

Probabilistic & computational

I

aspects of OT

Plan

I What is OT?

- i) The Monge Problem
- ii) Relaxat.^o to ~~the~~ the Kantorovich Problem.
- iii) Discrete setting

II Kantorovich duality

III Computational aspects: Sinkhorn

IV From Kantorovich to Monge:
Brenier's Thm.

I What is OT?

Let X & Y be Polish spaces.

$\mu \in \mathcal{M}_1(X) \mapsto \mathbb{P}$ -measures on X

$\nu \in \mathcal{M}_1(Y)$

(Often: X & Y will be

- open subsets of \mathbb{R}^n
- discrete

)

i) \rightarrow Monge's historical problem.
(1781) « ~~Thèse~~ Mémoire sur la théorie des déblais & remblais »

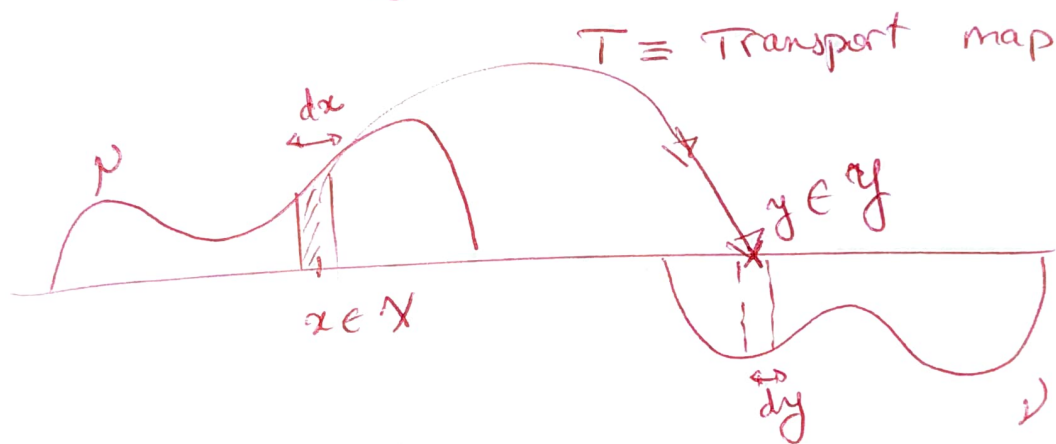
\approx « Dissertat.^o on the theory of excavat.^o & backfills »

where he asks the question

"What is the optimal way of moving a distribution μ (of dirt, to excavate) to a distribution ν (of dirt, to fill a hole) - ?"

To that end, let

$c: X \times Y \rightarrow \mathbb{R} \cup \{+\infty\}$ Cost funct^o



Monge's problem, formalized

Find a transport map $T: X \rightarrow Y$ realizing the infⁱⁿ

$$W_c^{\text{Monge}}(\mu, \nu) = \inf_{T: T_*\mu = \nu} \int_X \mu(dx) c(x, T(x))$$

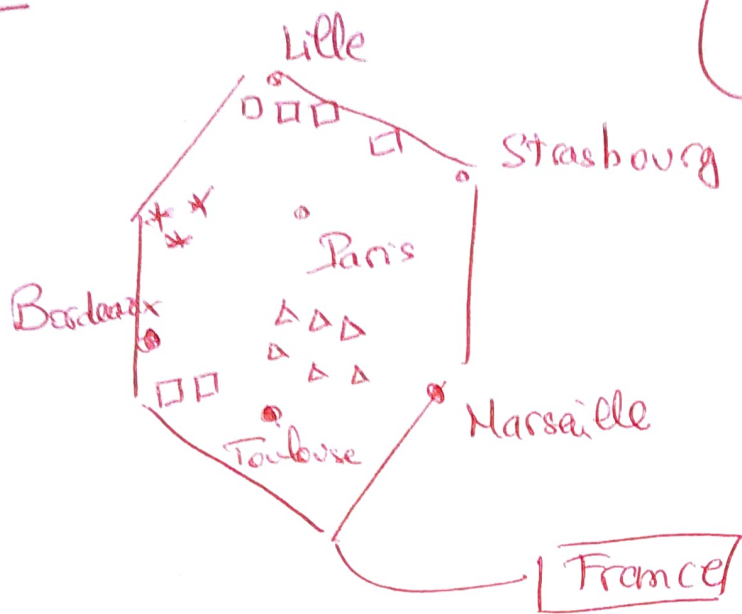
→ Difficult. Non-linear constraints on the space of maps Π .

ii) → Kantorovich relaxat^o

In 1942, Kantorovich proposed a relaxat^o of these constraints in a paper "On the translocat^o of masses"

Motivat°: Economics & Central Planning.

