

Regularized Optimal Transport is Ground Cost Adversarial

F-P. PATY

M. CUTURI



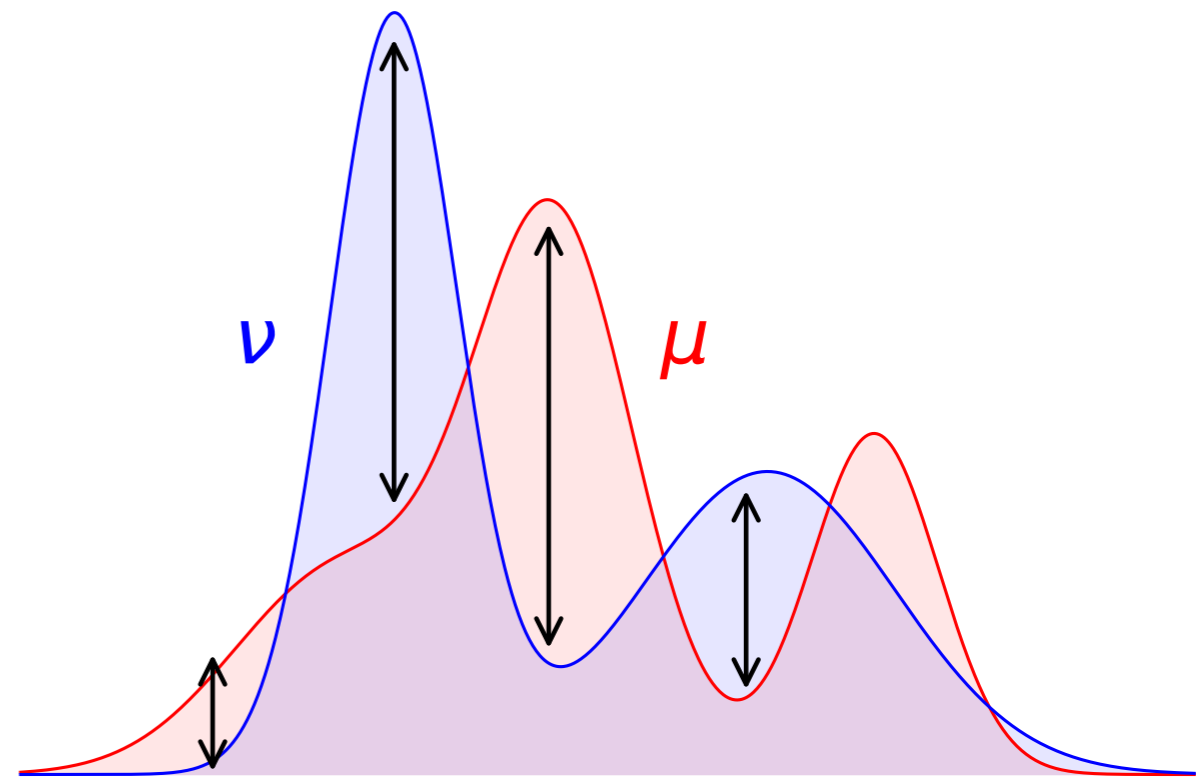
Google AI
Brain Team

COMPARING DISTRIBUTIONS

1. Vertical comparison

Look at the difference, or the ratio of the densities

e.g. Total Variation distance, Kullback Leibler divergence, etc.

