

# Optimal Transport

## Theory, Computation and Applications

Wenchong Huang

School of Mathematical Sciences,  
Zhejiang University.

Dec. 30th, 2024

# Overview

**Principal concern:** the distance between two probability measures.

**First introduced** in 1781 by Monge.

**Relative subjects:** probability theory, geometry, graph theory, machine learning...

**Applications:**

- Image registration and warping;
- Reflector design;
- Retrieving information from shadowgraphy and proton radiography;
- Seismic tomography and reflection seismology.

**Some well-known researchers:**

- Gaspard Monge (France);
- Leonid Kantorovich (Russia);
- Yann Brenier (France);
- Xianfeng Gu (顾险峰, China);

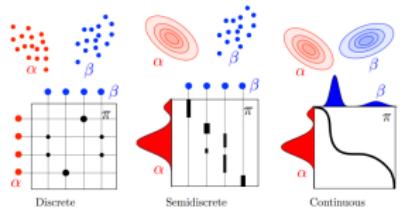


Fig. 1. Three main scenarios for Kantorovich OT

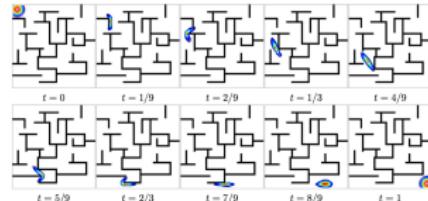


Fig. 2. Solving maze with OT

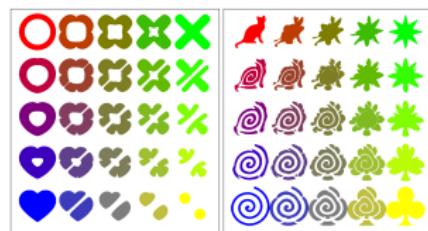


Fig. 3. 2D shape interpolation with OT

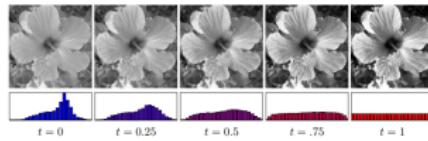


Fig. 4. Histogram equalization with OT

## ① Theory

## ② Computation

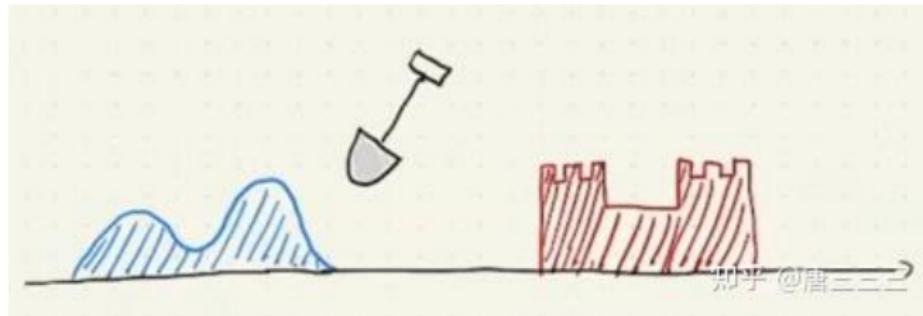
## ③ Applications

# The sand-moving problem

A child wants to make a pile of sand in the shape of a castle.

**Cost:** 1 kcal per shovel and per meter horizontally.

**Target:** Minimize the total cost.



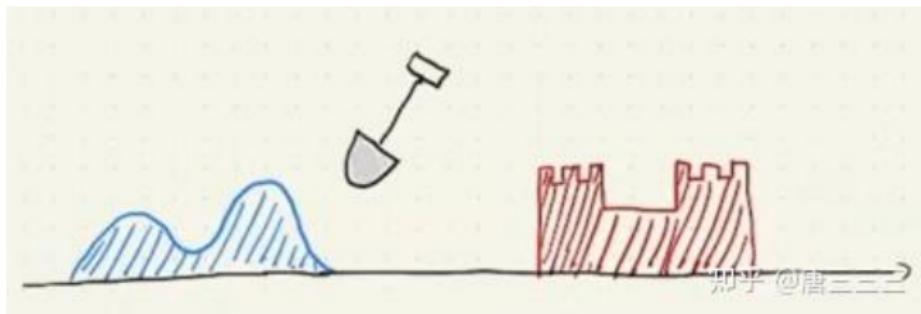
**Fig. 5.** The sand-moving problem.

# The sand-moving problem

A child wants to make a pile of sand in the shape of a castle.

