

# Computational Optimal Transport

## Basic Properties and Sinkhorn's Algorithm

WU Haoyu  
Supervisor: CAI Jianfeng

HKUST Math Department

SCIE2500 Presentation, 12 May 2022

# Table of Contents

## 1 Definition of Optimal Transport

- Discrete Case
- Continuous Case

## 2 Entropic Regularization

## 3 Sinkhorn's Algorithm

- Lagrangian Method
- Sinkhorn's Algorithm
- Hilbert Projective Metric
- Convergence of Sinkhorn's algorithm

## 4 Complexity and Application to Solve Original Problem

## 1 Definition of Optimal Transport

- Discrete Case
- Continuous Case

## 2 Entropic Regularization

## 3 Sinkhorn's Algorithm

- Lagrangian Method
- Sinkhorn's Algorithm
- Hilbert Projective Metric
- Convergence of Sinkhorn's algorithm

## 4 Complexity and Application to Solve Original Problem

# Definition of Optimal Transport

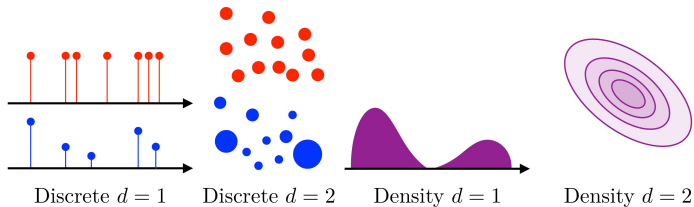


Figure: Mass Distribution

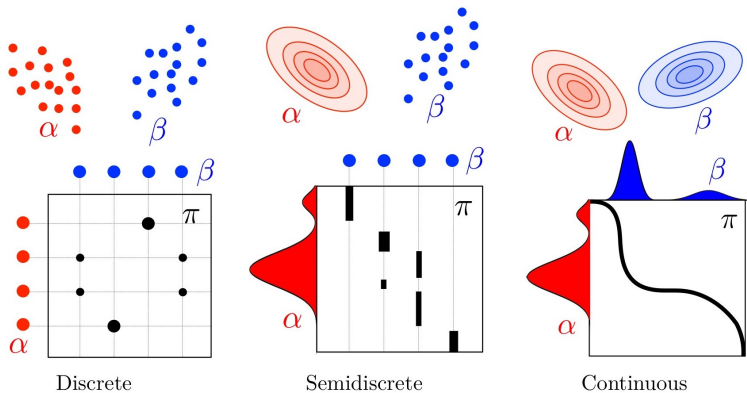


Figure: Coupling Method

# Discrete Case

Using vector to denote the mass distribution

$$a \in \mathbb{R}^m, \sum_{k=1}^m a_k = 1, a_k \geq 0, \quad b \in \mathbb{R}^n, \sum_{k=1}^n b_k = 1, b_k \geq 0$$

Then the coupling matrix set:

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

with given Cost Matrix  $C \in \mathbb{R}_+^{n \times m}$ , where  $C_{ij}$  represents the cost of 1 mass transports from  $j^{th}$  location (of  $a$ ) to  $i^{th}$  location (of  $b$ ).

Then optimal transport problem is:

$$L_C(a, b) \stackrel{\text{def.}}{=} \min_{P \in U(a, b)} \langle C, P \rangle \stackrel{\text{def.}}{=} \sum_{i,j} C_{i,j} P_{i,j}$$

# Continuous Case

Let  $(X, \mathcal{M}_\alpha, \alpha), (X, \mathcal{M}_\beta, \beta)$  be measure space satisfying

$$\int_X d\alpha = \int_X d\beta = 1$$

For any  $A \in \mathcal{M}_\alpha$ ,  $\alpha(A)$  represents the mass located in the set  $A$  (at the beginning), similar for  $\beta(B)$ . Then any coupling method could be represented by a product measure

$$\mu : \mathcal{M}_\alpha \times \mathcal{M}_\beta \rightarrow \mathbb{R}^+ \cup \{0\}$$

where input  $(A, B)$ , output the mass transport from  $A$  to  $B$ .

