

# Optimal Transport: Lecture #5

Last time: We established the form of the Kantorovich Problem in the setting of discrete measures. As usual, let  $\Pi(\mu, \nu)$

$$:= \left\{ \gamma: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \mid \begin{aligned} \gamma(A \times \mathbb{R}^d) &= \mu(A), \\ \gamma(\mathbb{R}^d \times B) &= \nu(B), \text{ for } \\ &\text{Borel } A, B \end{aligned} \right\}$$

be the space of couplings between measures  $\mu, \nu$ .

• When we consider the common squared Euclidean cost function  $c(x, y) = \|x - y\|_2^2$ , the Kantorovich problem is

$$\operatorname{argmin}_{\Pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|_2^2 d\Pi(x, y).$$

• When  $\mu, \nu$  are discrete and supported, respectively, on  $\{x_i\}_{i=1}^n$  and  $\{y_j\}_{j=1}^m$ , this reduces to a minimization over bistochastic matrices: let  $a = \sum_{i=1}^n \mu_i$  and  $b = \sum_{j=1}^m \nu_j$ .

$$\operatorname{argmin}_{P \in U(a, b)} \sum_{ij} C_{ij} \cdot P_{ij}, \quad \text{where: } \begin{aligned} \text{(i)} \quad C_{ij} &:= \|x_i - y_j\|_2^2 \\ \text{(ii)} \quad U(a, b) &:= \left\{ P \in \mathbb{R}_{\geq 0}^{n \times m} \mid \begin{aligned} P \mathbf{1}_m &= a, \quad P \mathbf{1}_n = b \end{aligned} \right\} \end{aligned}$$

• We saw last time this can be written as a linear program. This can be solved using off-the-shelf methods (e.g. simplex algorithms, interior point methods), though these scale quite poorly in the problem size. ②

• More precisely, if  $n=m$ , so we are optimizing over  $U(a,b) = \{P \in \mathbb{R}_{\geq 0}^{n \times n} \mid P \mathbb{1}_n = a, P^T \mathbb{1}_n = b\}$ , then typical algorithms will on average scale roughly like  $\mathcal{O}(n^3)$ , and could be worse in practice.

• This area is well-developed and there are bespoke algorithms available that achieve essentially cubic complexity in practice, but this is still bad. (A possible presentation could focus on some of these bespoke algorithms)

• A major innovation is due to M. Cuturi in 2013. The idea is to regularize the problem, allowing alternative, simpler approaches to be used for optimization.

• Let's establish some preliminaries. Let  $P \in \mathbb{R}_{\geq 0}^{n \times m}$  be a coupling matrix. We define the (discrete) entropy of  $P$  as

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

where  $H(P) = -\infty$  if any entry of  $P$  is 0.

Note that ~~the~~  $H: \mathbb{R} \rightarrow \mathbb{R}$ , and we can compute its gradient and Hessian:

$$\partial_{ij} H(P) = -[\log(p_{ij}) + 1 - 1]$$

$$\Rightarrow \partial_{ij, i'j'} H(P) = \begin{cases} 0, & ij \neq i'j' \\ -\frac{1}{p_{ij}}, & ij = i'j' \end{cases}$$

$$\Rightarrow \partial^2 H(P) = -\text{diag}\left(\frac{1}{p_{ij}}\right)$$

Now, by assumption,  $p_{ij} \leq 1 \Rightarrow \frac{1}{p_{ij}} \geq 1 \Rightarrow H$  is strongly concave. In particular,  $-\epsilon H(P)$  is  $\epsilon$ -strongly convex in  $P$ .

Now, the map  $P \rightarrow \sum_{ij} p_{ij} c_{ij}$  is linear in  $P$ , and in particular ~~the~~ its Hessian is 0. We conclude that the entropy-regularized cost functional

$$P \rightarrow \underbrace{\sum_{ij} p_{ij} c_{ij}}_{= \langle P, C \rangle} - \epsilon H(P)$$



is  $\epsilon$ -strongly convex.