**Cost:** 1 kcal per shovel and per meter horizontally.

**Target:** Minimize the total cost.



**Fig. 5.** The sand-moving problem.

Let's denote the source shape by  $f(x)$  and the target by  $g(x)$ . The sand-moving problem could be formulated as: find a **transport mapping**  $T : \mathbb{R} \rightarrow \mathbb{R}$  to minimize

$$\int_{\mathbb{R}} |T(x) - x| f(x) \, dx, \quad (1)$$

which satisfies

$$\int_{T(U)} g(x) \, dx = \int_U f(x) \, dx \text{ for all open interval } U \subset \mathbb{R}. \quad (2)$$

# The allocation problem

There are some steel coils to be transported from warehouses to factories. The transport cost is \$1 per coil and per kilometer. How to minimize the total cost?

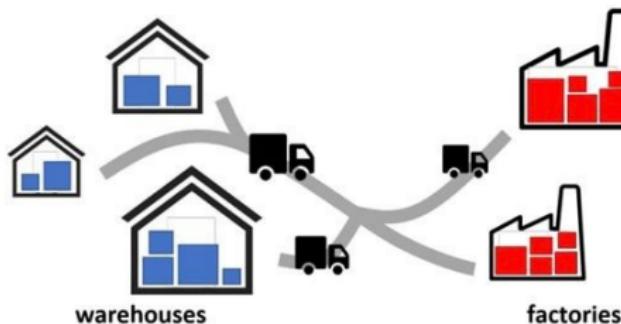


Fig. 6. The allocation problem.

# The allocation problem

There are some steel coils to be transported from warehouses to factories. The transport cost is \$1 per coil and per kilometer. How to minimize the total cost?

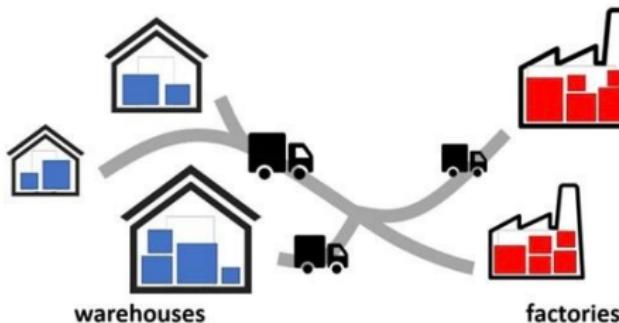


Fig. 6. The allocation problem.

Assume the  $i$ -th warehouse has  $a_i$  coils and the  $j$ -th factory needs  $b_j$  coils. And assume the distance between the  $i$ -th warehouse and the  $j$ -th factory is  $d_{ij}$ . The allocation problem could be formulated as: find a **transport matrix**  $v_{ij}$  to minimize

$$\sum_{i,j} d_{ij} v_{ij} \quad (3)$$

which satisfies

$$a_i = \sum_j v_{ij}, \quad \forall i, \quad \text{and} \quad b_j = \sum_i v_{ij}, \quad \forall j. \quad (4)$$

# The Monge formulation

Denote  $\mathcal{M}_+^1(\mathcal{X})$  the set of probability measures on  $\mathcal{X}$ .

## Definition (push-forward)

Suppose  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and a map  $T : \mathcal{X} \rightarrow \mathcal{Y}$ . Say  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$  is the push-forward of  $\mu$  by  $T$  if

$$\int_{\mathcal{Y}} h(y) d\nu(y) = \int_{\mathcal{X}} h(T(x)) d\mu(x), \quad \forall h \in \mathcal{C}(\mathcal{Y}). \quad (5)$$

Write  $T_{\#}\mu := \nu$ .

## Example (push-forward of a discrete measure)

Suppose  $\alpha$  is a discrete measure

$$\alpha = \sum_{i=1}^n a_i \delta_{x_i}.$$

Then the push-forward of  $\alpha$  by  $T$  is

$$T_{\#}\alpha = \sum_{i=1}^n a_i \delta_{T(x_i)}.$$

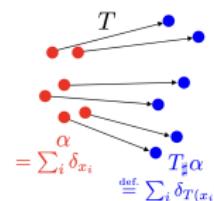


Fig. 7. push-forward of a discrete measure

<sup>1</sup> Gaspard Monge. "Mémoire sur la théorie des déblais et des remblais". In: *Histoire de l'Académie Royale des Sciences* (1781).



# The Monge formulation

Denote  $\mathcal{M}_+^1(\mathcal{X})$  the set of probability measures on  $\mathcal{X}$ .

## Definition (push-forward)

Suppose  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  and a map  $T : \mathcal{X} \rightarrow \mathcal{Y}$ . Say  $\nu \in \mathcal{M}_+^1(\mathcal{Y})$  is the push-forward of  $\mu$  by  $T$  if

$$\int_{\mathcal{Y}} h(y) d\nu(y) = \int_{\mathcal{X}} h(T(x)) d\mu(x), \quad \forall h \in \mathcal{C}(\mathcal{Y}). \quad (5)$$

Write  $T_{\#}\mu := \nu$ .

