

Optimal Transport

Makoto Yamada (Kyoto University)

- 1 Notation
- 2 Transportation problem
- 3 Wasserstein distance
- 4 Entropic regularized OT
- 5 Dual problem
- 6 Summary

- 1 Notation
- 2 Transportation problem
- 3 Wasserstein distance
- 4 Entropic regularized OT
- 5 Dual problem
- 6 Summary

Notation

Matrix operation

- scalar: $x \in \mathbb{R}$
- vector: $\mathbf{x} = (x_1, x_2, \dots, x_D)^\top \in \mathbb{R}^D$
- matrix: $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^\top \in \mathbb{R}^{D \times N}$
- transpose: $^\top$
- ℓ_2 norm: $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^D x_i^2}$ (Euclid norm)

Probability

- random variable: X
- probability density: $p(\cdot)$

- 1 Notation
- 2 Transportation problem
- 3 Wasserstein distance
- 4 Entropic regularized OT
- 5 Dual problem
- 6 Summary

Transportation problem

The optimal transport has been studied since the 17th century (many Nobel Prizes and Fields Medals recipients). In particular, in the last few years, a great deal of research has been reported in the field of machine learning.

- Wasserstein GAN (Machine learning)
- Earth Mover's Distance (Computer vision)
- Word Mover's Distance (NLP)

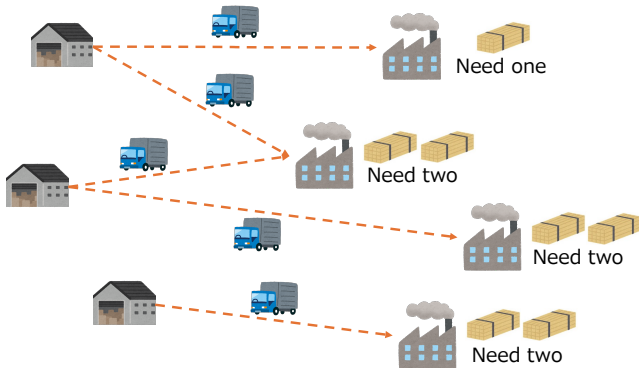
輸送問題

The problem of transporting goods from one location (warehouse) to another (factory). How is the best way to transport the goods?



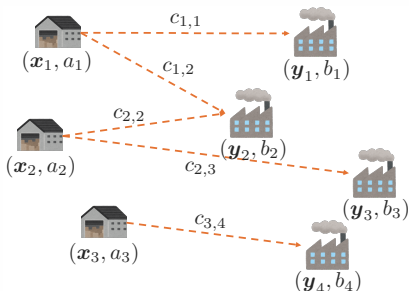
Transportation problem

Use trucks to transport goods (**costly**). Consider transportation that is less costly and satisfies the constraints.



Transportation problem

The percentage of the quantity of goods in the warehouse is $\mathbf{a} = (a_1, a_2, \dots, a_N)^\top$, the percentage of consumption in the factory is $\mathbf{b} = (b_1, b_2, \dots, b_N)^\top$, $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}^D$ are the locations of the warehouses and factories. transportation costs is $c_{i,j} = c(\mathbf{x}_i, \mathbf{y}_j)$. The quantity to be transported (by truck) is $[\mathbf{P}]_{ij} = p_{i,j}$.



Optimal transport (OT)

Optimal transport problem

$$\min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} c(\mathbf{x}_i, \mathbf{y}_j)$$

or

$$\min_{P \in \mathbb{R}_+^{N \times M}} \sum_{i,j} p_{ij} c(\mathbf{x}_i, \mathbf{y}_j) \quad \text{s.t.} \quad P \mathbf{1}_M = \mathbf{a}, P^\top \mathbf{1}_N = \mathbf{b}$$

- $\mu = \sum_{i=1}^N a_i \delta_{\mathbf{x}_i}$, $\nu = \sum_{j=1}^M b_j \delta_{\mathbf{y}_j}$ are discrete measures.
- $U(\mu, \nu) = \{P \in \mathbb{R}_+^{N \times M} : P \mathbf{1}_M = \mathbf{a}, P^\top \mathbf{1}_N = \mathbf{b}\}$
- $\mathbf{a}^\top \mathbf{1}_N = 1$ (probability)
- $\mathbf{b}^\top \mathbf{1}_M = 1$ (probability)
- $P \mathbf{1}_M = \mathbf{a}, P^\top \mathbf{1}_N = \mathbf{b}$.

