

CS480/680: Introduction to Machine Learning

Lec 18: Optimal Transport

Yaoliang Yu

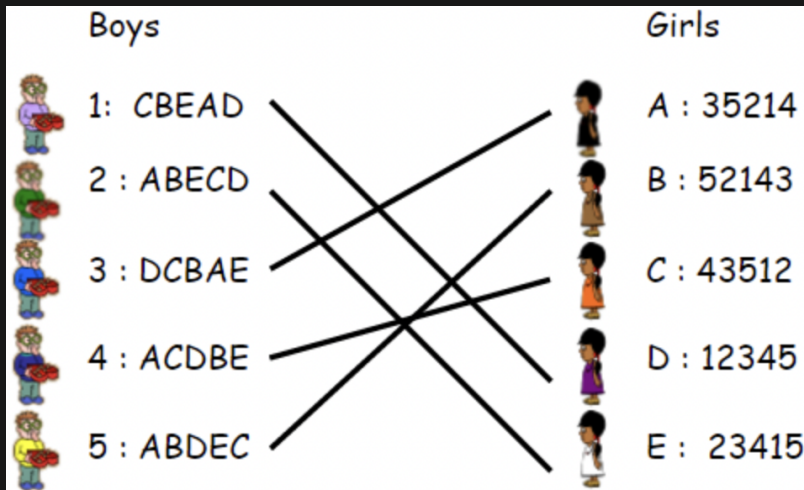


UNIVERSITY OF
WATERLOO

FACULTY OF MATHEMATICS
**DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE**

July 10, 2024

Matching Problem



<https://towardsdatascience.com/stable-matching-as-a-game-a68c279d70b>

- Matching co-op students/organ donors with companies/patients

Stable Matching

Definition: Blocking pair

A pair (i, j) and (i', j') where both i and j' would prefer to swap.

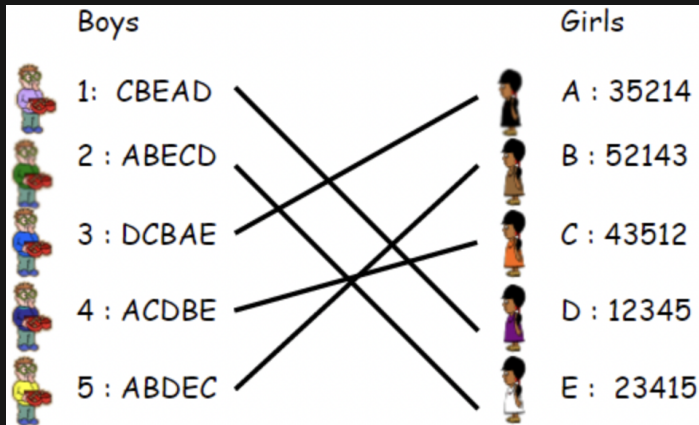
Example: Blocking pair

$(1, E)$ and $(2, C) \rightarrow (1, C)$ and $(2, E)$

- A stable matching is one when there is no blocking pair
 - Gale-Shapley algorithm
- More generally, can define a cost $c(i, j)$ for matching i -th boy with j -th girl
- Need $c(i, j) + c(i', j') \leq c(i, j') + c(i', j)$

Monge's Formulation

$$\min_{\mathbf{T}: [n] \rightarrow [n] \text{ bijective}} \sum_{i=1}^n c(i, \mathbf{T}(i))$$



From Discrete to Continuous

$$\min_{\mathbf{T}_{\#p=q}} \mathbb{E}[c(\mathbf{X}, \mathbf{T}(\mathbf{X}))]$$

- A distribution p of boys and a distribution q of girls
- Let $\mathbf{X} \sim p$ be a random boy
- $\mathbf{T}(\mathbf{X})$ is the matching for \mathbf{X}
- Require $\mathbf{T}_{\#p} = q$ to preserve mass

G. Monge. "Mémoire sur la théorie des déblais et des remblais". In: *Histoire de l'Académie royale des sciences avec les mémoires de mathématique et de physique tirés des registres de cette Académie*. 1781, pp. 666–705.