## Example (push-forward of a discrete measure)

Suppose  $\alpha$  is a discrete measure

$$\alpha = \sum_{i=1}^n a_i \delta_{x_i}.$$

Then the push-forward of  $\alpha$  by  $T$  is

$$T_{\#}\alpha = \sum_{i=1}^n a_i \delta_{T(x_i)}.$$

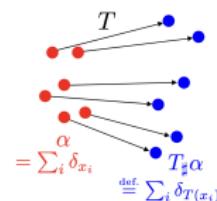


Fig. 7. push-forward of a discrete measure

Given two probability measures  $\mu$  on  $\mathcal{X}$  and  $\nu$  on  $\mathcal{Y}$ , and a cost function  $c(x, y)$ . Optimal transport could be generally formulated as the Monge problem:

$$\min_T \left\{ \int_{\mathcal{X}} c(x, T(x)) d\mu(x) : T_{\#}\mu = \nu \right\} \quad (6)$$

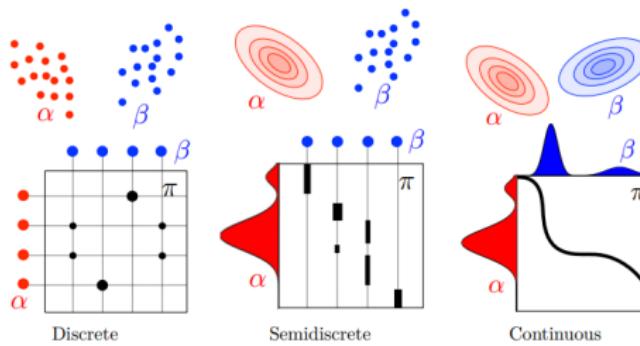
The Monge problem between discrete measures is introduced by Monge<sup>1</sup>.

<sup>1</sup> Gaspard Monge. "Mémoire sur la théorie des déblais et des remblais". In: *Histoire de l'Académie Royale des Sciences* (1781).



# The Kantorovich formulation

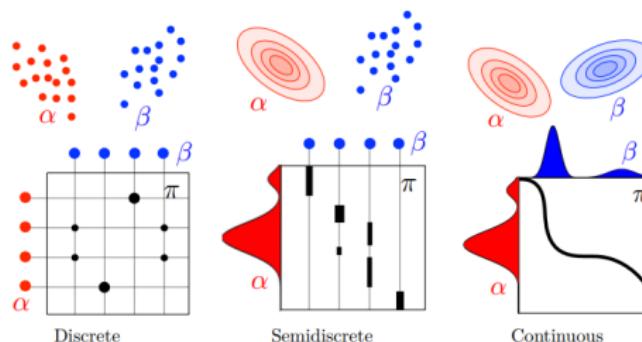
Here's another general formulation of OT, we first recall the three main scenarios for OT.



<sup>2</sup>Leonid Kantorovich. "On the transfer of masses". In: *Doklady Akademii Nauk* 37.2 (1942).

# The Kantorovich formulation

Here's another general formulation of OT, we first recall the three main scenarios for OT.



Given two probability measures  $\mu$  on  $\mathcal{X}$  and  $\nu$  on  $\mathcal{Y}$ , and a cost function  $c(x, y)$ . Optimal transport could be generally formulated as the Kantorovich problem<sup>2</sup>:

$$\mathcal{L}_c(\mu, \nu) = \min_{\pi} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y), \quad (7)$$

where  $\pi$  is a measure on  $\mathcal{X} \times \mathcal{Y}$ , whose marginals are  $\mu$  and  $\nu$ , that is,

$$\mu = \int_{\mathcal{Y}} \pi(\cdot, y) dy, \quad \nu = \int_{\mathcal{X}} \pi(x, \cdot) dx. \quad (8)$$

---

<sup>2</sup>Leonid Kantorovich. "On the transfer of masses". In: *Doklady Akademii Nauk* 37.2 (1942).

# Wasserstein disrtance

Here we suppose  $\mathcal{X} = \mathcal{Y}$  and  $c(x, y) = d(x, y)^p$  ( $p > 1$ ), where  $d$  is a distance on  $\mathcal{X}$ .

<sup>3</sup> Cédric Villani. *Optimal Transport: Old and New.* Vol. 338. Springer Verlag, 2009.

# Wasserstein distance

Here we suppose  $\mathcal{X} = \mathcal{Y}$  and  $c(x, y) = d(x, y)^p$  ( $p > 1$ ), where  $d$  is a distance on  $\mathcal{X}$ .

## Theorem (Wasserstein distance)

Under the above assumptions,  $\mathcal{L}_c(\mu, \nu)^{1/p}$  is a distance on  $\mathcal{M}_+^1(\mathcal{X})$ .

The distance  $\mathcal{W}_p(\mu, \nu) := \mathcal{L}_c(\mu, \nu)^{1/p}$  is called  $p$ -Wasserstein distance.

<sup>3</sup> Cédric Villani. *Optimal Transport: Old and New*. Vol. 338. Springer Verlag, 2009.