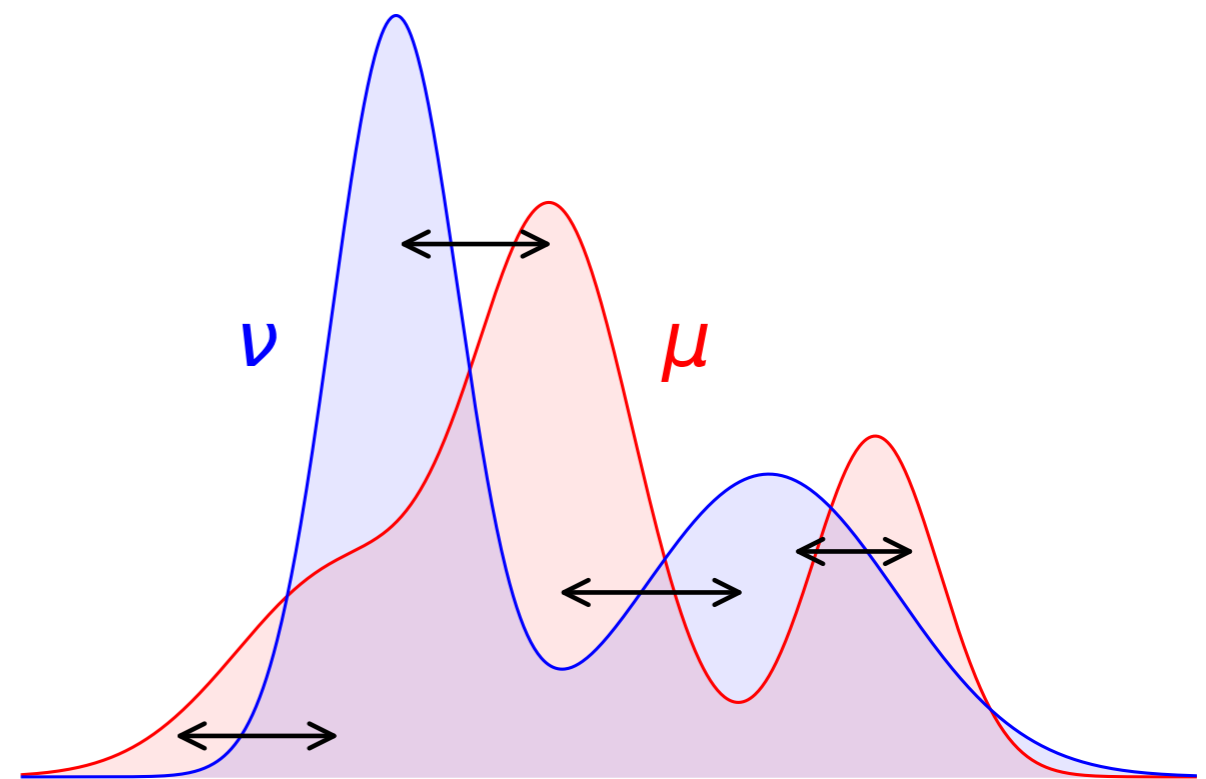


COMPARING DISTRIBUTIONS

2. Horizontal comparison aka Optimal Transport

Move the mass across the
ground space

⚠ Need for a notion of
displacement cost on the
ground space



SOME HISTORY

666. MÉMOIRES DE L'ACADÉMIE ROYALE

M É M O I R E

SUR LA

THÉORIE DES DÉBLAIS
ET DES REMBLAIS.

Par M. M O N G E.

LORSQU'ON doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.

1781



1939

SOME HISTORY



Brenier



Otto



McCann

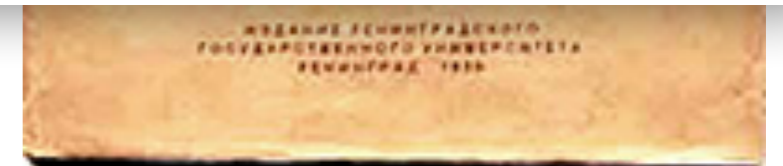


Villani



Figalli

Remblai à l'espace qu'elles doivent occuper après le transport.



1781

1939

A portrait of Leonid Kantorovich, a middle-aged man with glasses, wearing a grey pinstriped suit jacket, a green tie, and a patterned shirt. He is holding a newspaper in front of him, which has the word 'OPTIMUM' visible on its masthead. The background is a wood-paneled wall.

OPTIMAL TRANSPORT

Leonid Kantorovich

OPTIMAL TRANSPORT

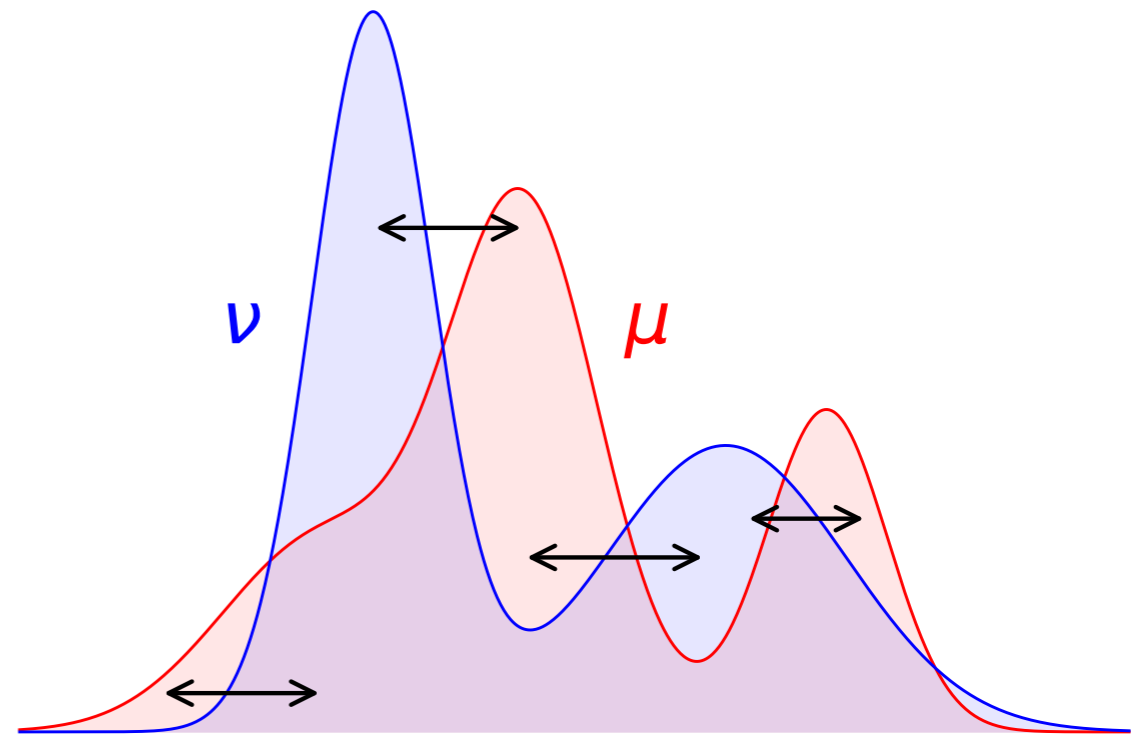
Data:

