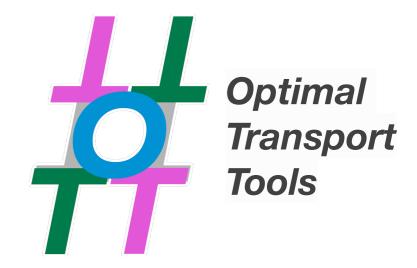
Differentiating through Optimal Transport

Marco Cuturi

EPFL - Zinal Jan. 2022







First, a big thanks!

To Prof. Kuhn and his lab for organising this event!

To all my recent collaborators on these topics,



M. Seigal

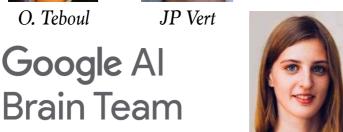


L. Papaxanthos





L. Petrini



C. Bunne



J. Niles-Weed







M. Scetbon



A. Gramfort

J. Josse



F.P. Paty



B. Muzellec



G. Peyré

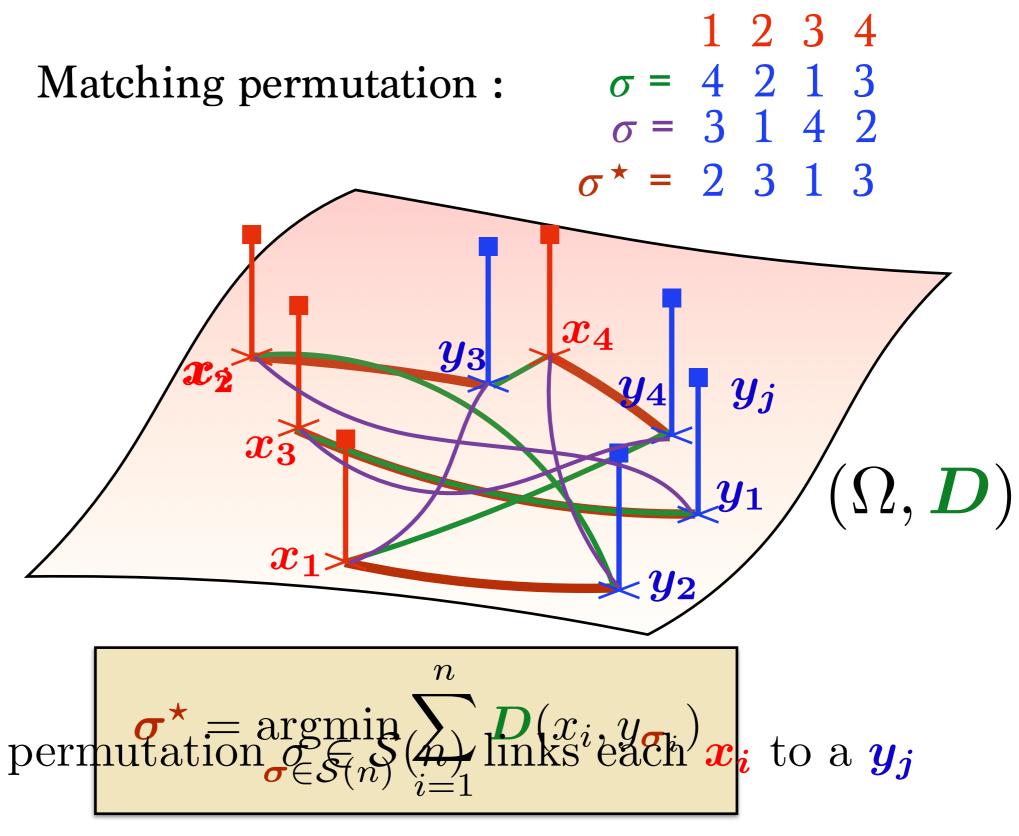






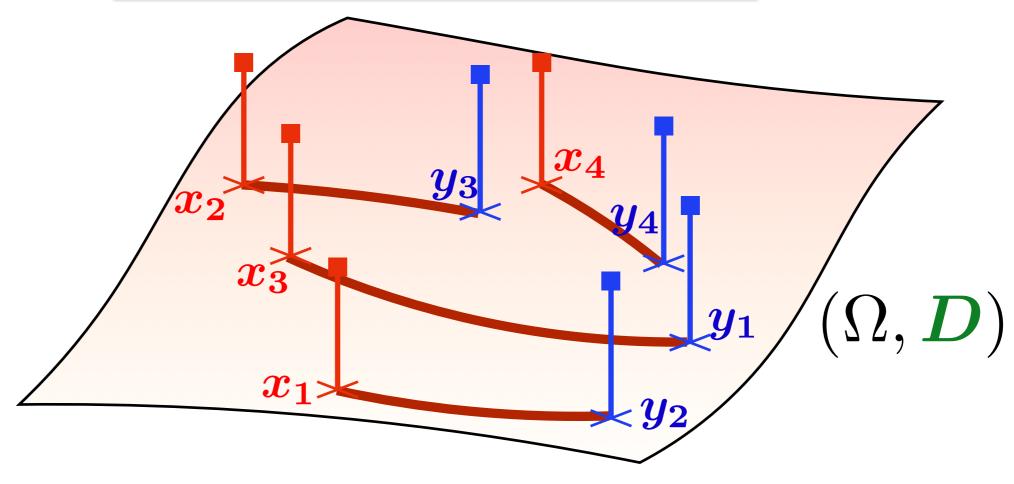


Warm-up: Optimal Matchings

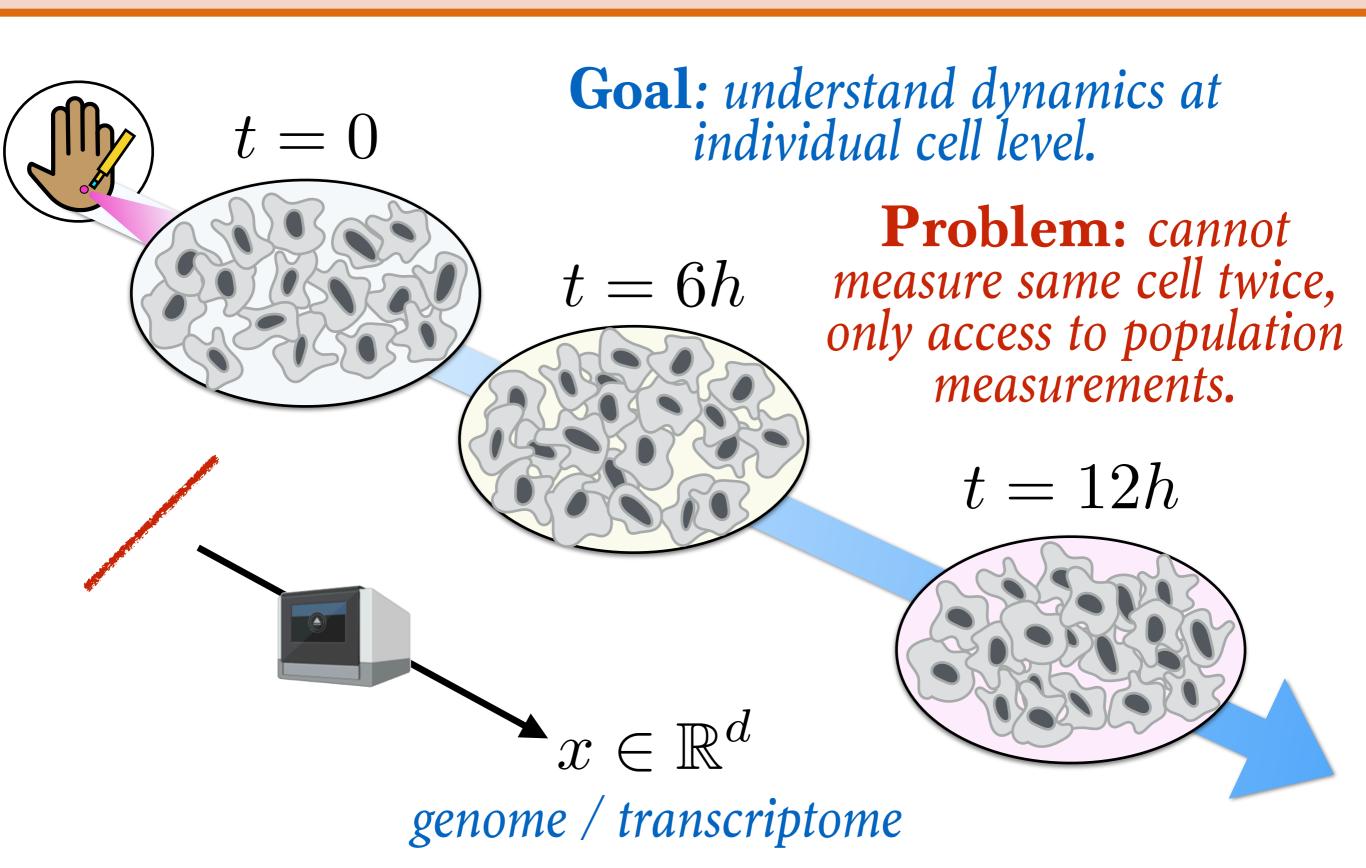


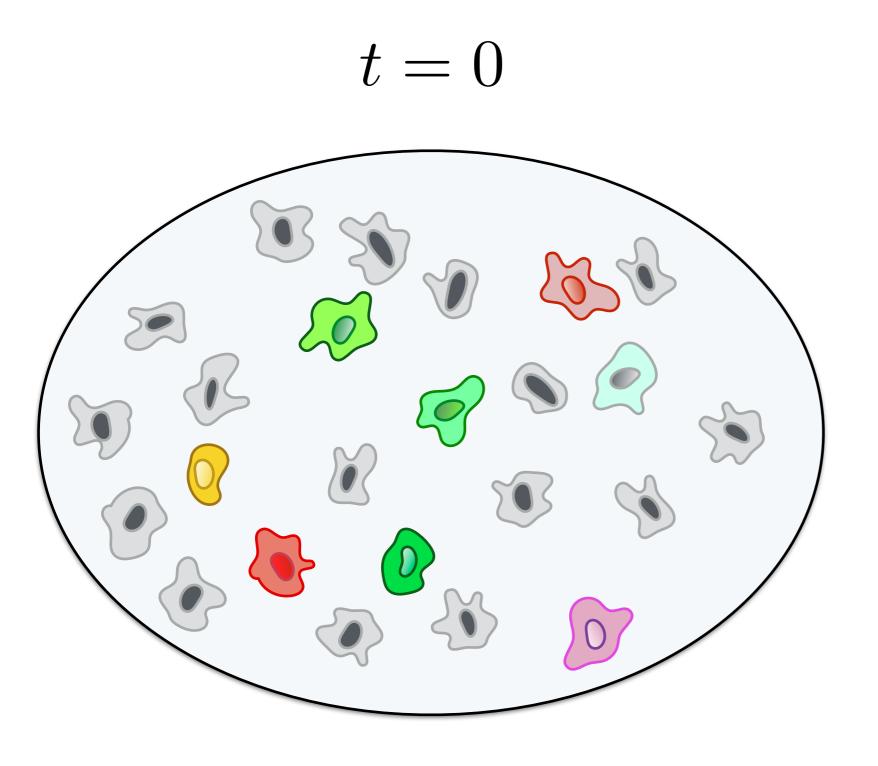
Warm-up: Optimal Matchings

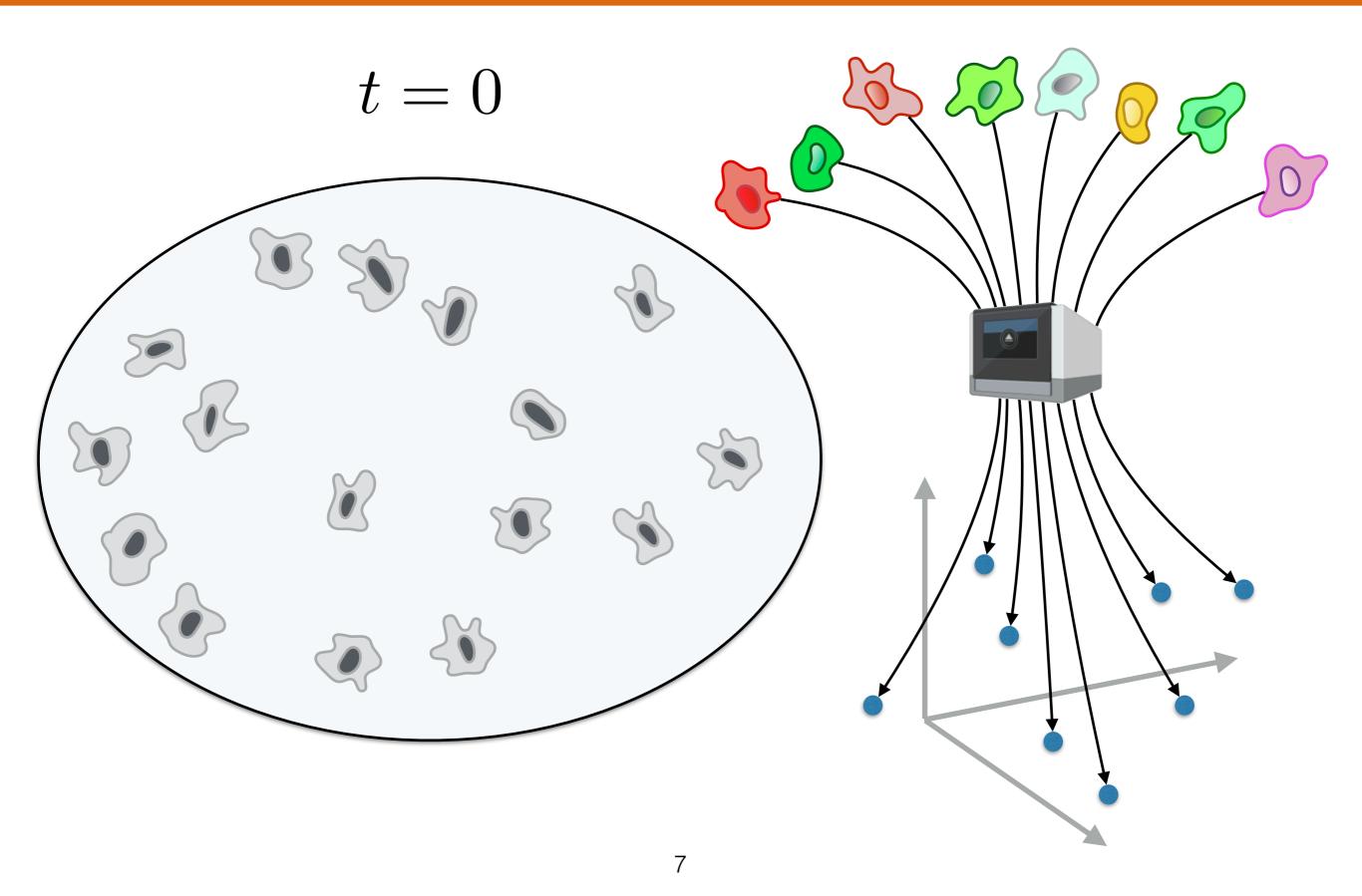
$$\boldsymbol{\sigma^{\star}} = \underset{\boldsymbol{\sigma} \in \mathcal{S}(n)}{\operatorname{argmin}} \sum_{i=1}^{n} \boldsymbol{D}(\boldsymbol{x_i}, \boldsymbol{y_{\sigma_i}})$$