Kantorovich's Relaxation

$$\min_{\mathbf{T}_{\#p=q}} \mathbb{E}[c(\mathbf{X}, \mathbf{T}(\mathbf{X}))] \geq \min_{\mathbf{X} \sim p, \mathbf{Y} \sim q} \mathbb{E}[c(\mathbf{X}, \mathbf{Y})]$$

Definition: Coupling

$(\mathbf{X}, \mathbf{Y}) \sim \pi$, where the joint coupling π has marginals p and q

- **Deterministic** pairing: \mathbf{x} is matched with **some** $\mathbf{y} = \mathbf{T}\mathbf{x}$
- **Stochastic** pairing: \mathbf{x} is matched to **every** \mathbf{y} with probability $\pi(\mathbf{y}|\mathbf{x})$
- Surprisingly, at optimality, $\pi(\mathbf{y}|\mathbf{x})$ could be deterministic anyway!

L. V. Kantorovich. "On the Translocation of Masses". *Journal of Mathematical Sciences*, vol. 133, no. 4 (2006). Originally published in Dokl. Akad. Nauk SSSR, vol. 37, No. 7–8, 227–229 (1942)., pp. 1381–1382, L. V. Kantorovich. "On a Problem of Monge". *Journal of Mathematical Sciences*, vol. 133, no. 4 (2006). Originally published in Uspekhi Mat. Nauk, vol. 3, No. 2, 225–226 (1948)., pp. 1383–1383.

Duality

$$\min_{\mathbf{X} \sim p, \mathbf{Y} \sim q, (\mathbf{X}, \mathbf{Y}) \sim \pi} \mathbb{E}[c(\mathbf{X}, \mathbf{Y})] = \min_{\pi \geq 0} \int c(\mathbf{x}, \mathbf{y}) \pi(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y}$$
$$\text{s.t.} \quad \int \pi(\mathbf{x}, \mathbf{y}) \, d\mathbf{y} = p(\mathbf{x}), \quad \int \pi(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} = q(\mathbf{y})$$

$$\begin{aligned} & \min_{\pi \geq 0} \max_{u(\cdot), v(\cdot)} \int [c(\mathbf{x}, \mathbf{y}) - u(\mathbf{x}) - v(\mathbf{y})] \pi(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} + \int u(\mathbf{x}) p(\mathbf{x}) \, d\mathbf{x} + \int v(\mathbf{y}) p(\mathbf{y}) \, d\mathbf{y} \\ &= \max_{u(\cdot), v(\cdot)} \min_{\pi \geq 0} \int [c(\mathbf{x}, \mathbf{y}) - u(\mathbf{x}) - v(\mathbf{y})] \pi(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} + \int u(\mathbf{x}) p(\mathbf{x}) \, d\mathbf{x} + \int v(\mathbf{y}) p(\mathbf{y}) \, d\mathbf{y} \\ &= \max_{u(\cdot), v(\cdot)} \int u(\mathbf{x}) p(\mathbf{x}) \, d\mathbf{x} + \int v(\mathbf{y}) p(\mathbf{y}) \, d\mathbf{y}, \quad \text{s.t.} \quad u(\mathbf{x}) + v(\mathbf{y}) \leq c(\mathbf{x}, \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \\ &= \max_{u, v} \mathbb{E}[u(\mathbf{X})] + \mathbb{E}[v(\mathbf{Y})], \quad \text{s.t.} \quad u(\mathbf{x}) + v(\mathbf{y}) \leq c(\mathbf{x}, \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \end{aligned}$$

Conjugacy

$$\forall \mathbf{x}, \forall \mathbf{y}, \quad u(\mathbf{x}) + v(\mathbf{y}) \leq c(\mathbf{x}, \mathbf{y})$$

- $u(\mathbf{x}) \leq [\inf_{\mathbf{y}} c(\mathbf{x}, \mathbf{y}) - v(\mathbf{y})] =: v^c(\mathbf{x})$
- $v(\mathbf{y}) \leq [\inf_{\mathbf{x}} c(\mathbf{x}, \mathbf{y}) - u(\mathbf{x})] =: u^c(\mathbf{y})$
- Since we are maximizing $\mathbb{E}[u(\mathbf{X})] + \mathbb{E}[v(\mathbf{Y})]$, at optimality:

$$u(\mathbf{x}) = v^c(\mathbf{x}), \quad v(\mathbf{y}) = u^c(\mathbf{y})$$

- $u^{cc} \geq u$ and $u^{ccc} = u^c$; similarly for v
- u is called c -concave iff $u = u^{cc}$ (or equivalently $u = v^c$ for some v)

