

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

1 Preliminaries.	3
1.1 Definitions and important theorems to remember.	3
1.1.1 Topology.	3
1.1.2 Functional Analysis.	5
1.1.3 Convex Analysis.	6
1.1.4 Measure Theory.	6
2 Linear Programming	9
2.1 Simplex Method.	10
2.2 Duality.	14
2.3 Complementary Slackness.	15
2.3.1 The Dual simplex Method.	15
2.4 Min-Max Theorems.	16
3 Optimal Transport Theory	17
3.0.1 Existence of a minimizer for Kantorovich's Problem.	22
3.0.2 Properties of Optimal plans	23
3.1 Kantorovich's formulation as relaxation of Monge's formulation.. . . .	24
3.2 Cyclical Monotonicity and Duality.	26
3.2.1 Duality.	26
3.2.2 Duality by convex analysis.	34
3.2.3 Brenier's theorem.	34
4 Computation of an Optimal Transport	35
4.1 Linear Programming Formulation.	35
4.1.1 Simplex Method Algorithm and Duality.	36
4.1.2 Sinkhorn-Knopp Algorithm.	36
4.2 Continuous Formulation.	36
4.2.1 Beckman Problem and Optimal Transport.	36
4.2.2 Proximal Splitting Algorithms.	36
5 Applications	37
5.1 Wasserstein's Distances.	37
5.1.1 Statistical Distances	37
5.2 Data Assimilation of a Dynamic.	37
5.3 Isoperimetric Inequality.	37
Bibliography	39

Notation Table.

\emptyset	Empty set
\mathbb{R}	Real numbers field.
$\overline{\mathbb{R}}$	$\mathbb{R} \cup \{\pm\infty\}$. That is $[-\infty, \infty]$
\mathbb{R}_+	The set of nonnegative real numbers, that is the interval $[0, \infty)$.
$\overline{\mathbb{R}}_+$	The set of nonnegative extended real numbers, that is the interval $[0, \infty]$
\mathbb{R}^d	The d -dimensional Euclidean space.
$B \setminus A$	Given a set B and a subset $A \subset B$, the set $B \setminus A$ is the complement of A in B .
$A \cup B$	Union of two sets.
$A \cap B$	Intersection of two sets.
proj_X	Projection of X . Given two sets X and Y the projection of X is the function $\text{proj}_X : X \times Y \rightarrow X$, defined by $\text{proj}_X(x, y) = x$.
proj_Y	Projection of Y . Given two sets X and Y the projection of Y is the function $\text{proj}_Y : X \times Y \rightarrow Y$, defined by $\text{proj}_Y(x, y) = y$.
$\mathbb{1}_\Omega$	Indicator function of a set $\Omega \subset X$. If $x \in \Omega$, then $\mathbb{1}_\Omega(x) = 1$. If $x \in X \setminus \Omega$, we have $\mathbb{1}_\Omega(x) = 0$.
$B(x; \epsilon)$	Open ball with radius ϵ centered at x .
$\overline{B}(x; \epsilon)$	Closed ball with radius ϵ centered at x .
$\text{supp}(f)$	Support of a continuous function f .
id	Identity map. That is $\text{id} : X \rightarrow X$, defined by $\text{id}(x) = x$.
l.s.c.	Lower semicontinuous.
$\frac{\delta F}{\delta \rho}$	First variation of $F : \mathcal{P}(X) \rightarrow \mathbb{R}$, that is $\left. \frac{d}{d\epsilon} F(\rho + \epsilon\chi) \right _{\epsilon=0} = \int \frac{\delta F}{\delta \rho} d\chi$
$DT(x)$	Jacobian matrix of a map $T(x)$.
$f _\Omega$	The restriction of a function f to a set Ω .
$\mathcal{M}(X)$	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 μ is absolutely continuous with respect to the measure ν .
$\mu \llcorner A$	A measure μ is restricted to a set A .
ω_d	The Measure of the unite ball in \mathbb{R}^d .
$T_\# \mu$	The image measure (or pushforward measure) of μ through the map T .
$\text{spt}(\mu)$	Support of a measure μ .
i.i.d.	Independent and identical probability distributions.
\mathcal{L}^p	Lebesgue measure on \mathbb{R}^p
$\mathcal{H} \llcorner A$	Hausdorff measure applied to some set $A \subset \mathbb{R}^d$.
δ_x	The Dirac mass at point x .
$\Pi(\mu, \nu)$	The set of transport plans from μ to ν .
W_p	Wasserstein distance of order p .
\mathbb{W}_p	Wasserstein space of order p .
γ_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 γ .
$\mu \otimes \nu$	The product measure of μ and ν 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 .

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 and convergence. We recall some important notions used during the development of the present text, for a deeper and better understanding of theory we refer to [4] for more details in topology.

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

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

If A is a subset of a topological space X , a point $x \in X$ is called a **limit point** of A if every neighborhood of x intersects A in some point different than x itself. A subset of topological space is closed if and only if contains all its limit points.

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

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

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

Continuity not only depends upon the function f , but also in the topologies specified for the range and the domain of f . The support $\text{supp}(f)$ of a continuous real valued function f is the closure of the set $\{x \in X : f(x) \neq 0\}$.

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

The notion of sequences is useful for tracking compactness in topological spaces satisfying the *first axiom of countability for topological spaces*¹. Although for arbitrary topological spaces we can use the notion of nets. A **net** is a generalization of sequences for arbitrary topological spaces. A net in X is a function f from a directed set² \mathcal{J} into X . If $\alpha \in \mathcal{J}$, we usually denote $f(\alpha) = x_\alpha$. We use the notation $(x_\alpha)_{\alpha \in \mathcal{J}}$ to refer a net. Every sequence is a net but the converse does not hold.