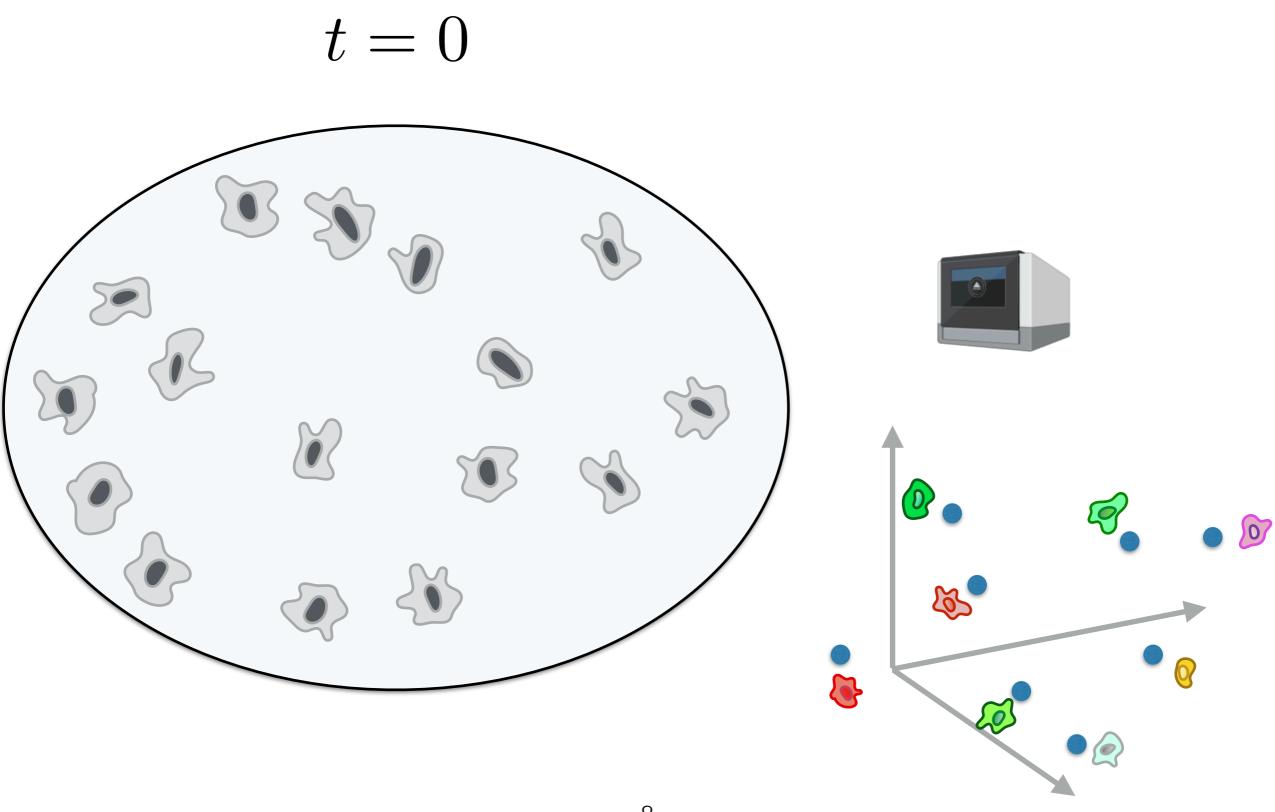


$$W(\{x_1 | y_{\boldsymbol{c}_i}, \{y_1, y_2, y_3\}, \{y_1, y_2, y_3\},$$

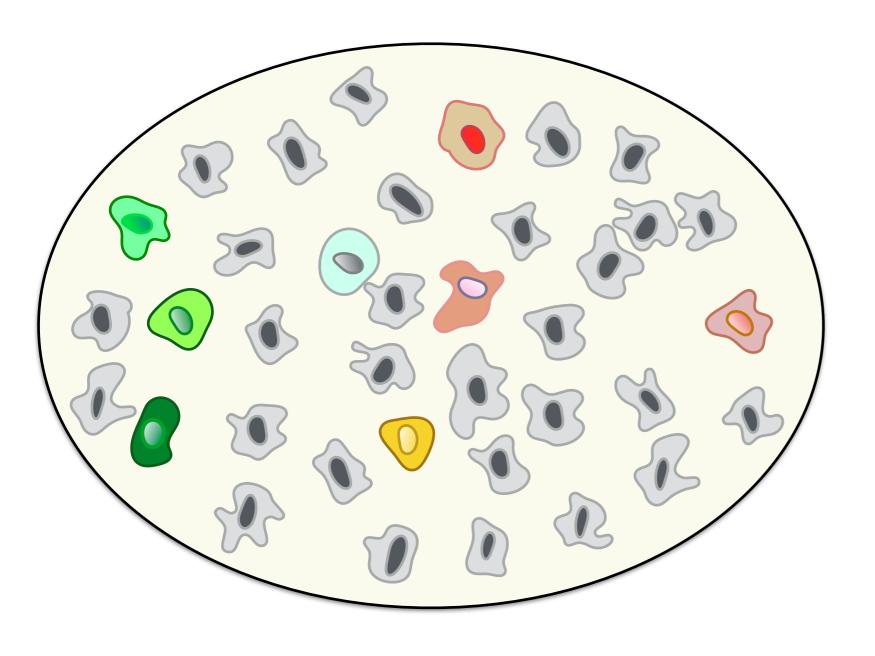




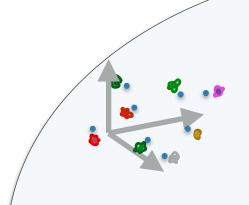


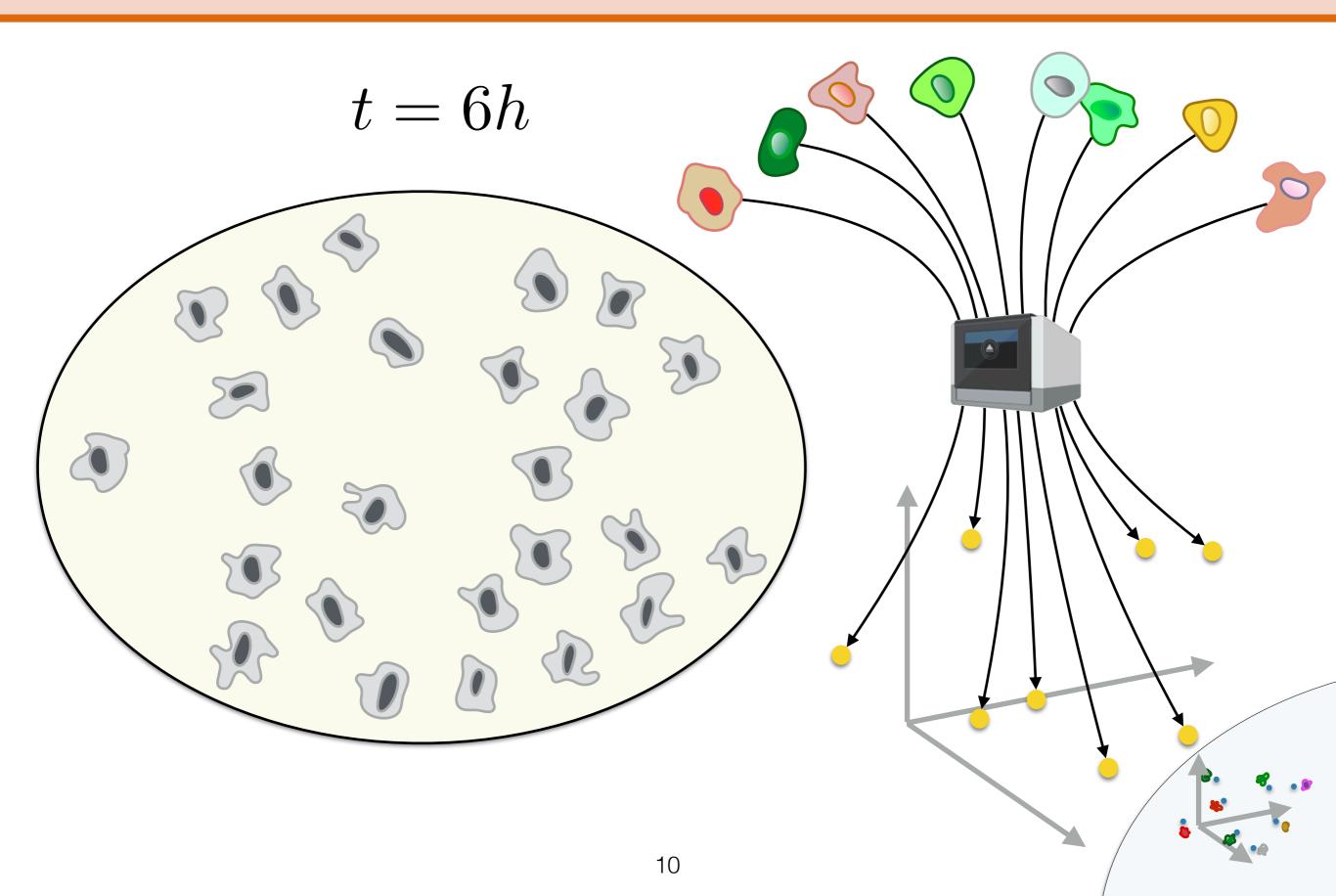


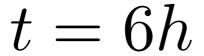
$$t = 6h$$

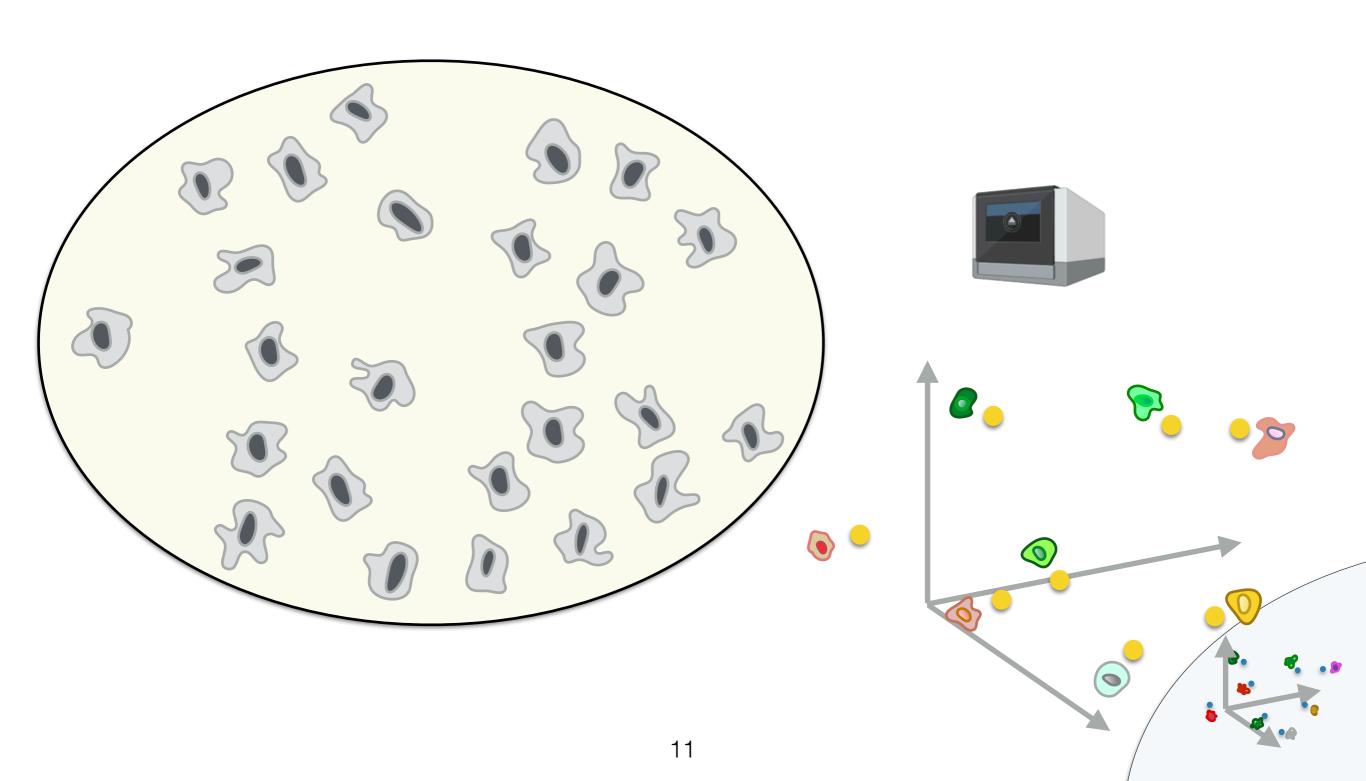




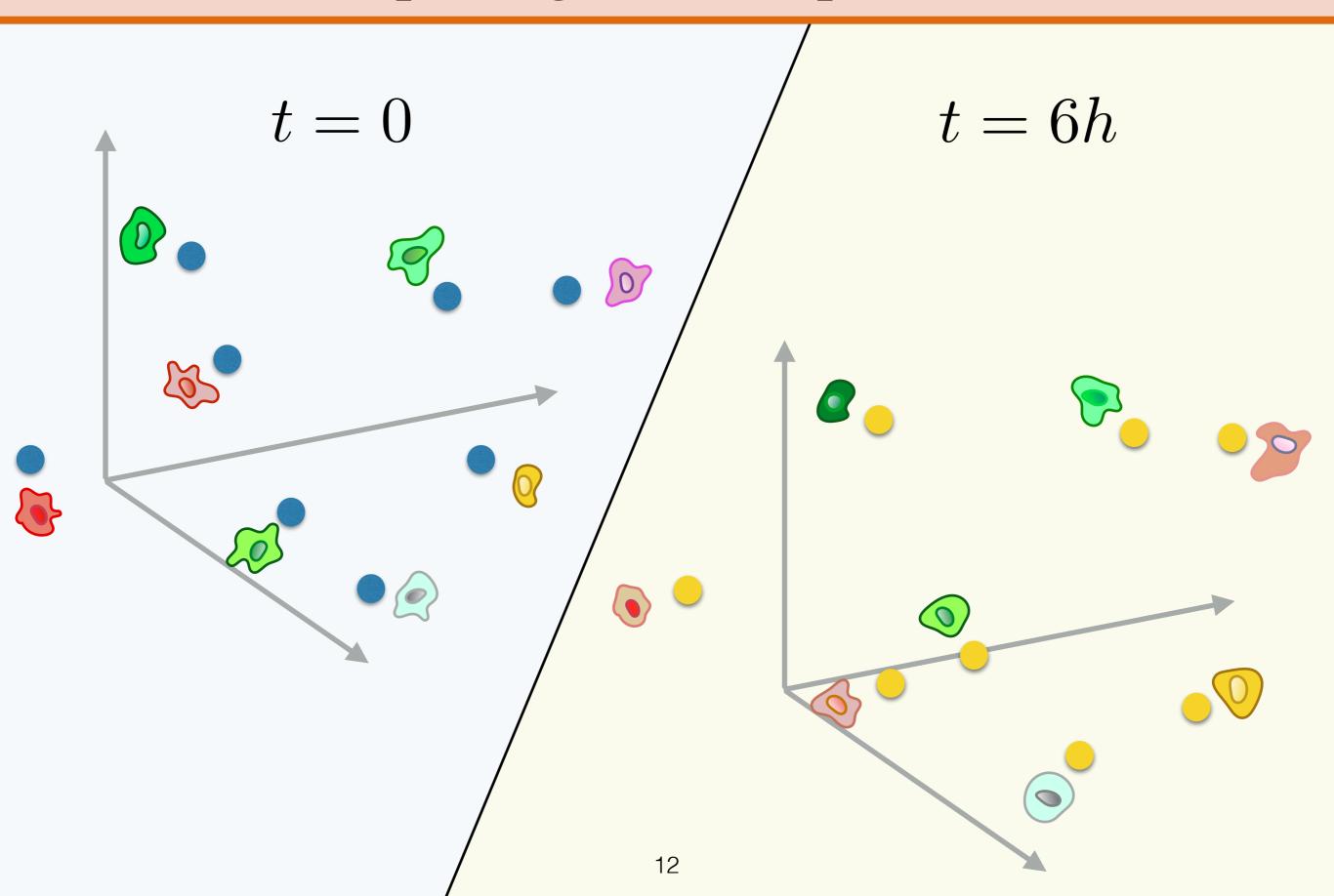




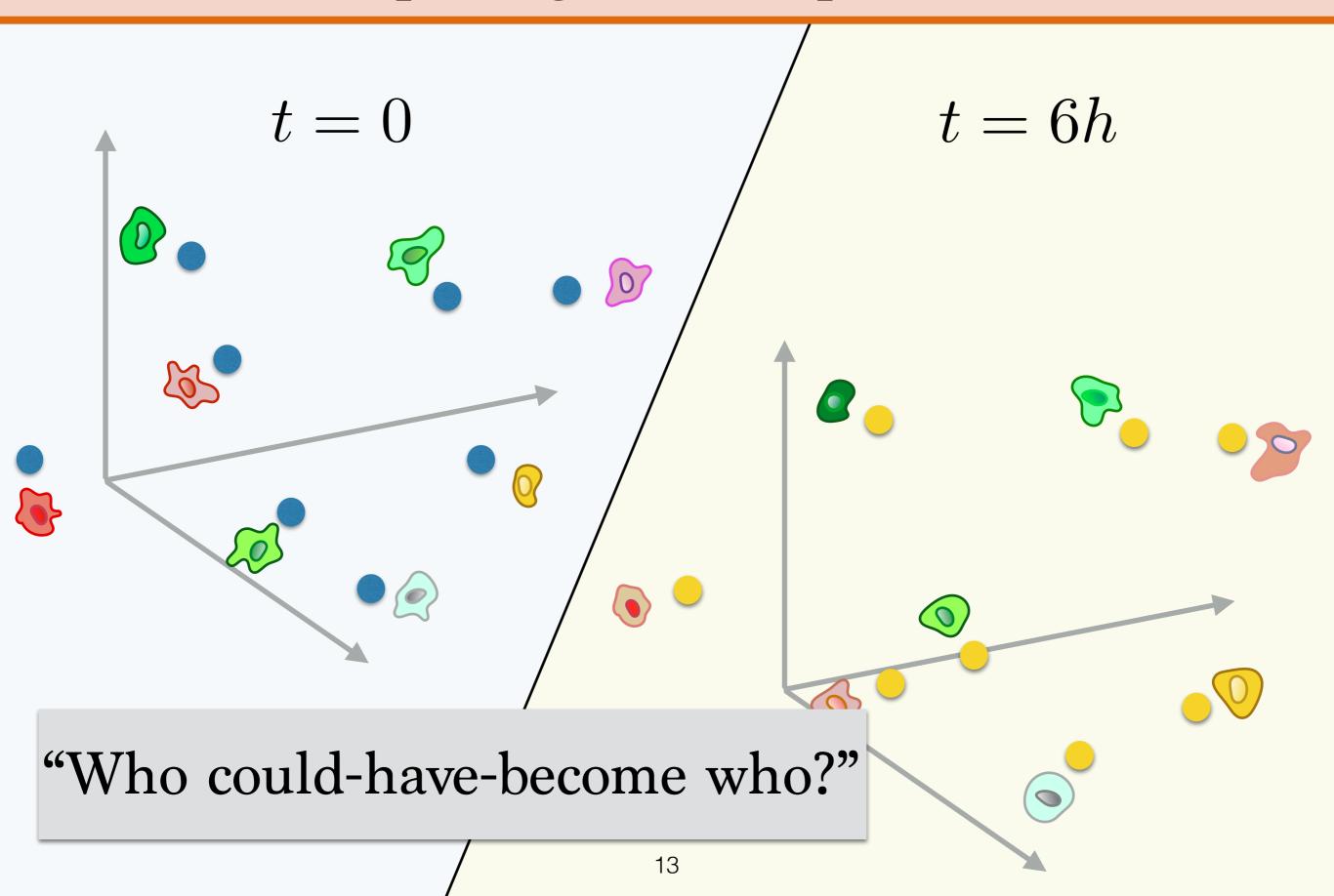




Comparing Two Populations

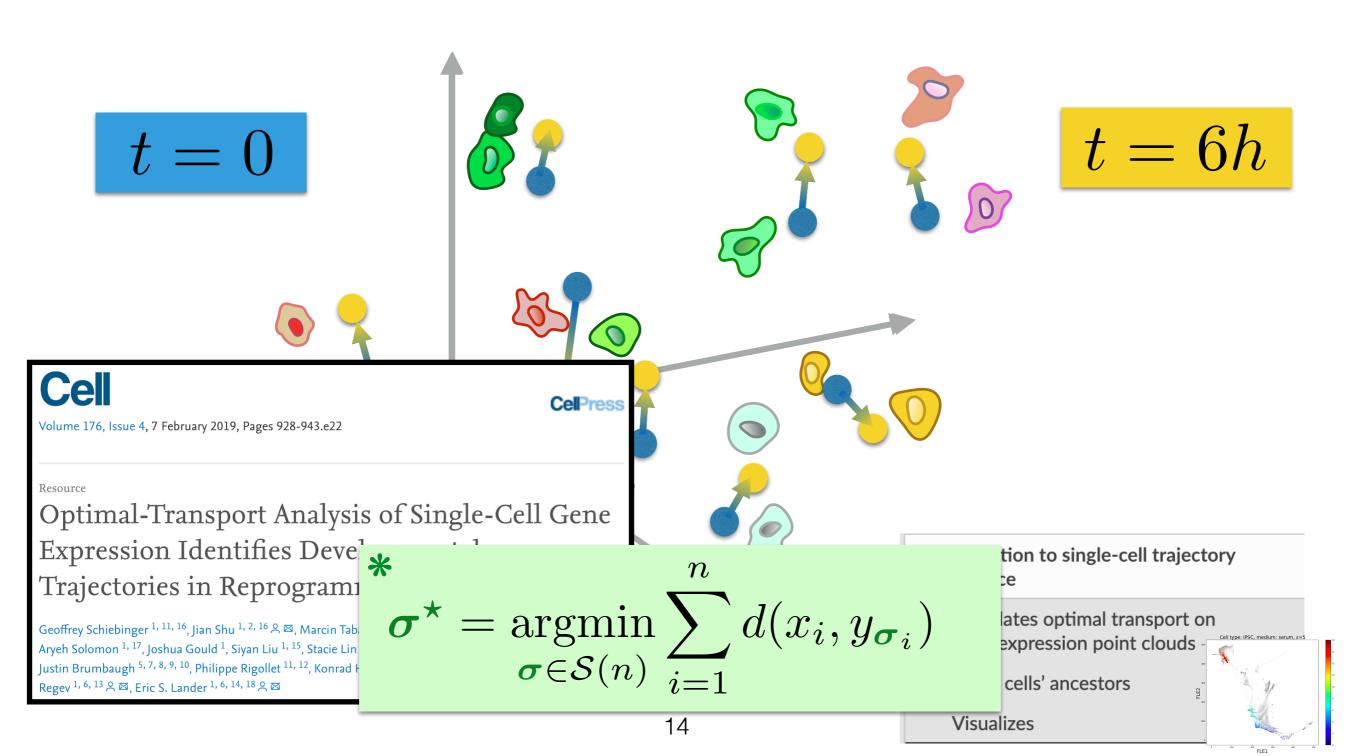


Comparing Two Populations



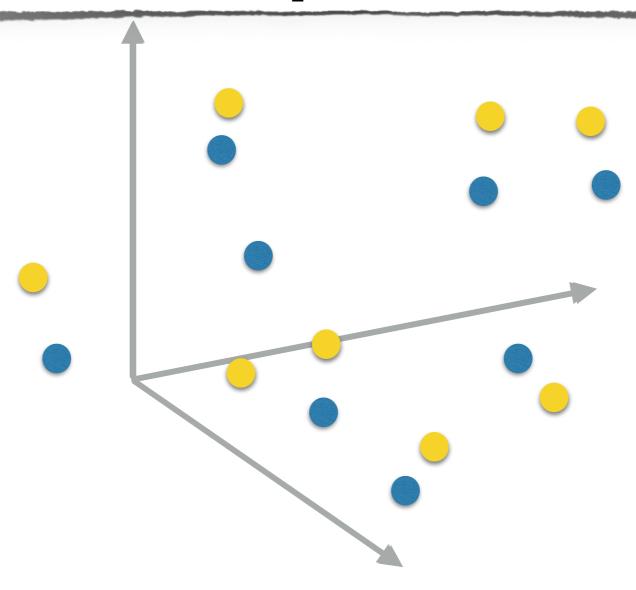
Matching Two Populations

to answer that question, use an optimal matching*

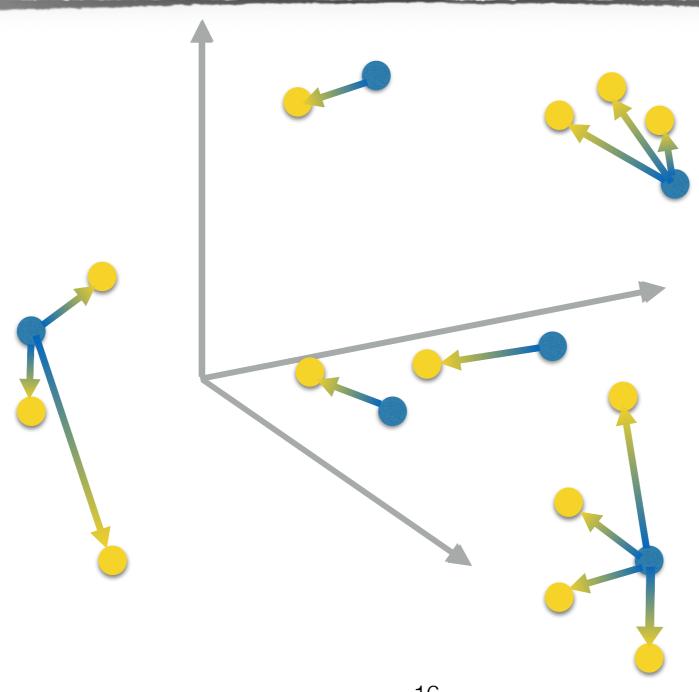


From there, many generalisations

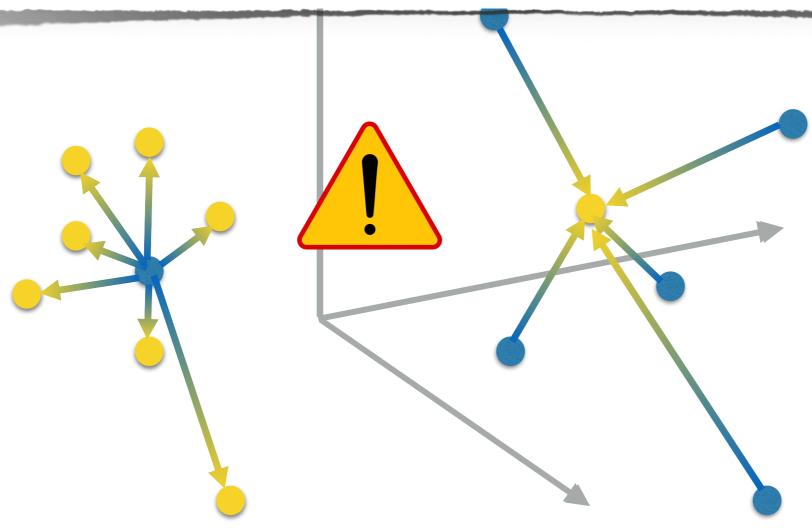
This story describes an ideal experimental setting, notably because this requires the same number of cells.



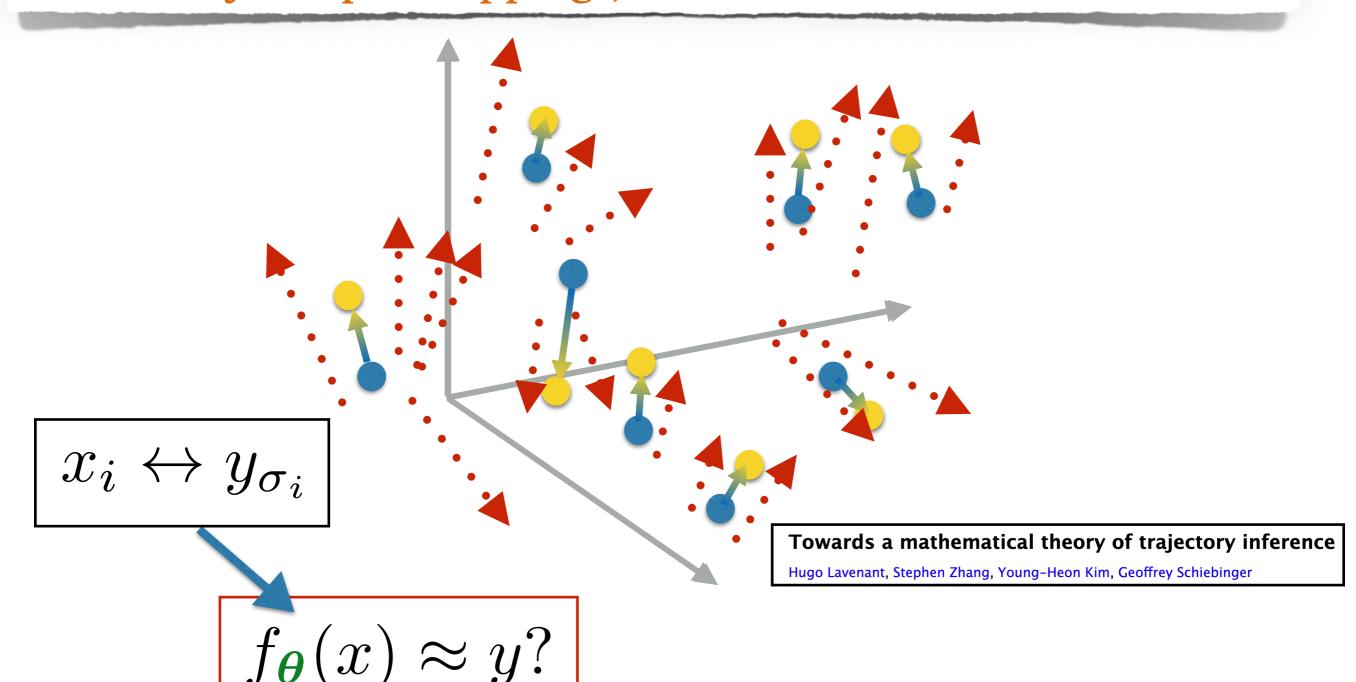
More likely to happen: variable number of cells. This requires a more versatile notion of matching,

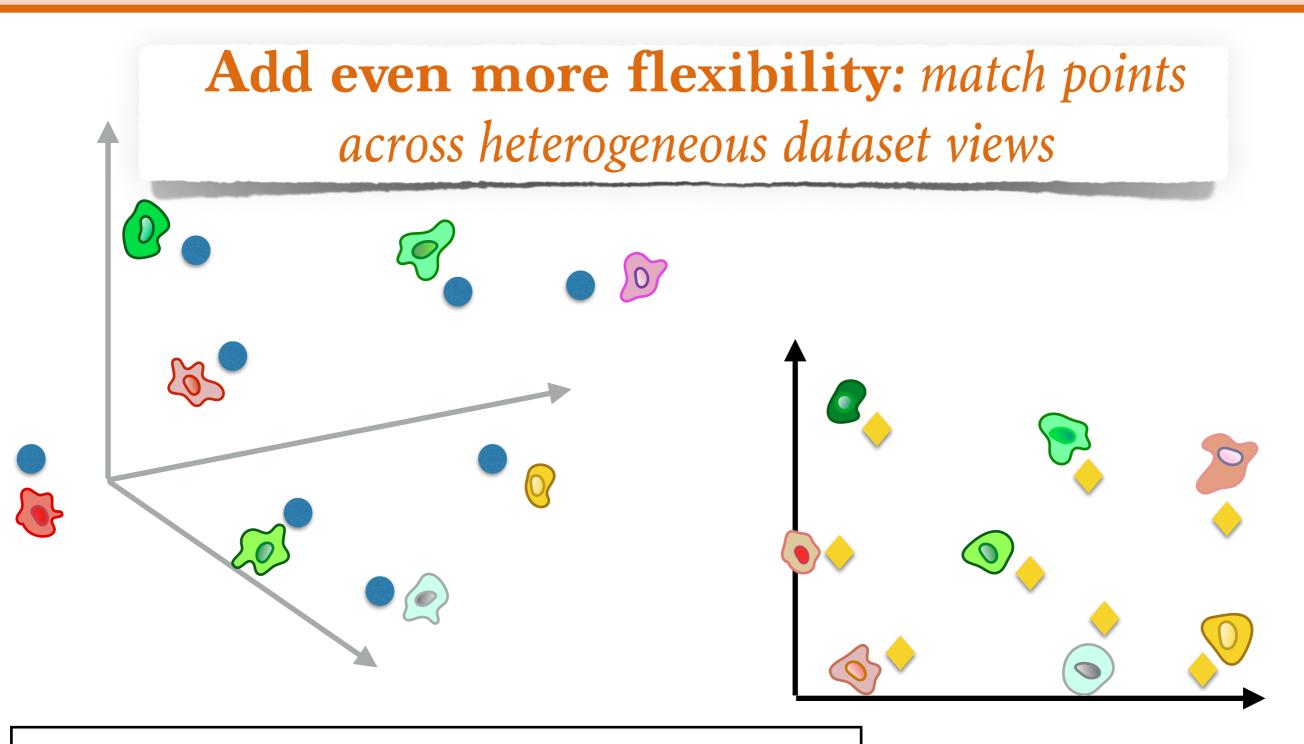


