

CHAPTER 3 NUMERICAL OPTIMAL TRANSPORT

OPTIMAL TRANSPORT studies the cost of moving quantities from one place to another and aims at finding the optimal way to do it, that is minimizing the cost of the displacement. Its applications can go from moving heaps of sand to holes (Monge, 1781), reorganizing military troops and cargo (Kantorovich, 1942), correcting an image histogram to a prescribed values (Haker et al., 2004), finding the origin of seismic events (Métivier et al., 2016), or transferring a learning model over a new data distribution (Courty et al., 2016). Its study dates back to Monge in 1781, had a renewed interest in the beginning of the XX-th century, and is still today a very active field of mathematics both pure (Villani, 2008) and applied (Santambrogio, 2015), notably in relation with machine learning (Peyré and Cuturi, 2019). The book *Computational optimal transport* by Gabriel Peyré and Marco Cuturi will serve as the main reference for this chapter, and is to be credited for some of the illustrations of this chapter.

3.1 INTRODUCTION

3.1.1 Measuring the mass

Let us consider a set X . To measure the mass, it is convenient to define a positive (so-called Radon) measure μ on X which associates at each point $x \in X$ a positive number $\mu(x)$.

Example 3.1 (Measure in continuous and discrete spaces).

3.1.2 Transporting the mass

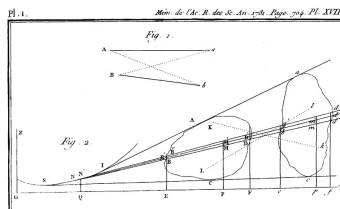
Example 3.2 (Original Monge problem).

M É M O I R E
S U R L A
T H É O R I E D E S D É B L A I S
E T D E S R E M B L A I S.

Par M. M O N G E.

Lorsqu'on doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.

Le prix du transport d'une molécule étant, toutes choses d'ailleurs égales, proportionnel à son poids & à l'espace qu'on lui fait parcourir, & par conséquent le prix du transport total devant être proportionnel à la somme des produits des molécules multipliées chacune par l'espace parcouru, il s'enluit que le déblai & le remblai étant donnés de figure & de position, il n'est pas indifférent que telle molécule du déblai soit transportée dans tel ou tel autre endroit du remblai, mais qu'il y a une certaine distribution à faire des molécules du premier dans le second, d'après laquelle la somme de ces produits fera la moindre possible, & le prix du transport total fera un *minimum*.



Let us define a starting set X and a target set Y , endowed with measures μ and ν . A transport operation is a *mapping* from X to Y

$$T : X \rightarrow Y$$

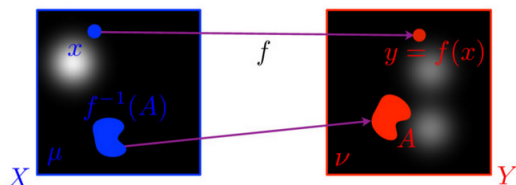
implying that $\nu(A) = \mu(T^{-1}(A))$ for all $A \subset Y$

We note $\nu = T_{\#}\mu$ with $T_{\#}$ called the *push-forward* operator.

Example 3.3.

$$\nu = f_{\#}\mu \text{ defined by: } \nu(A) \stackrel{\text{def.}}{=} \mu(f^{-1}(A))$$

$$\iff \int_Y g(y) d\nu(y) \stackrel{\text{def.}}{=} \int_X g(f(x)) d\mu(x)$$



But, intuitively, all transportation maps are not equivalent, we thus need define the *cost* $c(x, y)$ of moving (a unit of mass) from $x \in X$ to $y \in Y$.

With these definitions, we can formulate the *Monge problem* of minimizing the transportation cost:

$$\min_{T: \nu = T_{\#}\mu} \int_X c(x, T(x)) d\mu(x) \quad (\text{Monge problem})$$

We know from Brenier (Brenier, 1991) that this problem has a unique solution when $c(x, y) = \|x - y\|^2$ and μ, ν have densities. Furthermore, the optimal transport plan T^* is the gradient of a convex function.

3.1.3 The discrete Monge problem

Let us denote a discrete measure $\alpha = \sum_{i=1}^n a_i \delta_{x_i}$ as a sum of diracs at positions (x_i) weighted by non-negative coefficients (a_i) .