(↑ Actually same as Monge)



- Coal mines
- * Salt mines
- △ Quarries

& Kantorovich was thinking metal mines & factories.

Key idea: Allow the splitting of masses.

Not all the mass $\mu(dx)$ at x need to move to $T(x)$

but it can ~~move to~~ π_x
 be split over $\pi_x(dy)$ for $\neq y$ conditional probability

a unique point

$$\int \mu(dx) \int c(x, y) \pi_x(dy) = \int \pi(dx, dy) c(x, y)$$

* marginal.

Kantorovich Problem (Formalized)

$$W_c(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int_{X \times Y} \pi(dx, dy) c(x, y)$$

$$= \inf_{\substack{\mathcal{L}(X) = \mu \\ \mathcal{L}(Y) = \nu}} E c(X, Y)$$

Couplings

* marginals

Clearly, we have $W_C^{\text{Monge}}(\mu, \nu) \geq W_C(\mu, \nu)$ IV

because $\left\{ \left(\mu(dx) \otimes \frac{(dy)}{T(x)} \right); T_*\mu = \nu \right\} \neq \Pi(\mu, \nu)$

Rmk: Considerable conceptual gain; as Kantorovich's

problem is:

- linear in π

- ~~convex in $\pi \in \Pi(\mu, \nu)$~~

~~is~~ is convex + linear constraints

\Rightarrow Foundat^o of linear programming & operat^o research

iii) Discrete setup

$$X = \{x_1, x_2, \dots, x_m\} \simeq \{1, 2, \dots, m\} = [m]$$

$$Y = \{y_1, \dots, y_m\} \simeq [m]$$

$$\mu = (\mu_1, \dots, \mu_m) \in \mathbb{R}_+^m \text{ with } \sum \mu_i = 1$$

$$\nu = (\nu_1, \dots, \nu_m) \in \mathbb{R}_+^m \text{ with } \sum \nu_j = 1$$

$$\Pi(\mu, \nu) \ni \pi = (\pi_{ij}) \text{ with } \sum_j \pi_{ij} = \mu_i$$

$$\sum_i \pi_{ij} = \nu_j$$

\Rightarrow Always compact

$\pi_* = \operatorname{Argmin}_{\pi \in \Pi(\mu, \nu)} \int \pi(dx dy) c(x, y)$ reached for general Polish spaces

\Downarrow In general, need for Prokhorov!

Even more particular case: $m = n$

$\mu = \nu = \text{uniform}$

$\leadsto \tilde{\pi}$ is just a ~~big~~ bistochastic matrix

$$\leadsto W_C(\mu, \nu) = \inf_{\tilde{\pi} \in \mathcal{D}_m \leadsto \text{bistochastic matrices}} \left\{ \frac{1}{n} \sum c(x_i, y_j) \tilde{\pi}_{ij} \right\}$$

Fact: (Choquet) Optimum reached on $\text{ext}(\mathcal{D}_n)$

(Birkhoff) $\text{ext}(\mathcal{D}_n) = \mathcal{P}_n$.

\leadsto Matching Problem

$$W_C(\mu, \nu) = \inf_{\sigma \in \mathcal{P}_n} \left\{ \frac{1}{n} \sum_{i=1}^n c(x_i, y_{\sigma(i)}) \right\}$$

Solvable using the Hungarian algorithm
/ Kuhn-Munkres in $\mathcal{O}(n^3)$
(See Wikipedia).

II Kantorovich duality

VI

This nothing but convex duality which holds at a high level of generality

Then [Kantorovich (convex) duality]

Let $\begin{cases} X \& Y \text{ Polish spaces} \\ c: X \times Y \rightarrow \mathbb{R}_+ \cup \{+\infty\} \text{ lower semi-continuous} \end{cases}$

$$\begin{aligned} \text{Then } W_c(\mu, \nu) &= \inf_{\pi \in \Pi(\mu, \nu)} \int \pi(dx dy) c(x, y) \\ &= \sup_{(\varphi, \psi) \in \Phi_c} \int \varphi(x) \mu(dx) + \int \psi(y) \nu(dy) \end{aligned}$$

with $\Phi_c = \{ (\varphi, \psi) \in \mathcal{C}_b(X) \times \mathcal{C}_b(Y) \mid \varphi(x) + \psi(y) \leq c(x, y) \forall x \in X, y \in Y, \mu\text{-ac}, \nu\text{-ac} \}$

In the discrete setting $\begin{cases} X \approx [m] \\ Y \approx [m] \end{cases}$

the statement is much more simple:

Corollary (Kantorovich 1942)

$$W_c(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \sum_{i,j} c_{ij} \pi_{ij}$$

$$= \sup_{(f, g) \in \mathbb{R}^m \times \mathbb{R}^m} \sum_i f_i \mu_i + \sum_j g_j \nu_j$$

$f_i + g_j \leq c_{ij}$ better to saturate

Proof: Identify p, v to vectors
 π to a matrix, C as well. VII

Recall
$$\Pi(p, v) = \left\{ \pi \in \mathbb{R}_+^{m \times m} \mid \begin{array}{l} \pi \mathbf{1}_m = p \\ \pi^T \mathbf{1}_m = v \end{array} \right\}$$

and
$$W_C(p, v) = \min_{\pi \in \mathbb{R}_+^{m \times m}} \langle \pi, C \rangle$$

$$\begin{array}{l} \pi \mathbf{1}_m = p \quad \leftarrow n \text{ constraints} \\ \pi^T \mathbf{1}_m = v \quad \leftarrow m \text{ ---} \end{array}$$

Linear problem with constraints \Rightarrow Introduce Lagrangian.

Lagrange multipliers are $f \in \mathbb{R}^n$ and $g \in \mathbb{R}^m$

$$\mathcal{L}(\pi, f, g) = \langle \pi, C \rangle + \sum_i f_i (\pi \mathbf{1}_m - p)_i + \langle g, \pi^T \mathbf{1}_m - v \rangle$$

$$\frac{\partial \mathcal{L}}{\partial f_i} = 0 \Rightarrow p_i = (\pi \mathbf{1}_m)_i$$

$$\frac{\partial \mathcal{L}}{\partial g_j} = 0 \Rightarrow v_j = (\pi^T \mathbf{1}_m)_j$$

$$\frac{\partial \mathcal{L}}{\partial \pi_{ij}} = 0 \Rightarrow 0 = C_{ij} + (f_i + g_j)$$

$$\Leftrightarrow f_i + g_j = -C_{ij}$$

Now
$$W_C(p, v) \xrightarrow{\min} \xrightarrow{\max} \pi \in \mathbb{R}_+^{m \times m} \quad f, g$$

$$\begin{aligned} W_C(p, v) &= \langle \pi^*, C \rangle = \sum \pi_{ij}^* (f_i^* + g_j^*) \\ &= \sum_i f_i^* \left(\sum_j \pi_{ij}^* \right) + \sum_j g_j^* \left(\sum_i \pi_{ij}^* \right) = \langle f^*, p \rangle + \langle g^*, v \rangle \\ &= \sup_{f, g \in C} \langle f, p \rangle + \langle g, v \rangle \end{aligned}$$

Remarks & consequences of Kantorovich duality VIII

① Economical interpretation

$$* W_c(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \sum \pi_{ij} c_{ij} \quad \leftarrow \text{Minimising cost}$$

$$= \sup_{f_i + g_j \leq c_{ij}} \left[\sum f_i \mu_i + \sum g_j \nu_j \right] \quad \leftarrow \text{Maximizing gain}$$

f_i gain in buying one unit of i
 g_j gain in selling one unit of j

* Cost / Surplus

If $c(x, y) \leq a(x) + b(y)$
 $a \in L^1(\mu), b \in L^1(\nu)$

$s(x, y) = -c(x, y) + a(x) + b(y)$ "Surplus"

$$\inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) \pi(dx dy) = \mu(a) + \nu(b) - \sup_{\pi \in \Pi(\mu, \nu)} \int s(x, y) \pi(dx dy)$$

② c -concavity and the double convexification trick

Def: [c -concavity]

$c: X \times Y \rightarrow \mathbb{R} \cup \{-\infty\}$ given.

$\varphi: X \rightarrow \mathbb{R} \cup \{-\infty\}$ is called c -concave ~~iff~~ $\exists \psi$

if $\exists \psi: Y \rightarrow \mathbb{R} \cup \{-\infty\}, \psi \neq -\infty$

st $\varphi(x) = \inf_{y \in Y} c(x, y) - \psi(y)$. We write $\varphi = \varphi^c$.