# Wasserstein distance

Here we suppose  $\mathcal{X} = \mathcal{Y}$  and  $c(x, y) = d(x, y)^p$  ( $p > 1$ ), where  $d$  is a distance on  $\mathcal{X}$ .

## Theorem (Wasserstein distance)

Under the above assumptions,  $\mathcal{L}_c(\mu, \nu)^{1/p}$  is a distance on  $\mathcal{M}_+^1(\mathcal{X})$ .

The distance  $\mathcal{W}_p(\mu, \nu) := \mathcal{L}_c(\mu, \nu)^{1/p}$  is called  $p$ -Wasserstein distance.

## Definition (weak convergence)

Suppose  $\mathcal{X}$  is compact. Say  $(\mu_k)_{k \geq 1} \subset \mathcal{M}_+^1(\mathcal{X})$  converges weakly to  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  if

$$\int_{\mathcal{X}} g \, d\mu_k \rightarrow \int_{\mathcal{X}} g \, d\mu, \quad \forall g \in \mathcal{C}(\mathcal{X}). \quad (9)$$

## Theorem (Wasserstein distance and weak convergence<sup>3</sup>)

On a compact domain  $\mathcal{X}$ ,  $(\mu_k)_{k \geq 1} \subset \mathcal{M}_+^1(\mathcal{X})$  converges weakly to  $\mu \in \mathcal{M}_+^1(\mathcal{X})$  if and only if  $\mathcal{W}_p(\mu_k, \mu) \rightarrow 0$ .

<sup>3</sup> Cédric Villani. *Optimal Transport: Old and New*. Vol. 338. Springer Verlag, 2009.

# Equivalence between the Kantorovich and Monge problems

## Theorem (Kantorovich dual problem)

The Kantorovich problem can be solved in the dual space by

$$\mathcal{L}_c(\mu, \nu) = \sup_{(f,g) \in \mathcal{R}(c)} \int_{\mathcal{X}} f(x) \, d\mu(x) + \int_{\mathcal{Y}} g(y) \, d\nu(y), \quad (10)$$

where the set of admissible dual potential is

$$\mathcal{R}(c) := \{(f, g) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y}) : \forall (x, y), f(x) + g(y) \leq c(x, y)\}. \quad (11)$$

The pair  $(f, g)$  is called Kantorovich potentials.

<sup>4</sup>Yann Brenier. "Polar factorization and monotone rearrangement of vector-valued functions" In: *Communications on Pure and Applied Mathematics* 44.4 (1991).

# Equivalence between the Kantorovich and Monge problems

## Theorem (Kantorovich dual problem)

The Kantorovich problem can be solved in the dual space by

$$\mathcal{L}_c(\mu, \nu) = \sup_{(f, g) \in \mathcal{R}(c)} \int_{\mathcal{X}} f(x) d\mu(x) + \int_{\mathcal{Y}} g(y) d\nu(y), \quad (10)$$

where the set of admissible dual potential is

$$\mathcal{R}(c) := \{(f, g) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y}) : \forall (x, y), f(x) + g(y) \leq c(x, y)\}. \quad (11)$$

The pair  $(f, g)$  is called Kantorovich potentials.

## Theorem (Brenier<sup>4</sup>)

In the case  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$  and  $c(x, y) = \|x - y\|_2^2$ , if at least one of the two input measures (denoted  $\mu$ ) has a density  $\rho_\mu$  with respect to the Lebesgue measure, then the optimal  $\pi$  in the Kantorovich formulation is unique and is supported on the graph  $(x, T(x))$  of a Monge map  $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . This means that  $\pi = (\text{Id}, T)_\# \mu$ , i.e.

$$\int_{\mathcal{X} \times \mathcal{Y}} h(x, y) d\pi(x, y) = \int_{\mathcal{X}} h(x, T(x)) d\mu(x), \quad \forall h \in \mathcal{C}(\mathcal{X} \times \mathcal{Y}). \quad (12)$$

Furthermore, this map  $T$  is uniquely defined as the gradient of a convex function  $\varphi$ ,  $T(x) = \nabla \varphi(x)$ , where  $\varphi$  is the unique (up to an additive constant) convex function such that  $(\nabla \varphi)_\# \mu = \nu$ . This convex function is related to the dual potential  $f$  solving (10) as

$$\varphi(x) = \frac{\|x\|_2^2}{2} - f(x). \quad (13)$$

<sup>4</sup>Yann Brenier. "Polar factorization and monotone rearrangement of vector-valued functions" In: *Communications on Pure and Applied Mathematics* 44.4 (1991).