Two distributions μ and ν over \mathbb{R}^d

Parameter:

A (continuous) cost function

$$c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$



OPTIMAL TRANSPORT

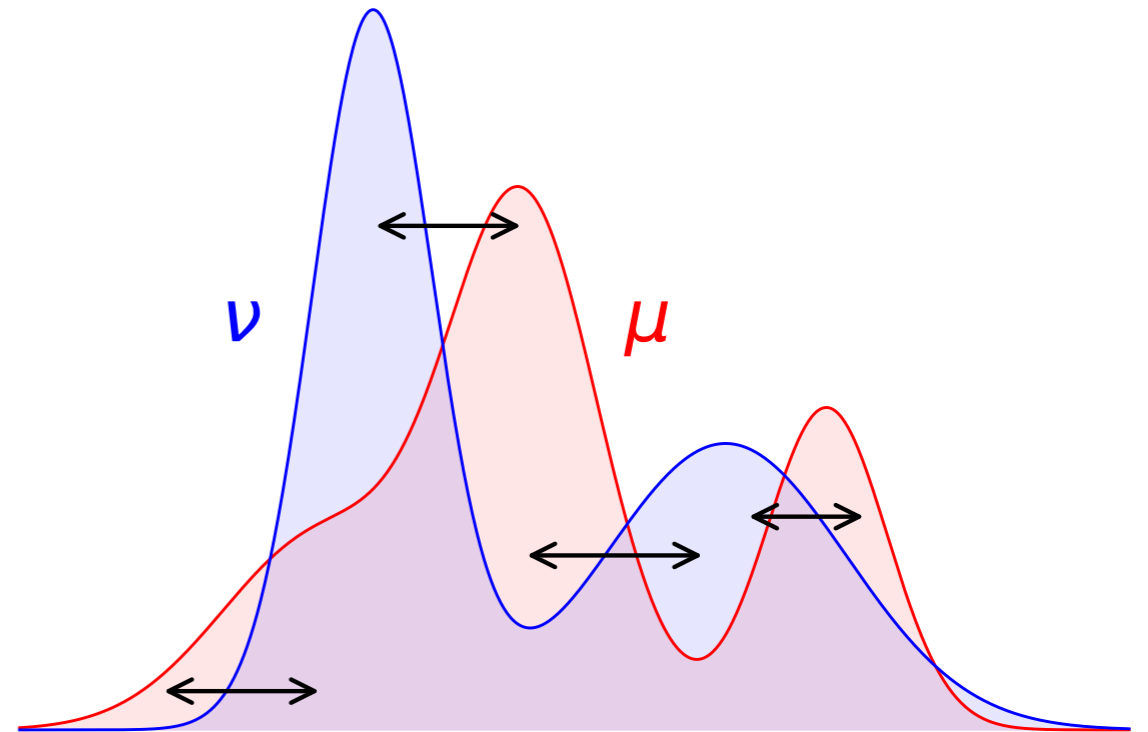
Data:

Two distributions μ and ν over \mathbb{R}^d

Parameter:

A (continuous) cost function

$$c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$



Definition of Optimal Transport (OT):