More likely to happen: variable number of cells. This requires a more versatile notion of matching, one that is more robust than simple nearest-neighbor.



A broader goal: infer functions able to provide out of sample mappings, i.e. a continuous extension.



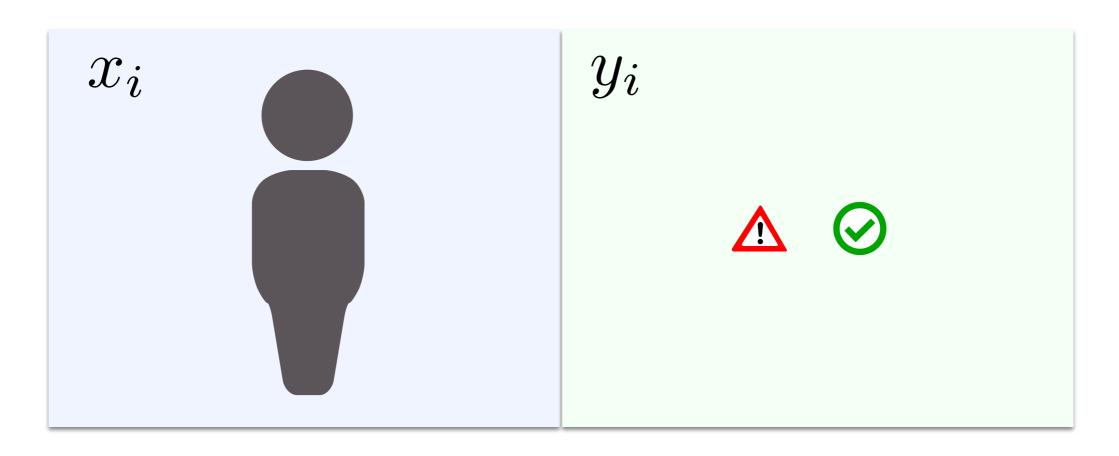


Gromov-Wasserstein optimal transport to align single-cell multi-omics data

Pinar Demetci, Rebecca Santorella, Björn Sandstede, D William Stafford Noble, Ritambhara Singh doi: https://doi.org/10.1101/2020.04.28.066787

This article is a preprint and has not been certified by peer review [what does this mean?].

Person x_i , caracteristic y_i we wish to predict.



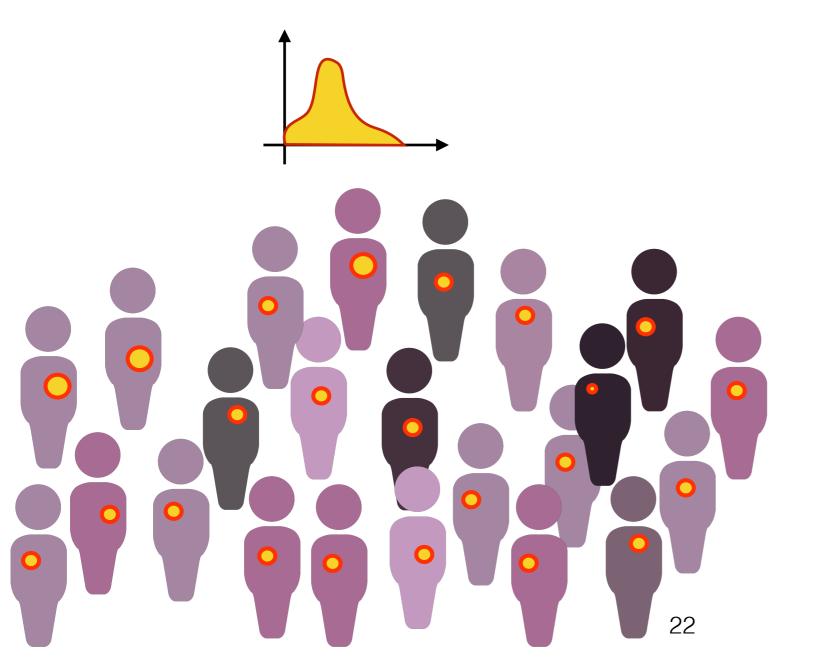
Find $\boldsymbol{\theta}$ such that $f_{\boldsymbol{\theta}}(x_i) \approx y_i$ typically $\min_{\boldsymbol{\theta}} \sum_i \ell(f_{\boldsymbol{\theta}}(x_i), y_i)$

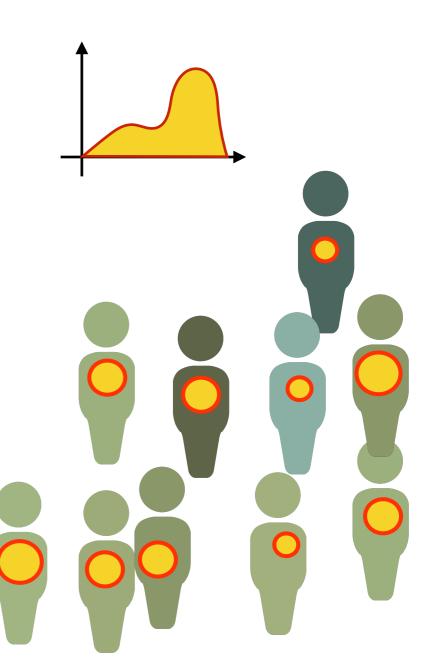
Found θ^* such that $f_{\theta^*}(x_i) \approx y_i$

For each individual i, a different error $\ell(f_{\theta^*}(x_i), y_i)$.



