

Recap:

Setting: given $\mu \& \nu \rightarrow$ distributions on \mathbb{R}^d

we want to transform μ to ν

minimizing some (specified) cost.

Monge: Define a valid transport map: $T_\# \mu = \nu$

(if $x \sim \mu$, then $T(x) \sim \nu$)

OT Map: $T_0 = \arg \min_{\substack{T_\# \mu = \nu}} \int c(x, T(x)) d\mu(x)$

Kantorovich: Define valid couplings τ
(joint distribution on $\mathbb{R}^d \times \mathbb{R}^d$ with
first marginal μ & second marginal ν)

OT Coupling:

$$\tau_0 = \arg \min_{\tau \in \Gamma_{\mu, \nu}} \int c(x, y) d\tau(x, y)$$

Wasserstein distances:

$$W_p(\mu, \nu) = \left[\min_{\tau \in \Gamma_{\mu, \nu}} \int \underbrace{\|x - y\|^p}_{d\tau(x, y)} \right]^{1/p}.$$

→ nice metric on distributions

→ More generally, OT gives a way to lift
 a metric on some space to a metric
 on distributions defined on that space.]

Brenier's Theorem: Squared Euclidean cost

μ, ν have at least 2 bdd moments

μ is absolutely cont. then:

1. Monge has a solution T_0

$T_0 := \nabla \varphi_0$. φ_0 is convex.

(ie) T_0 is monotonic.

$\nabla \varphi_0$ is unique μ a.e

2. Kantorovich has a solution supported on the

$\underline{\tau}_0 := (\underline{x}, \underline{T}_0(\underline{x}))$ graph of a function

3. If ν is also absolutely cont. then

$S_0 := \nabla \varphi_0^*$ is the OT map from ν to μ .

Duality: Discrete case:

$$\min_{\underline{\gamma}} \sum_{ij} c_{ij} \gamma_{ij} \leftarrow \quad]$$

$\gamma_{ij} \geq 0$ $\sum_i \gamma_{ij} = 1/n, \sum_j \gamma_{ij} = 1/n$]

Dual:

$$\max_{\alpha, \beta} \underbrace{\sum_{i=1}^n \alpha_i}_{\frac{n}{n}} + \underbrace{\sum_{j=1}^n \beta_j}_{\frac{n}{n}} \leftarrow$$
$$\alpha_i + \beta_j \leq c_{ij}$$

(Interpretation as a shipper's problem).

Outline

- (1) Duality & Polar Factorization Theorem
 - (2) Statistical optimal transport }
 - Basic questions
 - (3) Estimating the OT Map.]
-

Polar Factorization Theorem: (Brenier)

→ $\|\cdot\|^2$, μ, ν μ ac., $T_0 = \nabla \varphi_0$.
Any other valid transport map T

$$T = T_0 \circ S$$

S : measure preserving map $\mu \rightarrow \mu$.

$$S = \text{Id} \quad T = T_0$$

Matrices: any matrix M

$$M = P O$$

PSD

orthonormal

$$x \sim N(0, I_d)$$

$$Mx \sim N(0, MM^T)$$

$$\varphi(x) = \underbrace{x^T Px}_{2}$$

$$M = P \circ O \leftarrow$$

DT Map

$$T_d(x) = \underline{Px}$$

Duality:

$$\max_{\alpha, \beta} \left[\underbrace{\sum_i \alpha_i}_{\alpha_i + \beta_j \leq c_{ij}} + \underbrace{\sum_j \beta_j}_{\alpha_i + \beta_j \leq c_{ij}} \right]$$

→ general Kantorovich dual:

$$\begin{aligned} & \int |f| d\mu < \infty \quad \max_{f \in L^1(\mu)} \quad \int f d\mu \\ & \int g d\nu < \infty \quad \left[\max_{g \in L^1(\nu)} \quad \int g d\nu \right] \end{aligned}$$

$\int f d\mu + \int g d\nu$

$\rightarrow f(x) + g(y) \leq C(x, y) \quad \forall (x, y)$

→ Strong duality: C is lower semi-cont. & bdd below, then strong duality holds.