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

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

1. $d(x, y) = d(y, x)$, for all $x, y \in X$.
2. Symmetry: $d(x, y) = 0 \iff x = y$, for all $x, y \in X$.
3. Triangle inequality: $d(x, z) \leq d(x, y) + d(y, z)$ for all $x, y, z \in X$.

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

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

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

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

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

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

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

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

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

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

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

²Any set with a partial order J relation

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

In topology we can find different notions of compactness. We start with the usual notion, a topological space X is called **compact** if every open covering contains a finite subcover also covering X . A metrizable space X is compact if and only if it is sequentially compact. Every compact metric space is separable.

Every closed subspace of a compact space is compact. Every compact subspace of a Hausdorff space is closed. The image of a compact space under a continuous map is compact. The product of *finitely many* compact spaces is compact.

A metric space X is said to be **complete** if any Cauchy sequence³ has a limit in X . Every compact metric space is complete. We call **Polish space** to any topological space that is separable and completely metrizable.

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

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

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

$$d_X(x_0, x_1) < \delta \Rightarrow d_Y(f(x_0), f(x_1)) < \epsilon \quad (1.1)$$

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

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

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

Functional Analysis

Definition 1.4 (Lower Semicontinuous).

Theorem 1.2 (Arzela-Ascoli).

Definition 1.5 (Vector Space).

Definition 1.6 (Linear map).

Definition 1.7 (Affine maps).

Definition 1.8 (Duality).

Theorem 1.3 (Banach-Alaoglu). *Let X be an arbitrary normed space, and let X^* be its dual. A closed unit ball $B \subset X^*$ is weak-star compact.*

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

Theorem 1.4 (Banach-Alaoglu weak-star sequentially compactness). *Let X be a separable normed vector space, and let X^* be its dual. Then the weak-star topology on a closed ball $B \subset X^*$ is metrizable.*

Corollary 1.1. *Let X be a separable normed vector space and let X^* be its topological dual space. Then every bounded sequence $(\phi_n)_{n \in \mathbb{N}} \in X^*$ has a weak-star convergent subsequence.*

Theorem 1.5 (Weak Topology is not metrizable). *content...*

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

Theorem 1.6 (Direct Method of Calculus of Variations).

Convex Analysis.

Measure Theory.

All measures considered in this text are Borel measures on Polish spaces, equipped with their respective Borel σ -algebra. [3]

Theorem 1.7 (Regular Measure). *The outer measure coincides with the inner measure. That is*

Definition 1.9 (Support of a measure.). *Given a separable metric space X , the support of a measure μ is defined as the smallest closed set on which μ is concentrated, that is*

$$\text{spt}(\mu) := \bigcap_{\substack{\mu(X \setminus A) = 0 \\ A = \text{clo } A}} A \quad (1.2)$$

Let $\mathcal{M}^+(X)$ be the set of all nonnegative, finite, finitely additive and regular measures of \mathcal{A}_X . The set of generalized measures is denoted by $\mathcal{M}(X)$. It is a linear vector space; a norm can be introduced by the variation of a measure:

$$|m| = m^+(X) + m^-(X) \quad (1.3)$$

where,

$$\begin{aligned} m^+(X) &:= \sup \{m(B) : B \in \mathcal{A}\} \\ m^-(X) &:= -\inf \end{aligned}$$

Theorem 1.8 (Aleksandrov). *For an arbitrary topological space X any linear continuous functional on $C(X)$ is of the form,*

$$\phi(f) = \langle f, m \rangle = \int_X f(x) dm(x) \quad (1.4)$$

Moreover,

$$|m| = \sup_{\|f\|_X \leq 1} \left| \int_X f(x) dm(x) \right| \quad (1.5)$$

There is an isometrical, isomorphical and one to one mapping between the space of all continuous linear functionals on $C(X)$ and the space of $\mathcal{M}(X)$; We write $C(X) = \mathcal{M}(X)$

Theorem 1.9 (Luizin). *Let X be a locally compact Hausdorff space, let \mathcal{A} be a σ -algebra on X containing the Borel σ -algebra of X . Let μ a regular measure on (X, \mathcal{A}) and let $f : X \rightarrow \mathbb{R}$ be \mathcal{A} -measurable. If A belongs to \mathcal{A} and satisfies $\mu(A) < \infty$ and given a positive number $\epsilon > 0$, then there is a compact $K \subset X$ such that $\mu(X \setminus K) < \epsilon$ and the restriction $f|_K$ is continuous in K . Moreover, there is a function $g \in C_c(K)$ such that $g(x) = f(x)$ for each $x \in K$.*

It follows that the supremum of a collection of continuous (or lower semicontinuous) functions is lower semicontinuous and that each lower semicontinuous function on a Hausdorff space is Borel measurable.

The space of probabilities $\mathcal{P}(X)$ is not closed in $\mathcal{M}(X)$.

The structure of the second dual $\mathcal{C}(X)^{**}$ is too complex, but it is well known that $B(X)$ —the set of all bounded Borel-measurable functions—is a subset of $\mathcal{C}(X)^{**}$.

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

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

Theorem 1.10. *If X is a complete and separable topological space, then $\mathcal{P}(X)$ is tight.*

Theorem 1.11 (Prohorov). *Let X an arbitrary metric space and let Γ be a tight set of measures. Then Γ is relatively compact.*

Results for Image Measure.