OPTIMAL TRANSPORT

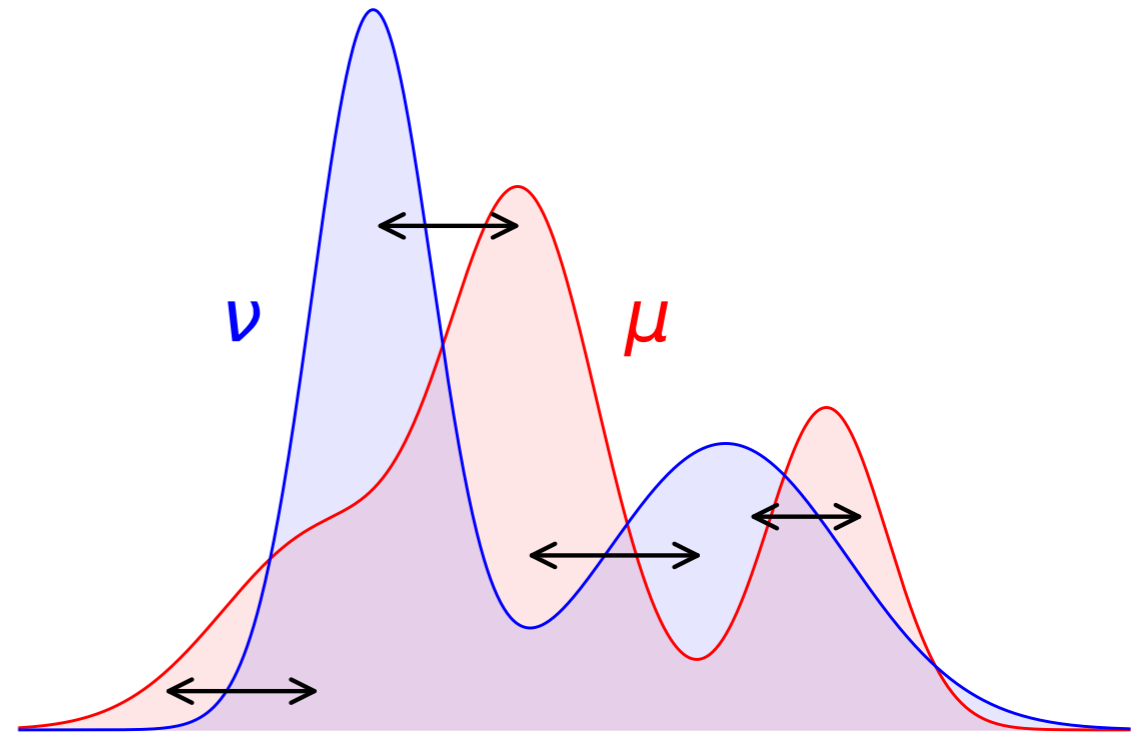
Data:

Two distributions μ and ν over \mathbb{R}^d

Parameter:

A (continuous) cost function

$$c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$



Definition of Optimal Transport (OT):

$$c(x, y) d\pi(x, y)$$

OPTIMAL TRANSPORT

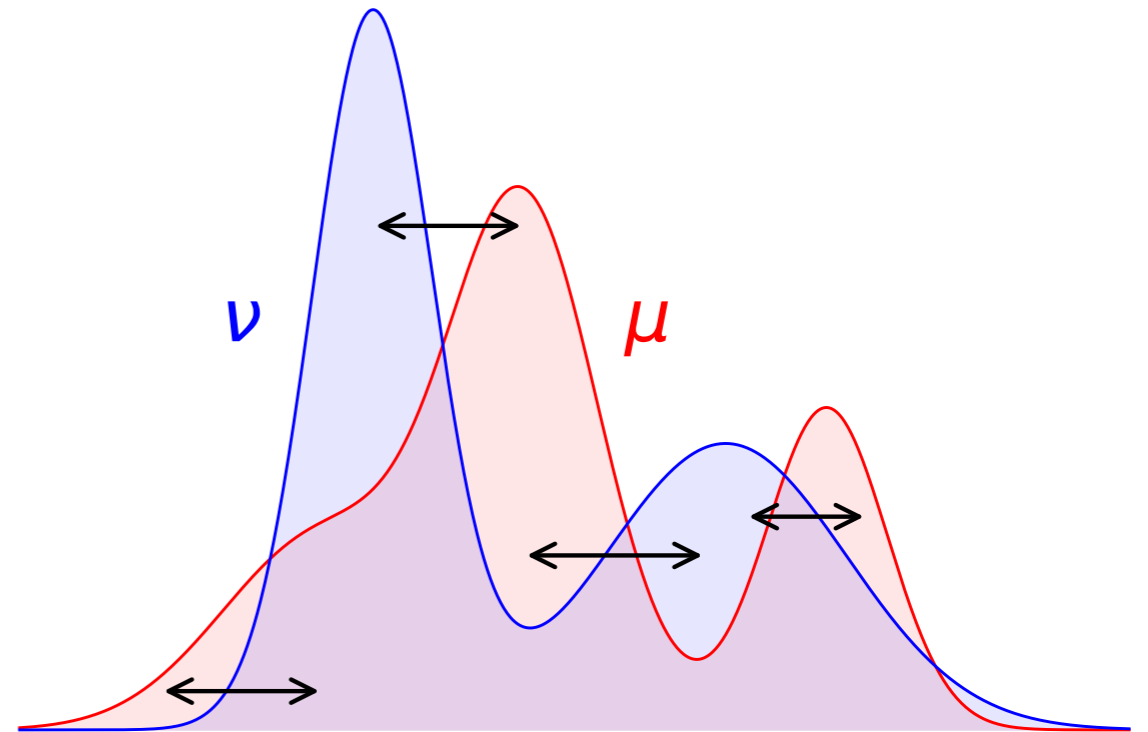
Data:

Two distributions μ and ν over \mathbb{R}^d

Parameter:

A (continuous) cost function

$$c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$



Definition of Optimal Transport (OT):

$$\iint c(x, y) d\pi(x, y)$$

OPTIMAL TRANSPORT

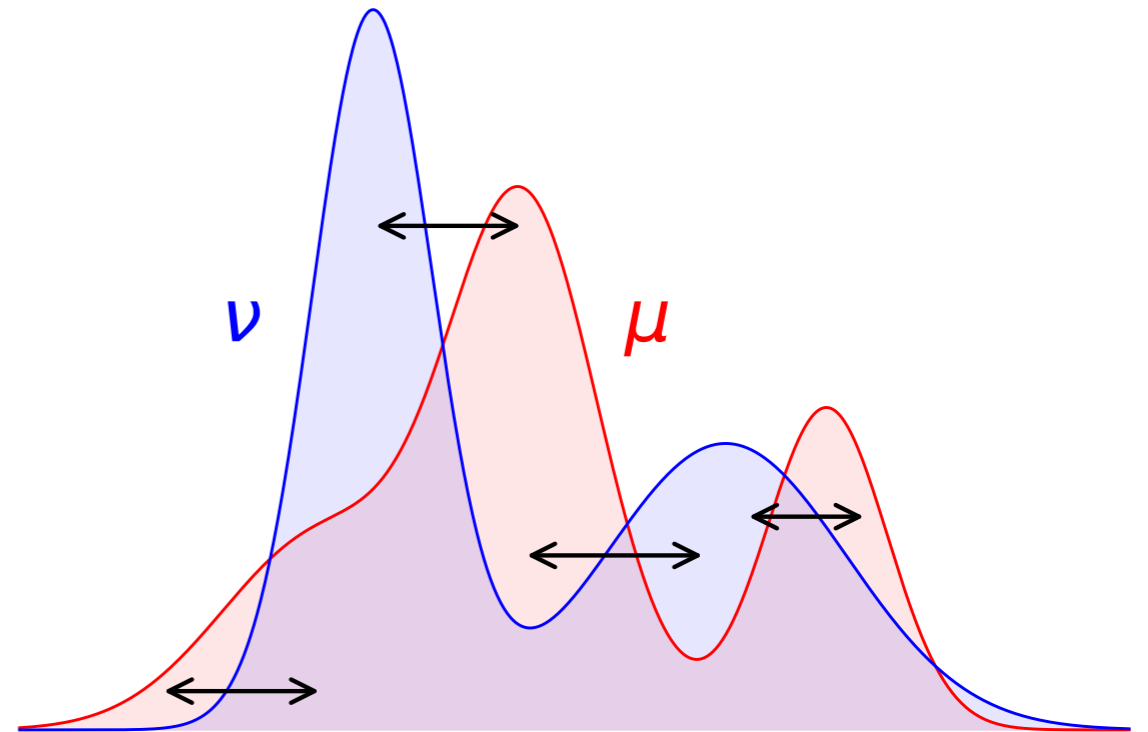
Data:

Two distributions μ and ν over \mathbb{R}^d

Parameter:

A (continuous) cost function

$$c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$



Definition of Optimal Transport (OT):

$$\inf_{\pi} \iint c(x, y) d\pi(x, y)$$

OPTIMAL TRANSPORT

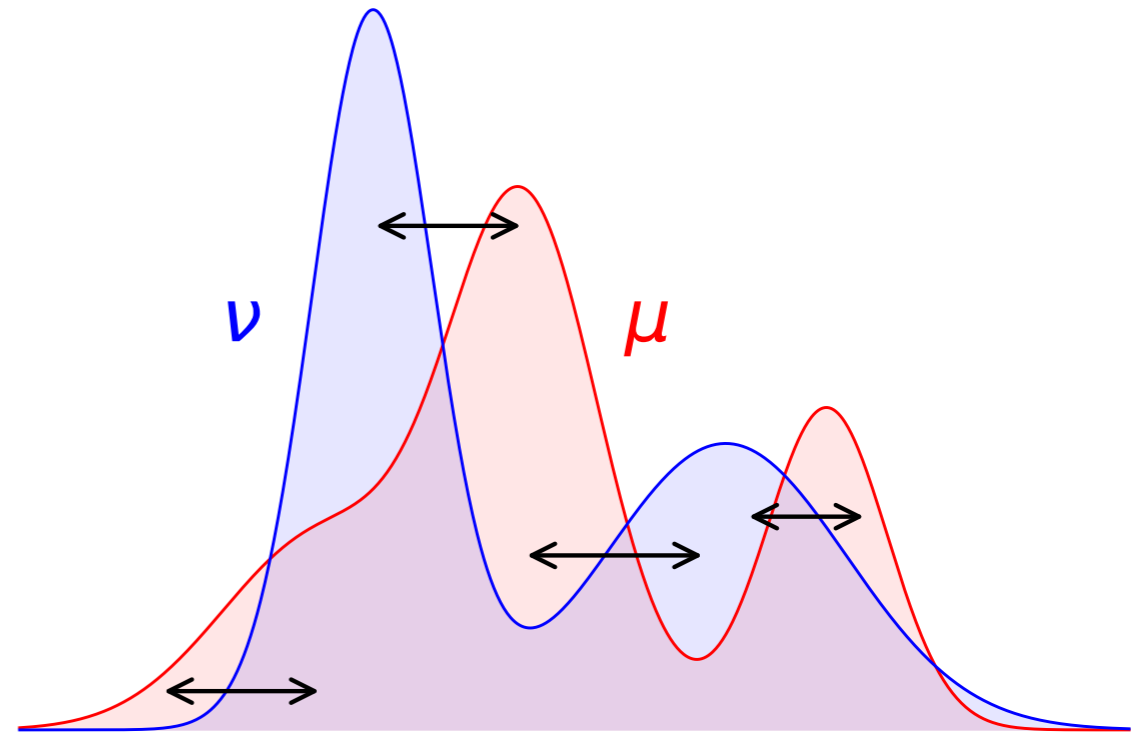
Data:

Two distributions μ and ν over \mathbb{R}^d

Parameter:

A (continuous) cost function

$$c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$



Definition of Optimal Transport (OT):

$$\mathcal{T}_c(\mu, \nu) = \inf_{\pi} \iint c(x, y) d\pi(x, y)$$

OPTIMAL TRANSPORT

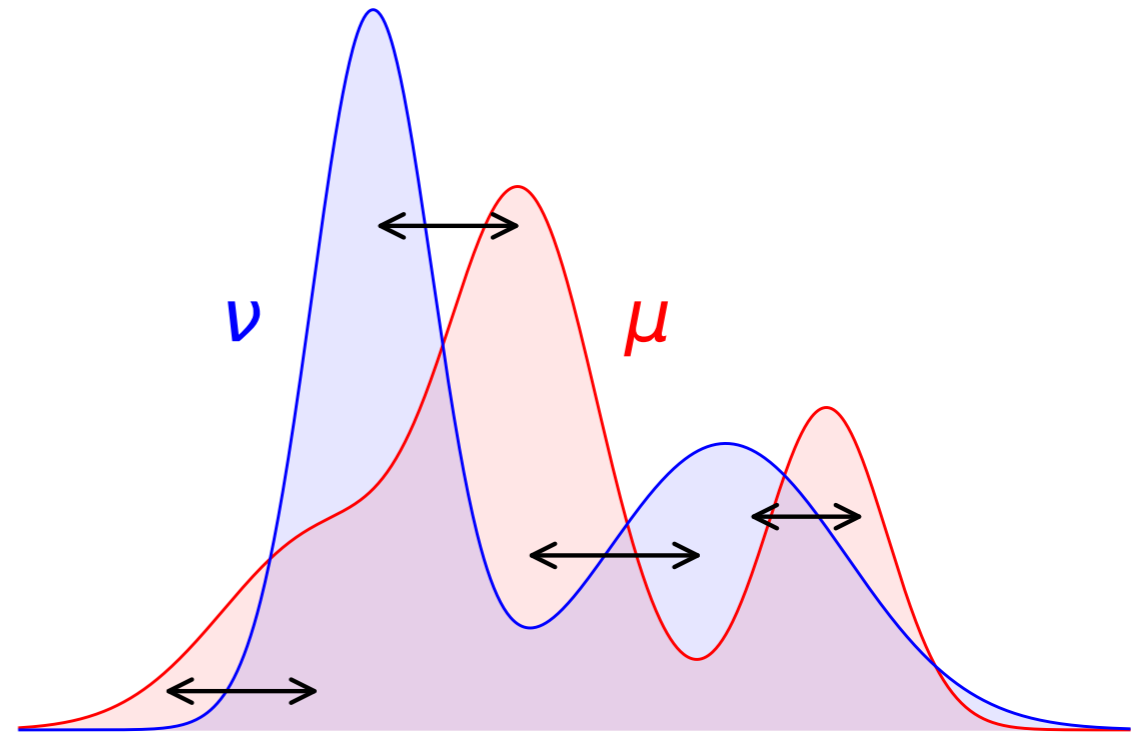
Data:

Two distributions μ and ν over \mathbb{R}^d

Parameter:

A (continuous) cost function

$$c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$



Definition of Optimal Transport (OT):

$$\mathcal{T}_c(\mu, \nu) = \inf_{\pi} \iint c(x, y) d\pi(x, y)$$

over all π such that

$$\begin{cases} \int d\pi(x, y) = d\mu(x) \quad \forall x \\ \int d\pi(x, y) = d\nu(y) \quad \forall y \end{cases}$$