Complementarity

$$\max_{u,v} \min_{\pi \geq 0} \int [c(\mathbf{x}, \mathbf{y}) - u(\mathbf{x}) - v(\mathbf{y})] \pi(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} + \int u(\mathbf{x}) p(\mathbf{x}) \, d\mathbf{x} + \int v(\mathbf{y}) p(\mathbf{y}) \, d\mathbf{y}$$

- $\pi(\mathbf{x}, \mathbf{y}) > 0 \implies u(\mathbf{x}) + v(\mathbf{y}) = c(\mathbf{x}, \mathbf{y})$
- Recall that $u^c = v$, we define the subdifferential:

$$\partial u(\mathbf{x}) := \operatorname{argmin}_{\mathbf{y}} [c(\mathbf{x}, \mathbf{y}) - u^c(\mathbf{y})] = \{\mathbf{y} : u(\mathbf{x}) + u^c(\mathbf{y}) = c(\mathbf{x}, \mathbf{y})\}$$

– for a c -concave u , $\mathbf{y} \in \partial u(\mathbf{x}) \iff \mathbf{x} \in \partial u^c(\mathbf{y})$

- Thus, $\operatorname{supp} \pi \subseteq \operatorname{gph} \partial u$

– in particular, if u is differentiable, π is deterministic and the Kantorovich relaxation is tight!

Cyclic Monotonicity

Definition: Cyclic monotonicity

We call a set $\Gamma \subseteq \mathbb{X} \times \mathbb{Y}$ c -cyclically monotone if for any n and $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n) \in \Gamma$, any (cyclic) permutation $\sigma : [n] \rightarrow [n]$, we always have

$$\sum_{i=1}^n c(\mathbf{x}_i, \mathbf{y}_i) \leq \sum_{i=1}^n c(\mathbf{x}_i, \mathbf{y}_{\sigma(i)})$$

Theorem: Optimal coupling

Suppose $\mathbb{X} = \mathbb{R}^{d_x}$ and $\mathbb{Y} = \mathbb{R}^{d_y}$, $c : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{R}$ is continuous, and p and q are probability densities. There exists an optimal coupling π whose support is c -cyclically monotone. Moreover, there exists a c -concave function u such that $\text{supp } \pi \subseteq \text{gph } \partial u$.

1-Wasserstein Distance

- $c(\mathbf{x}, \mathbf{y}) = d(\mathbf{x}, \mathbf{y})$ for some distance metric d :

$$\begin{aligned} \mathbb{W}_1(p, q) &:= \min_{(\mathbf{X}, \mathbf{Y}) \sim \pi, \mathbf{X} \sim p, \mathbf{Y} \sim q} \mathbb{E}[d(\mathbf{X}, \mathbf{Y})] \\ &= \max_{u, v} \mathbb{E}[u(\mathbf{X})] + \mathbb{E}[v(\mathbf{Y})], \quad \text{s.t.} \quad \forall \mathbf{x}, \mathbf{y}, \quad u(\mathbf{x}) + v(\mathbf{y}) \leq d(\mathbf{x}, \mathbf{y}) \end{aligned}$$

- Lipschitz envelope: $v^c(\mathbf{x}) := [\inf_{\mathbf{y}} d(\mathbf{x}, \mathbf{y}) - v(\mathbf{y})]$
 - v^c is Lipschitz continuous: $|v^c(\mathbf{x}) - v^c(\mathbf{z})| \leq d(\mathbf{x}, \mathbf{z})$
 - v^c is the largest Lipschitz continuous function majorized by $-v$
- Thus, $u = v^c = -v$ and hence

$$\mathbb{W}_1(p, q) = \max_u \mathbb{E}[u(\mathbf{X})] - \mathbb{E}[u(\mathbf{Y})], \quad \text{s.t.} \quad \forall \mathbf{x}, \mathbf{y}, \quad u(\mathbf{x}) - u(\mathbf{y}) \leq d(\mathbf{x}, \mathbf{y})$$