Obviously, θ^* achieved performance by skewing distribution of $\ell(f_{\theta}(x_i), y_i)$





Goal: enforce fairness mechanisms.

Problem: constraints are distributional, not individual



Goal: enforce fairness

fairness optimal transport

hor name Word matching

Show the calculator

About 72,400 results (0.05 sec; Showing 50 m

Problem: constraints are distributional, not individual

Fairness with Overlapping Groups

24

Forest Yang, Moustapha Cisse, Sanmi Koyejo

Obtaining fairness using optimal transport theory

P Gordaliza, E Del Barrio, G Fabrice - ... on Machine Learning, 2019 - proceedings.mlr.press
In the fair classification setup, we recast the links between fairness and predictability in terms of probability metrics. We analyze repair methods based on mapping conditional distributions to the Wasserstein barycenter. We propose a Random Repair which yields a ...

☆ 99 Cited by 55 Related articles >>>

Fliptest: fairness testing via optimal transport

<u>E Black, S Yeom, M Fredrikson, Fairness</u> - ... of the 2020 Conference on Fairness ..., 2020 We present FlipTest, a black-box technique for uncovering discrimination in classifiers. FlipTest is motivated by the intuitive question: had an individual been of a different protestatus, would the model have treated them differently? Rather than relying on causal ...

☆ 💯 Cited by 18 Related articles All 4 versions

Obtaining fairness using optimal transport theory

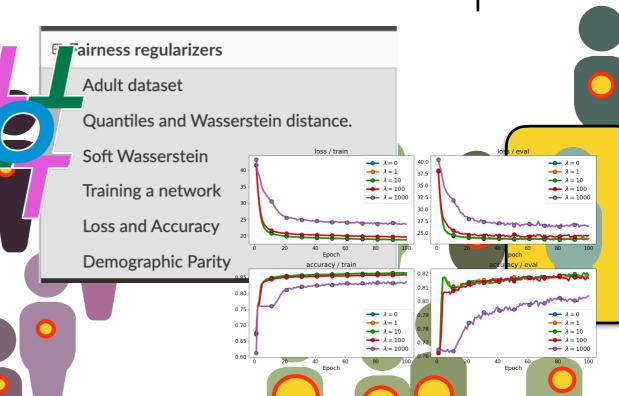
<u>E Del Barrio</u>, <u>F Gamboa</u>, <u>P Gordaliza</u> - arXiv preprint arXiv ..., 2018 - arxiv.org Statistical algorithms are usually helping in making decisions in many aspects of our lives. But, how do we know if these algorithms are biased and commit unfair discrimination of a particular group of people, typically a minority?\textit {**Fairness**} is generally studied in a ...

☆ 💯 Cited by 11 Related articles All 11 versions 🕸

A general approach to fairness with optimal transport

<u>C Silvia</u>, <u>J Ray</u>, <u>S Tom</u>, <u>P Aldo</u>, <u>J Heinrich</u> - Proceedings of the AAAI ..., 2020 - ojs.aaai.org We propose a general approach to **fairness** based on transporting distributions corresponding to different sensitive attributes to a common distribution. We use **optimal transport** theory to derive target distributions and methods that allow us to achieve **fairness**

☆ ワワ Cited by 7 Related articles All 2 versions ১৯

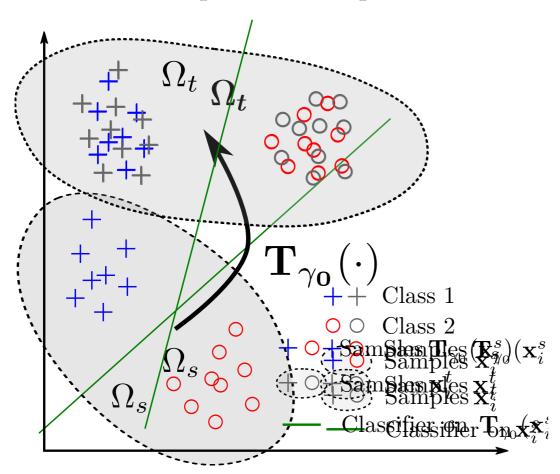


rly matched individuals

Optimal matchings appear everywhere!

Goal: domain adaptation

Classifica Coptinated states sported samples



Problem: train and test data must be matched first

EEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL. X, NO. X, JANUARY XX

Optimal Transport for Domain Adaptation

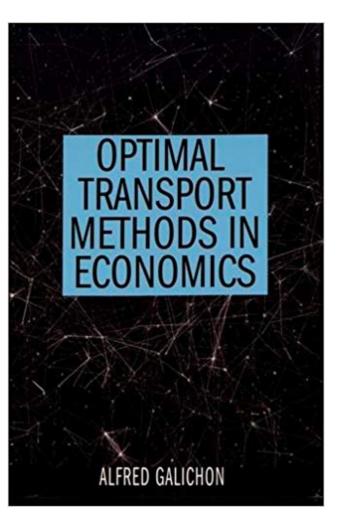
Nicolas Courty, Rémi Flamary, Devis Tuia, Senior Member, IEEE, Alain Rakotomamoniy, Member, IEEE

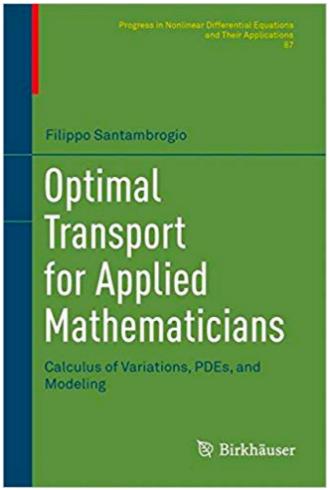
In all these problems...

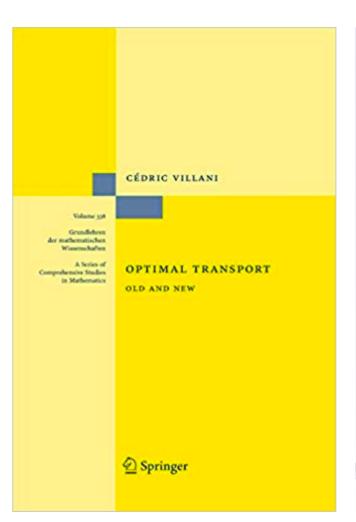
- •An intermediate matching/mapping step is needed.
- That step is only **intermediate** in the sense that other *advanced learning procedures* build on top of it.
- Wishlist for a *matching* framework:
 - guided by rich theory, to inform new algorithms,
 - versatile, to accommodate several settings,
 - scalable, to handle large dim/n datesets,
 - differentiable, for end-to-end learning,
 - robust, to perform well in most settings.

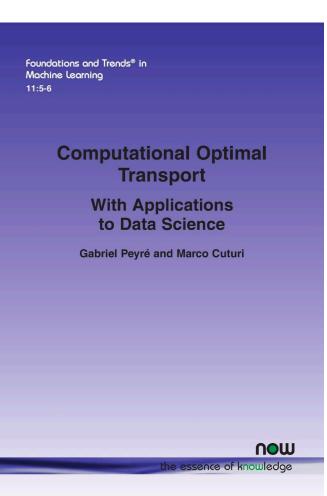
That theory is Optimal Transport

Matchings happen everywhere, so many people have (re)discovered it: economics, *applied* maths, *pure* maths, graphics, and, increasingly ML people!









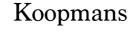
Why is it popular across so many fields?

It has a bit of everything! theory provides theorems; algorithms translate them in tangible results; applications guide new developments.









Nobel'75



Dantzig



McCann







Villani

Fields'10





Brenier





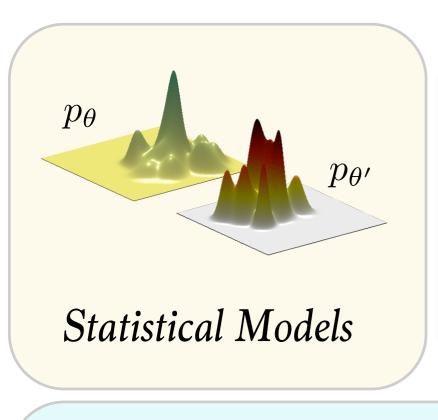




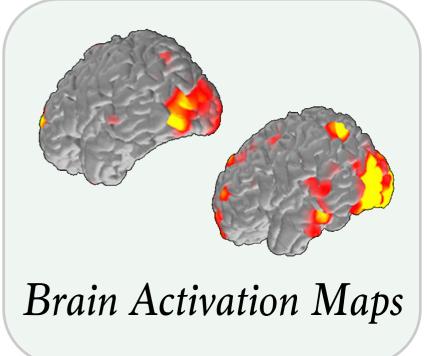


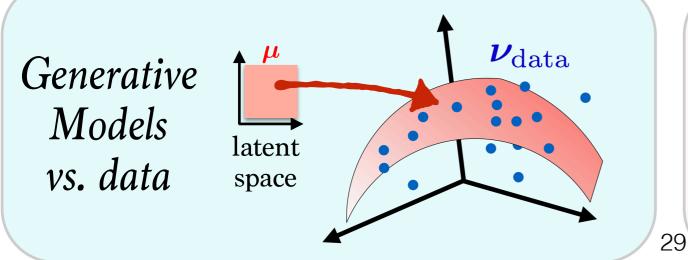
Optimal Transport in Data Sciences

Advertised as the natural way to define geometry, through matching, for probability measures, when supported on a geometric space.









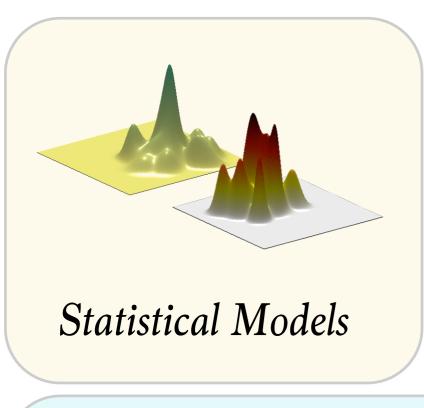




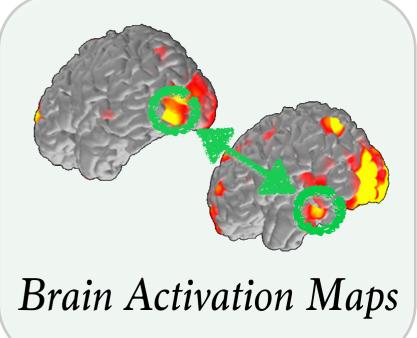
Color Histograms

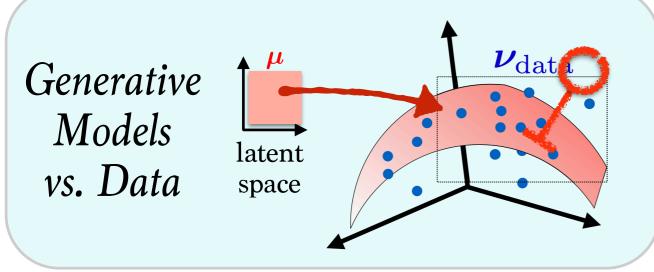
Optimal Transport in Data Sciences

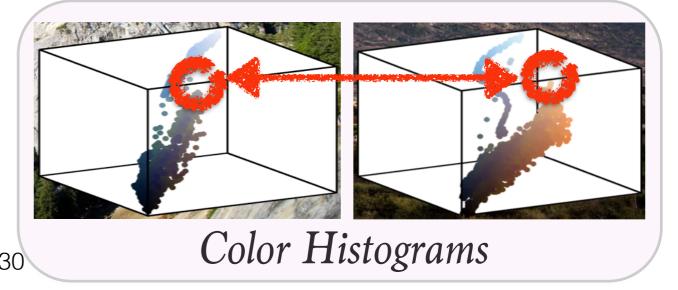
Advertised as the natural way to define geometry, through matching, for probability measures, when supported on a geometric space.











Short Course Outline

- 1. Introduction to optimal transport
- 2. Computing OT exactly
- 3. Computing OT for data science

Introduction to OT

- Two examples: moving earth & soldiers
- Monge problem, Kantorovich problem
- OT as geometry, OT as a loss function

Origins: Monge Problem (1781)

tripadvisor*

Paris ▼







c Profile

Join



Travel feed: Paris

Hotels

Things to do

Restaurants

Flights Holiday Homes

Shopping

Car Hire

•••

Europe > France > Ile-de-France > Paris > Things to do in Paris > Place Monge

Place Monge, Paris: Address, Place Monge Reviews: 4/5

Gaspard Monge (1746 - 1818)

Place Monge

84 Reviews

#323 of 2 272 things to do in Paris

Shopping, Flea & Street Markets

Pl. Monge, Paris, France

☐ Save A Share

Review Highlights

"Good local market."

Place Monge market is one of our regular haunts when staying in Paris, it has a good range of... read more



Reviewed 26 September 2018
johngl8492UH , Middle Park ☐ via mobile

"One of the most amazing streets we have be...

We loved place monge. All the shops slide their products out to the streets, it has a real Parisian... read more



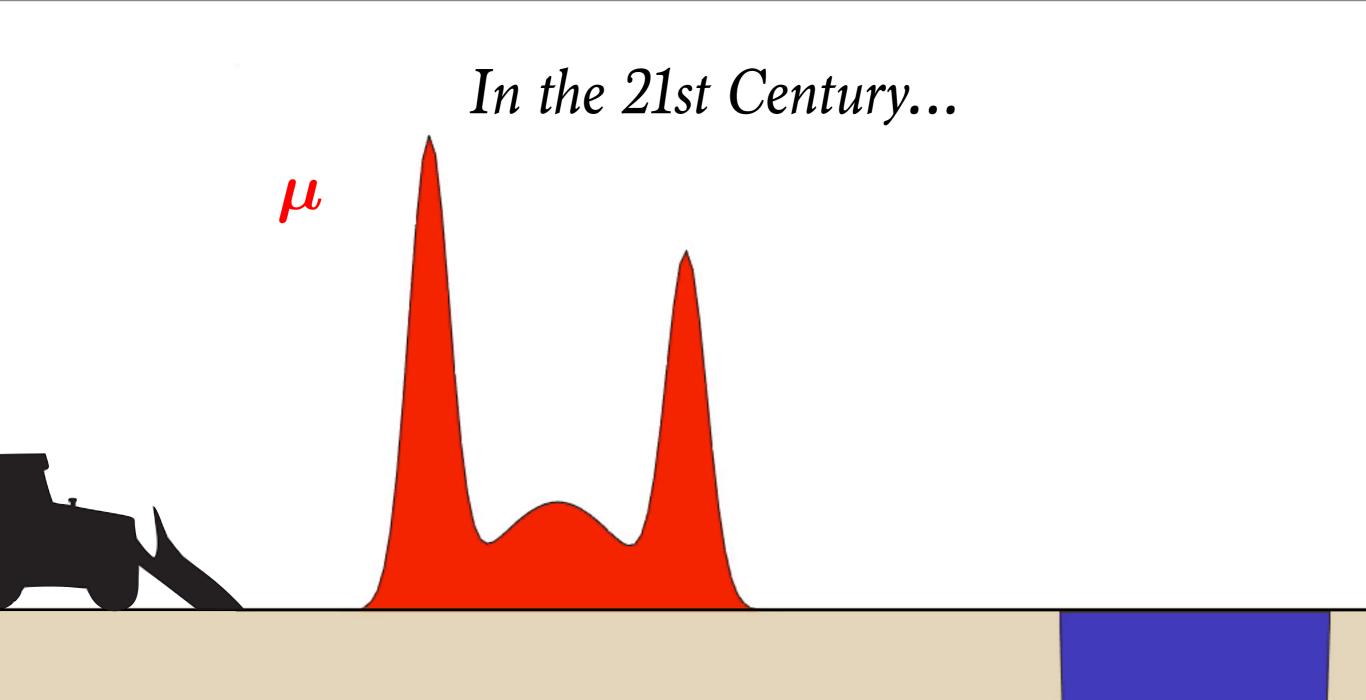
©©©© Reviewed 11 October 2018 Steve L , Pacific Coast Australia, Australia

Read all 84 reviews

via mobile

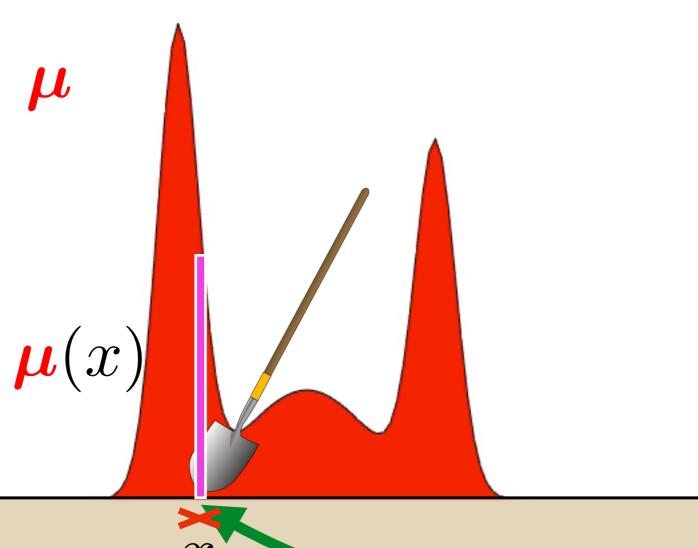


Origins: Monge Problem



Origins: Monge's Problem



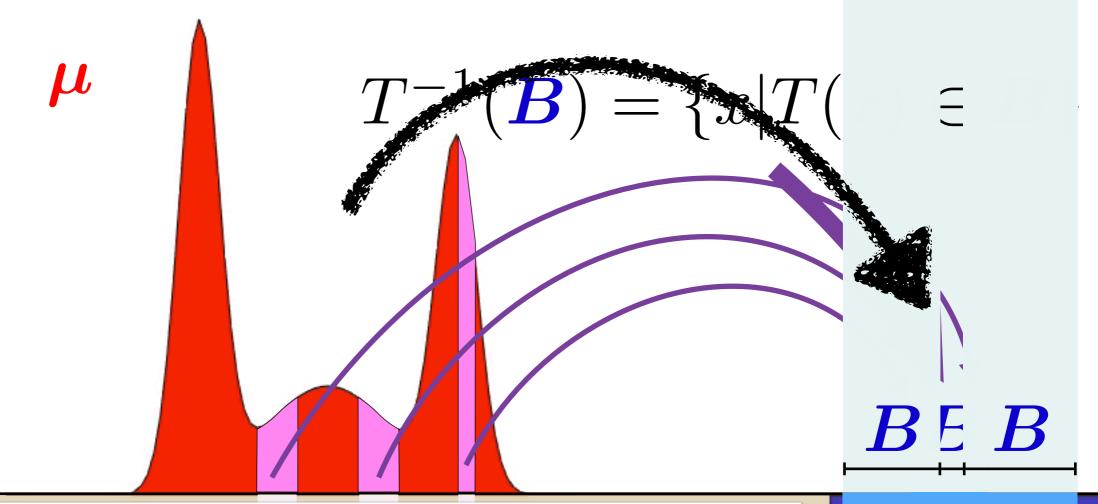


$$y = T(x)$$



Origins: Monge's Proble

T must push-forwandttheapedechtasblreetowards the blue



What T s.t. $T_{\sharp}\mu = \nu$ minimizes $\int D(x, T(x))\mu(dx)$?