Noticed

$$\mu(A, B) \leq \min\{\alpha(A), \beta(B)\}, \quad \mu(A, X) = \alpha(A), \mu(X, B) = \beta(B)$$

And the cost of  $\mu$  is

$$\int_X^2 C(x, y) d\mu$$

where  $C : X^2 \rightarrow \mathbb{R}^+ \cup \{0\}$  is the cost function, and the optimal problem is

$$L_C(\alpha, \beta) = \min_{\mu} \int_X^2 C(x, y) d\mu$$

where  $\mu$  is a product measure on  $X^2$  satisfying

$$\mu(A, X) = \alpha(A), \mu(X, B) = \beta(B)$$



Then we just focus on discrete case.

## 1 Definition of Optimal Transport

- Discrete Case
- Continuous Case

## 2 Entropic Regularization

## 3 Sinkhorn's Algorithm

- Lagrangian Method
- Sinkhorn's Algorithm
- Hilbert Projective Metric
- Convergence of Sinkhorn's algorithm

## 4 Complexity and Application to Solve Original Problem

# Entropic Regularization

## Definition(Entropic)

For a coupling matrix  $P$ , the discrete entropy is defined as:

$$H(P) \stackrel{\text{def.}}{=} \sum_{i,j} P_{i,j} (1 - \log(P_{i,j}))$$

## Regularization

it's always greater or equal 0 since  $P_{i,j} \in [0, 1]$ . The function  $-H$  is 1-strongly convex since by computing the Hessian,

$$\partial^2 -H(P) = \text{diag}(1/P_{i,j})$$

and  $P_{i,j} \leq 1$ .

## Regularization

The idea of the entropic regularization is to make  $-H$  the regularizing function ( $L_C^\varepsilon(a, b)$ ) to approach solutions or approximation of  $L_C(a, b)$ , the original problem:

$$L_C^\varepsilon(a, b) \stackrel{\text{def.}}{=} \min_{P \in \mathcal{U}(a, b)} \langle P, C \rangle - \varepsilon H(P)$$

Due to it is  $\varepsilon$ -strongly convex function, then  $L_C^\varepsilon(a, b)$  has a unique optimal solution.

For every different  $\varepsilon > 0$ , the solution  $P_\varepsilon$  of  $L_C^\varepsilon(a, b)$  is unique due to convexity. Then we will claim  $P_\varepsilon$  converges to the original optimal solution with maximal entropic, exactly:

$$P_\varepsilon \longrightarrow \arg \min_{P \in \mathcal{U}(a, b)} \{-H(P) : \langle P, C \rangle = L_C(a, b)\}, \text{ as } \varepsilon \rightarrow 0$$

In particular,

$$L_C^\varepsilon(a, b) \longrightarrow L_C(a, b), \text{ as } \varepsilon \rightarrow 0$$

# Proof of Convergence with $\varepsilon$

## Proof.

Consider a sequence  $(\varepsilon_l)_l$  s.t.  $\varepsilon_l \downarrow 0$ . Denote  $P_l$  the solution of  $L_C^{\varepsilon_l}(a, b)$ . Noticed  $U(a, b)$  is bounded, by Bazano-Weierstrass Theorem, there is a subsequence such that  $P_k \rightarrow P^*$ , and  $P^* \in U(a, b)$  because  $U(a, b)$  is closed. Consider any  $P$  as the solution of  $L_C(a, b)$ , due to optimality:

$$\langle C, P \rangle \leq \langle C, P_k \rangle, \quad \langle C, P_k \rangle + \varepsilon_k H(P) \leq \langle C, P_k \rangle + \varepsilon_k H(P_k)$$

$$\implies 0 \leq \langle C, P_k \rangle - \langle C, P \rangle \leq \varepsilon_k (H(P_k) - H(P))$$

Noticed  $H$  is continuous,  $\langle C, P^* \rangle = \langle C, P \rangle$  as  $k \rightarrow \infty$ . and  $H(P) \leq H(P^*)$ , which means that  $P^*$  is a solution in the set of all optimal solutions of  $L_C(a, b)$  with maximal entropy. By strictly convexity, the solution is unique and just  $P^*$ .



## 1 Definition of Optimal Transport

- Discrete Case
- Continuous Case

## 2 Entropic Regularization

## 3 Sinkhorn's Algorithm

- Lagrangian Method
- Sinkhorn's Algorithm
- Hilbert Projective Metric
- Convergence of Sinkhorn's algorithm

## 4 Complexity and Application to Solve Original Problem

We use Lagrangian to show that the unique solution of  $L_C^\varepsilon(a, b)$  has specific form.