Definition 1.10 (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 μ 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,*

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

Lemma 1.1. *If μ, ν are two probability measures on the real line \mathbb{R} and μ 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 1.12. *If μ and ν are two probability measures on \mathbb{R}^d and μ is atomless, then there exists at least a map T such that $T_{\#}\mu = \nu$.*

Theorem 1.13. *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 ρ_n on X such that, for every n and $i \in I_p$, we have $\rho_n(C_{i,n})$. Then $\rho_n \rightarrow \rho$.*

2

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. The literature in this branch is abundant and we can find theory. We dedicate one chapter to present important results for the finite dimensional part of this branch, since the discrete version of the optimal transport problem can be seen as linear program. The following results are based in [1].

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

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

$$\min \quad \mathbf{c}^\top \mathbf{x} \quad (2.1)$$

$$\text{subject to} \quad \mathbf{Ax} = \mathbf{b} \quad (2.2)$$

$$\mathbf{x} \geq 0 \quad (2.3)$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

Simplex Method.

The idea of the simplex method is to proceed from one basic feasible solution that 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,

$$\begin{array}{cccccccc}
 \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_m & \mathbf{a}_{m+1} & \mathbf{a}_{m+2} & \dots & \mathbf{a}_n & \mathbf{b} \\
 \hline
 1 & 0 & \dots & 0 & y_{1,m+1} & y_{1,m+2} & \dots & y_{1,n} & y_{1,0} \\
 0 & 1 & \dots & 0 & y_{2,m+1} & y_{2,m+2} & \dots & y_{2,n} & y_{2,0} \\
 \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\
 0 & 0 & \dots & 1 & y_{m,m+1} & y_{m,m+2} & \dots & y_{m,n} & y_{m,0}
 \end{array} \tag{2.5}$$

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

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

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

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,

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

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 non-degenerated. This assumption simplifies the description of the simplex method, however the arguments can be modified to include degenerated basic feasible solutions.

Suppose we have the basic feasible solution $\mathbf{x} = (x_1, x_2, \dots, x_m, 0, 0, \dots, 0)$. We are assuming non-degeneracy 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 ϵ ,

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

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,

$$\mathbf{A} = [\mathbf{B}, \mathbf{D}] \quad (2.7)$$

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

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

$$\mathbf{x} = (\mathbf{x}_B, \mathbf{x}_D), \quad \mathbf{c} = (\mathbf{c}_B, \mathbf{c}_D) \quad (2.8)$$

Where \mathbf{x}_B has m entries and \mathbf{x}_D has $n - m$ entries; in similar way for \mathbf{c}_B and \mathbf{c}_D . Then, our primal problem can be written as follows,

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

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,

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

where,

$$z_j = y_{1,j} c_1 + y_{2,j} c_2 + \dots + y_{m,j} c_m, \quad \text{for } m+1 \leq j \leq n. \quad (2.10)$$

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

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

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

Proof. Consider equations (2.9) and (2.10), if $c_j - z_j$ is negative for some j , $m + 1 \leq j \leq n$, then changing x_j from zero to a positive value decreases the total cost z . Let $(x_1, x_2, \dots, x_m, 0, \dots, 0)$ 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 $x'_{m+1} > 0$, substituting this new solution into equation (2.9) 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 2.4 (Optimality Condition Theorem). *If for some basic feasible solution $c_j - z_j \geq 0$ for all j , then that solution is optimal.*

Proof. This result comes from equation (2.9), since any other feasible solution must have $x_i \geq 0$ for all i , and hence the value z of the objective will satisfy $z - z_0 \geq 0$. \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 (2.5) corresponding to a basic feasible solution and compute the relative cost coefficients r_j . For this purpose we can use row reduction.
2. If $r_j \geq 0$ for all j , then the current basic feasible solution is optimal; stop.
3. Select l such that $r_l < 0$ to determine which nonbasic variable is to become basic.
4. Calculate the ratios $y_{i,0}/y_{i,l}$ for $y_{i,l} > 0$, $i = 1, 2, \dots, m$. If no $y_{i,l} > 0$, then problem is unbounded; stop. Otherwise, select k as the index i corresponding to the minimum ratio.
5. Apply the pivoting procedure to introduce \mathbf{a}_l substituting \mathbf{a}_k in the basis. Return to step 1.

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

Revised simplex method.

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

For any \mathbf{x}_D the necessary value of \mathbf{x}_B as follows,

$$\mathbf{x}_B = \mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{D}\mathbf{x}_D$$

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

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

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,

$$\max \quad \lambda^T \mathbf{b} \quad (2.11)$$

$$\text{subject to} \quad \lambda^T \mathbf{A} \leq \mathbf{c} \quad (2.12)$$

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

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

$$\lambda^T \mathbf{b} = \lambda^T (\mathbf{A}\mathbf{x}) \leq \mathbf{c}^T \mathbf{x}$$

□

Corollary 2.5. If \mathbf{x}_0 and λ_0 are feasible for the (2.2) and (2.12) respectively and $\mathbf{c}^T \mathbf{x}_0 = \lambda_0^T \mathbf{b}$, then \mathbf{x}_0 and λ_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 λ_0 for their respective problems, satisfying the equality means that each problem has reached its optimal value.

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

Proof. We see from corollary 2.5 that the first condition holds. If the primal is unbounded and λ is feasible for the dual we must have, $\lambda^T \mathbf{b} \leq -M$ 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

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

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 $\lambda \in \mathbb{R}^m$ 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 $\mathbf{w} = \mathbf{0}$, $\mathbf{w} = \alpha \mathbf{b} - \mathbf{A}\mathbf{x}_0 = \mathbf{0}$ with $\mathbf{x}_0 \geq \mathbf{0}$, then $\mathbf{x} = \mathbf{x}_0/\alpha$ is feasible for the primal problem.

