

Let us recall the Kantorovich problem. Given  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ , we let

$$\Pi(\mu, \nu) := \left\{ \gamma: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \mid \begin{array}{l} \forall A \subset \mathbb{R}^d, \gamma(A \times \mathbb{R}^d) = \mu(A) \\ \forall B \subset \mathbb{R}^d, \gamma(\mathbb{R}^d \times B) = \nu(B) \end{array} \right\}.$$

We call  $\Pi(\mu, \nu)$  the space of couplings between  $\mu, \nu$ . The Kantorovich formulation of the OT problem with quadratic cost solves

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

We saw last time that for general forms of this problem (i.e.  $\mathbb{R}^d$  replaced with abstract Polish spaces,  $\|x-y\|_2^2$  replaced with an abstract cost function satisfying certain properties) that (KP) admits a dual formulation:

$$\operatorname{min}_{\Pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x-y\|_2^2 d\Pi(x, y) = \sup_{\substack{(\ell, \psi) \in \mathcal{C}_b(\mathbb{R}^d \times \mathbb{R}^d) \\ \text{bounded}}} \int_{\mathbb{R}^d} \ell(x) d\mu(x) + \int_{\mathbb{R}^d} \psi(y) d\nu(y).$$

s.t.  $\ell(x) + \psi(y) \leq \|x-y\|_2^2$

The "dual" formulation (RHS) says: maximize a linear functional

$$(\ell, \psi) \mapsto \int_{\mathbb{R}^d} \ell(x) d\mu(x) + \int_{\mathbb{R}^d} \psi(y) d\nu(y) \quad \text{subject to the constraint}$$

$\ell(x) + \psi(y) \leq \|x-y\|_2^2$ . Note that the constraint is linear in  $\ell, \psi$ , so in fact we have a linear objective subject to linear constraints.

**Q:** (Can we address this with practical algorithms in the discrete case?)

- First, let's see how this works out in the discrete setting.
- Suppose instead of arbitrary measures  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$ , we instead have discrete measures, i.e. those supported on a finite number of points in  $\mathbb{R}^d$ , i.e. sums of (potentially non-uniformly weighted) Dirac masses:

$$\mu = \sum_{i=1}^n a_i \delta_{x_i}, \quad \nu = \sum_{j=1}^m b_j \delta_{y_j} \quad \text{s.t.} \quad \sum_{i=1}^n a_i = \sum_{j=1}^m b_j = 1.$$

In this case, (KP) shakes out as follows:

~~min~~  $\arg \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x-y\|_2^2 d\pi(x,y)$  can be reduced to a matrix problem by noting that any  $\pi \in \Pi(\mu, \nu)$  as above must be supported on points in the Cartesian product  $\{x_i\}_{i=1}^n \times \{y_j\}_{j=1}^m = \{(x,y) \mid x=x_i \text{ for some } i=1,\dots,n \text{ and } y=y_j \text{ for some } j=1,\dots,m\}$ .

So, any  $\pi \in \Pi(\mu, \nu)$  is supported on only (at most!)  $m \cdot n$  points in  $\mathbb{R}^d \times \mathbb{R}^d$ . This means

$$\int_{\mathbb{R}^d \times \mathbb{R}^d} \|x-y\|_2^2 d\pi(x,y) = \sum_{i,j} \|x_i - y_j\|_2^2 \cdot \hat{\pi}_{ij},$$

where  $\hat{\pi}_{ij} \in [0,1]$  are how much mass gets sent from  $x_i$  to  $y_j$  (intuitively). This allows us to write (KP) in a more familiar form for those

with background in discrete optimization:

$$\text{arg min}_{P \in U(a,b)} \sum_{i,j} C_{ij} \cdot P_{ij}, \quad \text{where: } C_{ij} := \|x_i - y_j\|_2^2 \text{ are costs}$$

$$P \text{ is an "assignment" matrix}$$

$$U(a,b) = \{ P \in \mathbb{R}_{\geq 0}^{n \times m} \mid P \mathbb{1}_m = a, P^T \mathbb{1}_n = b \}$$

is the space of admissible transport matrices

$$\mathbb{1}_n^T = (1, 1, \dots, 1) \quad n \text{ times}$$

This is evidently a linear program, but we can write it in standard form (linear objective, equality constraints defined with a matrix and constant vector, and non-negativity constraints on the variables).

Let  $I_n = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$   $n$  times

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & \dots & a_{1m}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1}B & \dots & \dots & a_{nm}B \end{bmatrix} \in \mathbb{R}^{n \times p \times m \times q}$$

be the Kronecker product of A, B

Let  $A = \begin{bmatrix} \mathbb{1}_n^T \otimes I_m & \\ I_n \otimes \mathbb{1}_m^T \end{bmatrix} \in \mathbb{R}^{(n+m) \times (nm)}$

Exercise:  $P \in U(a,b) \Leftrightarrow p \in \mathbb{R}_{\geq 0}^{nm}$  and  $A_p = \begin{bmatrix} a \\ b \end{bmatrix}$ , where

$p$  is the "vectorized" form of  $P$  with  $p_{i+n(j-1)} = P_{ij}$ .