Structural remarks when $p=1$

1. Fix f , optimize over g .

$$g(y) \leq \inf_x [c(x, y) - f(x)] := \underline{f^c(y)}$$

Optimal \underline{g} is $\underline{f^c}$.

2. $p=1, c(x, y) = \|x - y\|$

- \rightarrow c -transform of f is 1-Lipschitz.
- \rightarrow c -transform of 1-Lip f , $-f$.

3. Assume they are true:

$$\max_{f \in L^1(\mu)} \int_{\sim} (f^c) d\mu + \int_{\sim} (f^c)^c d\nu.$$

\downarrow 1-Lipschitz -ve of f^c

$$\max_{f \in 1\text{-Lipschitz}} \int f d\mu - \int f d\nu$$

\downarrow Kantorovich-Rubenstein duality.

$$W_1(\mu, \nu) = \sup_{f \in 1\text{-Lip}} \left[\int f d\mu - \int f d\nu \right].$$

\mathbb{W}_1 is an Integral Probability Metric.

$$S(\mu, \nu) = \sup_{f \in \mathcal{F}} \left| \int f d\mu - \int f d\nu \right|$$

TV distance is an IPM

$$\text{TV}(\mu, \nu) = \sup_{f \in \{f : \|f\|_\infty \leq 1\}} \left| \int f d\mu - \int f d\nu \right|.$$

Maximum Mean Discrepancy

$$\rightarrow \text{MMD}(\mu, \nu) = \sup_{f, \|f\|_K \leq 1} \left| \int f d\mu - \int f d\nu \right|.$$

RKHS unit ball {not very big}.

Property 1: f .

$$g(y) = \inf_x \left[\|x - y\| - f(x) \right] \xrightarrow{\inf_x f(x)} \begin{cases} \inf_x f(x) \\ \downarrow \\ \text{1-Lip.} \end{cases}$$

$$\begin{aligned} \tilde{g}(y) &= \|x_0 - y\| - f(x_0) \xleftarrow{\text{1-Lip}} \\ \tilde{g}(y) - \tilde{g}(y') &= \|x_0 - y\| - \|x_0 - y'\| \\ &\leq \|y - y'\|. \end{aligned}$$

Conclusion: $g = f^*$ is 1-Lipschitz.

Property 2: If f is 1-Lip, $(-f)^c = f$.

$$g(y) = \inf_x [\|x - y\| + f(x)].$$

$$\underline{g(y) \leq f(y)}$$

take $x = y$

$$\begin{aligned} f(y) &\leq \|x - y\| + f(x) \\ f(y) &\leq g(y). \end{aligned}$$

$$f = g. \text{, (ie) } (-f)^c = f.$$

Duality for squared Euclidean cost.

Starting pt:

$$\sup_{\substack{f \in L^1(\mu) \\ g \in L^1(\nu)}} \int f d\mu + \int g d\nu$$

$$f(x) + g(y) \leq \|x - y\|^2$$

Optimal f & g - Kantorovich potentials.

$$f(x) = \|x\|^2 - 2\varphi(x)$$

$$g(y) = \|y\|^2 - 2\psi(y)$$

$$\sup_{\varphi, \psi} \underbrace{\int \|x\|^2 d\mu}_{\int \|y\|^2 d\nu} + \underbrace{\int \|y\|^2 d\nu}_{-2 \int \varphi(x) d\mu} - 2 \underbrace{\int \psi(y) d\nu}_{-2 \int \varphi(y) d\mu}$$

$$\varphi(x) + \varphi(y) \geq x^T y.$$

Mod. Dual: $\inf_{\varphi, \psi} \int \varphi d\mu + \int \psi d\nu$

$$\varphi(x) + \psi(y) \geq x^T y.$$

$$\psi(y) \geq \sup_x \left[x^T y - \varphi(x) \right] := \varphi^*(y)$$

$\rightarrow \inf_{\substack{\varphi \in L^1(\mu), \\ \varphi \text{ is convex}}} \int \varphi d\mu + \int \varphi^* d\nu$ ←

↑ Semi-dual

If you satisfy condns of Brenier's thm:

$$T_0 := \nabla \varphi_0.$$

where $\underline{\varphi}_0 = \arg \inf_{\varphi \in L^1(\mu)} \int \varphi d\mu + \int \varphi^* d\nu.$

$$(\varphi(x) + \varphi^*(\nabla \varphi(x))) = \langle x, \nabla \varphi(x) \rangle.$$

Statistical Optimal Transport.