Hence $r/\alpha = z_0 - \mathbf{c}^T \mathbf{x} \leq 0$, implying that $r \leq 0$. For $\alpha = 0$, we have $\mathbf{w} = -\mathbf{A}\mathbf{x}_0 = \mathbf{0}$ with $\mathbf{x}_0 \geq \mathbf{0}$ and $\mathbf{c}^T \mathbf{x}_0 = -1$. 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 $(s, \lambda) \in \mathbb{R}^{m+1}$ and constant c satisfying

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

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

$$-r + \lambda^T \mathbf{w} \geq 0 \quad (2.13)$$

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

$$(\mathbf{c} - \lambda^T \mathbf{A}) - \alpha z_0 + \alpha \lambda^T \mathbf{b} \geq 0 \quad (2.14)$$

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

Complementary Slackness.

Theorem 2.6. Let \mathbf{x} and λ 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

- $x_i > 0 \Rightarrow \lambda^T \mathbf{a}_i = c_i$,
- $x_i = 0 \Leftarrow \lambda^T \mathbf{a}_i < c_i$.

Proof. If the above conditions hold, then $(\lambda^T \mathbf{A} - \mathbf{c}^T) \mathbf{x} = 0$. By the means of the Weak Duality lemma and corollary 2.5, $\lambda^T \mathbf{b} = \mathbf{c}^T \mathbf{x}$ implies two solutions are optimal. Conversely, if the two solutions are optimal, by the means of the Duality Theorem $\lambda^T \mathbf{b} = \mathbf{c}^T \mathbf{x}$. Since each component of \mathbf{x} is nonnegative and each component of $\lambda^T \mathbf{A} - \mathbf{c}^T$ 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, since it is more efficient to work with the dual tableau instead, making use of the complementary slackness conditions to recover the primal solution.

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

$$\lambda^T \mathbf{a}_j = c_j, \quad \text{for } j = 1, \dots, m \quad (2.15)$$

Assuming non-degeneracy there is an inequality,

$$\boldsymbol{\lambda}^\top \mathbf{a}_j < c_j, \quad \text{for } j = m+1, \dots, n \quad (2.16)$$

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

$$\bar{\boldsymbol{\lambda}}^\top = \boldsymbol{\lambda}^\top - \epsilon \boldsymbol{\beta}^i \quad (2.17)$$

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

$$\bar{\boldsymbol{\lambda}}^\top \mathbf{a}_j = c_j \quad \text{for } j = 1, 2, \dots, m \text{ and } i \neq j \quad (2.18)$$

$$\bar{\boldsymbol{\lambda}}^\top \mathbf{a}_j = c_j - \epsilon y_{i,j} \quad \text{for } j = m+1, m+2, \dots, n \quad (2.19)$$

$$\bar{\boldsymbol{\lambda}}^\top \mathbf{a}_i = c_i - \epsilon \quad (2.20)$$

In the other hand,

$$\bar{\boldsymbol{\lambda}}^\top = \boldsymbol{\lambda}^\top \mathbf{b} - \epsilon (\mathbf{x}_\mathbf{B})_i \quad (2.21)$$

The idea behind the algorithm is

1. We start with a dual feasible solution $\mathbf{x}_\mathbf{B}$, if $\mathbf{x}_\mathbf{B} \geq 0$ the solution is optimal. $\mathbf{x}_\mathbf{B}$ is not nonnegative we can find i such that the i -th component of $\mathbf{x}_\mathbf{B}$ is less than zero, i.e. $(\mathbf{x}_\mathbf{B})_i < 0$.
2. If all $y_{i,j} \geq 0$, for $j = 1, 2, \dots, n$, then the dual has no maximum, due to the fact we have feasibility of $\bar{\boldsymbol{\lambda}}$ for any choice of $\epsilon > 0$.
If $y_{i,j} < 0$ for some j , we set

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

3. Form a new basis \mathbf{B} by pivoting \mathbf{a}_i and \mathbf{a}_l . Using this basis determine the new $\mathbf{x}_\mathbf{B}$ and return to Step 1.

Min-Max Theorems

$$\max_{\mathbf{y} \in Y} \left(\min_{\mathbf{x} \in X} (\mathbf{x}^\top \mathbf{A} \mathbf{y} + \mathbf{B} \mathbf{x} + \mathbf{C} \mathbf{y}) \right) = \min_{\mathbf{x} \in X} \left(\max_{\mathbf{y} \in Y} (\mathbf{x}^\top \mathbf{A} \mathbf{y} + \mathbf{B} \mathbf{x} + \mathbf{C} \mathbf{y}) \right) \quad (2.22)$$

3

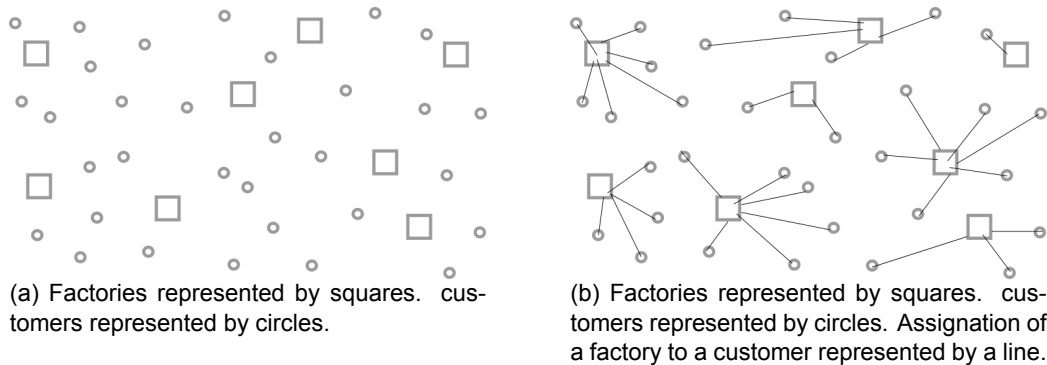
Optimal Transport Theory

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 3.1: 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,

