

# Lesson 3: Numerical aspects of optimal transport

I

## I - Reminders / Crash course

1/  $\triangle!$  \* This section is intended as a quick reminder of some elements from the class of Max Fathi.

\* In order to simplify exposition, we discard all growth / integrability hypothesis.

Let  $\Omega, \Omega' \subset \mathbb{R}^d$  be measurable subsets

$$\begin{aligned} \mathcal{M}_1(\Omega) &:= \{ \text{Probability measures on } \Omega \} \\ &= \{ \mathcal{L}(X) \mid X \text{ is an } \Omega\text{-valued r.v.} \} \end{aligned}$$

Given  $(\mu, \nu) \in \mathcal{M}_1(\Omega) \times \mathcal{M}_1(\Omega')$ , then set

$$\Pi(\mu, \nu) := \{ \pi \in \mathcal{M}_1(\Omega \times \Omega') \mid (\pi_1)_* \pi = \mu; (\pi_2)_* \pi = \nu \}$$

$$\simeq \left\{ \mathcal{L}(X, Y) \mid \begin{array}{l} (X, Y) \text{ r.v. on } \Omega \times \Omega' \\ \mathcal{L}(X) = \mu \\ \mathcal{L}(Y) = \nu \end{array} \right\}$$

"Transport plans" (PDE)  
or "Couplings" (P)

$c: \Omega \times \Omega' \rightarrow \mathbb{R}$  a cost function

Optimal transport (stated in the form of  $\Pi$  the Kantorovich problem) is the study of the optimization problem

$$\inf_{\pi \in \Pi(\mu, \nu)} \int_{\Omega \times \Omega'} c(x, y) \pi(dx, dy)$$

Thm [Brenier]

If  $\left\{ \begin{array}{l} \mu \text{ has a density } \mu(dx) = dx \rho(x) \\ \mu \text{ \& } \nu \text{ have } n \text{ atoms, each of size } 1/n \end{array} \right.$

Then there exists a unique transport plan  $\pi$  of the form

$$\pi = (\text{id}, T)_\# \mu \in \Pi(\mu, \nu)$$

where  $T: \Omega \rightarrow \Omega'$  is called the Brenier map

As such, in these favorable cases:

$$\begin{aligned} \inf_{\pi \in \Pi(\mu, \nu)} \int_{\Omega \times \Omega'} c(x, y) \pi(dx, dy) \\ = \inf_{\substack{T: \Omega \rightarrow \Omega' \\ T_\# \mu = \nu}} \int_{\Omega} c(x, T(x)) \mu(dx) \end{aligned}$$

This form is the Monge problem

(older & less general).

Examples:

1/  $\Omega = \Omega' = \mathbb{R}$

Then  $T = F_\nu^{\leftarrow} \circ F_\mu$  (Left inverse)

where  $\begin{cases} F_\nu^{\leftarrow}(x) := \inf \{t \geq 0 \mid F_\nu(x) \geq t\} \\ F_\mu(t) = \mu([-\infty, t]) \end{cases}$

2/ If  $\begin{cases} \mu = \frac{1}{n} \sum_{i=1}^m \delta_{x_i} \\ \nu = \frac{1}{n} \sum_{j=1}^m \delta_{y_j} \end{cases}$

Then the Brenier map identifies to a permutation

$T \simeq \sigma: \{1, 2, \dots, m\} \rightarrow \{1, 2, \dots, n\}$

via  $T(x_i) = y_{\sigma(i)}$ . Also, upon writing

$C_{ij} = C(x_i, y_j)$ , the transport problem

reduces to

$\inf_{\sigma \in \mathcal{P}_m} \frac{1}{n} \sum_{i=1}^m C_{i, \sigma(i)}$

The Assignment/Matching problem.

This is the most relevant example in ML: Computers work with discrete masses!

I.2 / Topology & Geometry:  
Geodesics

Here  $c(x, y) = |x - y|^p$

$$W_p(\mu, \nu) = \left( \inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) \pi(dx, dy) \right)^{1/p}$$

Thm: If  $\{\Omega, \Omega\}$  compact then  
 $\{p \geq 1\}$

$W_p$  is a distance over  $\mathcal{M}_1(\Omega)$   
 & the induced topology is equivalent to  
 the topology of weak convergence

Better than that,  $(\mathcal{M}_1(\Omega), W_p)$  is an  
 arc length space: curves have a good  
 notion of speed & geodesics can be defined

Thm [McCann interpolation]

If  $\mu, \nu$  belong to  $(\mathcal{M}_1(\Omega), W_p)$ , then ~~there~~

- there is a geodesic curve between them

$$\begin{cases} \mu_t : [0, \frac{1}{p}] \rightarrow \mathcal{M}_1(\Omega) \\ \mu_0 = \mu; \mu_1 = \nu \end{cases}$$