# Dynamic formulation

In the case  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ , and  $c(x, y) = \|x - y\|_2$ , the optimal transport distance  $\mathcal{W}_2^2(\mu, \nu) = \mathcal{L}_c(\mu, \nu)$  can be defined as

$$\mathcal{W}_2^2(\mu, \nu) = \min_{(\alpha_t, v_t)_t} \int_0^1 \int_{\mathbb{R}^d} \|v_t(x)\|^2 d\alpha_t(x) dt, \quad (14)$$

where  $\alpha_t$  is a scalar-valued measure and  $v_t$  a vector-valued measure which satisfy the conservation of mass formula,

$$\frac{\partial \alpha_t}{\partial t} + \nabla \cdot (\alpha_t, v_t) = 0, \quad (15)$$

and the boundary conditions  $\alpha_0 = \mu$  and  $\alpha_1 = \nu$ .

## 1 Theory

## 2 Computation

## 3 Applications

# 1-D discrete case

Here  $\mathcal{X} = \mathcal{Y} = \mathbb{R}$ . Suppose  $\alpha = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  and  $\beta = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$  where  $x_1 \leq \dots \leq x_n$  and  $y_1 \leq \dots \leq y_n$ .

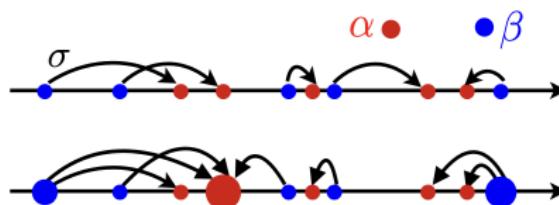


Fig. 8. 1-D optimal transport in discrete case

Then the  $p$ -Wasserstein distance can be simply computed by

$$\mathcal{W}_p(\alpha, \beta)^p = \frac{1}{n} \sum_{i=1}^n |x_i - y_i|^p. \quad (16)$$

It's in fact a greedy algorithm.

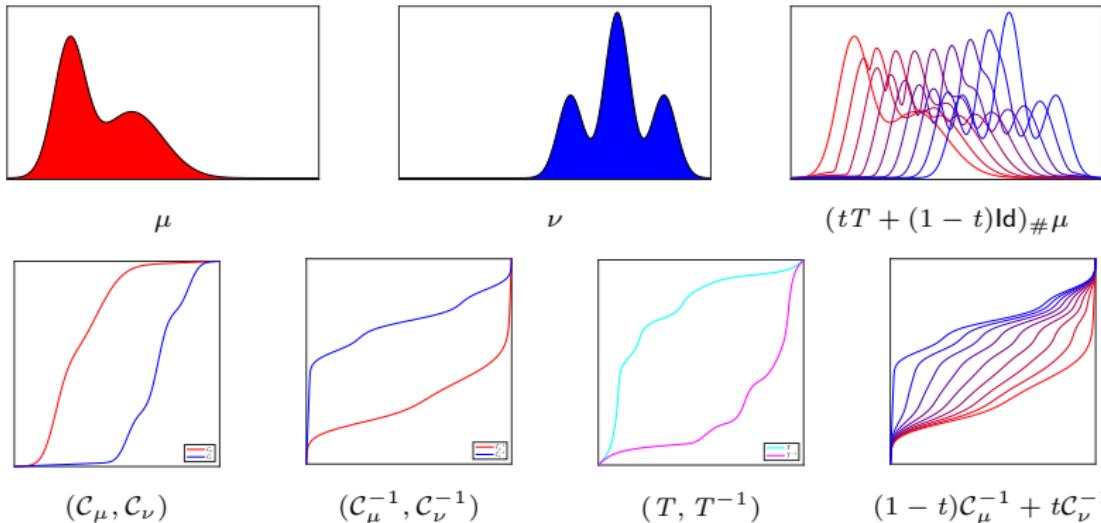
# 1-D continuous case

If  $\mu, \nu$  are 1-D measures with densities. Suppose their cumulative distribution functions are  $\mathcal{C}_\mu$  and  $\mathcal{C}_\nu$ , respectively. Then the  $\mathcal{W}_1$  distance could be computed by

$$\mathcal{W}_1(\mu, \nu) = \int_{\mathbb{R}} |\mathcal{C}_\mu(x) - \mathcal{C}_\nu(x)| dx = \int_{\mathbb{R}} \left| \int_{-\infty}^x d(\mu - \nu) \right| dx. \quad (17)$$