- 1 Notation
- 2 Transportation problem
- 3 Wasserstein distance**
- 4 Entropic regularized OT
- 5 Dual problem
- 6 Summary

Wasserstein distance

p-Wasserstein distance

$$W_p(\mu, \nu) = \left(\min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} d(\mathbf{x}_i, \mathbf{y}_j)^p \right)^{1/p}$$

$d(\mathbf{x}, \mathbf{y})$ is a distance. (e.g., Euclid distance

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2)$$

1-Wasserstein distance

$$W_1(\mu, \nu) = \min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} d(\mathbf{x}_i, \mathbf{y}_j)$$

2-Wasserstein distance

$$W_2(\mu, \nu) = \left(\min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} d(\mathbf{x}_i, \mathbf{y}_j)^2 \right)^{1/2}$$

Wasserstein distance

2-Wasserstein distance

$$W_2(\mu, \nu) = \left(\min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} d(\mathbf{x}_i, \mathbf{y}_j)^2 \right)^{1/2}$$

2-Wasserstein distance If $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2$ is a Euclidean distance, we have

$$W_2(\mu, \nu) = \left(\min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} \|\mathbf{x}_i - \mathbf{y}_j\|^2 \right)^{1/2}$$

$\min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} \|\mathbf{x}_i - \mathbf{y}_j\|^2 = (W_2(\mu, \nu))^2$ is not a distance!

- 1 Notation
- 2 Transportation problem
- 3 Wasserstein distance
- 4 Entropic regularized OT**
- 5 Dual problem
- 6 Summary

Entropic regularized OT

Optimization problem

$$\min_{P \in U(\mu, \nu)} \sum_{i,j} p_{ij} c(\mathbf{x}_i, \mathbf{y}_j) + \lambda \sum_{i,j} p_{ij} (\log(p_{ij}) - 1)$$

$\lambda \geq 0$ is a regularization parameter.

- $c(\mathbf{x}, \mathbf{y})$ does **not need to be a distance**.
- When $\lambda = 0$, it is Earth Mover's Distance (EMD).
- **Strongly convex function** ($0 < p_k \leq 1$).

$$\frac{\partial}{\partial p_k \partial p_{k'}} \sum_i p_i (\log(p_i) - 1) = \frac{\partial}{\partial p_{k'}} \log(p_k) = \begin{cases} \frac{1}{p_k} & k = k' \\ 0 & k \neq k' \end{cases}$$

For $p \in [0, 1]$, $p \log(p)$ is convex & The Hessian is positive definite \rightarrow Strongly convex.

Sinkhorn Algorithm

Optimization problem

$$\begin{aligned} \min_{\mathbf{P}} \quad & \sum_{i,j} p_{ij} c(\mathbf{x}_i, \mathbf{y}_j) + \lambda \sum_{i,j} p_{ij} (\log(p_{ij}) - 1) \\ \text{s.t.} \quad & \mathbf{P} \mathbf{1}_M = \mathbf{a}, \mathbf{P}^\top \mathbf{1}_N = \mathbf{b} \end{aligned}$$

Using Lagrange multipliers $\tilde{\mathbf{u}} \in \mathbb{R}^N, \tilde{\mathbf{v}} \in \mathbb{R}^M$

$$\begin{aligned} J(\mathbf{P}, \mathbf{u}, \mathbf{v}) = & \sum_{i,j} p_{ij} c(\mathbf{x}_i, \mathbf{y}_j) + \lambda \sum_{i,j} p_{ij} (\log(p_{ij}) - 1) \\ & + \tilde{\mathbf{u}}^\top (\mathbf{a} - \mathbf{P} \mathbf{1}_M) + \tilde{\mathbf{v}}^\top (\mathbf{b} - \mathbf{P}^\top \mathbf{1}_N) \end{aligned}$$