Tolstoi 1930

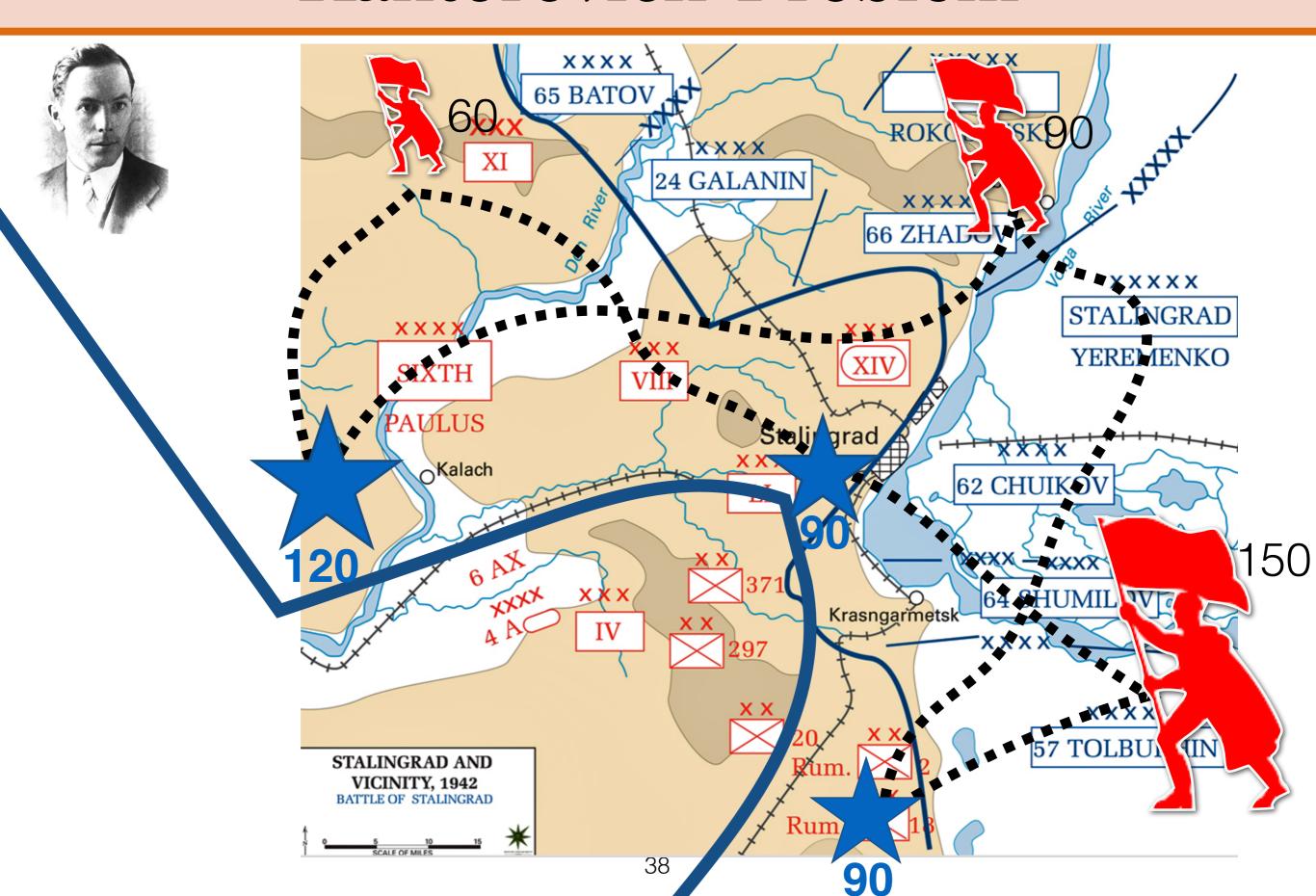


Hitchcock

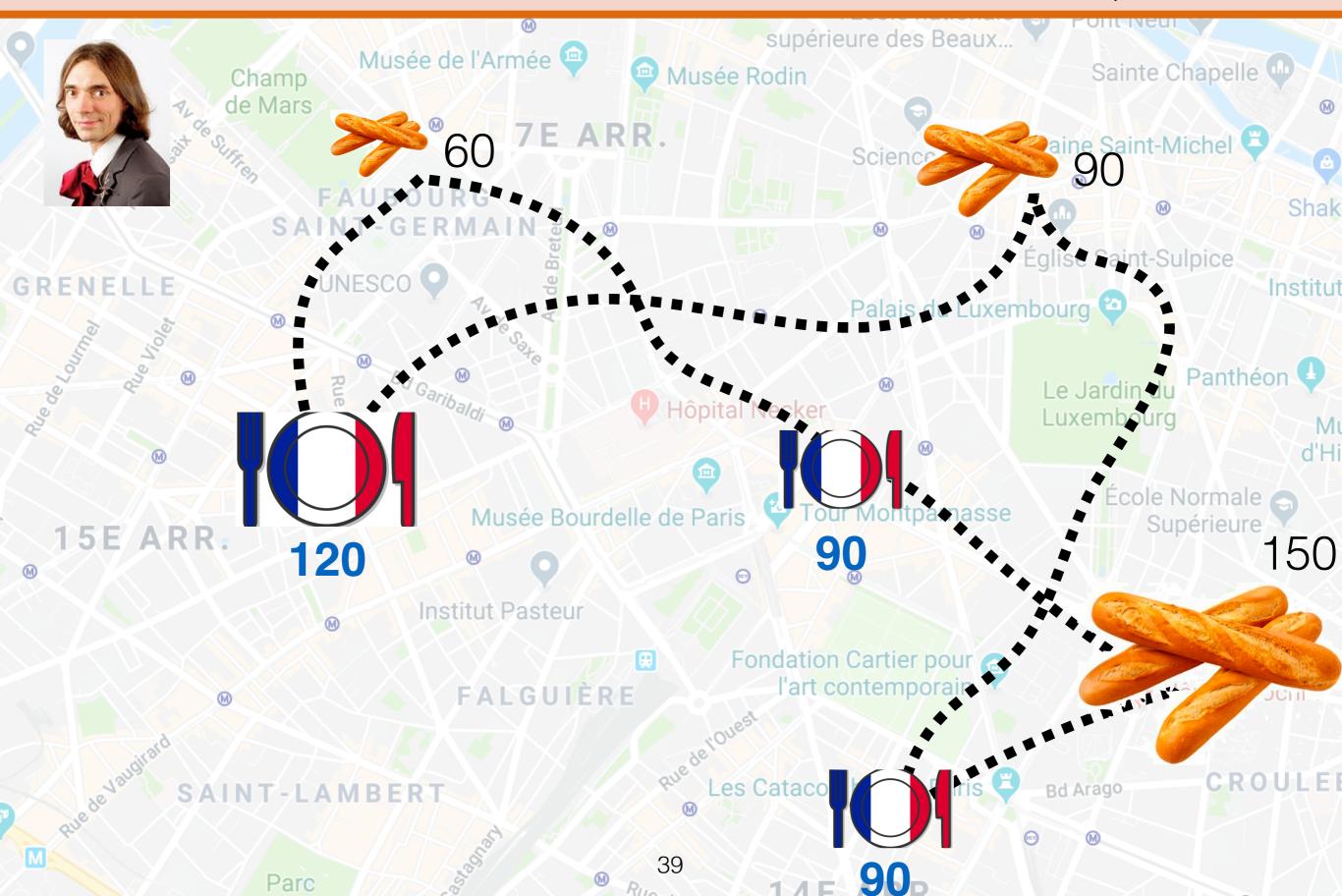
THE DISTRIBUTION OF A PRODUCT FROM SEVERAL SOURCES TO NUMEROUS LOCALITIES

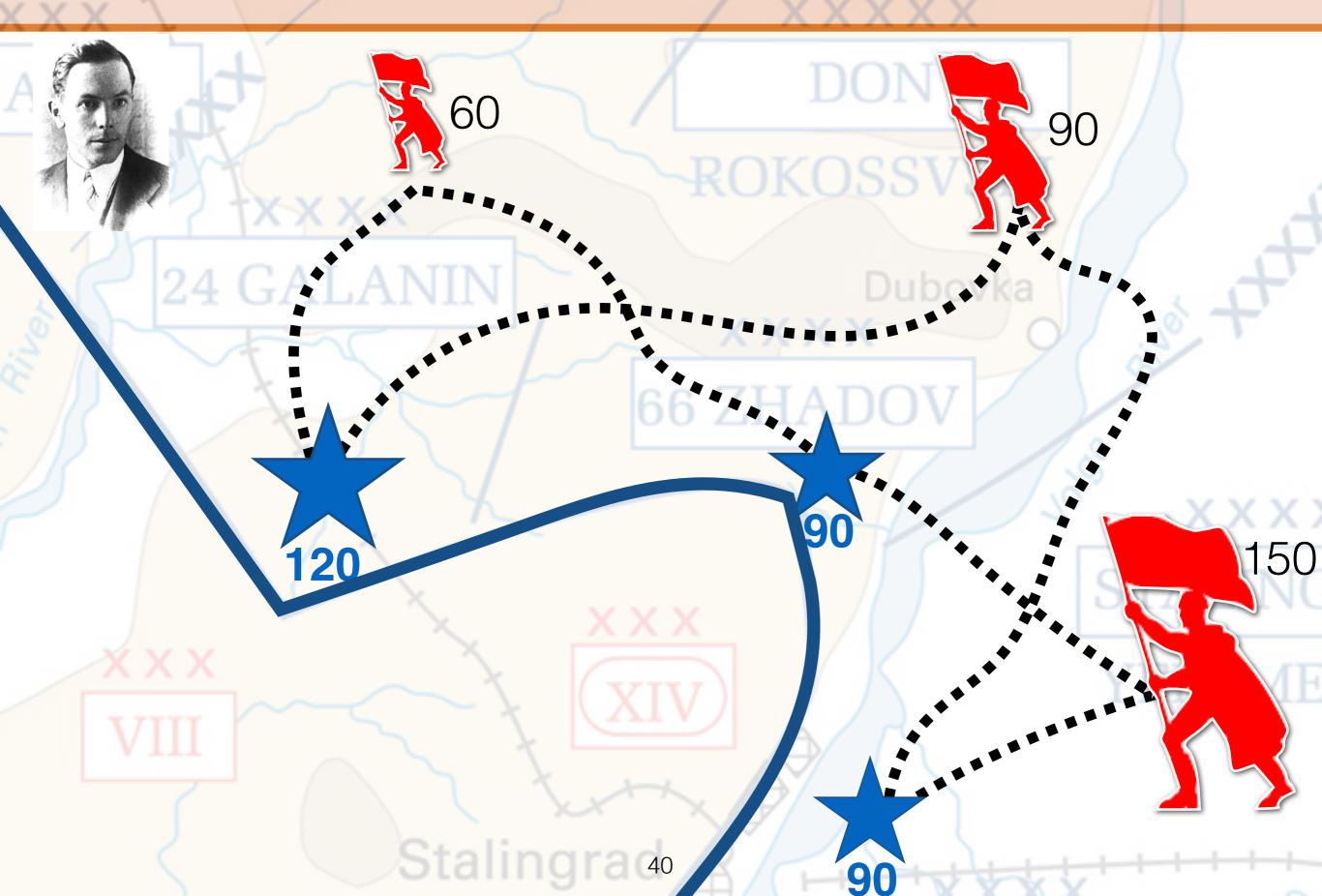
By Frank L. Hitchcock

1. Statement of the problem. When several factories supply a product to a number of cities we desire the least costly manner of distribution. Due to freight rates and other matters the cost of a ton of product to a particular city will vary according to which factory supplies it, and will also vary from city to city.