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

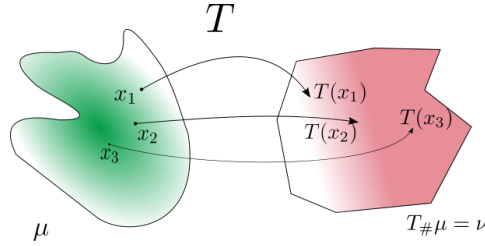
Therefore, we need to search for the optimum in the set of measurable maps $T : X \rightarrow Y$ such that the condition (3) is translated to,

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

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

$$g(T(x)) \det(DT(x)) = f(x) \quad (3.2)$$

Figure 3.2: Monge problem. Finding a map.



The equation (3.2) 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 $f_n : [0, 1] \rightarrow [0, 1]$,

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

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, this is sequence of oscillating function converging weakly¹ to its mean value $f_n \rightarrow f = \frac{1}{2}$, which makes $f_{\#}\mu \neq \mathcal{L}^1 \llcorner [0, 1]$.

¹A sequence of bounded periodic functions converges weakly to its mean value.

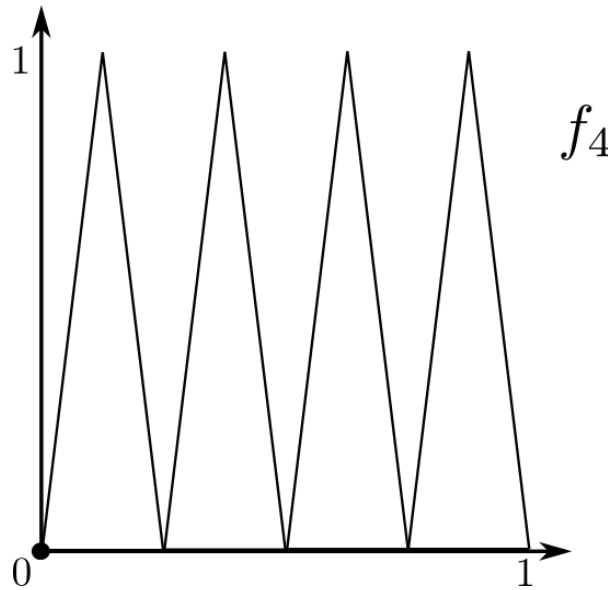


Figure 3.3: f_n constructed using hat functions. The picture shows the case $n = 4$.

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 for 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 the 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 μ and ν . We introduce the following notation to give the necessary formalism to this approach.

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

The above definition is equivalently to say that coupling two measures means to find a probability measure γ , 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 ϕ, ψ 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 3.1 is given for measures on probabilistic spaces, we can rephrase it in terms of stochastic variables. Let (X, μ) and (Y, ν) be two probability spaces. Coupling μ and ν 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 (μ, ν) .

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

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

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

Proof. Take $\gamma = \mu \otimes \nu$. □

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 $\left(\text{proj}_x \right)_\# \gamma = \mu$ and $\left(\text{proj}_y \right)_\# \gamma = \nu$. This inspires a definition for a coupling between two measures generated by a transport map.

Definition 3.2 (Deterministic Coupling). *Let (X, μ) and (Y, ν) 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 μ and ν .*

For the sake of simplicity, we refer as γ_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 μ and ν , 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 μ to ν , for all ν -integrable (nonnegative measurable) function ϕ ,

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

The increasing rearrangement on \mathbb{R} is an example of a coupling between two probability measures over one dimensional euclidean space. Let μ, ν 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,

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

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

Then, we set the map T as $T = G^{-1} \circ F$. If μ 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 μ and ν be two probability measures on \mathbb{R}^n , such that μ 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 μ_1 being atomless. Then define $y_1 = T_1(x_1)$ by the composition of the pseudo-inverse functions of the increasing rearrangement, with F and G considered as they are in (3.5) and (3.6) respectively.
2. Now take the marginal on the first two variables and disintegrate it with respect to the first variable. This gives probability measures $\mu_2(dx_1dx_2) = \mu_1(dx_1)\mu_2(dx_2|x_1)$, $\nu_2(dy_1dy_2) = \nu_1(dy_1)\nu_2(dy_2|y_1)$. 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 3.2 (Gluing lemma). *If \mathcal{Z} is a function of \mathcal{Y} and \mathcal{Y} is a function of \mathcal{X} , then \mathcal{Z} is a function of \mathcal{X} . Let (X_i, μ_i) , $i = 1, 2, 3$, be Polish probability spaces. If (X_1, X_2) is a coupling of (μ_1, μ_2) and (Y_2, Y_3) is a coupling of (μ_2, μ_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

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

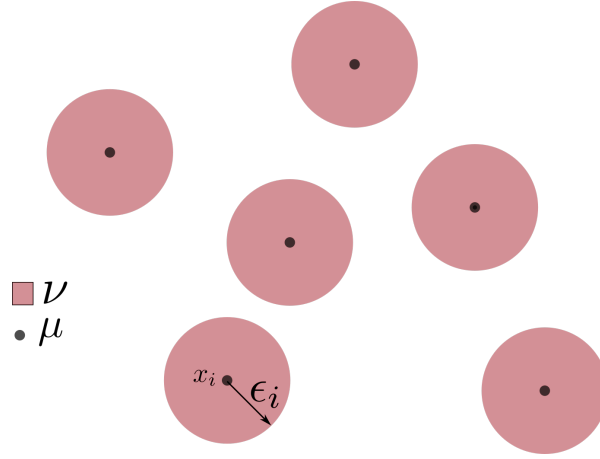
where $\Pi(\mu, \nu)$ is the set of transport plans.