Important example: $X = Y = \mathbb{R}^n$

IX

before $c(x, y) = \frac{\|x-y\|^2}{2} = \frac{\|x\|^2}{2} + \frac{\|y\|^2}{2} - \underbrace{\langle x, y \rangle}_{s(x, y)}$

Given $\varphi: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{-\infty\}$, $\|$

Then $\varphi^c(y) = \inf_{x \in \mathbb{R}^n} \frac{1}{2} \|x-y\|^2 - \varphi(x)$

$$\Leftrightarrow \underbrace{\left(\varphi^c - \frac{\|x\|^2}{2}\right)}_{-f^*}(y) = \inf_{x \in \mathbb{R}^n} -\langle x, y \rangle - \underbrace{\left(\varphi - \frac{\|x\|^2}{2}\right)}_{-f}(x)$$

$$\Leftrightarrow f^*(y) = \sup_{x \in \mathbb{R}^n} \langle x, y \rangle - f(x)$$

usual Legendre-Fenchel duality.

$$\triangleleft W_2(p, \nu) = \sqrt{\inf \int \frac{\|x-y\|^2}{2} \pi(dx, dy)}$$

very geometric with strong links to Legendre-Fenchel (no Part IV).

Now

Fact 1: $(\varphi^{cc})^c = \varphi^c$

Fact 2: $p(\varphi) + \nu(\varphi)$ with $\varphi(x) + \varphi(y) \leq c(x, y)$

$$\leq p(\varphi) + \nu(\varphi^c)$$

$$\leq p(\varphi^{cc}) + \nu(\varphi^c)$$

Double convexification trick:

$$\Rightarrow W_c(p, \nu) = \inf_{\pi} \langle \pi, C \rangle = \sup_{\varphi} p(\varphi^c) + \nu(\varphi^c)$$

③ Kantorovich-Rubinstein ($c=d$ is a distance)

Thm [Kantorovich-Rubinstein]

Let $X=Y$, $c=d$ is a distance.

$$\begin{aligned} \text{Then } W_1(\mu, \nu) &:= \inf_{\pi \in \Pi(\mu, \nu)} \int d(x, y) \pi(dx dy) \\ &= \sup_{\|\varphi\|_{\text{Lip}} \leq 1} |\mu(\varphi) - \nu(\varphi)| \end{aligned}$$

Proof: $W_1(\mu, \nu) \stackrel{\text{K.duality}}{=} \sup_{\varphi(x) + \varphi(y) \leq d(x, y)} \mu(\varphi) + \nu(\varphi)$
 $\stackrel{2x \text{ convexificato}}{=} \sup_{\varphi \text{ bounded}} \mu(\varphi^d) + \nu(\varphi^{dd}) \quad (\star)$

But: Fact 1: $\varphi^d(y_1) = \inf_{x \in X} \underbrace{d(x, y_1) - \varphi(x)}_{\leq d(x, y_2) + d(y_1, y_2)}$
 is 1-lip
Fact 2: $-\varphi^d = \varphi^{dd} \leq d(y_1, y_2) + \varphi^d(y_2)$

Indeed $-\varphi^d(x) \leq \varphi^{dd}(x) = \inf_y \underbrace{d(x, y) - \varphi^d(y)}_{\geq \varphi^d(x)}$
 $\leq -\varphi^d(x)$
 $x=y$

from (\star)

Therefore $W_1(\mu, \nu) = \sup_{\varphi \text{ bounded}} \mu(\varphi^d) - \nu(\varphi^d)$
 $\stackrel{\varphi^d \text{ is 1-lip}}{\leq} \sup_{\|\varphi\|_{\text{Lip}} \leq 1} \mu(\varphi) - \nu(\varphi)$
 $\stackrel{\varphi \text{ 1-lip}}{\Rightarrow} (\varphi, -\varphi) \in \Phi_c \text{ with } \varphi(y) \leq d(x, y) \text{ and } \varphi(x) \leq d(x, y)$
 $\sup_{(\varphi, \psi) \in \Phi_c} \mu(\varphi) + \nu(\psi) = W_1(\mu, \nu) \quad \square$

III Computational aspects of OT: Sinkhorn.

Ref: Cuturi, Peyre. Computational OT

- 1) Idea: * Strongly convex optim. problem are better.
- * Regularize using entropy / KL.

$$\forall \pi \in \mathbb{R}_+^{n \times m}, \quad H(\pi) = \begin{cases} + \sum \pi_{ij} (\log \pi_{ij} - 1) & \text{if } \pi_{ij} > 0 \\ -\infty & \text{otherwise} \end{cases}$$

Fact: $\nabla^2 H \preceq -Id$. Δ -concave.