Letting  $L_c(a, b) = \min_{P \in U(a, b)} \sum_{\substack{i=1, \dots, n \\ j=1, \dots, m}} C_{ij} P_{ij}$  be the ~~minimum~~ OT distance (4)

(using that word informally), ~~what~~ we thus have

$$L_c(a, b) = \min_{\substack{p \in \mathbb{R}^{nm} \\ p \geq 0}} C^T p, \quad (\text{KP-LP})$$

$$\max_{\pi} \text{ s.t. } A p = \begin{bmatrix} a \\ b \end{bmatrix}$$

where again  $C$  is the vectorized form of the cost matrix  $C \in \mathbb{R}^{nm}$ , namely

~~$C_{i+n \cdot (j-1)} = C_{ij}$~~ . This has the standard form of a linear program,

so there are many off-the-shelf tools to handle it.

Remark: Formulating the transport as a (discrete) linear program is what Kantorovich's original work and applications centered on. Note, Kantorovich won the Nobel (Memorial) Prize in Economics, indicating the huge impact his work had on the applied side of things.

As we might hope (KP-LP) admits a dual formulation that matches what we saw in the general setting:

$$L_c(a, b) = \max_{h \in \mathbb{R}^{n+b}} \begin{bmatrix} a \\ b \end{bmatrix}^T h.$$

$$A^T h \leq c$$

Remark: Suppose  $m = n$ , so that there are the same number of support

points in  $\mu$  and  $\nu$ , and suppose moreover ~~uniform~~  $d \equiv b \equiv \frac{1}{n} \mathbb{1}_n$ , i.e. the weights are uniform on both measures. Then

$$U(a, b) = U\left(\frac{1}{n} \mathbb{1}_n, \frac{1}{n} \mathbb{1}_n\right) \\ = \left\{ P \in \mathbb{R}_{\geq 0}^{n \times n} \text{ s.t. } P \mathbb{1}_n = P^T \mathbb{1}_n = \frac{1}{n} \mathbb{1}_n \right\}$$

So,  $U(a, b)$  is the set of non-negative (entrywise) matrices with rows and columns that sum up to  $\frac{1}{n} \mathbb{1}_n$ .

An important point is that, ~~the~~  $P \in U\left(\frac{1}{n} \mathbb{1}_n, \frac{1}{n} \mathbb{1}_n\right)$  could be very diffuse, (e.g.,  $P = \begin{pmatrix} \frac{1}{n^2} & \frac{1}{n^2} & \dots & \frac{1}{n^2} \\ \frac{1}{n^2} & \frac{1}{n^2} & \dots & \frac{1}{n^2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{n^2} & \dots & \frac{1}{n^2} \end{pmatrix} \in U$ ). This is because KP allows us to "spread mass"

for and wide, at least in principle. In contrast, the Monge problem specifies a map from  $\text{supp}(\mu)$  to  $\text{supp}(\nu)$ . In the discrete setting, we can capture this idea through permutation matrices: let  $\sigma = \{1, \dots, n\}$  be a permutation. Then we have an associated matrix

$$[P_\sigma]_{ij} = \begin{cases} \frac{1}{n}, & j = \sigma(i) \\ 0, & \text{else} \end{cases}$$

Note  $P_\sigma \in U\left(\frac{1}{n} \mathbb{1}_n, \frac{1}{n} \mathbb{1}_n\right)$  for all permutations  $\sigma$ , and in fact

$\sum_{i,j=1}^n [P_\sigma]_{ij} C_{ij} = \sum_{i=1}^n C_{i\sigma(i)}$ , which is the discrete analogue of the objective in the Monge problem,  $\int_{\mathbb{R}^d} \|x - T(x)\|_2 \mu(x)$ .

• So, the discrete Marge problem involves optimizing the usual cost functional over permutation matrices rather than  $n \times n$  bi-stochastic matrices: ⑥

$$\min_{\substack{\sigma \in \{1, \dots, n\} \\ \text{permutation}}} \sum_{i=1}^n C_i \sigma(i) \quad (\text{discrete MP}).$$

• Note, since  $(\text{discrete KP}) = \min_{U \in \Pi(\mathbf{1}_n, \mathbf{1}_n)} \sum_{i,j=1}^n C_{ij} U_{ij}$  optimization

over a larger set,  $(\text{discrete KP}) \leq (\text{discrete MP})$ . But, (discrete MP) is a combinatorial optimization (hard to work with permutation matrices directly), so (discrete KP) is preferable. In this sense, KP relaxes MP.

Note: A focus as the course goes on will be to understand in both narrow and general settings when the relaxed KP coincides with MP.

Next time: More on (KP-LP) and discussion of approximations that improve computational complexity.