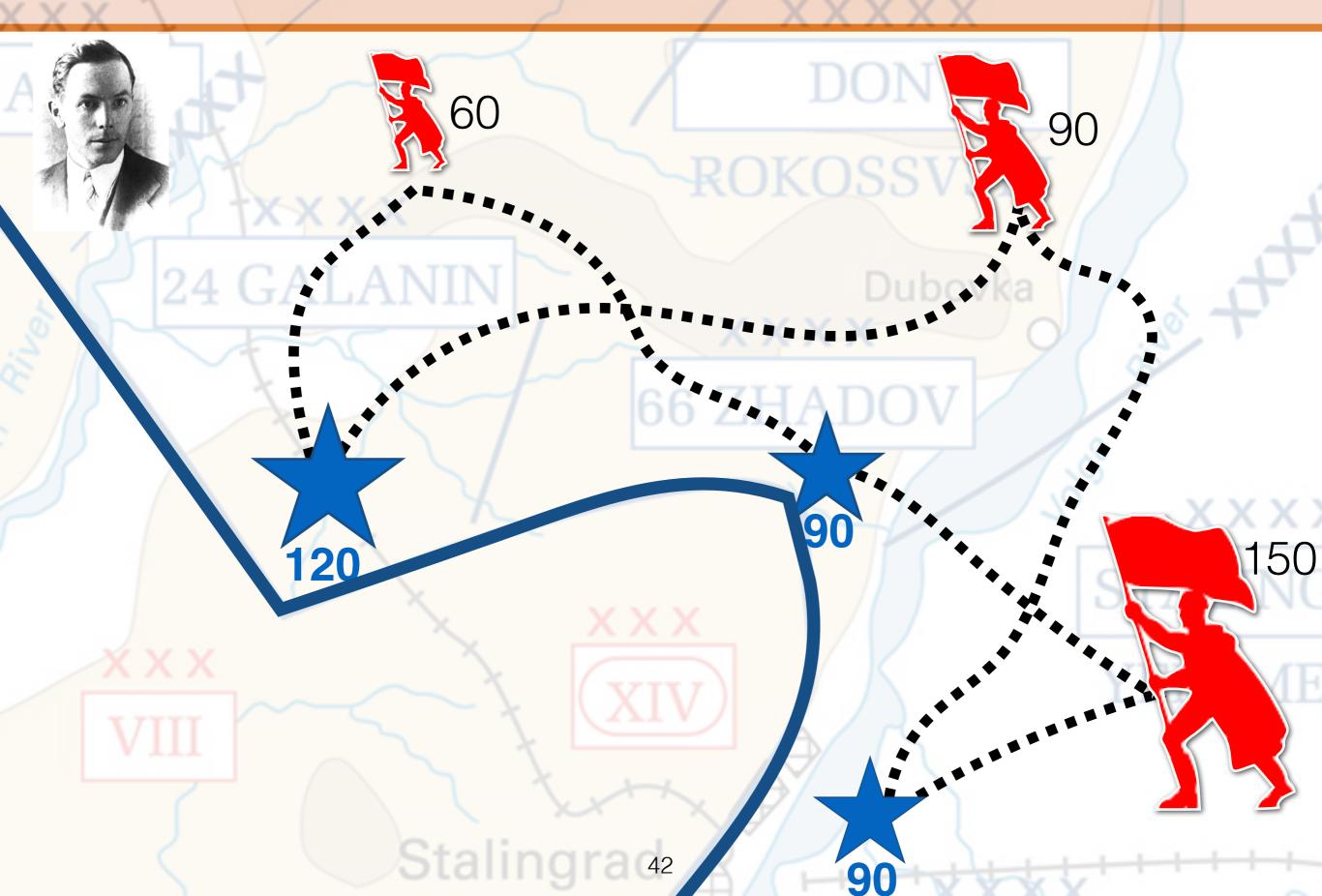


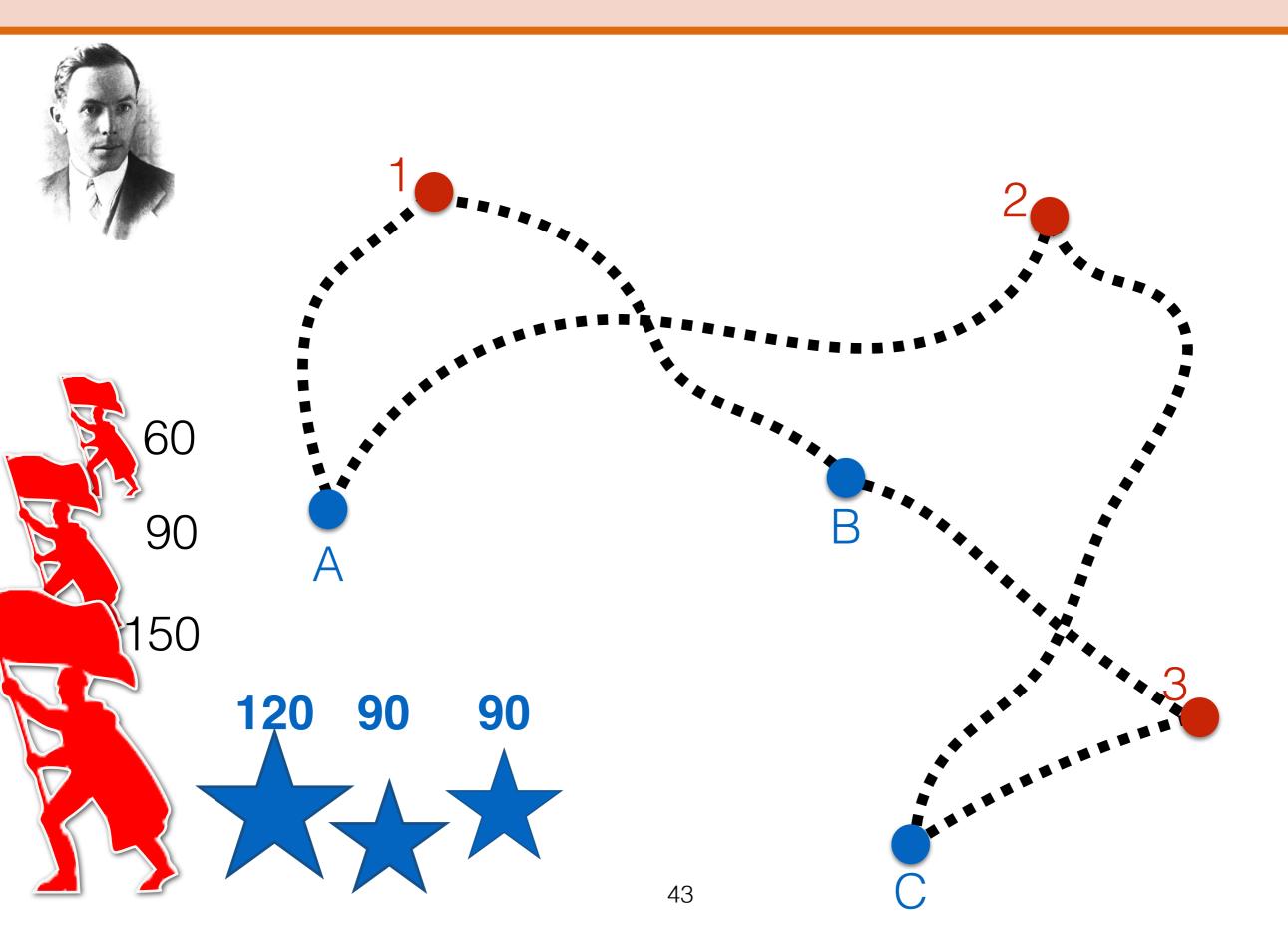
Kantorovich Problem à la française

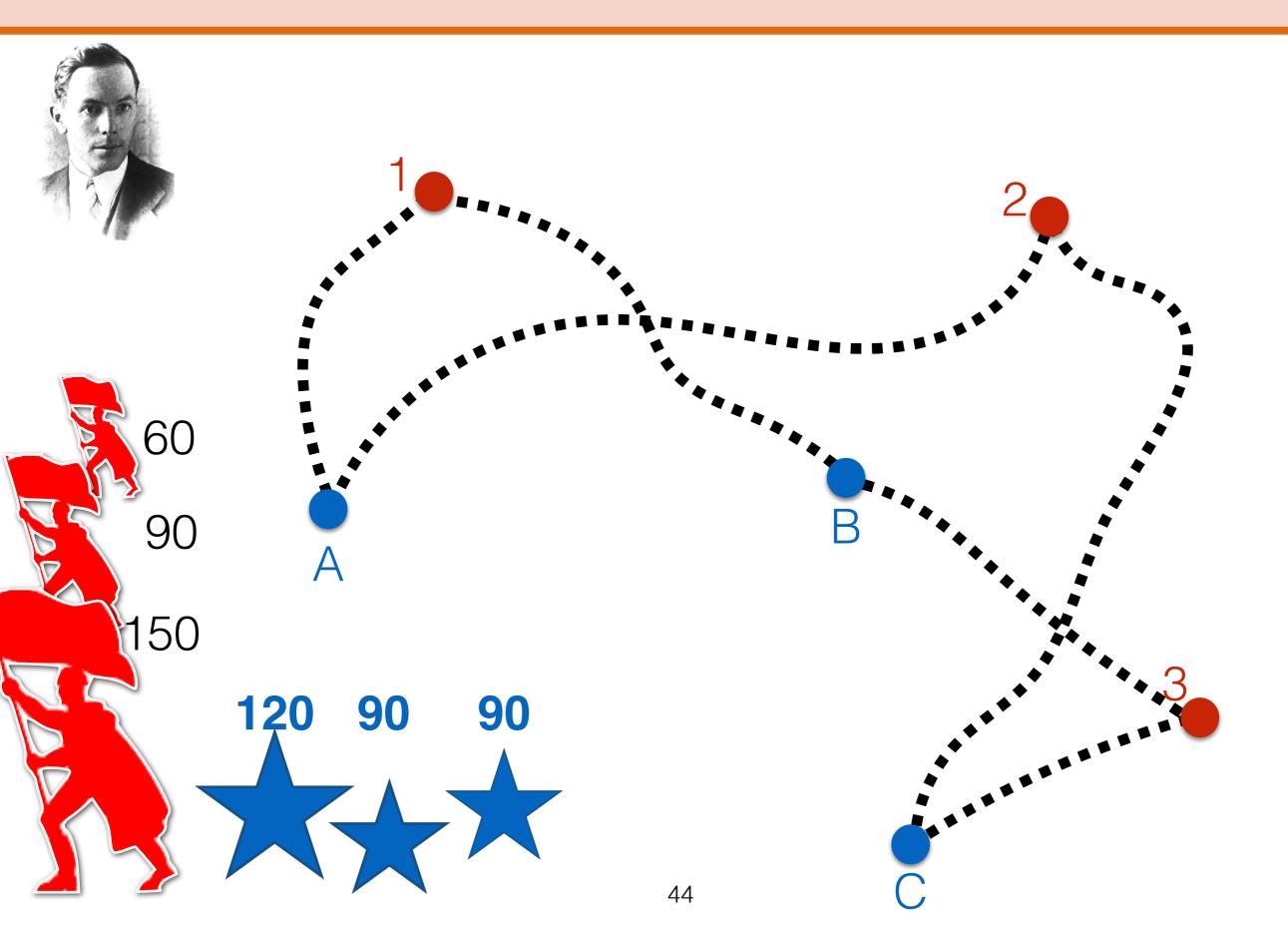








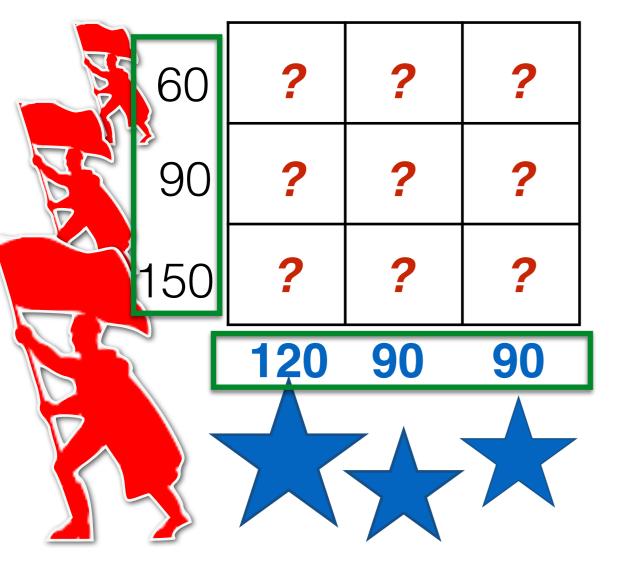




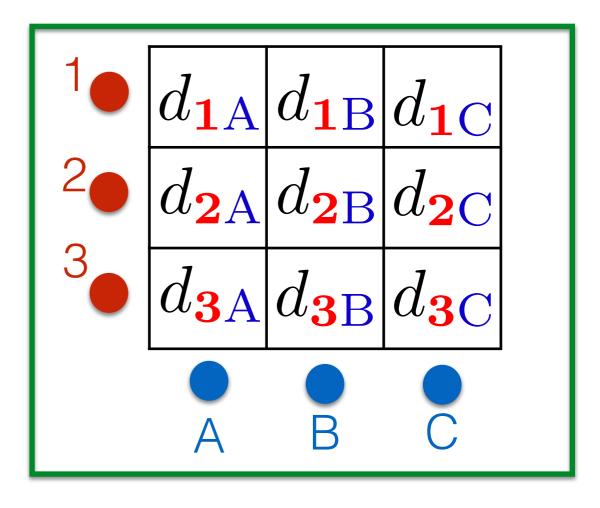


The problem is entirely described by counts and a cost/distance matrix

Transportation matrix



Distance matrix



Transportation matrix

Constraints

$$orall i \in \{1,2,3\}, \sum_{oldsymbol{j} \in \{ ext{A,B,C}\}} p_{ioldsymbol{j}} = a_i$$

$$\forall j \in \{A, B, C\}, \sum_{i \in \{1,2,3\}} p_{ij} = b_j$$

$$p_{ij} \ge 0$$

Distance matrix

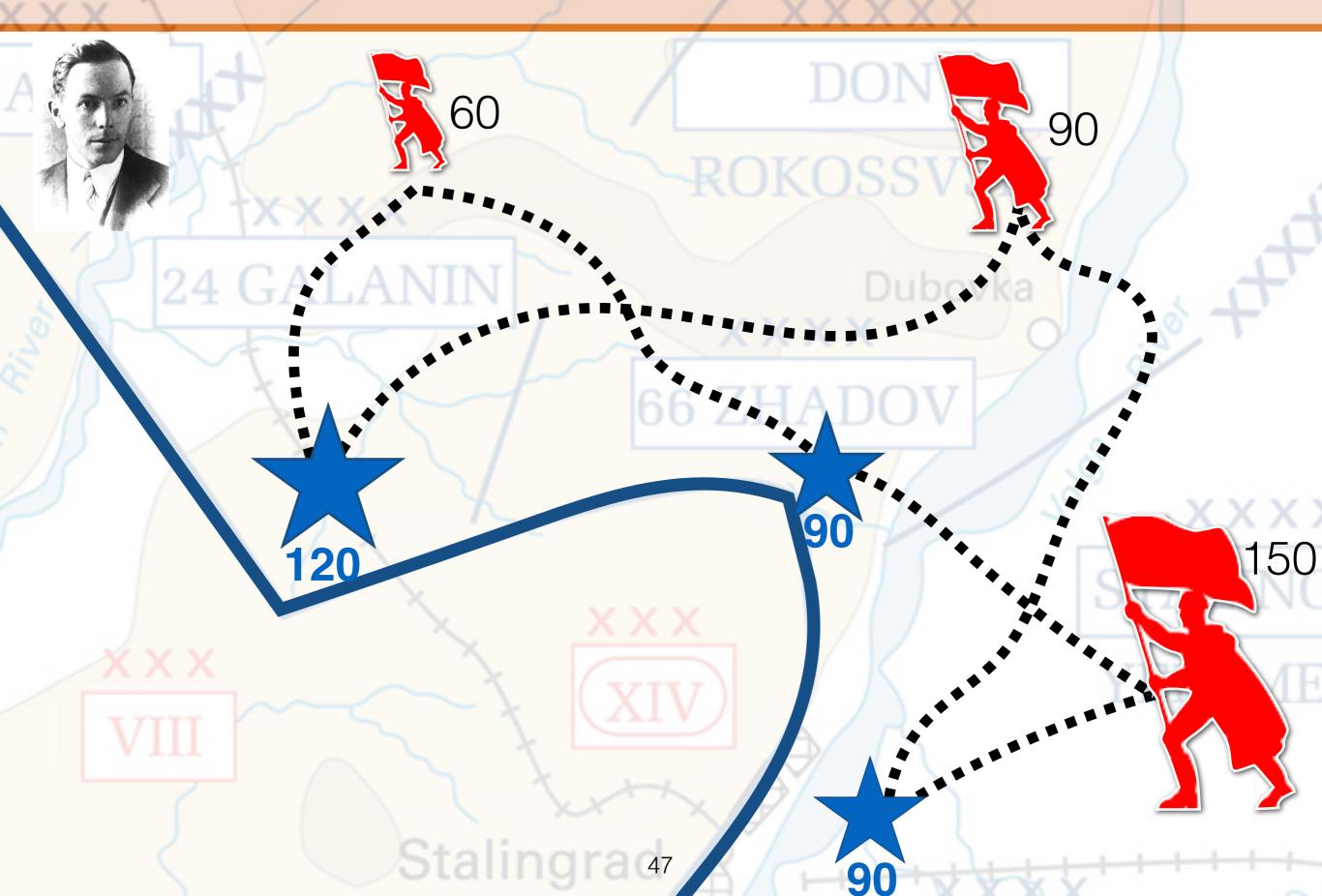
1
$$d_{1A} d_{1B} d_{1C}$$
2 $d_{2A} d_{2B} d_{2C}$
3 $d_{3A} d_{3B} d_{3C}$
A B C

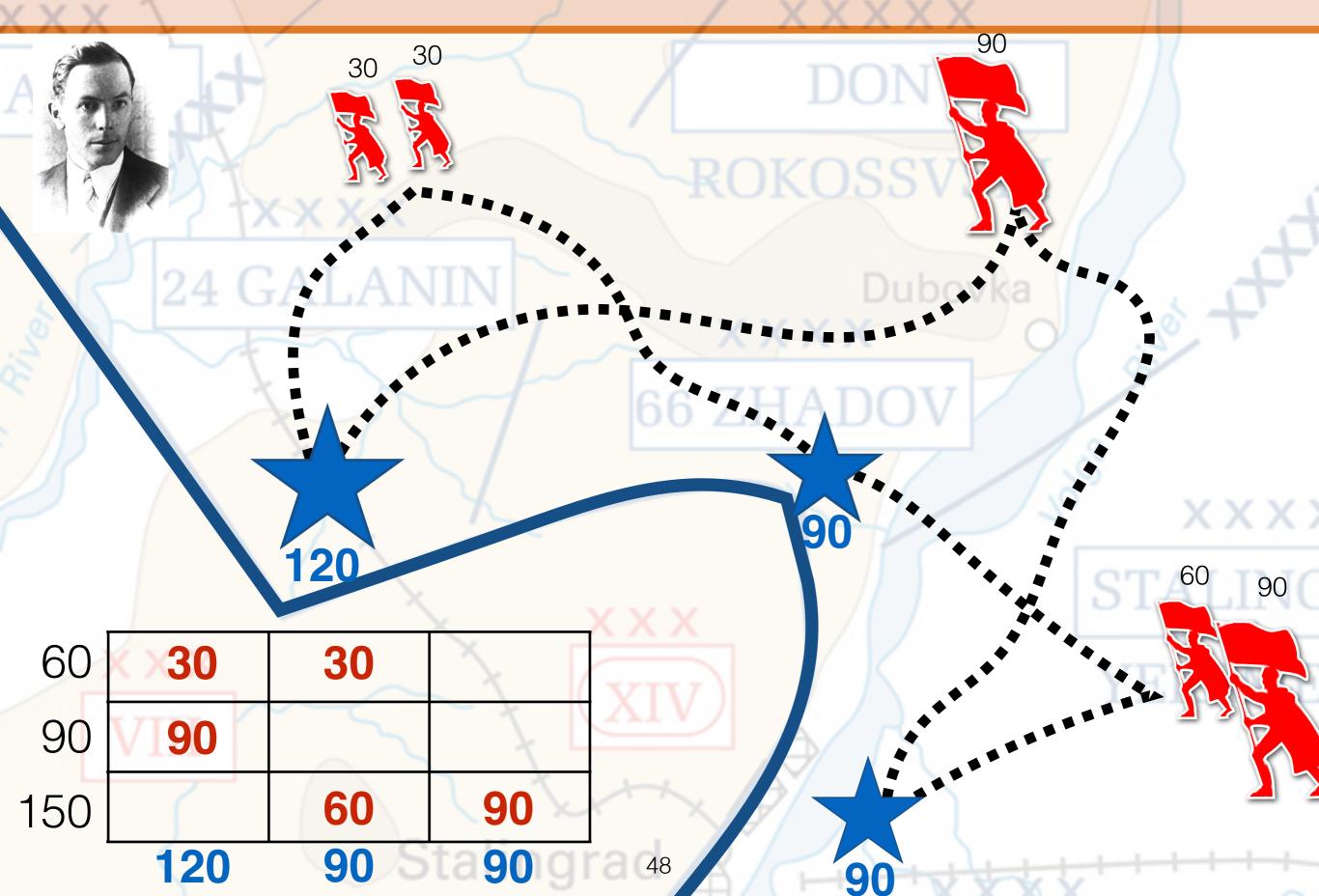
Cost function

$$C(\boldsymbol{P}) = \sum_{\boldsymbol{j} \in \{A,B,C\}} \sum_{\boldsymbol{i} \in \{1,2,3\}} \boldsymbol{p_{ij}} d_{\boldsymbol{ij}}$$