## Proposition

There exist scaling variable  $u \in \mathbb{R}_+^n, v \in \mathbb{R}_+^m$  such that the solution  $P$  of  $L_C^\varepsilon(a, b)$  satisfying:

$$P_{i,j} = u_i K_{i,j} v_j, \quad 1 \leq i \leq n, 1 \leq j \leq m, i, j \in \mathbb{N}$$

where  $K \in \mathbb{R}^{n \times m}$  with  $K_{i,j} = e^{-\frac{c_{i,j}}{\varepsilon}}$ .

# Proof of Proposition

## Proof.

Consider the Lagrangian of  $L_C^\varepsilon(a, b)$  with variable  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$  for all marginal constraints:

$$L(P, x, y) = \langle C, P \rangle - \varepsilon H(P) - \langle x, P \mathbb{1}_m - a \rangle - \langle y, P^T \mathbb{1}_n - b \rangle$$

By the first order condition,

$$\frac{\partial L(P, x, y)}{\partial P_{i,j}} = C_{i,j} + \varepsilon \log(P_{i,j}) - x_i - y_j = 0$$

$$\implies P_{i,j} = e^{\frac{x_i}{\varepsilon}} e^{\frac{-C_{i,j}}{\varepsilon}} e^{\frac{y_j}{\varepsilon}}, \quad (u)_i = e^{\frac{x_i}{\varepsilon}}, (v)_j = e^{\frac{y_j}{\varepsilon}}$$





# Sinkhorn's Algorithm

Noticed for the corresponding  $(u, v)$ ,  $P = \text{diag}(u)K\text{diag}(v)$ . Therefore the variable  $(u, v)$  must satisfy the restriction of  $U(a, b)$ :

$$u \odot (Kv) = a, v \odot (K^T u) = b$$

where  $\odot$  is the vector product under entrywise  $((a \odot b)_i = a_i b_i)$ . The naturally way to find  $(u, v)$  is solve them iteratively, and the Sinkhorn's algorithm is defined by the idea.

## Definition(Sinkhorn's Algorithm)

$$u^{(k+1)} \stackrel{\text{def.}}{=} \frac{a}{Kv^{(k)}}, \quad v^{(k+1)} \stackrel{\text{def.}}{=} \frac{b}{Ku^{(k+1)}}$$

where the division is also entrywise and initialized with an arbitrary positive vector  $v^{(0)} = \mathbb{1}_m$ .

## Remark

Noticed that if Sinkhorn's algorithm converges for any initialization (will be proved later), different initialization may obtain different limit up to a multiplicative constant since if  $(u, v)$  satisfies the  $U(a, b)$  restriction, then  $(\lambda u, \lambda^{-1} v)$  also.

Then we prove the global convergence of Sinkhorn's algorithm. It's first proved by [Franklin:1989], using the property of Hilbert projective metric introduced from [Birkhoff:1957].

# Hilbert Projective Metric

## Definition

The Hilbert metric is defined on  $\mathbb{R}_+^n$ :

$$\forall (u, u') \in (\mathbb{R}_+^n)^2, \quad d_H(u, u') \stackrel{\text{def.}}{=} \log \max_{i,j} \frac{u_i u'_j}{u_j u'_i}$$

It's actually a distance on the projective cone  $\mathbb{R}_+^n / \sim$  with  $u \sim v$  iff  $\exists r > 0, v = ru$ .

## Important Property

For  $K \in \mathbb{R}_+^{n \times m}$ ,  $(u, u') \in (\mathbb{R}_+^m)^2$ :

$$d_H(Ku, Ku') \leq \lambda(K) d_H(u, u')$$

$$\text{where } \lambda(K) = \frac{\sqrt{\mu(K)} - 1}{\sqrt{\mu(K)} + 1} < 1, \mu(K) = \max_{i,j,k,l} \frac{K_{i,j} K_{k,l}}{K_{i,l} K_{j,k}}.$$

# Important Property

For  $K \in \mathbb{R}_+^{n \times m}$ ,  $(u, u') \in (\mathbb{R}_+^m)^2$ :

$$d_H(Ku, Ku') \leq \lambda(K) d_H(u, u'), \text{ where } \begin{cases} \lambda(K) = \frac{\sqrt{\mu(K)} - 1}{\sqrt{\mu(K)} + 1} < 1 \\ \mu(K) = \max_{i,j,k,l} \frac{K_{i,j}K_{k,l}}{K_{i,l}K_{j,k}} \end{cases}$$

## Remark

The inequality is also used to prove the Perron-Frobenius theorem, known as a matrix cases of contraction mapping theorem.