And the Monge map is then defined by

$$T = \mathcal{C}_\nu^{-1} \circ \mathcal{C}_\mu. \quad (18)$$



**Fig. 9.** Computation of OT and displacement interpolation between two 1-D measures.

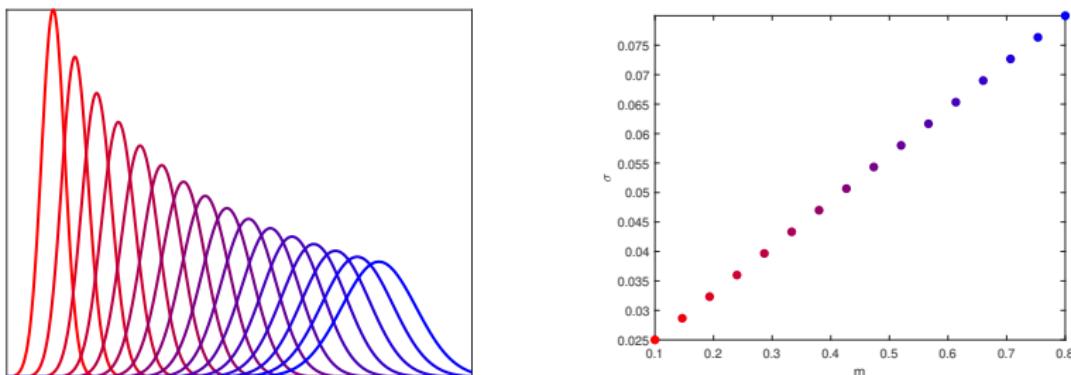


# 1-D Gaussian

If  $\mu = \mathcal{N}(m_1, \sigma_1^2)$ ,  $\nu = \mathcal{N}(m_2, \sigma_2^2)$  are 1-D Gaussians. Then the  $\mathcal{W}_2$  distance can be directly computed by

$$\mathcal{W}_2(\mu, \nu) = \sqrt{|m_1 - m_2|^2 + |\sigma_1 - \sigma_2|^2}, \quad (19)$$

which is thus the Euclidean distance on the 2-D plane plotting the mean and the standard deviation of a Gaussian  $\mathcal{N}(m, \sigma)$ .



**Fig. 10.** Computation of displacement interpolation between two 1-D Gaussians.

Learn more in [Takatsu, 2011]<sup>5</sup>.

<sup>5</sup> Asuka Takatsu. "Wasserstein geometry of Gaussian measures". In: *Osaka Journal of Mathematics* 48, 3 (2011).

# Discretization

Suppose  $\mu$  is a measure with density  $\rho$ , supported on  $[0, 1]$ . Let

$$\tilde{\mu} = \sum_{i=0}^N u_i \delta_{x_i}, \quad (20)$$

where

$$u_i = \frac{\rho(x_i)}{N+1}, \quad x_i = \frac{i}{N}, \quad i = 0, \dots, N. \quad (21)$$

We call  $\tilde{\mu}$  the *discretization* of  $\mu$ . This technique can also be used in  $\mathbb{R}^d$ .

# Discretization

Suppose  $\mu$  is a measure with density  $\rho$ , supported on  $[0, 1]$ . Let

$$\tilde{\mu} = \sum_{i=0}^N u_i \delta_{x_i}, \quad (20)$$

where

$$u_i = \frac{\rho(x_i)}{N+1}, \quad x_i = \frac{i}{N}, \quad i = 0, \dots, N. \quad (21)$$

We call  $\tilde{\mu}$  the *discretization* of  $\mu$ . This technique can also be used in  $\mathbb{R}^d$ .

Let  $\tilde{\nu} = \sum_{i=0}^M v_i \delta_{y_i}$  and  $(C)_{ij}$  be the cost matrix. The Kantorovich problem then becomes

$$L_C(\mathbf{u}, \mathbf{v}) := \min_{\mathbf{P} \in U(\mathbf{u}, \mathbf{v})} \langle \mathbf{P}, \mathbf{C} \rangle := \min_{\mathbf{P} \in U(\mathbf{u}, \mathbf{v})} \sum_{i,j} \mathbf{P}_{ij} C_{ij}, \quad (22)$$

where

$$U(\mathbf{u}, \mathbf{v}) := \left\{ \mathbf{P} \left| \begin{array}{l} \sum_j \mathbf{P}_{ij} = u_i, \forall i, \\ \text{and} \\ \sum_i \mathbf{P}_{ij} = v_j, \forall j \end{array} \right. \right\}. \quad (23)$$

# Entropy regularization

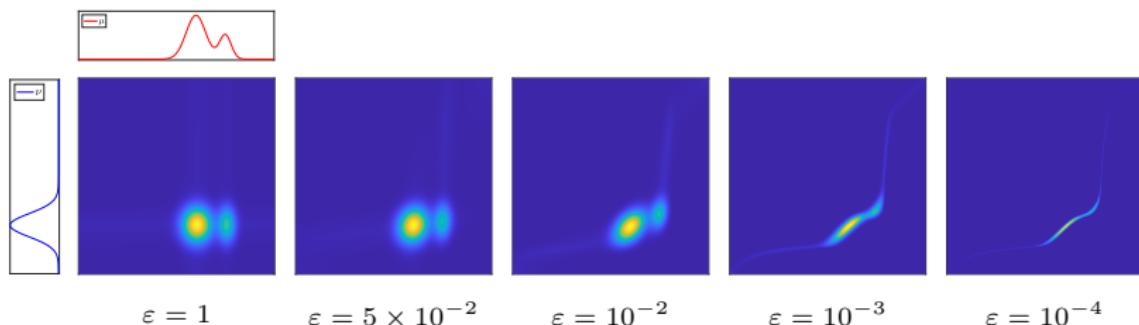
Define the entropy

$$H(\mathbf{P}) := - \sum_{i,j} \mathbf{P}_{ij} (\log(\mathbf{P}_{ij}) - 1). \quad (24)$$