- when  $\mu_t$  is parametrized by arc length

&  $\pi$  optimal coupling then

$$\mu_t = (\mu_t)_* \pi \quad \& \quad \mu_t : \Omega \times \Omega \rightarrow \Omega^{\#}$$

(Linear interpolation)  $(x, y) \mapsto x + t(y-x)$

In particular, if  $\pi$  is given by a Brenier map,  $\underline{V}$

$$\mu_t = \left[ (1-s) \text{id} + sT \right]_* \mu$$

And in the even more particular case of the assignment problem:

$$\mu = \frac{1}{n} \sum_i \delta_{x_i}$$

$$\nu = \frac{1}{m} \sum_j \delta_{y_j}$$

Then  $\mu_t = \frac{1}{n} \sum_{i=1}^m \delta_{\alpha_i(t)}$  where  $\alpha_i(t) = (1-t)x_i + tT(x_i)$   
 $= (1-t)x_i + t(y_{\sigma(i)})$

↑ In the tutorial, there will be a drawing of that!

### I-3/ Duality

As a particular case of convex duality, which transforms a primal problem to its convex dual.

Thm: Kantorovich Duality:

$$\inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) \pi(dx, dy) = \sup_{\substack{f, g \\ f(x) + g(y) \leq c(x, y)}} \int f d\mu + \int g d\nu$$

In particular for the assignment problem:

$$\inf_{\sigma \in \mathcal{A}_m} \sum_{i=1}^m c_{i, \sigma(i)} = \sup_{f_i + g_j \leq c_{ij}} \sum_{i=1}^m f_i + \sum_{j=1}^m g_j$$

## II The Hungarian algorithm.

VI

If the complexity of testing all combinations is  $O(m!)$ , the above duality is the basis of the so-called  $\left\{ \begin{array}{l} \text{Hungarian algorithm} \\ \text{Munkres algorithm} \end{array} \right.$  whose complexity is  $O(m^3)$ .

The name "Hungarian algorithm" is in honor of Dénes König & Jenő Egerváry; as the validity of the ~~the~~ algorithm relies on graph theory due to them.

### Description of the algorithm:

Input: The Cost Matrix  
 $C = (C_{ij})_{1 \leq i, j \leq m}$

Ex:  $C_4 = \begin{pmatrix} 5 & 9 & \textcircled{1} \\ 10 & 3 & \textcircled{2} \\ 8 & 7 & \textcircled{4} \end{pmatrix}$

Step 1: Remove smallest entry from each row.

• Same for columns

$\rightarrow$  Each row & column contain a zero.

$$\begin{pmatrix} \textcircled{4} & 8 & 0 \\ 8 & \textcircled{1} & 0 \\ \textcircled{4} & 3 & 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 & 7 & 0 \\ 4 & 0 & 0 \\ \textcircled{0} & 2 & 0 \end{pmatrix}$$

### Step 2:

• Draw lines through row & columns that have 0 entries with ~~the~~ fewest possible lines

$$\begin{pmatrix} 0 & 7 & 0 \\ 4 & 0 & 0 \\ 0 & 2 & 0 \end{pmatrix} \text{ Here } 3 \text{ is smallest}$$

Step 3:

VII

If  $n$  lines are drawn, then we can find an assignment of 0.

DONE

here for  $C_1$   
 $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$  or  $\begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}$   
 work

If  $k < n$  lines are drawn,

call  $V$  the indices  $(i,j)$  covered at least by column or line  
 $\Lambda$  the indices  $(i,j)$  covered by both column & line.

$$\Delta = \min_{(i,j) \notin V} c_{ij}$$

Transform  $c_{ij} \leftarrow \begin{cases} c_{ij} - \Delta, & (i,j) \notin V \\ c_{ij}, & (i,j) \in V, (i,j) \notin \Lambda \\ c_{ij} + \Delta, & (i,j) \in \Lambda \end{cases}$   
 Loop to step 2.

Example:

$$C_2 = \begin{pmatrix} 1 & 4 & 5 \\ 5 & 7 & 6 \\ 5 & 8 & 8 \end{pmatrix}$$

Step 1  $\rightsquigarrow$   $\begin{pmatrix} 0 & 3 & 4 \\ 0 & 2 & 1 \\ 0 & 3 & 3 \end{pmatrix}$

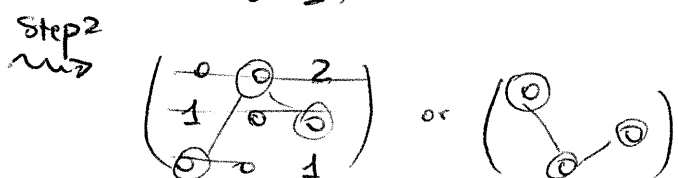
Step 1  $\rightsquigarrow$   $\begin{pmatrix} 0 & 1 & 3 \\ 0 & 0 & 0 \\ 0 & 1 & 2 \end{pmatrix}$

Step 2  $\rightsquigarrow$   $\begin{pmatrix} 0 & 1 & 3 \\ 0 & 0 & 0 \\ 0 & 1 & 2 \end{pmatrix}$

Step 3  $\rightsquigarrow$   $\begin{pmatrix} 0 & 0 & 2 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \Delta = 1$

Cost = 15

= 6 + 4 + 5  
 = 6 + 8 + 1



### III - Sinkhorn algorithm

Again, in 2013, another little revolution happened in the numerical solving of the optimal transport problem.