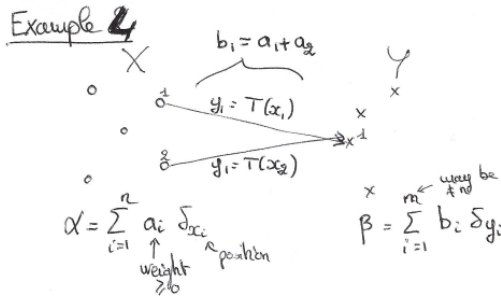
The problem of transporting $\alpha = \sum_{i=1}^n a_i \delta_{x_i}$ to $\beta = \sum_{j=1}^m b_j \delta_{y_j}$ amounts to finding a *Monge transport map* T that associates to each point x_i , a single point y_j so that

$$\text{for all } j \in \{1, \dots, m\}, \quad b_j = \sum_{i: y_j = T(x_i)} a_i.$$

This equation, sometimes called *mass transportation*, defines the set of valid transport maps from α to β by imposing that the mass of a target point y_j (i.e. b_j) is equal to the mass that is transported from all x_i such that $y_j = T(x_i)$.

An important point is that in this problem, the mass of point x_i cannot be split: even though two input points can go to the same target point, the mass a_i of an input point cannot be split into several target points¹³. This means that *there may not exist a Monge transport plan*.

¹³ following this remark, we can say that the discrete Monge problem is actually an *assignment problem*



When Monge transport maps exist, it is possible to evaluate their cost defined as the sum of the costs of transport for all input point, that is $\sum_{i=1}^n c(x_i, T(x_i))$. The associated optimal transport problem thus writes:

$$\min_T \sum_{i=1}^n c(x_i, T(x_i)) \quad \text{s.t. } \forall j, \quad b_j = \sum_{i: y_j = T(x_i)} a_i.$$

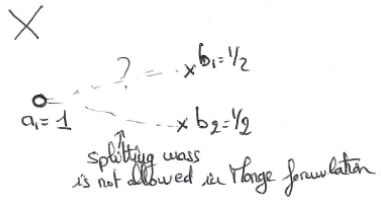
In this case, we notice that the transport plan T can be simply rewritten as an $n \times m$ matrix T with $T_{ij} = 1$ if $T(x_i) = y_j$ and 0 elsewhere; we can also define a cost matrix C as $C_{ij} = c(x_i, y_j)$. Then:

- A transport matrix T must have i) exactly one 1 per row (all others coefficients are null); and ii) verify the mass transportation equality which rewrites $b_j = \sum_{i=1}^n T_{ij} a_i$;
- The transport cost is equal to $\langle T; C \rangle$ where $\langle A; B \rangle = \sum_{i,j} A_{ij} B_{ij}$ is called the Frobenius scalar product.

It is thus a highly combinatorial problem (maybe with no solution).

3.1.4 Kantorovitch's relaxation

Example 5



Monge’s problem may not have a feasible solution due to the impossibility of splitting mass. Allowing such a mass splitting is at the core of Kantorovitch’s relaxation (Kantorovich, 1942). Instead of considering a mapping transport matrix T as in Monge problem (see above), we consider a coupling matrix P where $P_{ij} \geq 0$ represents the quantity *NOT* the proportion of mass going from x_i to y_j . In order for the transport to be valid, one must have for all j that $a_j = \sum_{i=1}^m P_{ji}$ and $b_j = \sum_{i=1}^n P_{ij}$.

Monge transport	Kantorovitch relaxation
T is a surjective mapping $X \rightarrow Y$	T is a coupling matrix $P \in \mathbb{R}_+^{n \times m}$
$\forall j, b_j = \sum_{i=1}^n T_{ij} a_i$	$\forall j, a_j = \sum_{i=1}^m P_{ji}$ and $b_j = \sum_{i=1}^n P_{ij}$

In order to work properly on such transport couplings, it is interesting to define the set of *admissible couplings*

$$U(a, b) = \{P \in \mathbb{R}_+^{n \times m} : P \mathbf{1}_m = a, \mathbf{1}_n^T P = b^T\}$$

where $\mathbf{1}_d$ is the size- d vector with unit unit entries.

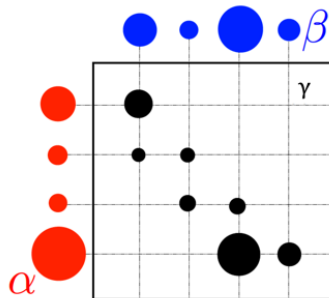
Lemma 3.4. For any pair of probability vectors $a \in \Delta_n, b \in \Delta_m, U(a, b)$ is a convex non-empty linear polytope.