Then the regularized Kantorovich problem<sup>6</sup> is defined by

$$L_C^\varepsilon(\mathbf{u}, \mathbf{v}) := \min_{\mathbf{P} \in U(\mathbf{u}, \mathbf{v})} \langle \mathbf{P}, \mathbf{C} \rangle - \varepsilon H(\mathbf{P}). \quad (25)$$

It can be shown that  $L_C^\varepsilon(\mathbf{u}, \mathbf{v}) = L_C(\mathbf{u}, \mathbf{v}) + O(\varepsilon)$ .



**Fig. 11.** Graphs of optimal  $\mathbf{P}$ s when choose different  $\varepsilon$ . Set  $\mathbf{C}_{ij} = |x_i - x_j|^2$ .

<sup>6</sup> Alan G. Wilson. "The use of entropy maximizing models, in the theory of trip distribution, mode split and route split". In: *Journal of Transport Economics and Policy* (1969), pp. 108–126.

# Sinkhorn iteration

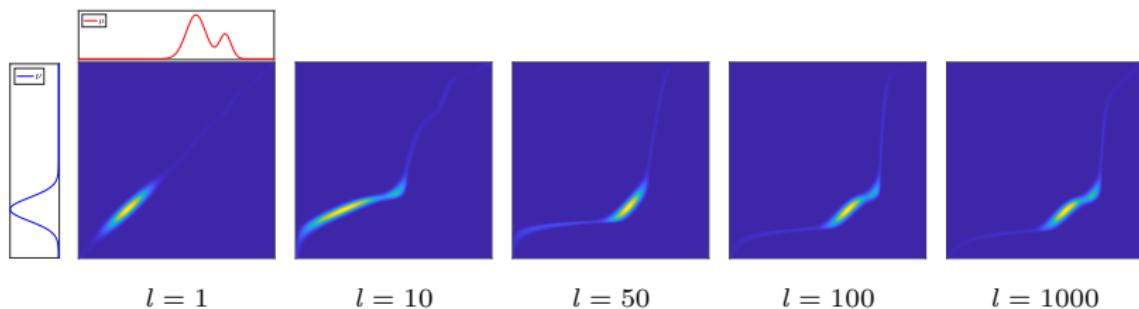
Let  $K_{ij} = e^{-\frac{c_{ij}}{\varepsilon}}$ . Sinkhorn iteration writes

$$\mathbf{a}^{(l+1)} \leftarrow \frac{\mathbf{u}}{\mathbf{K}\mathbf{b}^{(l)}}, \quad \text{and} \quad \mathbf{b}^{(l+1)} \leftarrow \frac{\mathbf{v}}{\mathbf{K}^T \mathbf{a}^{(l+1)}}, \quad \text{for } l = 0, 1, \dots \quad (26)$$

which starts with an arbitrary  $\mathbf{b}^{(0)}$ . The transport matrix  $\mathbf{P}$  can be rebuilt by

$$\mathbf{P}^{(l)} = \text{diag}(\mathbf{b}^{(l)}) \cdot \mathbf{K} \cdot \text{diag}(\mathbf{a}^{(l)}). \quad (27)$$

The convergence is proved by Sinkhorn<sup>7</sup>. And Altschuler et al<sup>8</sup> give an analysis of the computational complexity.



**Fig. 12.** Graphs of  $\mathbf{P}^{(l)}$ . Set  $C_{ij} = |x_i - x_j|^2$  and  $\varepsilon = 10^{-3}$ .

<sup>7</sup>Richard Sinkhorn. "A relationship between arbitrary positive matrices and doubly stochastic matrices". In: *Annals of Mathematical Statistics* 35 (1964).

<sup>8</sup>Jason Altschuler, Jonathan Weed, and Philippe Rigollet. "Near-linear time approximation algorithms for optimal transport via Sinkhorn iteration". In: *Advances in Neural Information Processing Systems* (2017).

## ① Theory

## ② Computation

## ③ Applications

## 2-D shape interpolation



Fig. 13. From Kunkun to chicken. Top: color interpolation. Bottom: shape interpolation.

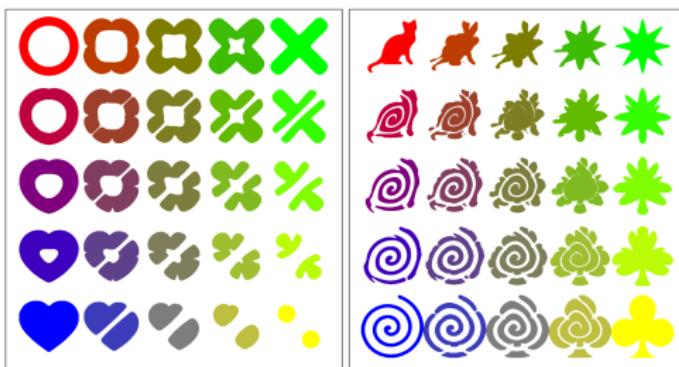


Fig. 14. Barycenter of four shapes<sup>9</sup>.