Wasserstein GAN

$$\begin{aligned}\min_{\mathbf{T}} W_1(q, \mathbf{T}_{\#}r) &= \min_{\mathbf{T}} \max_u \mathbb{E}_{\mathbf{X} \sim q}[u(\mathbf{X})] - \mathbb{E}_{\mathbf{Z} \sim r}[u(\mathbf{T}(\mathbf{Z}))], \quad \text{s.t. } u \text{ is Lipschitz} \\ &\approx \min_{\mathbf{T}} \max_u \hat{\mathbb{E}}_{\mathbf{X} \sim q}[u(\mathbf{X})] - \hat{\mathbb{E}}_{\mathbf{Z} \sim r}[u(\mathbf{T}(\mathbf{Z}))], \quad \text{s.t. } u \text{ is Lipschitz}\end{aligned}$$

- r is the noise density, e.g., standard normal
- q is the data density: only a training sample is available
- \mathbf{T} is the generator network: maps noise to data
- u is the discriminator network: maps data to a real scalar

2-Wasserstein Distance

- $c(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ (note the square):

$$\begin{aligned} W_2^2(p, q) &:= \min_{(\mathbf{X}, \mathbf{Y}) \sim \pi, \mathbf{X} \sim p, \mathbf{Y} \sim q} \mathbb{E}[\tfrac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_2^2] \\ &= \max_{u, v} \mathbb{E}[u(\mathbf{X})] + \mathbb{E}[v(\mathbf{Y})], \quad \text{s.t. } \forall \mathbf{x}, \mathbf{y}, \quad u(\mathbf{x}) + v(\mathbf{y}) \leq \tfrac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \end{aligned}$$

- Conjugate: $v^c(\mathbf{x}) := [\inf_{\mathbf{y}} \tfrac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 - v(\mathbf{y})]$
– $\tfrac{1}{2} \|\mathbf{x}\|_2^2 - v^c(\mathbf{x}) = \sup_{\mathbf{y}} \langle \mathbf{x}, \mathbf{y} \rangle - [\tfrac{1}{2} \|\mathbf{y}\|_2^2 - v(\mathbf{y})]$: convex conjugate
- Thus, $u = v^c = \tfrac{1}{2} \|\cdot\|_2^2 - (\tfrac{1}{2} \|\cdot\|_2^2 - v)^*$ and hence

$$W_2^2(p, q) = \max_f \mathbb{E}[\tfrac{1}{2} \|\mathbf{X}\|_2^2 - f(\mathbf{X})] + \mathbb{E}[\tfrac{1}{2} \|\mathbf{Y}\|_2^2 - f^*(\mathbf{Y})], \quad \text{s.t. } f \text{ is convex}$$

$$\pi = (\text{Id} \times \partial f)_{\#} p, \quad \text{in particular } q = (\partial f)_{\#} p$$

Potential GAN

$$\begin{aligned}\min_{\mathbf{T}} W_2^2(q, \mathbf{T}_{\#}r) &= \min_{\mathbf{T}} \max_f \mathbb{E}_{\mathbf{X} \sim q} [\tfrac{1}{2} \|\mathbf{X}\|_2^2 - f(\mathbf{X})] + \mathbb{E}_{\mathbf{Z} \sim r} [\tfrac{1}{2} \|\mathbf{T}(\mathbf{Z})\|_2^2 - f^*(\mathbf{T}(\mathbf{Z}))], \\ &\approx \min_{\mathbf{T}} \max_f \hat{\mathbb{E}}_{\mathbf{X} \sim q} [-f(\mathbf{X})] + \hat{\mathbb{E}}_{\mathbf{Z} \sim r} [\tfrac{1}{2} \|\mathbf{T}(\mathbf{Z})\|_2^2 - f^*(\mathbf{T}(\mathbf{Z}))], \text{ s.t. } f \text{ is convex}\end{aligned}$$

- r is the noise density, e.g., standard normal
- q is the data density: only a training sample is available
- \mathbf{T} is the generator network: maps noise to data
- f is the discriminator network: maps data to a real scalar

T. Salimans, H. Zhang, A. Radford, and D. Metaxas. "Improving GANs Using Optimal Transport". In: *International Conference on Learning Representations*. 2018, H. Liu, X. Gu, and D. Samaras. "Wasserstein GAN With Quadratic Transport Cost". In: *IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019, pp. 4831–4840.