1/ Marco Cuturi (2013) had the idea of introducing ~~another~~ a penalized problem, which has better convexity properties:

$$\inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) \pi(dx, dy) + \varepsilon \text{Ent}(\pi)$$

$\xrightarrow{\quad} \text{Entropy}$

Written for a discrete problem:

$$\inf_{\pi \in \Pi(\mu, \nu)} \sum_{i,j=1}^m c_{ij} \pi_{ij} + \varepsilon \pi_{ij} \log \pi_{ij}$$

$$\Leftrightarrow \begin{cases} \pi \mathbf{1} = \mu \\ \pi^T \mathbf{1} = \nu \end{cases}$$

$$= \inf_{\substack{\pi \in M_m(\mathbb{R}_+), \\ \pi \mathbf{1} = \mu \\ \pi^T \mathbf{1} = \nu}} \sum_{i,j=1}^m c_{ij} \pi_{ij} + \underbrace{\varepsilon \pi_{ij} \log \pi_{ij}}_{\text{gives convexity}}$$



Method of Lagrangian Multipliers

$$L(\pi, \alpha, \beta) = \sum_{i,j=1}^m c_{ij} \pi_{ij} + \varepsilon \sum_{i,j} \pi_{ij} \log \pi_{ij} - \sum_i \alpha_i (\sum_j \pi_{ij} - \mu) - \sum_j \beta_j (\sum_i \pi_{ij} - \nu)$$

$(\pi, \alpha, \beta)$  critical point

$$\Leftrightarrow \begin{cases} \sum_j \pi_{ij} = \mu; \quad \sum_i \pi_{ij} = \nu \\ 0 = \frac{\partial L}{\partial \pi_{ij}} = c_{ij} + \varepsilon (\log \pi_{ij} + 1) - \alpha_i - \beta_j \end{cases}$$

$$\Leftrightarrow \begin{cases} \sum_j \pi_{ij} = \mu; \quad \sum_i \pi_{ij} = \nu \\ \log \pi_{ij} = -1 + \frac{1}{\varepsilon} (\alpha_i + \beta_j - c_{ij}) \end{cases}$$

$$\Leftrightarrow \begin{cases} \pi_{ij} = e^{-1} e^{-\frac{c_{ij}}{\varepsilon}} e^{\frac{1}{\varepsilon} \alpha_i} e^{\frac{1}{\varepsilon} \beta_j} \\ \sum_j \pi_{ij} = \mu; \quad \sum_i \pi_{ij} = \nu \end{cases}$$

$$\Leftrightarrow \begin{cases} \pi_{ij} = \cancel{e^{-1}} u_i e^{-\frac{c_{ij}}{\varepsilon}} v_j \\ u_i, v_j > 0 \\ \sum_j \pi_{ij} = \mu; \quad \sum_i \pi_{ij} = \nu \end{cases}$$

Hence Proposition: The  $\varepsilon$ -penalized problem

has a unique solution of the form

$$\pi_{ij} = u_i e^{-\frac{c_{ij}}{\varepsilon}} v_j, \quad u_i, v_j > 0$$

$$\sum_j \pi_{ij} = \mu; \quad \sum_i \pi_{ij} = \nu$$

2) Sinkhorn & Knopp (1967)

had the idea much earlier of solving this algorithm via a fixed-point iteration

Lemma:  $K = (k_{ij} = e^{-\frac{c_{ij}}{\epsilon}})$   $\pi_{ij} = u_i k_{ij} v_j$

$$\left\{ \begin{array}{l} \pi \mathbf{1} = \mu \quad \Rightarrow \quad u_i = \frac{\mu_i}{(Kv)_i} \\ \pi^T \mathbf{1} = \nu \quad \Rightarrow \quad v_j = \frac{\nu_j}{(K^T u)_j} \end{array} \right.$$

Proof: Easy.

Hence the idea:

while  $|\pi \mathbf{1} - \mu|, |\pi^T \mathbf{1} - \nu| \geq \text{Error}$

$$u \leftarrow \left( \frac{\mu_i}{(Kv)_i} \right)_i$$

$$\pi \leftarrow \left( u_i k_{ij} v_j \right)_{i,j}$$

$$v \leftarrow \left( \frac{\nu_j}{(K^T u)_j} \right)_j$$

End while

$$\pi \leftarrow \left( u_i k_{ij} v_j \right)_{i,j \leq m}$$