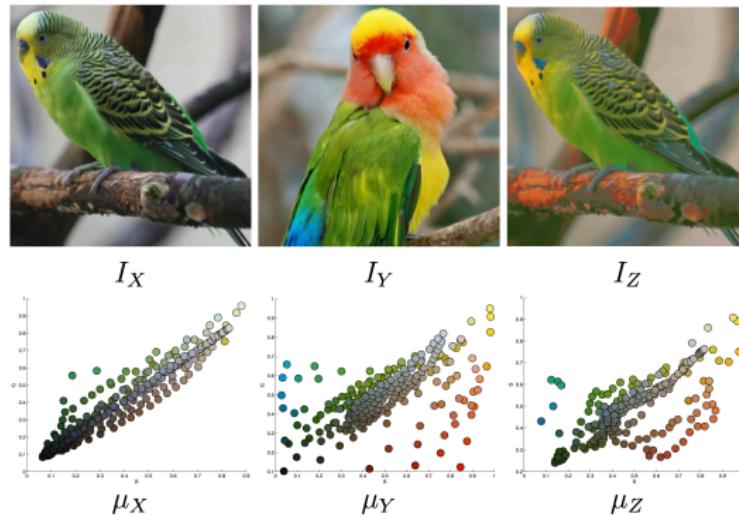
<sup>9</sup> Gabriel Peyré, Marco Cuturi, et al. "Computational optimal transport: With applications to data science". In: *Foundations and Trends® in Machine Learning* 11.5-6 (2019), pp. 355–607.

# Color transfer

Compute a transformation  $T$  such that

$$I_Z(x) = T(I_X(x)), \quad \text{for all pixel } x, \quad (28)$$

where the new color distribution  $\mu_Z$  is close or equal to  $\mu_Y$ .



**Fig. 15.** Example of color transfer<sup>10</sup>. The second row represents RGB color distributions using the 2-D projection of every pixel in the RG plane

<sup>10</sup> Nicolas Papadakis. "Optimal transport for image processing". PhD thesis. Université de Bordeaux; Habilitation à thèse, 2015.

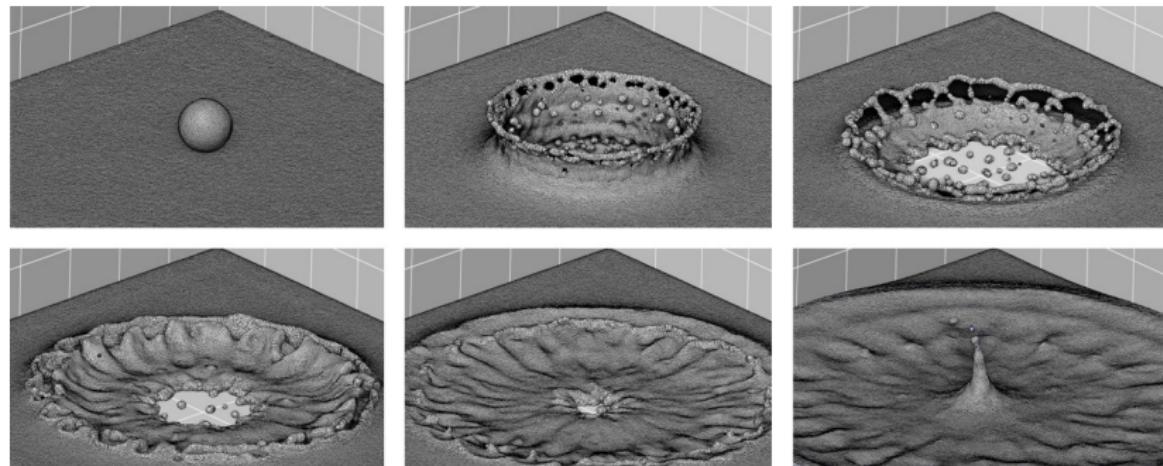
# Color transfer



Fig. 16. Another example of color transfer<sup>11</sup>.

<sup>11</sup> Nicolas Bonneel and Julie Digne. "A survey of optimal transport for computer graphics and computer vision". In: *Computer Graphics Forum*. Vol. 42. 2. Wiley Online Library. 2023, pp. 439–460.

# Fluid dynamics



**Fig. 17.** Simulation of the free boundary problem using partial OT<sup>12</sup>.

<sup>12</sup> Bruno Lévy. "Partial optimal transport for a constant-volume Lagrangian mesh with free boundaries". In: *Journal of Computational Physics* 451 (2022), p. 110838.

*Thank You*