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 $\gamma = \mu \otimes \nu$, as shown in lemma 3.1. By contrast, it is not always possible to find transportation maps (deterministic couplings).

Consider a measure μ 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 δ_{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 ν 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 μ and ν deterministically; since there is no map T , such that $T_{\#}\mu = \nu$.

Figure 3.4: Transportation maps. There is no deterministic coupling for μ and ν , but there is a transportation plan.

Existence of a minimizer for Kantorovich's Problem.

Properties of transportation plans.

Lemma 3.3. *Let X be a metric space. If $f : X \rightarrow \overline{\mathbb{R}}$ 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 μ . 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 that J is lower semicontinuous. \square

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

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

is lower semicontinuous.

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

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

Theorem 3.2. *Let X and Y be compact metric spaces, $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$ and a cost function $c : X \times Y \rightarrow \mathbb{R}$ 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 ν , 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 μ and ν is sequentially compact. \square

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

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

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

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

Proof. \square

Theorem 3.4. *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 the Kantorovich's problem (KP) 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,

$$\begin{aligned} \gamma_n((X \times Y) \setminus (K_X \times K_Y)) &\leq \gamma_n((X \setminus K_X) \times Y) + \gamma_n(X \times (Y \setminus K_Y)) \\ &= \mu(X \setminus K_X) + \nu(Y \setminus K_Y) \\ &= 2\epsilon \end{aligned}$$

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

Properties of Optimal plans

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

Proof. We see immediately that if γ_1 and γ_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

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

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

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

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

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

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

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

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

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

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

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

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

There are situations in which is possible to find a deterministic coupling between two measures, but not an optimal one for a cost function $c : X \times Y \rightarrow \overline{\mathbb{R}}$. 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 $D = \{(0, t)^\top \in \mathbb{R}^2 : \text{for } t \in [0, 1]\}$.

Let D_1 and D_2 be the segments given by,

$$\begin{aligned} D_1 &= \{(-1, t)^\top \in \mathbb{R}^2 : \text{for } t \in [0, 1]\} \\ D_2 &= \{(+1, t)^\top \in \mathbb{R}^2 : \text{for } t \in [0, 1]\} \end{aligned}$$

And we set the measure ν as follows,

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

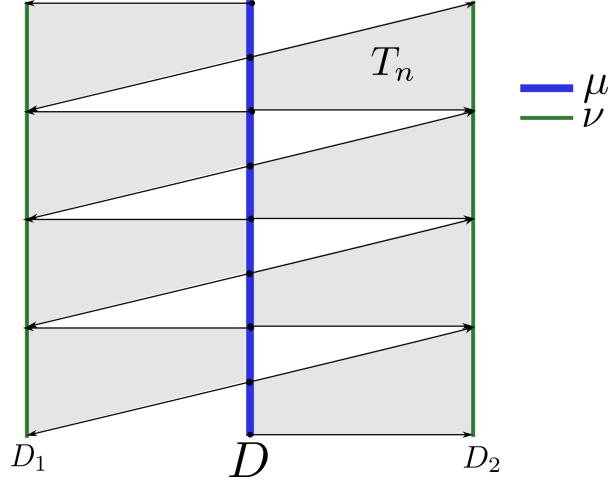


Figure 3.5: There is a deterministic coupling for μ and ν , but no optimal one. The map T_n shown in this picture with $n = 4$.

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 assigns 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 hand 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 $C(T_n)$,

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

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

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

Proof. Since K is continuous, then it is lower semicontinuous, and since we have $K \leq J$, then K is necessarily smaller than the relaxation of J . We only need to prove that, for each γ , 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 $\gamma_{T_n} \rightarrow \gamma$. It is possible to do the last step due to the density of transport plans generated by a map γ_{T_n} in the set of transport plans $\Pi(\mu, \nu)$. \square

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

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

Cyclical Monotonicity and Duality.

Imagine that the consortium changed its policy and it has decided not to be responsible any longer for the transportation of the goods, letting the customers to solve this problem by themselves (assume that the consortium has the monopoly of the goods and the customers have no choice but to adhere to this policy). An entrepreneur feeling that he can ship the goods more efficiently than the consortium did, he intend to buy the goods at the factories and 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.

Duality

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

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

where the supremum is taken among all bounded and continuous functions ϕ, ψ .

Note that the result of this problem is 0 if γ satisfies the constrain of being a probability measure over $X \times Y$ with given marginals μ and ν , and we obtain ∞ if γ does not satisfy this constrain.

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

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

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

Note that,

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

In the other hand, equation (3.10) it is not really useful if we are not able to exchange the min and sup.