Sinkhorn Algorithm

Taking derivative w.r.t. p_{kl} , and solving for p_{kl} .

$$\frac{\partial}{\partial p_{kl}} J(\mathbf{P}, \mathbf{u}, \mathbf{v}) = c(\mathbf{x}_k, \mathbf{y}_l) + \lambda \log(p_{kl}) - u_k - v_l = 0$$

The optimal solution of p_{kl} is given as

$$\begin{aligned}\log(p_{kl}) &= -\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l) + \frac{1}{\lambda} \tilde{u}_k + \frac{1}{\lambda} \tilde{v}_l \\ p_{kl} &= \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l)\right) \exp\left(\frac{1}{\lambda} \tilde{u}_k\right) \exp\left(\frac{1}{\lambda} \tilde{v}_l\right) \\ p_{kl} &= \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l)\right) \exp\left(\frac{1}{\lambda} \tilde{u}_k\right) \exp\left(\frac{1}{\lambda} \tilde{v}_l\right) \\ &= \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l)\right) u_k v_l\end{aligned}$$

Sinkhorn Algorithm

Updating v_l :

$$v_l \sum_{k=1}^N \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l)\right) u_k = \sum_{k=1}^N p_{kl}$$

$$v_l = \frac{b_l}{\sum_{k=1}^N \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l)\right) u_k}$$

Matrix version:

$$\mathbf{v}^{(\ell+1)} = \mathbf{b} / (\mathbf{K}^\top \mathbf{u}^{(\ell)})$$

Note: $\mathbf{K} \in \mathbb{R}^{n \times m}$, $[\mathbf{K}]_{ij} = \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_i, \mathbf{y}_j)\right)$.

Sinkhorn Algorithm

Updating u_k :

$$v_k \sum_{l=1}^N \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l)\right) v_l = \sum_{l=1}^M p_{kl}$$

$$u_k = \frac{a_k}{\sum_{l=1}^M \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_k, \mathbf{y}_l)\right) v_l}$$

Matrix version:

$$\mathbf{u}^{(\ell+1)} = \mathbf{a} / (\mathbf{K} \mathbf{v}^{(\ell)})$$

Note: $\mathbf{K} \in \mathbb{R}^{n \times m}$, $[\mathbf{K}]_{ij} = \exp\left(-\frac{1}{\lambda} c(\mathbf{x}_i, \mathbf{y}_j)\right)$.

- 1 Notation
- 2 Transportation problem
- 3 Wasserstein distance
- 4 Entropic regularized OT
- 5 Dual problem**
- 6 Summary

Dual problem

Optimization problem

$$\min_P \sum_{i,j} p_{ij} c(\mathbf{x}_i, \mathbf{y}_j) \text{ s.t. } P\mathbf{1}_M = \mathbf{a}, P^\top \mathbf{1}_N = \mathbf{b}$$

Using the Lagrange multipliers $\mathbf{u} \in \mathbb{R}^N, \mathbf{v} \in \mathbb{R}^M$ We denote $\langle P, C \rangle = \sum_{i,j} p_{ij} c(\mathbf{x}_i, \mathbf{y}_j)$, then we have

$$\min_{P \geq 0} \max_{(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^N \times \mathbb{R}^M} \langle P, C \rangle + \mathbf{u}^\top (\mathbf{a} - P\mathbf{1}_M) + \mathbf{v}^\top (\mathbf{b} - P^\top \mathbf{1}_N)$$