$\rightarrow \mu$ & ν not known.

$\rightarrow X_1, \dots, X_n \sim \mu$ & $Y_1, \dots, Y_n \sim \nu$.
 μ absolute cont.

Target of inference:

1. Estimate T_0 .
 $\hat{T}(x) \approx T_0$? \nearrow how many samples
do I need?

Applications:

* Generative modeling.
 $X_1, \dots, X_n \sim N(0, I_d)$
 Y_1, \dots, Y_n - real images.

* Waddington OT:

cell concentrations

cell
concentrations.

$t=0$

$t=1$

"broken correspondence".

OT map is a way to recover correspondence.

2. Estimate $W_2(\mu, \nu)$ or $W_1(\mu, \nu)$.

Hypothesis testing:

$$x_1, \dots, x_n \sim \mu$$

$$y_1, \dots, y_n \sim \nu$$

Two-sample: $H_0: \mu = \nu$]
 $H_1: W_p(\mu, \nu) \geq \epsilon$]

3. Density estimation: $x_1, \dots, x_n \sim \mu$.

construct $\hat{\mu}$, so that

$W_p(\hat{\mu}, \mu)$ is small.

Sample complexity.

~ Maybe $\hat{\mu} = \bar{x}_n$ is good/good enough.

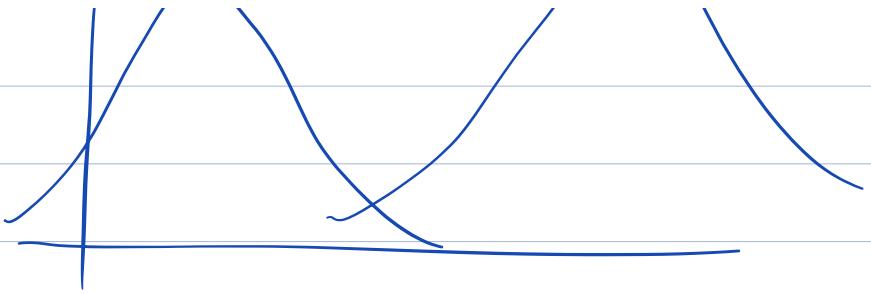
4. Barycenter: $\underbrace{\mu_1, \dots, \mu_k}_{\bar{\mu}}.$

$$\bar{\mu} = \underset{\mu}{\operatorname{argmin}} \sum_i W_2(\mu, \mu_i).$$

"shape preserving".

$\mu_1 \sim N(m_1, I_d)$

$\mu_2 \sim N(m_2, I_d)$



one defn. of "mid point"

$$= \frac{1}{2} N(m_1, I_d) + \frac{1}{2} N(m_2, I_d)$$

Wasserstein mid-point will be
 $N(\bar{m}, I_d)$

Estimating the OT map:

Dual Estimators

Primal Estimators

$$\rightarrow \varphi_0 = \arg \min_{\varphi \in \Phi} \underbrace{\int \varphi d\mu}_{\text{Brenier potential}} + \underbrace{\int \varphi^* d\nu}_{L^1(\mu) \text{ & convex.}}$$

$$T_0 := \nabla \varphi_0.$$

ERM - empirical risk min.

$$\hat{\varphi} = \underset{\varphi \in \Omega}{\operatorname{arg\,min}} \quad \int \varphi d\mu_n + \int \varphi^* d\nu_n.$$

$$\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad \nu_n = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$$

$$\left\{ \hat{T} = \nabla \hat{\varphi} \right\}.$$

→ Makkula et al. (2020)

Input Convex NN (Amos et al.).