Problem

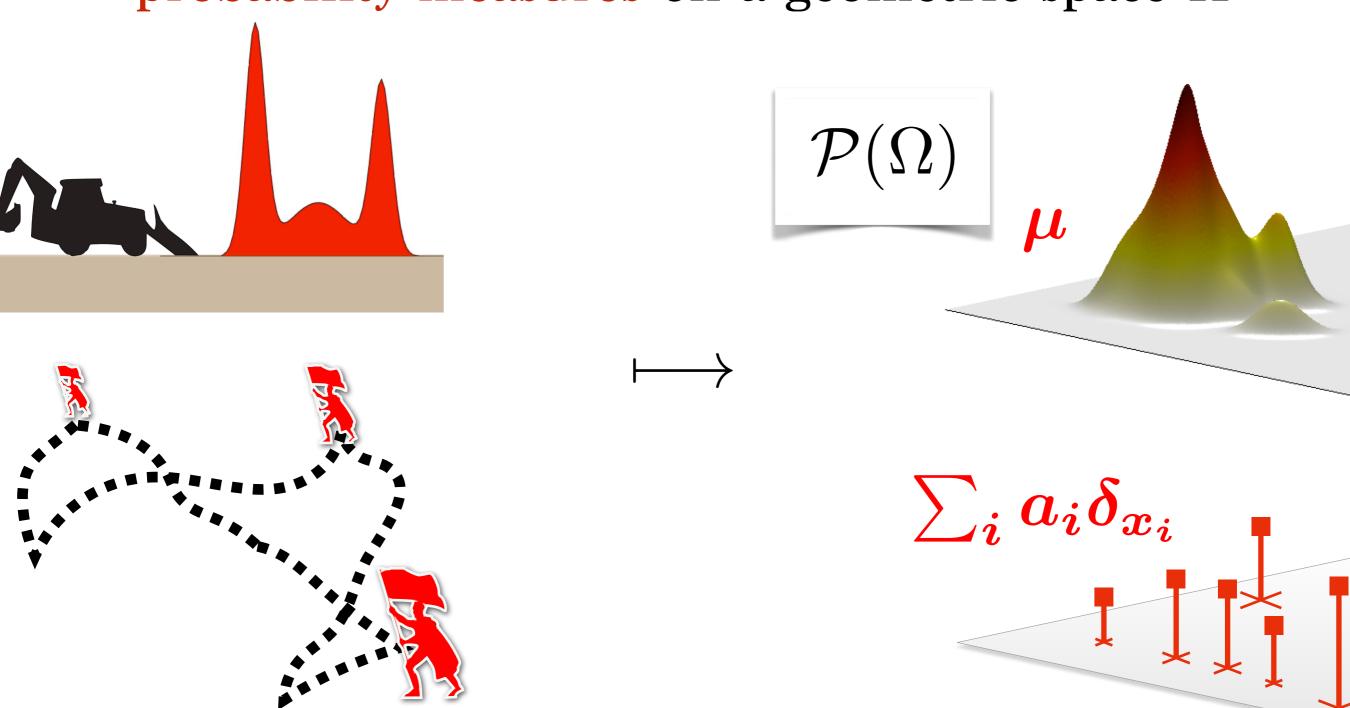
 $\min_{\text{all valid } \boldsymbol{P}} C(\boldsymbol{P})$



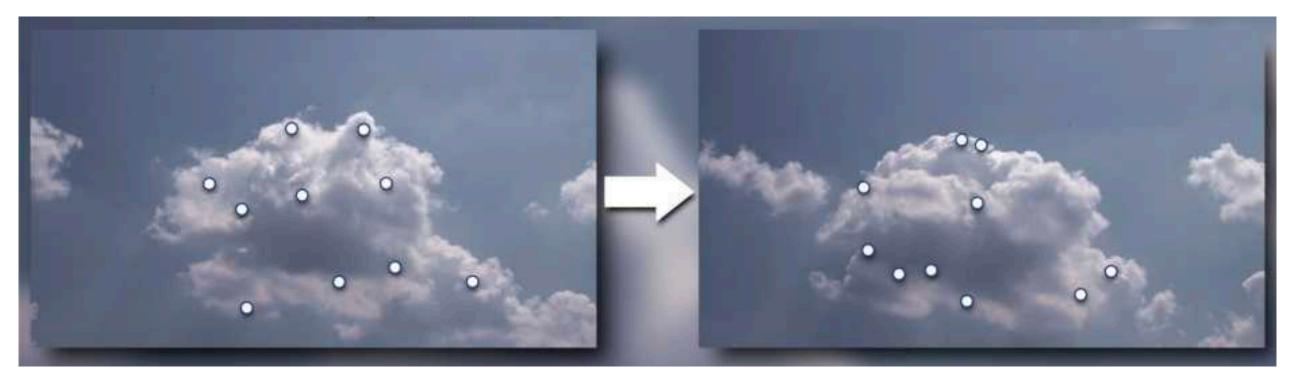


Mathematical Formalism

These problems involve discrete and continuous probability measures on a geometric space Ω



OT: Nature's way to move particles



screenshot video: ETH Zurich / Michael Steiner, vistory GmbH



Figalli

Fields'18

Monge Problem

 Ω a measurable space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T: \Omega \to \Omega$ [Brenier'87] If $\Omega = \mathbb{R}^d$, $T = \mathbb{R}^d$,

[Brenier'87]: For any u convex, ∇u is the OT Monge map between μ and $\nabla u_{\sharp}\mu$.

Links: Monge-Ampère Equation

If $\Omega = \mathbb{R}^d$, $\boldsymbol{c} = \|\cdot - \cdot\|^2$, $\boldsymbol{\mu}, \boldsymbol{\nu}$ have densities $\boldsymbol{p}, \boldsymbol{q}$, then $\boldsymbol{T}_{\sharp}\boldsymbol{\mu} = \boldsymbol{\nu}$ is equivalent to

$$p(x) = q(T(x)) |\det J_T(x)|$$

Monge-Ampère: find convex f such that

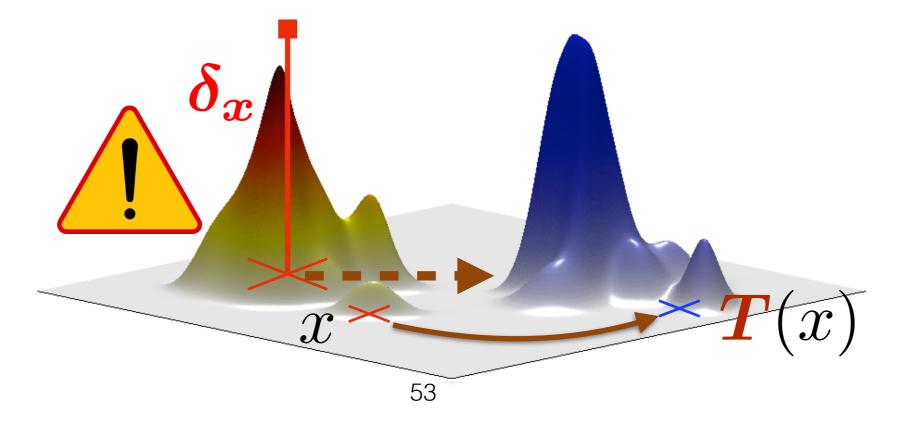
$$|\nabla^2 \boldsymbol{f}(\boldsymbol{x})| = \frac{\boldsymbol{p}(\boldsymbol{x})}{\boldsymbol{q}(\nabla \boldsymbol{f}(\boldsymbol{x}))}$$

Monge Problem

 Ω a measurable space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T: \Omega \to \Omega$

$$\inf_{\boldsymbol{T} \neq \boldsymbol{\mu} = \boldsymbol{\nu}} \int_{\Omega} \boldsymbol{c}(x, \boldsymbol{T}(x)) \boldsymbol{\mu}(dx)$$



Kantorovich Relaxation

Instead of maps $T: \Omega \to \Omega$, consider probabilistic maps, i.e. couplings $P \in \mathcal{P}(\Omega \times \Omega)$: P(Y|X=x)

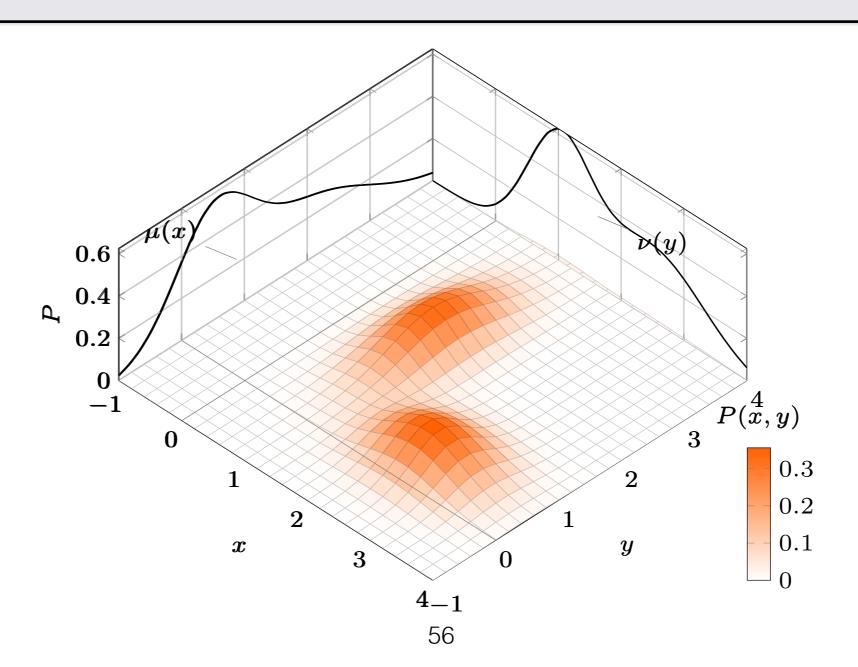
Kantorovich Relaxation

Instead of maps $T:\Omega\to\Omega$, consider probabilistic maps, i.e. couplings $P\in\mathcal{P}(\Omega\times\Omega)$:

$$\Pi(\mu, \nu) \stackrel{\mathrm{def}}{=} \{ P \in \mathcal{P}(\Omega \times \Omega) | \forall A, B \subset \Omega, \ P(A \times \Omega) = \mu(A), \ P(\Omega \times B) = \nu(B) \}$$

Kantorovich Relaxation

$$\Pi(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathcal{P}(\Omega \times \Omega) | \forall \boldsymbol{A}, \boldsymbol{B} \subset \Omega, \\ \boldsymbol{P}(\boldsymbol{A} \times \Omega) = \boldsymbol{\mu}(\boldsymbol{A}), \boldsymbol{P}(\Omega \times \boldsymbol{B}) = \boldsymbol{\nu}(\boldsymbol{B}) \}$$



$$\inf_{\boldsymbol{T} \not = \boldsymbol{\nu}} \int_{\Omega} \boldsymbol{c}(x, \boldsymbol{T}(x)) \boldsymbol{\mu}(dx) \quad \text{MONGE}$$

Def. Given μ, ν in $\mathcal{P}(\Omega)$; a cost function \boldsymbol{c} on $\Omega \times \Omega$, the Kantorovich problem is

$$\inf_{\boldsymbol{P}\in\Pi(\boldsymbol{\mu},\boldsymbol{\nu})}\int\int \boldsymbol{c}(x,y)\boldsymbol{P}(dx,dy).$$

PRIMAL

Def. Given μ, ν in $\mathcal{P}(\Omega)$; a cost function \boldsymbol{c} on $\Omega \times \Omega$, the Kantorovich problem is

$$\inf_{\mathbf{P}\in\Pi(\boldsymbol{\mu},\boldsymbol{\nu})} \iint \boldsymbol{c}(x,y) \mathbf{P}(dx,dy).$$

PRIMAL

$$\sup_{\substack{\boldsymbol{\varphi} \in L_1(\boldsymbol{\mu}), \boldsymbol{\psi} \in L_1(\boldsymbol{\nu}) \\ \boldsymbol{\varphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{c}}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu}.$$

Let
$$\varphi, \psi : \Omega \to \mathbb{R}$$
, and $P \in \Pi(\mu, \nu)$.

$$\int \varphi d\mu + \int \psi d\nu - \int \int \varphi \oplus \psi dP = \emptyset$$
Let $\varphi, \psi : \Omega \to \mathbb{R}$ and $P \in \mathcal{P}_{(x,y)} + \mathcal{P}_{(y)dP(x,y)}$

$$\int \varphi d\mu + \int \psi d\nu - \int \varphi \oplus \psi dP = \emptyset$$

$$\int \varphi d(\mu - P_X) + \int \psi d(\nu - P_Y)$$

$$\neq 0 \quad \text{and } / \text{ or } \neq 0$$

$$\iota_{\Pi}(\mathbf{P}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \left[\int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu} - \iint \boldsymbol{\varphi} \oplus \boldsymbol{\psi} d\mathbf{P} \right]$$

$$= \begin{cases} 0 & \text{if } P \in \Pi(\mu, \nu), \\ +\infty & \text{otherwise.} \end{cases}$$

$$\inf_{\boldsymbol{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \iint \boldsymbol{c} \, d\boldsymbol{P}$$

$$\inf_{\boldsymbol{P}\in\mathcal{P}_{+}(\Omega^{2})}\int\int \boldsymbol{c}\,d\boldsymbol{P}+\boldsymbol{\iota}_{\Pi}(\boldsymbol{P})$$

$$\inf_{\boldsymbol{P}\in\mathcal{P}_{+}(\Omega^{2})} \iint \boldsymbol{c} \, d\boldsymbol{P} + \iota_{\Pi}(\boldsymbol{P})$$

$$\inf_{\boldsymbol{P}\in\mathcal{P}_{+}(\Omega^{2})} \int \int (\boldsymbol{c}-\boldsymbol{\varphi}\oplus\boldsymbol{\psi})d\boldsymbol{P} \int \boldsymbol{\psi}d\boldsymbol{\psi} - \int \int \boldsymbol{\varphi}\oplus\boldsymbol{\psi}d\boldsymbol{P}$$

$$\inf_{\mathbf{P}\in\mathcal{P}_{+}(\Omega)} \iint (\mathbf{c} - \boldsymbol{\varphi} \oplus \boldsymbol{\psi}) d\mathbf{P} = \begin{cases} 0 & \text{if } \mathbf{c} - \boldsymbol{\varphi} \oplus \boldsymbol{\psi} \geq 0. \\ -\infty & \text{otherwise} \end{cases}$$

$$\sup_{\boldsymbol{\varphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{c}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu}.$$

$$\inf_{\boldsymbol{P}\in\mathcal{P}_{+}(\Omega^{2})} \iint \boldsymbol{c} \, d\boldsymbol{P} + \iota_{\Pi}(\boldsymbol{P})$$

$$\sup_{\varphi \oplus \psi \leq c} \int \varphi d\mu + \int \psi d\nu.$$
DUAL

Prop: Primal-dual relationship:

$$\mathbf{P}^{\star}(x,y) > 0 \Leftrightarrow \boldsymbol{\varphi}^{\star}(\mathbf{x}) + \boldsymbol{\psi}^{\star}(\mathbf{y}) = \mathbf{c}(x,y).$$

Wasserstein Distances

Let
$$p \ge 1$$
. Let $c(x,y) := D^p(x,y)$, a metric.

Def. The p-Wasserstein distance between μ, ν in $\mathcal{P}(\Omega)$ is

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \left(\inf_{\boldsymbol{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \iint \boldsymbol{D}(\boldsymbol{x}, \boldsymbol{y})^p \boldsymbol{P}(dx, dy) \right)^{\frac{1}{2/P}}.$$

Kantorovich Duality

$$W_{\boldsymbol{c}}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{c}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu}.$$

- Kantorovich duality is computationally easier: easier to store 2 functions than an entire coupling.
- *c*-transforms: useful machinery to consider only **two** instead of **one** dual potential, For instance, when p = 1

$$W_1(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi} \text{ 1-Lipschitz}} \int_{\boldsymbol{\varphi}} \boldsymbol{\varphi}(d\boldsymbol{\mu} - d\boldsymbol{\nu}).$$

W1

D transforms

$$W_{\boldsymbol{c}}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{c}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu}.$$

Imagine we choose a φ . Can we find a good ψ ? We need that ψ satisfies for all x, y

$$\varphi(x) + \psi(y) \le c(x, y)$$

$$\psi(y) \leq c(x, y) - \varphi(x)$$

$$\psi(y) \le \inf_{\boldsymbol{x}} c(\boldsymbol{x}, y) - \varphi(\boldsymbol{x})$$

D transforms

$$W_{\boldsymbol{c}}(\boldsymbol{\mu}, \boldsymbol{
u}) = \sup_{\boldsymbol{arphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{c}} \int \boldsymbol{arphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu}.$$

For given φ , cannot get a better ψ than

$$\overline{\boldsymbol{\varphi}}(\boldsymbol{y}) \stackrel{\mathrm{def}}{=} \inf_{\boldsymbol{x}} \boldsymbol{c}(\boldsymbol{x}, \boldsymbol{y}) - \boldsymbol{\varphi}(\boldsymbol{x}).$$

$$W_p^p(\mu, \nu) = \sup_{oldsymbol{arphi}} \int_{oldsymbol{arphi}} oldsymbol{arphi} d\mu + \int_{oldsymbol{\square}} \overline{oldsymbol{arphi}} d
u.$$

D transforms

$$\overline{\boldsymbol{\varphi}}(\boldsymbol{y}) \stackrel{\text{def}}{=} \inf_{\boldsymbol{x}} \boldsymbol{c}(\boldsymbol{x}, \boldsymbol{y}) - \boldsymbol{\varphi}(\boldsymbol{x}).$$

$$\overline{\boldsymbol{\psi}}(\boldsymbol{x}) = \inf_{\boldsymbol{y}} \boldsymbol{c}(\boldsymbol{x}, \boldsymbol{y}) - \boldsymbol{\psi}(\boldsymbol{y}).$$

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}} \int \overline{\overline{\boldsymbol{\varphi}}} d\boldsymbol{\mu} + \int \overline{\boldsymbol{\varphi}} d\boldsymbol{\nu}.$$

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi} \text{ is } \boldsymbol{c}\text{-concave}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \overline{\boldsymbol{\varphi}} d\boldsymbol{\nu}.$$

D transforms, W_1

Prop. If c = D, namely p = 1, then φ is D-concave $\Leftrightarrow \overline{\varphi} = -\varphi$, φ is 1-Lipschitz

For given
$$\mathbf{x}$$
, $\overline{\varphi}_{\mathbf{x}}(\mathbf{y}) \stackrel{\text{def}}{=} \mathbf{D}(\mathbf{x}, \mathbf{y}) - \varphi(\mathbf{x})$ is 1-Lipschitz.
 $\overline{\varphi}_{\mathbf{x}}(\mathbf{y}) \stackrel{\text{def}}{=} \mathbf{D}(\mathbf{y}, \mathbf{y})$ is 1-Lipschitz.
 $\Rightarrow \overline{\varphi}(\mathbf{y}) - \overline{\varphi}(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{D}(\mathbf{x}, \mathbf{y})$
 $\Rightarrow -\overline{\varphi}(\mathbf{x}) \leq \mathbf{D}(\mathbf{x}, \mathbf{y}) - \overline{\varphi}(\mathbf{y})$
 $\Rightarrow -\overline{\varphi}(\mathbf{x}) \leq \inf_{\mathbf{y}} \mathbf{D}(\mathbf{x}, \mathbf{y}) - \overline{\varphi}(\mathbf{y})$
 $\Rightarrow -\overline{\varphi}(\mathbf{x}) \leq \inf_{\mathbf{y}} \mathbf{D}(\mathbf{x}, \mathbf{y}) - \overline{\varphi}(\mathbf{y}) \leq -\overline{\varphi}(\mathbf{x})$
 $\Rightarrow -\overline{\varphi}(\mathbf{x}) \leq \overline{\varphi}(\mathbf{x}) \leq -\overline{\varphi}(\mathbf{x})$ and $\overline{\varphi}(\mathbf{x}) = -\varphi(\mathbf{x})$

D transforms, W_1

$$W_1(\mu, \nu) = \sup_{\varphi \text{ is } D\text{-concave}} \int \varphi d\mu + \int_{\text{SEMI-DUAL}} \overline{\varphi} d\nu.$$

Prop. If
$$c=D$$
, then φ is D -concave $\Leftrightarrow \overline{\varphi} = -\varphi$, φ is 1-Lipschitz

$$W_1(\mu, \nu) = \sup_{\varphi \text{ 1-Lipschitz}} \int_{\varphi} \varphi(d\mu - d\nu).$$

Links between Monge & Kantorovich

Prop. For "well behaved" costs c, if μ has a density then an *optimal* Monge map T^* between μ and ν must exist.

