma_{i,j} \log(\gamma_{i,j}) \quad (5.11)$$

### Continuous Formulation.

**Beckman Problem and Optimal Transport.**

**Proximal Splitting Algorithms.**

# 6

## Applications

In this section we describe some applications developed under the frame of optimal transportation.

**Domain Adaptation.**

**Isoperimetric Inequality.**

**Wasserstein's Distance**



# Bibliography

- [1] Vladimir I. Bogachev. *Measure theory*. Springer-Verlag, Berlin Heidelberg, 2007.
- [2] L. Kantorovich. On the translocation of masses. *Dokl. Acad. Nauk. USSR* 37, pages 7–8, 1942.
- [3] J.R. Munkres. *Topology*. Topology. Prentice-Hall, 2000. ISBN 9780131784499.