(4)

This implies

(A)

$$L_c^\epsilon(a, b) := \min_{P \in U(a, b)} \left\{ \langle P, C \rangle - \epsilon H(P) \right\}$$

has a unique optimum.

One can see immediately that introducing  $-\epsilon H(P)$  kills any sparsity in the optimizer  $P_\epsilon^* = \arg \min_{P \in U(a, b)} \left\{ \langle P, C \rangle - \epsilon H(P) \right\}$ . Indeed, if

$$[P_\epsilon^*]_{ij} = 0 \text{ for any } (ij) \text{ index, then } \langle P, C \rangle - \epsilon H(P) = +\infty.$$

In this sense, regularizing with entropy ensures we have a "diffuse" coupling, rather than a sparse one as in the classical Monge problem.

So, in some sense we are "giving up" on the (frankly, desirable) property of sparse maps/couplings. So, why do this? What's the gain?

The key point is, Sinkhorn's algorithm allows (A) to be solved much faster than its unregularized cousin.

Let  $K_\epsilon \in \mathbb{R}^{n \times m}$  be the Gibbs/heat kernel associated to the cost matrix  $C \in \mathbb{R}^{n \times m}$ ;  $[K_\epsilon]_{ij} = \exp(-C_{ij}/\epsilon)$ . So, if we have

quadratic cost, then

$$[K_\epsilon]_{ij} = \exp(-\|x_i - x_j\|^2 / \epsilon).$$

Theorem (structure of solution to  $\star$ ): The solution to  $\star$  has the form

$$[P_\epsilon^\star]_{ij} = u_i [K_\epsilon]_{ij} v_j \quad \text{for two (unknown) vectors } (u, v) \in \mathbb{R}_+^n \times \mathbb{R}^m$$

Proof: Let  $f_i \in \mathbb{R}^n$ ,  $g_j \in \mathbb{R}^m$  be dual variables for the marginal constraint enforcing  $P \in U(a, b)$ . Then the Lagrangian is

$$\mathcal{L}(P, f, g) = \underbrace{\langle P, C \rangle - \epsilon H(P)}_{\text{original objective}} - \underbrace{\langle f, P \mathbb{1}_{m \times a} \rangle - \langle g, P^T \mathbb{1}_{n \times b} \rangle}_{\text{constraints}}$$

Differentiating in  $P$  yields

$$\frac{\partial \mathcal{L}(P, f, g)}{\partial P_{ij}} = C_{ij} - \epsilon \log(P_{ij}) - f_i - g_j.$$

Setting equal to 0 and using first order optimality,

we have

$$C_{ij} - \epsilon \log(P_{ij}) - f_i - g_j = 0$$

$\Leftrightarrow$

$$P_{ij} = \exp(f_i / \epsilon) \exp(-C_{ij} / \epsilon) \exp(-g_j / \epsilon)$$

$$= U_i [K_\varepsilon]_{ij} V_j.$$

(6)

• One can find the vectors  $(u, v)$  through an alternating scheme as follows

• Note that  $P_{ij}^* = U_i [K_\varepsilon]_{ij} V_j \Leftrightarrow P = \text{diag}(u) K_\varepsilon \text{diag}(v)$ , and the

constraints shake out to

$$\begin{cases} \text{diag}(u) K_\varepsilon \text{diag}(v) \mathbb{1}_m = a, \\ \text{diag}(v) K_\varepsilon^T \text{diag}(u) \mathbb{1}_n = b \end{cases}$$

• Let ~~we~~ as define a sequence of iterates as  $V^{(0)} = \mathbb{1}_m$  and

$$\begin{aligned} U^{(1)} &= \frac{a}{K_\varepsilon V^{(0)}}, \\ V^{(2)} &= \frac{b}{K_\varepsilon^T U^{(1)}}. \end{aligned}$$

• One can show (see Pegré's (2018)) that this sequence of iterates converges to the optimal  $u, v$ . Note moreover that an update requires only a matrix multiplication  $\rightarrow O(mn)/O(n^2)$ . If there are not too many iterates, this gives us an essentially quadratic algorithm.