For a moment suppose that the conditions that allow to exchange them do exist, we can rewrite the equation (3.10) as follows,

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

If it exists $(x, y) \in X \times Y$ such that $\phi(x) + \psi(y) > c$, we can find measures γ 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),

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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,

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

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

$$\begin{aligned}\phi^{c\bar{c}}(x) &= \inf_{y \in Y} (c(x, y) - \phi^c(y)) \\ &\leq \inf_{y \in Y} (c(x, y) - \xi) \\ &= \xi^{\bar{c}} = \phi\end{aligned}$$

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

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

$$\begin{aligned}\chi^{\bar{c}c}(y) &= \inf_{x \in X} (c(x, y) - \chi^{\bar{c}}(x)) \\ &= \inf_{x \in X} (c(x, y) - \bar{\phi}(x)) \\ &\leq \inf_{x \in X} (c(x, y) - \phi(x)) \\ &= \phi^c.\end{aligned}$$

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

$$\begin{aligned}\phi^{c\bar{c}}(x) &= \inf_{y \in Y} (c(x, y) - \phi^c(y)) \\ &\leq \inf_{y \in Y} (c(x, y) - \chi) \\ &\leq \chi^{\bar{c}} = \bar{\phi}\end{aligned}$$

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

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

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

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

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

Having as result that for any ϕ , $\phi^{c\bar{c}c} = \phi^c$. We can use this fact to prove the theorem 3.8, as we can find it in [5].

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

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

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

Proof. 1. We see that

$$\begin{aligned}c(x, y) &= \phi(x) + c(x, y) - \phi(x) \\ &\geq \phi(x) + \inf_{z \in X} (c(z, y) - \phi(z)) = \phi(x) + \phi^c(y)\end{aligned}$$

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

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

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

Similarly $\forall y \in Y$,

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

□

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

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

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

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

admits a feasible solution (ϕ, ψ) . Moreover $\psi = \phi^c$, then the problem is equivalent to restrict the search to $\phi \in c - \text{conc}(A)$. That is our problem is reduced to,

$$\max \left\{ \int_A \phi d\mu + \int_B \phi^c d\nu : \phi \in c - \text{conc}(A) \right\} \quad (3.16)$$

Proof. □

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

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

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

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

Therefore, at the end if we have,

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

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

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

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

considering $N + 1 = 1$.

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

$$\sum_{i=1}^N c(x_i, y_i) \leq \sum_{i=1}^N c(x_i, y_{\sigma(i)}) \quad (3.17)$$

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

$$\partial^{\bar{c}} \xi = \{(x, y) \in X \times Y; \quad \xi(y) + \xi^{\bar{c}}(x) = c(x, y)\} \quad (3.18)$$

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

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

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

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

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

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

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

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

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

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

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

$$\Gamma \subset \partial^c \phi = \{(x, y) \in X \times Y; \quad \phi(x) + \phi^c(y) = c(x, y)\} \quad (3.19)$$

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

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

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

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

$$0 < k\epsilon < \sum_{i=1}^k c(x_i, y_i) - c(x_i, y_{\sigma(i)}).$$

Note that ϵ satisfying the above condition allows to write the inequality as follows,

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

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

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

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

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

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

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

$$\begin{aligned} (\text{proj}_X)_\# \hat{\gamma} &= \mu - \epsilon_0 \sum_{i=1}^k \mu_i + \epsilon_0 \sum_{i=1}^k \mu_i = \mu \\ (\text{proj}_Y)_\# \hat{\gamma} &= \nu - \epsilon_0 \sum_{i=1}^k \nu_i + \epsilon_0 \sum_{i=1}^k \nu_{\sigma(i)} = \nu \end{aligned}$$

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

$$\begin{aligned} \int c d\gamma - \int c d\hat{\gamma} &= \int c d\gamma - \int c d\gamma + \epsilon_0 \sum_{i=1}^k \int c d\gamma_i - \epsilon_0 \sum_{i=1}^k \int c d\tilde{\gamma}_i \\ &= \epsilon_0 \sum_{i=1}^k \int c d\gamma_i - \epsilon_0 \sum_{i=1}^k \int c d\tilde{\gamma}_i \\ &\geq \epsilon_0 \sum_{i=1}^k \int \left(c(x_i, y_i) - \frac{\epsilon}{2} \right) d\gamma_i - \epsilon_0 \sum_{i=1}^k \int \left(c(x_i, y_{\sigma(i)}) + \frac{\epsilon}{2} \right) d\tilde{\gamma}_i \\ &= \epsilon_0 \left(\sum_{i=1}^k c(x_i, y_i) - \sum_{i=1}^k c(x_i, y_{\sigma(i)}) - k\epsilon \right) > 0 \end{aligned}$$

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

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

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

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

$$\begin{aligned} \phi^c(y) &= \inf_{x \in X} (c(x, y) - \phi(x)) \leq \|c\|_\infty - \phi(x_0) \quad \forall y \in Y \\ \phi(x) &= \phi^{c^c}(x) = \inf_{y \in Y} (c(x, y) - \phi^c(y)) \leq \|c\|_\infty - \phi^c(y_0) \quad \forall x \in X \end{aligned}$$

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

$$\begin{aligned}\phi(x) &= \inf_{y \in Y} (c(x, y) - \phi^c(y)) \\ &\geq \inf_{y \in Y} (c(x, y) - (\|c\|_\infty - \phi(x_0))) \\ &\geq -\|c\|_\infty + \phi(x_0) + \inf_{x \in X} \inf_{y \in Y} c(x, y) \\ &\geq -2\|c\|_\infty + \phi(x_0)\end{aligned}$$

and all $y \in Y$,

$$\begin{aligned}\phi^c(y) &= \inf_{x \in X} (c(x, y) - \phi(x)) \\ &\geq \inf_{x \in X} (c(x, y) - (\|c\|_\infty - \phi^c(y_0))) \\ &\geq -\|c\|_\infty + \phi^c(y_0) + \inf_{y \in Y} \inf_{x \in X} c(x, y) \\ &\geq -2\|c\|_\infty + \phi^c(y_0).\end{aligned}$$

Having upper and lower bounds for ϕ and ϕ^c , we can integrate ϕ and ϕ^c with respect to μ and ν respectively,

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

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

By definition of Δ in the dual problem (DP),

$$\sup(\text{DP}) = \Delta \geq \int_X \phi d\mu + \int_Y \phi^c d\nu = \int_{X \times Y} c d\gamma = \min(\text{KP}) \quad (3.20)$$

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

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

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

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

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

Proving in this way that γ is optimal. \square

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

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

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

and let γ be a solution for the Kantorovich's problem with c as cost function. Then,

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

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

Duality by convex analysis.