Proof. As an exercise. □

Using a cost matrix C (defined as above as $C_{ij} = c(x_i, y_j)$) and the Frobenius scalar product, Kantorovitch’s optimal transport problem writes

$$\min_{P \in U(a, b)} \langle C; P \rangle$$

which is a linear program!



Remark 3.5 (Continuous version). For two measures α, β over \mathcal{X}, \mathcal{Y} , Kantorovitch's optimal transport problem writes

$$\min_{\gamma \in \Gamma(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\gamma(x, y)$$

where $\Gamma(\alpha, \beta)$ is the set of measures on $\mathcal{X} \times \mathcal{Y}$ admitting α and β as marginals. \blacktriangleleft

3.2 COMPUTING THE OPTIMAL TRANSPORT

In this section, we will be looking into the numerical computation of Kantorovitch's discrete optimal transport problem:

$$\min_{P \in \mathcal{U}(a, b)} \langle \mathbf{C}; P \rangle \quad (\mathcal{K})$$

3.2.1 Primal problem

Since $\mathcal{U}(a, b)$ is a convex non-empty linear polytope (See Lemma 3.4) and (\mathcal{K}) is a linear program, we have some information about the localization of the solutions.

Theorem 3.6. *There is a solution P^* of (\mathcal{K}) which is an extremal point of $\mathcal{U}(a, b)$.*

Proof. $\mathcal{U}(a, b)$ is a non-empty convex polytope; thus the solution of a Linear Program on such a set is necessary on the boundary by the maximum principle (see e.g. Chap. 32 in (Rockafellar, 1970)). \square

In terms of optimization:

- Kantorovitch's problem and Dantzig's simplex algorithm are concomitant;
- Direct LP may be hard due to the polytope constraints;
- When $m = n$ and $a = b = \mathbf{1}/n$, the Hungarian/Auction algorithm is in $\mathcal{O}(n^3)$;
- In 1D, sorting is in $\mathcal{O}(n \log(n))$.

3.2.2 Dual Problem

Let us dualize of the linear program (\mathcal{K}) :

$$\begin{aligned} & \min_{P \in \mathcal{U}(a, b)} \langle \mathbf{C}; P \rangle && (\mathcal{K}) \\ \Leftrightarrow & \min_{P \in \mathbb{R}_+^{n \times m}, P \mathbf{1}_m = a, \mathbf{1}_n^T P = b^T} \langle \mathbf{C}; P \rangle \\ \text{(Lagrange)} \Leftrightarrow & \min_{P \in \mathbb{R}_+^{n \times m}} \max_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle \mathbf{C}; P \rangle - \langle f; P \mathbf{1}_m - a \rangle - \langle g; \mathbf{1}_n^T P - b^T \rangle \\ \text{(Strong duality)} \Leftrightarrow & \max_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \min_{P \in \mathbb{R}_+^{n \times m}} \langle \mathbf{C}; P \rangle - \langle f; P \mathbf{1}_m - a \rangle - \langle g; \mathbf{1}_n^T P - b^T \rangle \\ \Leftrightarrow & \max_{f \in \mathbb{R}^n, g \in \mathbb{R}^m} \langle f; a \rangle + \langle g; b \rangle + \min_{P \in \mathbb{R}_+^{n \times m}} \langle \mathbf{C} - f \mathbf{1}^T - \mathbf{1} g^T; P \rangle. \end{aligned}$$

Since $P \in \mathbb{R}_+^{n \times m}$, the solution of the right part is attained if and only if

$$\mathbf{C} - f \mathbf{1}^T - \mathbf{1} g^T \geq 0$$

where the inequality is meant elementwise.

In this case, $\langle \mathbf{C} - f\mathbf{1}^\top - \mathbf{1}g^\top; P^\star \rangle = 0$ and we have

$$\begin{aligned} & \min_{P \in \mathcal{U}(a,b)} \langle \mathbf{C}; P \rangle & (\mathcal{K}) \\ \Leftrightarrow & \max_{f \in \mathbb{R}^n, g \in \mathbb{R}^m; f\mathbf{1}^\top + \mathbf{1}g^\top \leq \mathbf{C}} \langle f; a \rangle + \langle g; b \rangle. & (\mathcal{D}) \end{aligned}$$

Remark 3.7 (Interpretation). Consider m warehouses producing a and n factories needing b .

Primal: Find P^\star and pay $\langle \mathbf{C}; P^\star \rangle$ to transport.