For a linear problem, min and max can be exchangeable

$$\max_{(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^N \times \mathbb{R}^M} \langle \mathbf{u}, \mathbf{a} \rangle + \langle \mathbf{v}, \mathbf{b} \rangle + \min_{P \geq 0} \langle P, C \rangle - \mathbf{u}^\top P\mathbf{1}_M - \mathbf{v}^\top P^\top \mathbf{1}_N$$

Dual problem

We use $\mathbf{u}^\top \mathbf{P} \mathbf{1}_M = \text{tr}(\mathbf{u}^\top \mathbf{P} \mathbf{1}_M) = \text{tr}(\mathbf{P} \mathbf{1}_M \mathbf{u}^\top) = \text{tr}(\mathbf{u} \mathbf{1}_M^\top \mathbf{P}^\top) = \langle \mathbf{P}, \mathbf{u} \mathbf{1}_M^\top \rangle$ and $\mathbf{v}^\top \mathbf{P}^\top \mathbf{1}_N = \langle \mathbf{P}, \mathbf{1}_N \mathbf{v}^\top \rangle$.

$$\max_{(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^N \times \mathbb{R}^M} \langle \mathbf{u}, \mathbf{a} \rangle + \langle \mathbf{v}, \mathbf{b} \rangle + \min_{\mathbf{P} \geq 0} \langle \mathbf{P}, \mathbf{C} \rangle - \langle \mathbf{P}, \mathbf{u} \mathbf{1}_M^\top \rangle - \langle \mathbf{P}, \mathbf{1}_N \mathbf{v}^\top \rangle$$

$$\max_{(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^N \times \mathbb{R}^M} \langle \mathbf{u}, \mathbf{a} \rangle + \langle \mathbf{v}, \mathbf{b} \rangle + \min_{\mathbf{P} \geq 0} \langle \mathbf{P}, \mathbf{C} - \mathbf{u} \mathbf{1}_M^\top - \mathbf{1}_N \mathbf{v}^\top \rangle$$

$$\min_{\mathbf{P} \geq 0} \langle \mathbf{P}, \mathbf{C} - \mathbf{u} \mathbf{1}_M^\top - \mathbf{1}_N \mathbf{v}^\top \rangle = \begin{cases} 0 & \mathbf{C} - \mathbf{u} \oplus \mathbf{v} \geq 0 \\ -\infty & \text{Otherwise} \end{cases}$$

$\mathbf{u} \oplus \mathbf{v} = \mathbf{u} \mathbf{1}_M^\top - \mathbf{1}_N \mathbf{v}^\top$. If the following condition satisfied, the second term is zero.

$$\mathbf{u} \oplus \mathbf{v} \leq \mathbf{C}$$

Dual problem

Dual problem

$$\max_{(\mathbf{u}, \mathbf{v}) \in R(\mathbf{C})} \langle \mathbf{u}, \mathbf{a} \rangle + \langle \mathbf{v}, \mathbf{b} \rangle$$

$$R(\mathbf{C}) = \{(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^N \times \mathbb{R}^M : \forall (i, j) \in \llbracket N \rrbracket \times \llbracket M \rrbracket, \mathbf{u} \oplus \mathbf{v} \leq \mathbf{C}\}.$$

$$\llbracket N \rrbracket = \{1, 2, \dots, N\}.$$

The relationship between primal and dual problems

$$\sum_{i,j} p_{ij} c_{ij} \geq \sum_{i,j} p_{ij} (u_i + v_j) = \sum_{i=1}^N a_i u_i + \sum_{j=1}^M b_j v_j = \langle \mathbf{u}, \mathbf{a} \rangle + \langle \mathbf{v}, \mathbf{b} \rangle$$

Other OT based distances

- Sliced Wasserstein distance
- Tree-Wasserstein distance
- Gromov Wasserstein distance

- 1 Notation
- 2 Transportation problem
- 3 Wasserstein distance
- 4 Entropic regularized OT
- 5 Dual problem
- 6 Summary**

Summary

- Optimal transport (OT)
- Wasserstein distance
- Entropic regularized Optimal transport
- Dual problem