## Remark

The proof is basically first shown on  $\mathbb{R}_+^2$  then generalize to common case.

# Convergence of Sinkhorn's algorithm

We first has the unique solution  $P^*$  of  $L_C^\varepsilon(a, b)$ , by [Franklin;1989], there exist  $(u^*, v^*) \in \mathbb{R}^m \times \mathbb{R}^n$  such that:

$$P^* = \text{diag}(u^*)K\text{diag}(v^*), \quad u^* \odot (Kv^*) = a, v^* \odot (K^T u^*) = b$$

## Theorem

We have  $(u^{(k)}, v^{(k)}) \rightarrow (u^*, v^*)$  under Sinkhorn's algorithm, and:

$$d_H(u^{(k)}, u^*) = O(\lambda(K)^{2k}), \quad d_H(v^{(k)}, v^*) = O(\lambda(K)^{2k}) \quad (1)$$

$$d_H(u^{(k)}, u^*) \leq \frac{d_H(P^{(k)} \mathbb{1}_m, a)}{1 - \lambda(K)^2}, \quad d_H(v^{(k)}, v^*) \leq \frac{d_H(P^{(k),T} \mathbb{1}_n, b)}{1 - \lambda(K)^2} \quad (2)$$

$$\|\log(P^{(k)}) - \log(P^*)\|_\infty \leq d_H(u^{(k)}, u^*) + d_H(v^{(k)}, v^*) \quad (3)$$

where  $P^{(k)} \stackrel{\text{def.}}{=} \text{diag}(u^{(k)})K\text{diag}(v^{(k)})$ .

# Convergence of Sinkhorn's Algorithm

## Proof of (1)

Noticed by definition of Hilbert Metric,

$$\forall (u, u') \in (\mathbb{R}_+^m)^2, \quad d_H(u, u') = d_H(u/u', \mathbb{1}_m) = d_H\left(\frac{\mathbb{1}_m}{u}, \frac{\mathbb{1}_m}{u'}\right)$$

That means,

$$\begin{aligned} d_H(u^{(k+1)}, u^*) &= d_H\left(\frac{a}{Kv^{(k)}}, \frac{a}{Kv^*}\right) = d_H\left(\frac{\mathbb{1}_m}{Kv^{(k)}}, \frac{\mathbb{1}_m}{Kv^*}\right) \\ &= d_H(Kv^{(k)}, Kv^*) \leq \lambda(K) d_H(v^{(k)}, v^*) \end{aligned}$$

the last step used inequality proved above, and mutatis mutandis,

$$d_H(v^{(k)}, v^*) \leq \lambda(K) d_H(u^{(k)}, u^*) \implies d_H(u^{(k+1)}, u^*) \leq \lambda^2(K) d_H(u^{(k)}, u^*)$$

which proves (1).

# Convergence of Sinkhorn's Algorithm

## Proof of (2)

Consider Triangular inequality,

$$\begin{aligned}d_H(u^{(k)}, u^*) &\leq d_H(u^{(k+1)}, u^{(k)}) + d_H(u^{(k+1)}, u^*) \\&\leq d_H\left(\frac{a}{K_V^{(k)}}, u^{(l)}\right) + \lambda^2(K) d_H(u^{(k)}, u^*) \\&= d_H\left(a, u^{(k)} \odot K_V^{(k)}\right) + \lambda^2(K) d_H(u^{(k)}, u^*) \\&= d_H\left(a, P^{(k)} \mathbb{1}_m\right) + \lambda^2(K) d_H(u^{(k)}, u^*) \\&\implies d_H(u^{(k)}, u^*) \leq \frac{d_H(P^{(k)} \mathbb{1}_m, a)}{1 - \lambda(K)^2}\end{aligned}$$

This proves the first part of (2), the latter part is similar.