2) Entropy regularized OT. linear Δ concave

Define $W_C^\epsilon(p, \nu) = \inf_{\substack{\pi \in \mathbb{R}_+^{n \times m} \\ \sum \pi_{ij} = 1 \\ \pi \mathbb{1}_m = p \\ \pi^T \mathbb{1}_n = \nu}} \underbrace{\langle C, \pi \rangle}_{\text{linear}} - \underbrace{\epsilon H(\pi)}_{\epsilon\text{-concave}}$

$\implies \exists!$ minimize π_ϵ^*

Theorem (Properties)

- Duality: $W_C^\epsilon(p, \nu) = \inf_{\pi \in \Pi(p, \nu)} \langle C, \pi \rangle - \epsilon H(\pi)$
- = $\sup_{(f, g) \in \mathbb{R}^n \times \mathbb{R}^m} p(f) + \nu(g) - \epsilon \sum_{i,j} c_{ij} \frac{f_i g_j}{\epsilon}$
- Unique optimizers f^*, g^*, π^*

• Simple relat° between primal and dual.

(*) $(\pi_\epsilon^*)_{ij} = \exp\left(\frac{f_i^* + g_j^* - C_{ij}}{\epsilon}\right)$

Proof: write Lagrange multipliers & do like when $\epsilon = 0$.

More over; $\lim_{\epsilon \rightarrow 0} \pi_\epsilon^*$ exists and solves the K. problem.

Lemma: $\pi_\epsilon^* \rightarrow \underset{\pi}{\text{Argmin}} \{ -H(\pi) \mid W_c(\mu, \nu) = \langle C, \pi \rangle \}$

~~Proof: write first order condit°~~

Proof: ~~compare π_ϵ^*~~ Take a CV subsequence π_ϵ^* and how limit behaves when plugged into to minimizat° problem.

(3) Sinkhorn alg. or proportional fitting

write $K = e^{-\frac{C_{ij}}{\epsilon}}$ (Gibbs Kernel).

From (*) π_ϵ^* is the only rescaling of columns & rows such that $\pi_\epsilon^* \in \Pi(\mu, \nu)$

⇒ Iteratively rescale columns, then rows in order to have the correct marginals.

Algorithm:

$$P \in \mathbb{R}_+^{n \times m}$$

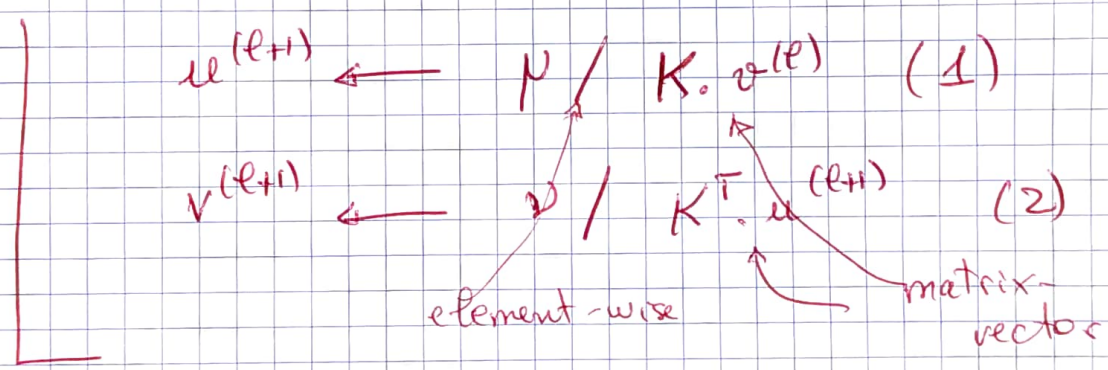
$$u^{(0)} \in \mathbb{R}^m$$

$$v^{(0)} \in \mathbb{R}^m$$

$$u^{(0)} \leftarrow 1; \quad v^{(0)} \leftarrow 1$$

while $P = \frac{K}{\sum u^{(l)}} \text{diag}(u^{(l)}) K \text{diag}(v^{(l)})$

satisfies $\|P\|_m - \mu + \|P^T\|_m - \nu \geq \text{Error}$



IV

Brenier's theorem:

• $(\mu, \nu) \in \mathcal{M}_+(\mathbb{R}^n) \times \mathcal{M}_+(\mathbb{R}^m)$ with 2nd moments

• $c(x, y) = \frac{\|x - y\|^2}{2} \quad (W_2!)$

• $\mathbb{R}^n \xrightarrow{\mu}$ ~~has a density~~

~~$\mu(A) = 0 \quad \forall A \text{ with } \dim_H(A) \leq n-1$~~

$\iff \mu$ has density.

Then $\exists!$ T optimal, $T(dx dy) = \mu(dx) \otimes T(x) (dy)$

with $T_* \mu = \nu$

Here $T = \nabla \left(\varphi - \frac{\|x\|^2}{2} \right)$ gradient of convex funct^o.