We can prove duality introducing a perturbation as follows,

Proposition 3.1. *Suppose X and Y are compact sets and c is continuous. Then there exists a solution (ϕ, ψ) for the*

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

$$H_\gamma(p) = -\max \left\{ \int_X \phi d\mu + \int_Y \psi d\nu : \phi(x) + \psi(y) \leq c(x, y) - p(x, y) \right\} \quad (3.21)$$

Lemma 3.9. *Using the above definition terminology, H_γ is convex.*

Proof. Take $p_0, p_1 \in C(X \times Y)$, □

Brenier's theorem

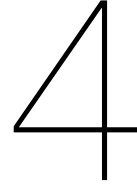
Theorem 3.16 (Brenier's theorem). *Let μ, ν probability measures on \mathbb{R}^d . Assume that μ is absolutely continuous with respect the Lebesgue measure, $\mu \ll \mathcal{L}^d$. Therefore,*

1. *There exists a unique transport plan.*
2. *The optimal transport plan is induced by a map.*
3. *The map is the gradient of a convex function.*

Theorem 3.17. *Let μ, ν be probabilities over \mathbb{R}^d and $c(x, y) = \frac{1}{2} |x - y|^2$. Suppose $\int |x|^2 dx < \infty$ and $\int |y| dy < \infty$. Consider*

Theorem 3.18 (Distance between two Gaussians). *The optimal transport between two gaussians for cost function $c(x, y) = |x - y|^2$ is given by a translation map.*

Proof. □



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 μ a probability measure defined over X ,

$$\mu = \sum_{i=1}^n a_i \delta_{x_i}, \quad (4.1)$$

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 ν a probability measure defined over Y ,

$$\nu = \sum_{i=1}^m b_i \delta_{y_i}, \quad (4.2)$$

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 $(\gamma)_{i,j} = \gamma_{i,j}$ be a joint probability distribution with marginals given by μ and ν . That is,

$$\sum_{j=1}^m \gamma_{i,j} = a_i \quad (4.3)$$

$$\sum_{i=1}^n \gamma_{i,j} = b_j \quad (4.4)$$

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

$$(\mathbf{C})_{i,j} = c_{i,j} = c(x_i, y_j) \quad (4.5)$$

Given the matrix nature of the optimal transport problem for discrete measures we can compute the total cost by the Frobenius inner product¹ of the two matrices, since \mathbf{C} and γ have the same dimensions.

$$\langle \mathbf{C}, \gamma \rangle = \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}} c_{i,j} \gamma_{i,j} = \text{tr}(\mathbf{C}^\top \gamma) \quad (4.6)$$

¹For matrices with real entries the product is defined as $\langle A, B \rangle = \text{tr}(A^\top B)$.

Simplex Method Algorithm and Duality.

We can solve the problem using the simplex method.

$$\Theta(\mu, \nu) = \min_{\gamma \in \Pi(\mu, \nu)} \langle \mathbf{C}, \gamma \rangle \quad (4.7)$$

The simplex method is Not polynomial time.

Consider the following example

$$\max \sum_{j=1}^n 10^{n-j} x_j \quad (4.8)$$

$$\text{subject to } 2 \sum_{j=1}^{i-1} 10^{i-j} x_j + x_i \leq \quad (4.9)$$

Sinkhorn-Knopp Algorithm.

Consider the problem adding a regularization.

$$\Theta_\epsilon(\mu, \nu) = \min_{\gamma \in \Pi(\mu, \nu)} \langle \mathbf{C}, \gamma \rangle + \epsilon H(\gamma) \quad (4.10)$$

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

$$H(\gamma) = - \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}} \gamma_{i,j} \log(\gamma_{i,j}) \quad (4.11)$$

Theorem 4.1. *The Linear programming program after adding the regularization, becomes a strictly convex program.*

Theorem 4.2. $\Theta_\epsilon(\mu, \nu)$ converges to the solution with maximum entropy as $\epsilon \rightarrow 0$

Continuous Formulation.

Beckman Problem and Optimal Transport.

Proximal Splitting Algorithms.

5

Applications

The applications of the optimal transport are many. We can find researches using this formulation in many fields.

Wasserstein's Distances.

Optimal couplings can be used to metrize the space of probability measures.

Statistical Distances

Data Assimilation of a Dynamic.

Isoperimetric Inequality.

Bibliography

- [1] Yinyu Ye David G. Luenberger. *Linear and Nonlinear Programming*. Springer, 2007.
- [2] L. Kantorovich. On the translocation of masses. *Dokl. Acad. Nauk. USSR*, 37(7-8):227–229, 1942.
- [3] Milan Merkle. Topics in weak convergence of probability measures. *Zb. radova Mat. Inst. Beograd*, 9(17):235–274, 2000.
- [4] J.R. Munkres. *Introduction to Topology*. Topology. Prentice-Hall, 2nd. edition, 2000. ISBN 9780131784499.
- [5] Cédric Villani. *Optimal Transport: Old and New*, volume 338 of Grundlehren der mathematischen Wissenschaften. Springer Science and Business Media, 2008.