OPTIMAL TRANSPORT

Two main questions in practice

OPTIMAL TRANSPORT

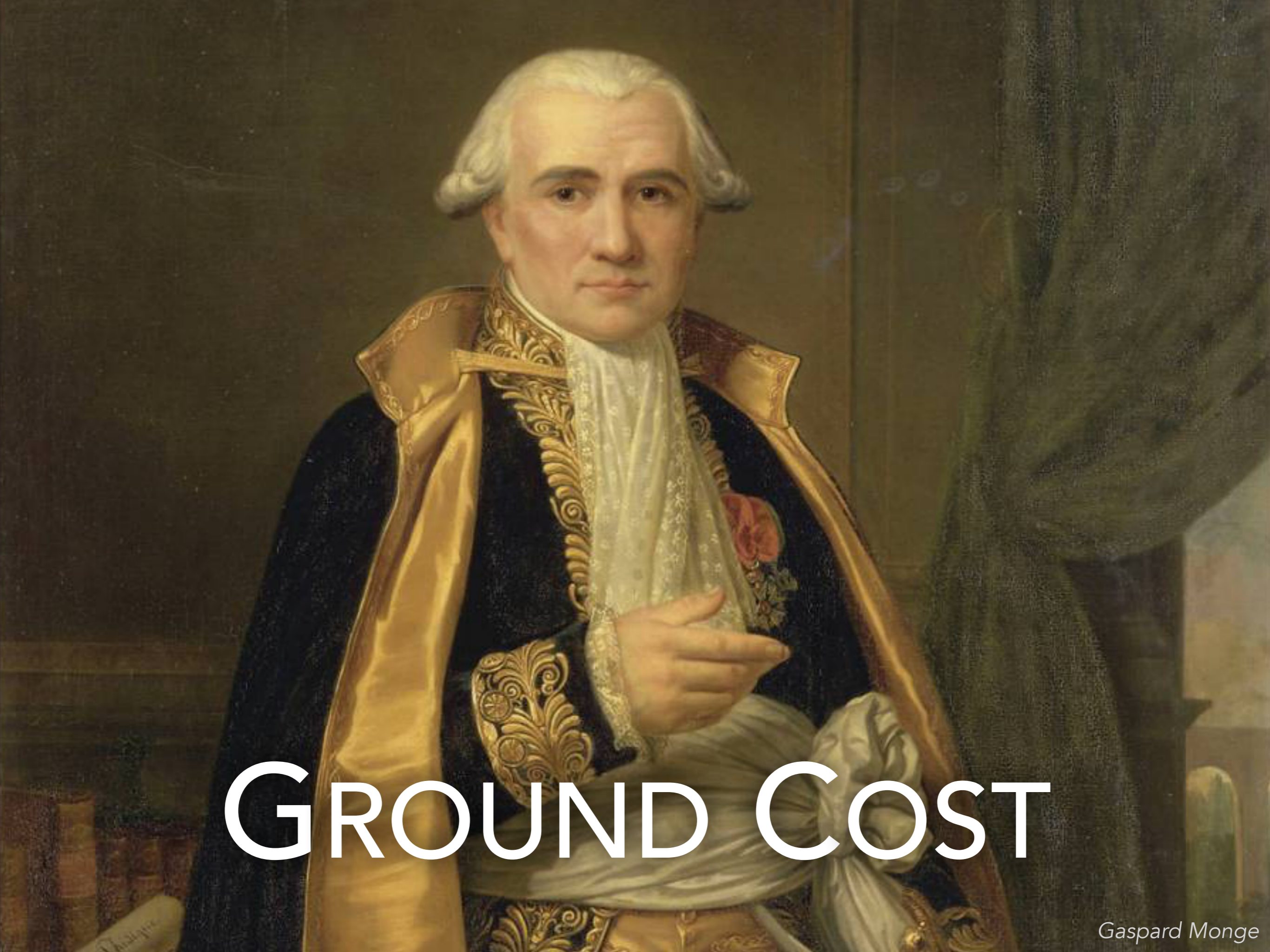
Two main questions in practice

1. How to choose the ground cost c in a way that makes sense for the data distributions μ and ν ?

OPTIMAL TRANSPORT

Two main questions in practice

1. How to choose the ground cost c in a way that makes sense for the data distributions μ and ν ?
2. How to compute/approximate the OT cost $\mathcal{T}_c(\mu, \nu)$, at least when the measures are discrete (*i.e. are finite sums of Dirac masses*) in a scalable way?



GROUND COST

GROUND COST

1. How to choose the ground cost c in a way that makes sense for the data distributions μ and ν ?

GROUND COST

1. How to choose the ground cost c in a way that makes sense for the data distributions μ and ν ?

1. Monge initially proposed $c(x, y) = \|x - y\|$

GROUND COST

1. How to choose the ground cost c in a way that makes sense for the data distributions μ and ν ?

1. Monge initially proposed $c(x, y) = \|x - y\|$

2. This was generalized to cost functions of the form

$$c(x, y) = \|x - y\|^p \quad \text{where } p \geq 1$$

(in this case, we say that \mathcal{T}_c is the p -Wasserstein distance)

GROUND COST

1. How to choose the ground cost c in a way that makes sense for the data distributions μ and ν ?

1. Monge initially proposed $c(x, y) = \|x - y\|$

2. This was generalized to cost functions of the form

$$c(x, y) = \|x - y\|^p \quad \text{where } p \geq 1$$

(in this case, we say that \mathcal{T}_c is the p -Wasserstein distance)

But does it make sense when the ground space is high-dimensional



GROUND COST

1. How to choose the ground cost c in a way that makes sense for the data distributions μ and ν ?

1. Monge initially proposed $c(x, y) = \|x - y\|$

2. This was generalized to cost functions of the form

$$c(x, y) = \|x - y\|^p \quad \text{where } p \geq 1$$

(in this case, we say that \mathcal{T}_c is the p -Wasserstein distance)