# Convergence of Sinkhorn's Algorithm

## Proff of (3)

Denote  $M_k = \exp(d_H(u^{(k)}, u^*) + d_H(v^{(k)}, v^*)) > 1$ . Noticed  $P^{(k)} = \text{diag}(u^{(k)})K\text{diag}(v^{(k)})$ ,  $P^{(k)} = \text{diag}(u^*)K\text{diag}(v^*)$ :

$$P^* = \text{diag}(u^*/u^{(k)})P^{(k)}\text{diag}(v^*/v^{(k)})$$

$$d_H(u^*, u^{(k)}) = d_H(u^*/u^{(k)}, \mathbb{1}_m), d_H(u^*, u^{(k)}) = d_H(v^*/v^{(k)}, \mathbb{1}_n)$$

Let  $(u^*/u^{(k)})_i$  be normalized by dividing the smallest entry among  $(u^*/u^{(k)})_i$  to obtain  $(u^*/u^{(k)})'_i$ .

$$1 \leq (u^*/u^{(k)})'_i \leq M_k, \forall i$$

And times the value divided to  $(v^*/v^{(k)})$  to obtain  $(v^*/v^{(k)})'_i$ .



# Convergence of Sinkhorn's Algorithm

## Proof of (3) (Continuous)

By definition of  $P^*, P^{(k)}$ ,

$$P^{(k),T} \mathbb{1}_n = P^{*,T} \mathbb{1}_m = b$$

By the diagonalization:

$$\text{diag} \left( v^*/v^{(k)} \right)' P^{(k),T} = P^{*,T} \text{diag} \left( u^*/u^{(k)} \right)'^{-1}$$

By the first equality,

$$\text{diag} \left( v^*/v^{(k)} \right)' b = \text{diag} \left( v^*/v^{(k)} \right) P^{(k),T} \mathbb{1}_m$$

By the second equality,

$$\text{diag} \left( v^*/v^{(k)} \right)' P^{(k),T} \mathbb{1}_m = P^{*,T} \text{diag} \left( u^*/u^{(k)} \right)'^{-1} \mathbb{1}_m$$

# Convergence of Sinkhorn's Algorithm

## Proof of (3)(Continuous)

Recall the range of  $(u^*/u^{(k)})$ :

$$M_k^{-1} \leq \left(u^*/u^{(k)}\right)_i^{-1} \leq 1, \forall i$$

Consider "times"  $P^{*,T}$ :

$$\begin{aligned} M_k^{-1} b_j &= M_k^{-1} \left(P^{*,T} \mathbb{1}_m\right)_j \leq \left(P^{*,T} \text{diag} \left(u^*/u^{(k)}\right)^{-1} \mathbb{1}_m\right)_j \\ &\leq \left(P^{*,T} \mathbb{1}_m\right)_j = b_j \end{aligned}$$

By the result from second equality and  $P^{(k),T} \mathbb{1}_m = b$ :

$$M_k^{-1} b_j \leq \left[\left(v^*/v^{(k)}\right)' b\right]_j \leq b_j, \forall j \implies 1 \leq \left(v^*/v^{(k)}\right)'_j \leq M_k$$

# Convergence of Sinkhorn's Algorithm

## Proof of (3)(Continuous)

$$M_k^{-1} \leq \left(u^*/u^{(k)}\right)'_i \left(v^*/v^{(k)}\right)'_j = \left(u^*/u^{(k)}\right)_i \left(v^*/v^{(k)}\right)_j = P_{ij}^*/P_{ij}^{(k)} \leq M_k$$

which can be represented to the form (3).

## 1 Definition of Optimal Transport

- Discrete Case
- Continuous Case

## 2 Entropic Regularization

## 3 Sinkhorn's Algorithm

- Lagrangian Method
- Sinkhorn's Algorithm
- Hilbert Projective Metric
- Convergence of Sinkhorn's algorithm

## 4 Complexity and Application to Solve Original Problem

# Complexity and Application to Solve Original Problem

It's mainly the result from [Altschuler:2017].

## Main Theorem

Sinkhorn's algorithm with a rounding step returns a point  $\hat{P} \in U(a, b)$  satisfying

$$\langle \hat{P}, C \rangle \leq \min_{P \in U(a, b)} \langle P, C \rangle + \varepsilon$$

in time  $O(n^2 L^3 (\log n) \varepsilon^{-3})$ , where  $L$  is the upper bound of entry of  $C$ .

## "Conclusion"

As the development of computer science, people can use more stronger method to compute some tradictional question. By using the entropic as the regularization, the orginal optimal problem can be evaluated more easily due the convexity.