Dual: Find f^\star, g^\star , f_i^\star is the price to take resource from warehouse i , g_j^\star is the price to deliver resource at factory j , thus the price is $\langle f; a \rangle$ (to take) + $\langle g; b \rangle$ (to deliver). \blacktriangleleft

Remark 3.8 (Complementary Slackness). $\langle \mathbf{C} - f^\star\mathbf{1}^\top - \mathbf{1}g^{\star\top}; P^\star \rangle = 0$ and thus for all (i, j)

$$\begin{cases} \text{either } P_{ij}^\star > 0 & \text{and } f_i^\star + g_j^\star = C_{ij} \\ \text{or } P_{ij}^\star = 0 & \text{and } f_i^\star + g_j^\star < C_{ij} \end{cases}$$

\blacktriangleleft

3.2.3 Associated Metric

The cost of moving from a distribution to another distribution naturally defines a distance between them when they are defined on the same space.

¹⁴that is: **Proposition 3.9.** Let $n = m$. Take $p \geq 1$ and let $\mathbf{C} = D^p$ where D defines a distance¹⁴ on $\{1, \dots, n\}$. Then,

i) D is symmetric;

ii) $D_{ij} = 0$ if and only if $i = j$;

iii) $D_{ik} \leq D_{ij} + D_{jk}$

$$W_p^p(a, b) := \min_{P \in \mathcal{U}(a,b)} \langle D^p; P \rangle$$

defines the (p -th power of the) p -Wasserstein distance on the simplex of size n .

$W_p(a, b)$ is a distance (without the power p) and thus for all $a, b, c \in \Delta_n$, $W_p(a, b) = 0$ if and only if $a = b$, $W_p(a, c) \leq W_p(a, b) + W_p(b, c)$.

Applications:

- bag of words distance for text classification;
- histogram distance.

3.3 ENTROPIC REGULARIZATION

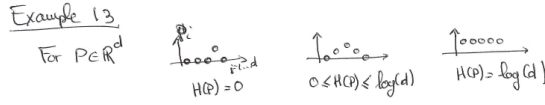
The problems we just saw are typically hard to compute numerically. There was a renewed interest towards these problems (especially in machine learning) following the introduction of an entropy-based regularization leading to more efficient computations.

Example 3.10 (Regularization leads to more stable solutions).

3.3.1 Entropy

The entropy function for a matrix $P \in \mathbb{R}_+^{m \times n}$ writes

$$H(P) = - \sum_{i,j} P_{ij} (\log(P_{ij}) - 1).$$



The *entropy-regularized* optimal transport problem (Cuturi, 2013; Wilson, 1969) then writes for some $\varepsilon > 0$

$$\min_{P \in \mathcal{U}(a,b)} \langle \mathbf{C}; P \rangle - \varepsilon H(P) \quad (\mathcal{P}_\varepsilon)$$

and promotes more “uniform/smoothed” transport maps. This means that now every point is transported to every other point (with potentially very small values), which allows the transport plan to vary smoothly whenever the weights or the cost is evolving, which is very interesting in practice.

3.3.2 Regularized Transport

Proposition 3.11. *The problem $(\mathcal{P}_\varepsilon)$ has a unique solution P_ε^* which verifies*

- $P_\varepsilon^* \xrightarrow{\varepsilon \rightarrow 0} \operatorname{argmin}_{P \text{ sol. of } (\mathcal{K})} -H(P)$
- $P_\varepsilon^* \xrightarrow{\varepsilon \rightarrow +\infty} ab^\top$

$$\begin{aligned} P_\varepsilon^* &= \operatorname{argmin}_{P \in \mathcal{U}(a,b)} \langle \mathbf{C}; P \rangle - \varepsilon H(P) \\ &= \operatorname{argmin}_{P \in \mathcal{U}(a,b)} \langle \mathbf{C}; P \rangle - \varepsilon \sum_{i,j} P_{ij} \log(P_{ij}) - \varepsilon \underbrace{\sum_{i,j} P_{ij}}_{\text{constant in } \mathcal{U}(a,b)} \\ &= \operatorname{argmin}_{P \in \mathcal{U}(a,b)} -\varepsilon \sum_{i,j} P_{ij} \frac{1}{\varepsilon} \log(\exp(-\mathbf{C}_{ij})) - \varepsilon \sum_{i,j} P_{ij} \log(P_{ij}) \\ &= \operatorname{argmin}_{P \in \mathcal{U}(a,b)} \sum_{i,j} P_{ij} \log\left(\frac{P_{ij}}{K_{ij}}\right) \quad \text{with } K_{ij} = \exp(-\mathbf{C}_{ij}/\varepsilon) \text{ called the Gibbs Kernel} \\ &= \operatorname{argmin}_{P \in \mathcal{U}(a,b)} \operatorname{KL}(P|K) \quad \text{with KL called the Kullback-Liebler divergence} \end{aligned}$$

3.3.3 Computational Interest

Proposition 3.12. *The problem $(\mathcal{P}_\varepsilon)$ has a unique solution P_ε^* and this solution writes*

$$P_{ij,\varepsilon}^* = u_i K_{ij} v_j$$

with $K_{ij} = \exp(-\mathbf{C}_{ij}/\varepsilon)$ called the Gibbs Kernel and two unknown vectors u, v .