But does it make sense when the ground space is high-dimensional



But does it make sense when the data lives on a low-dimensional manifold



GROUND COST

Idea: Find a ground cost c that is adversarial, *i.e.* that best separates the two distributions by maximizing the OT cost

$$\max_{c \in \mathcal{C}} \mathcal{T}_c(\mu, \nu) \quad \text{where } \mathcal{C} \text{ is a class of functions}$$

GROUND COST

Idea: Find a ground cost c that is adversarial, *i.e.* that best separates the two distributions by maximizing the OT cost

$$\max_{c \in \mathcal{C}} \mathcal{T}_c(\mu, \nu) \quad \text{where } \mathcal{C} \text{ is a class of functions}$$



$$\max_c \mathcal{T}_c(\mu, \nu) - f(c) \quad \text{for some convex } f$$

$$f(c) = \begin{cases} 0 & \text{if } c \in \mathcal{C} \\ +\infty & \text{if } c \notin \mathcal{C} \end{cases}$$

GROUND COST

Idea: Find a ground cost c that is adversarial, *i.e.* that best separates the two distributions by maximizing the OT cost

$$\max_c \mathcal{T}_c(\mu, \nu) - f(c) \quad \text{for some convex } f$$

- Links with the Robust Optimization literature
- Links with the matchings literature in Economics
- Initially proposed by Genevay *et al.* in 2017 to learn generative models
- When \mathcal{C} is the set of Mahalanobis distances, it defines the Subspace Robust Wasserstein distances (ICML 2019)



REGULARIZATION

REGULARIZATION

2. How to compute/approximate the OT cost $\mathcal{T}_c(\mu, \nu)$?

1. This is a Linear Program $\rightarrow \mathcal{O}(n^3)$ complexity
2. Entropic regularization $\rightarrow \mathcal{O}(n^2)$ Sinkhorn algorithm, GPU-friendly, differentiable...

$$\inf_{\pi} \iint c(x, y) d\pi(x, y) + \varepsilon R(\pi)$$

where $R(\pi) = \text{KL}(\pi || \mu \otimes \nu)$

Other regularizations have been proposed: e.g. quadratic, group-lasso, capacity constraints, with different algorithms and effects on the OT plan / value

REGULARIZATION

2. How to compute/approximate the OT cost $\mathcal{T}_c(\mu, \nu)$?

1. This is a Linear Program $\rightarrow \mathcal{O}(n^3)$ complexity
2. Entropic regularization $\rightarrow \mathcal{O}(n^2)$ Sinkhorn algorithm, GPU-friendly, differentiable...

$$\inf_{\pi} \iint c(x, y) d\pi(x, y) + \varepsilon R(\pi)$$

How can we interpret the effect of the regularization ?

TWO VIEWS OF THE SAME PHENOMENON



GROUND COST ROBUSTNESS \Leftrightarrow REGULARIZATION

Theorem: Regularized OT is ground cost adversarial in the following sense

$$\begin{aligned} & \inf_{\pi} \iint c_0(x, y) d\pi(x, y) + \varepsilon R(\pi) \\ &= \sup_c \mathcal{J}_c(\mu, \nu) - \varepsilon R^* \left(\frac{c - c_0}{\varepsilon} \right) \end{aligned}$$

where R is a convex regularizer

and R^* is the convex conjugate of R :

$$R^*(c) = \sup_{\pi} \int c d\pi - R(\pi)$$

GROUND COST ROBUSTNESS \Leftrightarrow REGULARIZATION

Theorem: Regularized OT is ground cost adversarial in the following sense

$$\begin{aligned} & \inf_{\pi} \iint c_0(x, y) d\pi(x, y) + \varepsilon R(\pi) \\ &= \sup_c \mathcal{T}_c(\mu, \nu) - \varepsilon R^* \left(\frac{c - c_0}{\varepsilon} \right) \end{aligned}$$

Is the adversarial cost c_* an interesting dissimilarity measure on the ground space



GROUND COST ROBUSTNESS \Leftrightarrow REGULARIZATION

Is the adversarial cost c_* an interesting
dissimilarity measure on the ground space



Short answer: In a sense, no.

GROUND COST ROBUSTNESS \Leftrightarrow REGULARIZATION

Is the adversarial cost c_* an interesting dissimilarity measure on the ground space



Short answer: In a sense, no.

Theorem: Under some technical assumption on R (verified for the entropic or quadratic regularizations), there exists functions ϕ and ψ such that

$$c : (x, y) \mapsto \phi(x) + \psi(y)$$

is an optimal adversarial cost, i.e. is solution to

$$\sup_c \mathcal{J}_c(\mu, \nu) - \varepsilon R^* \left(\frac{c - c_0}{\varepsilon} \right)$$

WHAT I COULD NOT TALK ABOUT

- Restriction to nonnegative adversarial costs $\sup_{c \geq 0} \dots$
- General duality result for regularized OT
- Extension to several measures

Thank you

[francoispierrepaty.github.io](https://github.com/francoispierrepaty)