Potential Flow

$$\min_f \mathbb{D}(q, (\nabla f)_\# r), \quad \text{s.t. } f \text{ is convex}$$

- r is the noise density, e.g., standard normal
- q is the data density: only a training sample is available
- ∇f is the generator network: maps noise to data
 - e.g., f is a Relu network with nonnegative weights
- \mathbb{D} is some “distance” function, e.g., the KL divergence

Triangular vs. Potential

- $\mathbf{T} : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $\mathbf{T}_{\#}p = q$
- \mathbf{T} is autoregressive
- $\nabla \mathbf{T}$ is always triangular
- composition holds
- no rotational equivariance

- $\mathbf{T} : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $\mathbf{T}_{\#}p = q$
- $\mathbf{T} = \nabla f$ for convex $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- $\nabla \mathbf{T} = \nabla^2 f$ is symmetric PSD
- composition fails
- rotationally equivariant

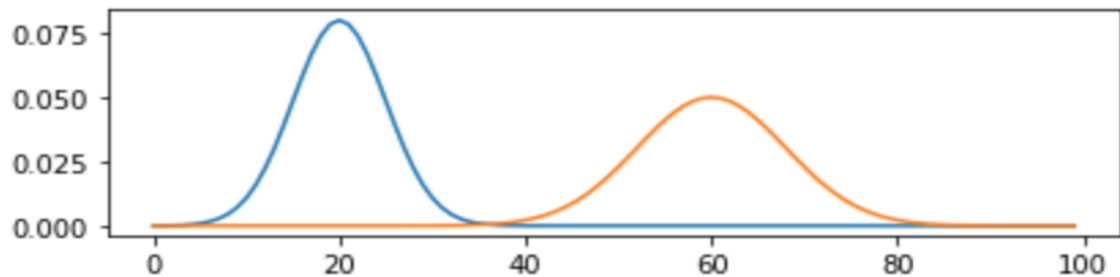
The two are equivalent iff \mathbf{T} is diagonal, in particular, if $d = 1$

Wasserstein Barycenter

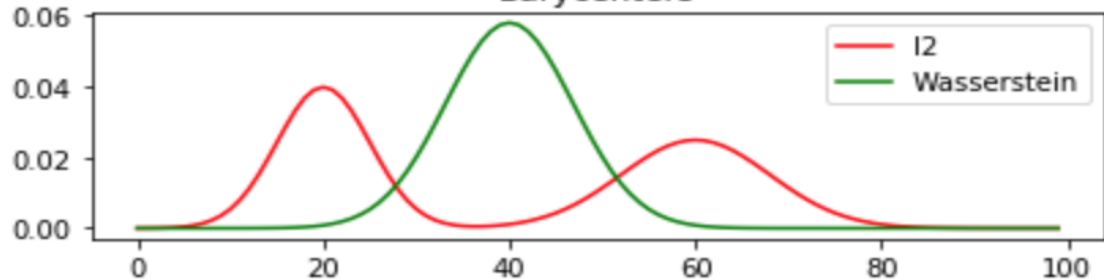
- Consider densities p_0 and p_1 , say two Gaussians with different mean and variance
- How to interpolate between them?
- Exists convex f such that $p_1 = (\nabla f)_\# p_0$
- Obviously $p_0 = (\text{Id})_\# p_0$
- Interpolate the push-forward maps!

$$p_t = [(1-t)\text{Id} + t\nabla f]_\# p_0 = \underset{p}{\operatorname{argmin}} (1-t)\mathbb{W}_2^2(p, p_0) + t\mathbb{W}_2^2(p, p_1)$$

Distributions



Barycenters



Topics in Optimal Transportation

Cédric Villani

Graduate Studies
in Mathematics
Volume 58



American Mathematical Society

Progress in Nonlinear Differential Equations
and Their Applications
87

Filippo Santambrogio

Optimal
Transport
for Applied
Mathematicians

Calculus of Variations, PDEs, and
Modeling

 Birkhäuser

Foundations and Trends® in
Machine Learning
11:5-6

Computational Optimal
Transport

With Applications
to Data Science

Gabriel Peyré and Marco Cuturi

now

the essence of knowledge

EMS TEXTBOOKS IN MATHEMATICS

Alessio Figalli
Federico Glaudo

An Invitation to
Optimal Transport,
Wasserstein Distances,
and Gradient Flows

==

EMS
PRESS