Proof. The solution is unique since the entropy is strictly concave.

Dualizing the constraints as in Section 3.2.2, the optimal P is obtained as the minimum of $\langle \mathbf{C} - f\mathbf{1}^\top - \mathbf{1}g^\top; P \rangle - \varepsilon H(p)$. Taking the first order optimality conditions, we obtain that for all i, j

$$\begin{aligned} \mathbf{C}_{ij} - f_i - g_j + \varepsilon \log(P_{ij,\varepsilon}^*) &= 0 \\ \Leftrightarrow P_{ij,\varepsilon}^* &= \underbrace{\exp(f_i/\varepsilon)}_{:=u_i} \underbrace{\exp(-\mathbf{C}_{ij}/\varepsilon)}_{:=K_{ij}} \underbrace{\exp(g_j/\varepsilon)}_{:=v_j} \end{aligned}$$

or, rewriting that in matrix form

$$P_\varepsilon^\star = \text{diag}(u)K \text{diag}(v).$$

□

¹⁵They depend on f and g which are the solutions to the dual problem, so no computational gain there.

Unfortunately, u and v are not explicit¹⁵ but since $P_\varepsilon^\star \in \mathcal{U}(a, b)$ we have

$$\begin{aligned} P_\varepsilon^\star \mathbf{1} &= \text{diag}(u)K \text{diag}(v) \mathbf{1} = \text{diag}(u)Kv = u \odot Kv = a \\ \text{and } \mathbf{1}^\top P_\varepsilon^\star &= \mathbf{1}^\top \text{diag}(u)K \text{diag}(v) = u^\top K \text{diag}(v) = (K^\top u \odot v)^\top = b^\top \end{aligned}$$

where \odot represents the Hadamard (entrywise) product.

Thus, we have to scale the matrix K to prescribed row and column sums, ie to get

$$\begin{cases} u \odot Kv = a \\ v \odot K^\top v = b \end{cases} .$$

Sinkhorn's algorithm solves this problem by alternating

$$u_{k+1} = \frac{a}{Kv_k} \quad v_{k+1} = \frac{b}{K^\top u_{k+1}}$$

BIBLIOGRAPHY

- Yann Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Communications on pure and applied mathematics*, 44(4):375–417, 1991.
- Laurent Condat. Fast projection onto the simplex and the ℓ_1 ball. *Mathematical Programming*, 158(1):575–585, 2016.
- Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2016.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pages 2292–2300, 2013.
- Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional variational inequalities and complementarity problems*. Springer Science & Business Media, 2007.
- Steven Haker, Lei Zhu, Allen Tannenbaum, and Sigurd Angenent. Optimal mass transport for registration and warping. *International Journal of computer vision*, 60(3):225–240, 2004.
- Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- Leonid V Kantorovich. On the translocation of masses. In *Dokl. Akad. Nauk. USSR (NS)*, volume 37, pages 199–201, 1942.
- Galina M Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756, 1976.
- Casimir Kuratowski. Sur le probleme des courbes gauches en topologie. *Fundamenta mathematicae*, 15(1):271–283, 1930.
- Ludovic Métivier, Romain Brossier, Quentin Mérigot, Edouard Oudet, and Jean Virieux. Measuring the misfit between seismograms using an optimal transport distance: Application to full waveform inversion. *Geophysical Supplements to the Monthly Notices of the Royal Astronomical Society*, 205(1):345–377, 2016.
- Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Histoire de l'Académie Royale des Sciences de Paris*, 1781.
- Martin J Osborne and Ariel Rubinstein. *A course in game theory*. MIT press, 1994.
- Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.

- R Tyrrell Rockafellar. *Convex analysis*. Number 28. Princeton university press, 1970.
- Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94, 2015.
- Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- John von Neumann. Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100 (1):295–320, 1928.
- Alan Geoffrey Wilson. The use of entropy maximising models, in the theory of trip distribution, mode split and route split. *Journal of transport economics and policy*, pages 108–126, 1969.