Prop. In that case

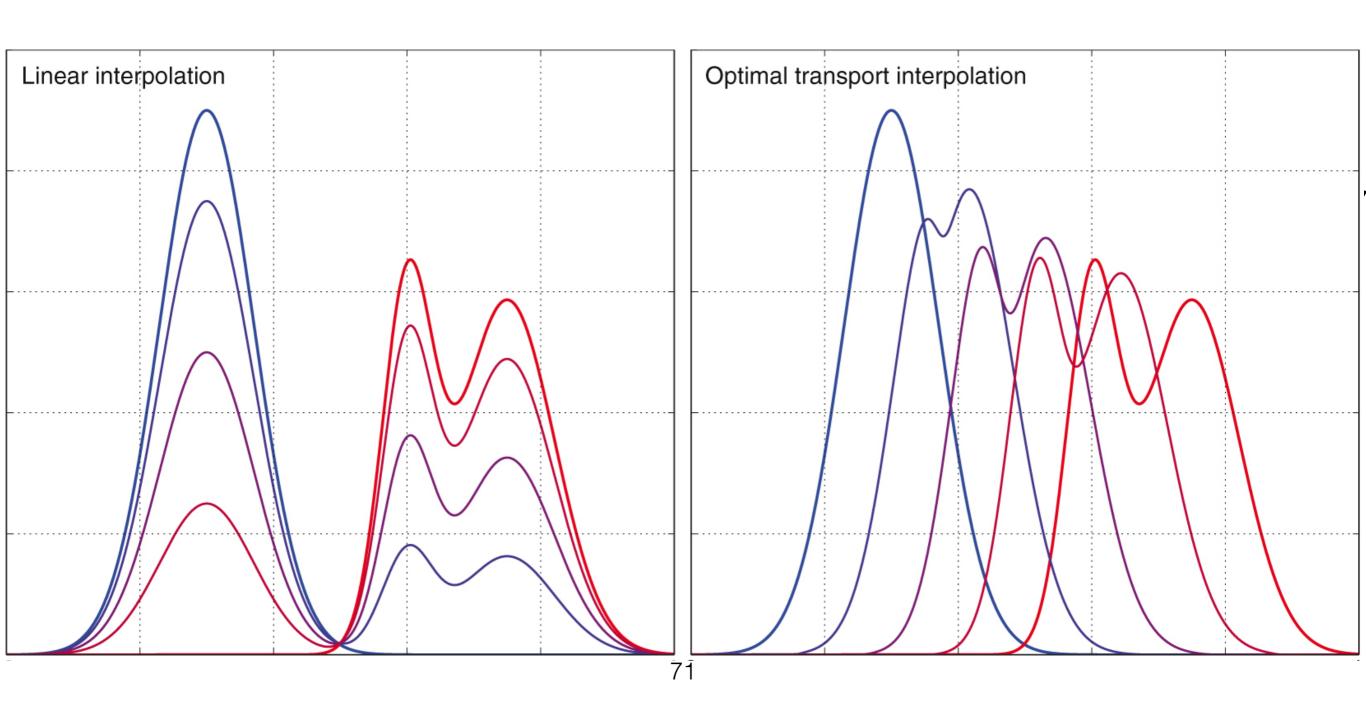
$$P^* := (\mathrm{Id}, T^*)_{\sharp} \mu \in \Pi(\mu, \nu)$$

is also optimal for the Kantorovich problem.

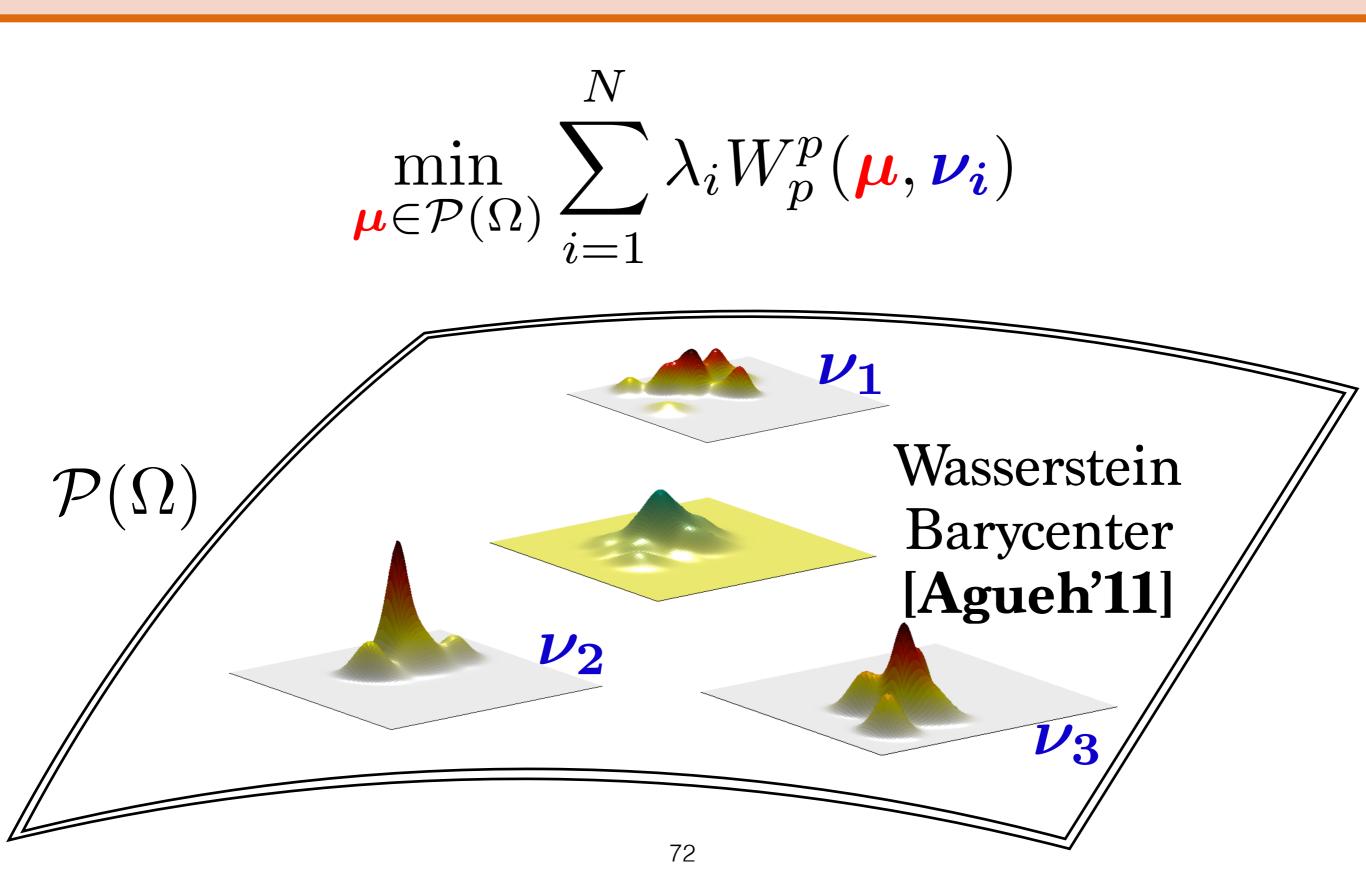
[Brenier'91] [Smith&Knott'87] [McCann'01]

Optimal Transport Geometry

Very different geometry than standard information divergences (*KL*, Euclidean)



Optimal Transport Geometry



Optimal Transport Geometry

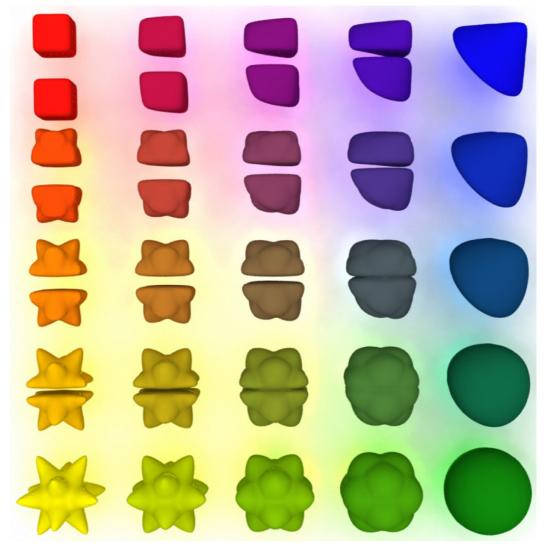
Very different geometry than standard information divergences (*KL*, Euclidean)



[SDPC..'15]

Optimal Transport Geometry

Very different geometry than standard information divergences (*KL*, Euclidean)



[SDPC..'15]

Variational OT Problems in ML

Up to 2010: OT solvers used mostly for retrieval in databases of histograms

$$W_p(\mu, \nu) = ?$$

$$W_p(\mu, \nu) \le \cdots ?$$

OT is now used as a loss or fidelity term:

$$\underset{\boldsymbol{\mu} \in \mathcal{P}(\Omega)}{\operatorname{argmin}} F(W_p(\boldsymbol{\mu}, \boldsymbol{\nu_1}), W_p(\boldsymbol{\mu}, \boldsymbol{\nu_2}), \dots, \boldsymbol{\mu}) = ?$$

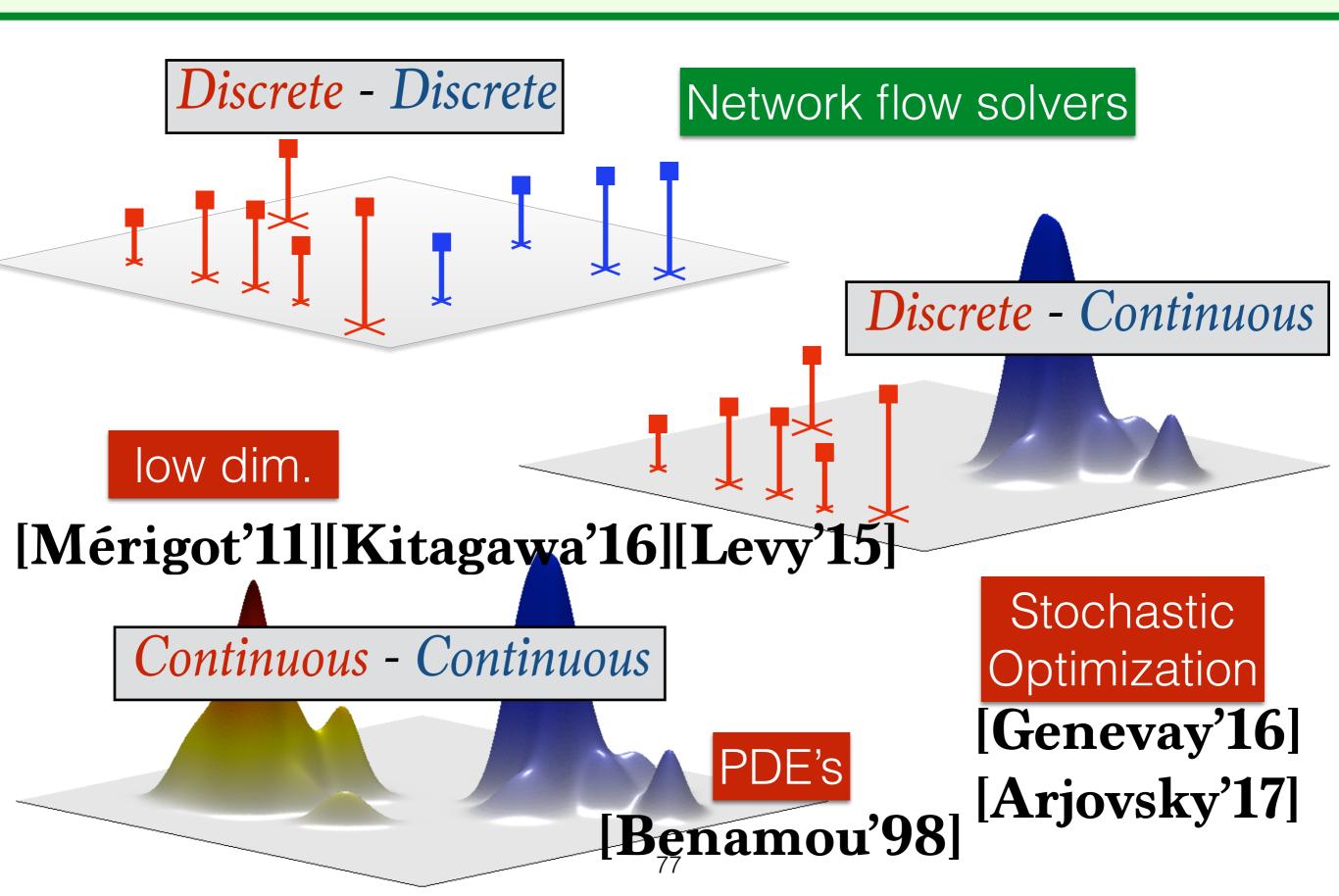
"
$$\nabla_{\boldsymbol{\mu}}$$
" $W_p(\boldsymbol{\mu}, \boldsymbol{\nu}) = ?$

[Jordan Kinderlehrer Ofto Bobsio Gigli Savaré'05]

2. Computing OT exactly

- Typology: discrete/continuous problems
- Easy cases and exact solvers for discrete measures.

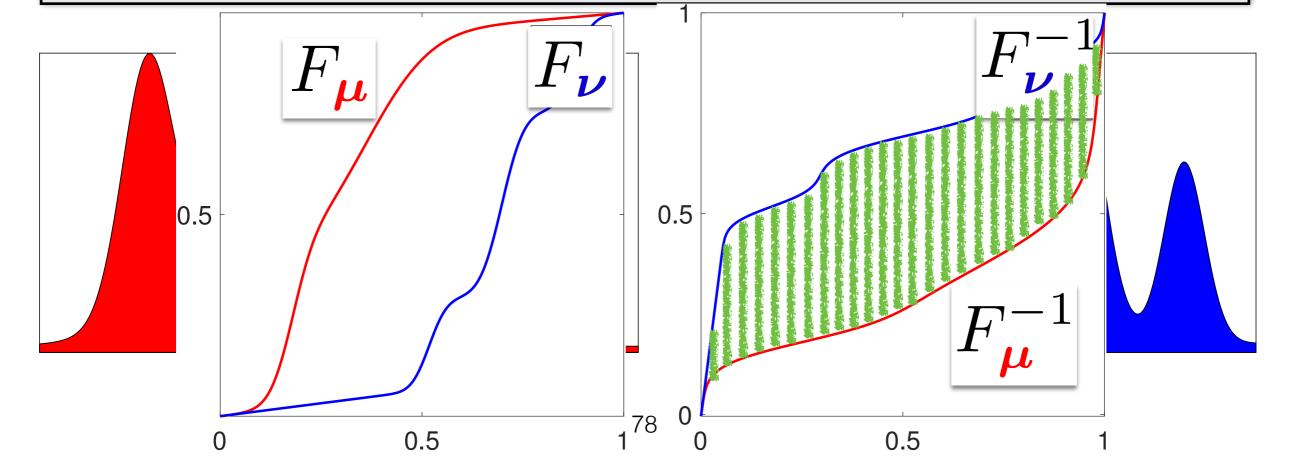
When can we compute OT?



Easy (1): Univariate Measures

Remark. If $\Omega = \mathbb{R}$, c(x,y) = c(|x-y|), c convex, F_{μ}^{-1} , F_{ν}^{-1} quantile functions,

$$W(\boldsymbol{\mu}, \boldsymbol{\nu}) = \int_0^1 c(|F_{\boldsymbol{\mu}}^{-1}(x) - F_{\boldsymbol{\nu}}^{-1}(x)|) dx$$



Easy (2): Gaussian Measures

Remark. If
$$\Omega = \mathbb{R}^d$$
, $\boldsymbol{c}(x,y) = \|x - y\|^2$, and $\boldsymbol{\mu} = \mathcal{N}(\mathbf{m}_{\boldsymbol{\mu}}, \boldsymbol{\Sigma}_{\boldsymbol{\mu}}), \boldsymbol{\nu} = \mathcal{N}(\mathbf{m}_{\boldsymbol{\nu}}, \boldsymbol{\Sigma}_{\boldsymbol{\nu}})$ then
$$W_2^2(\boldsymbol{\mu}, \boldsymbol{\nu}) = \|\mathbf{m}_{\boldsymbol{\mu}} - \mathbf{m}_{\boldsymbol{\nu}}\|^2 + B(\boldsymbol{\Sigma}_{\boldsymbol{\mu}}, \boldsymbol{\Sigma}_{\boldsymbol{\nu}})^2$$

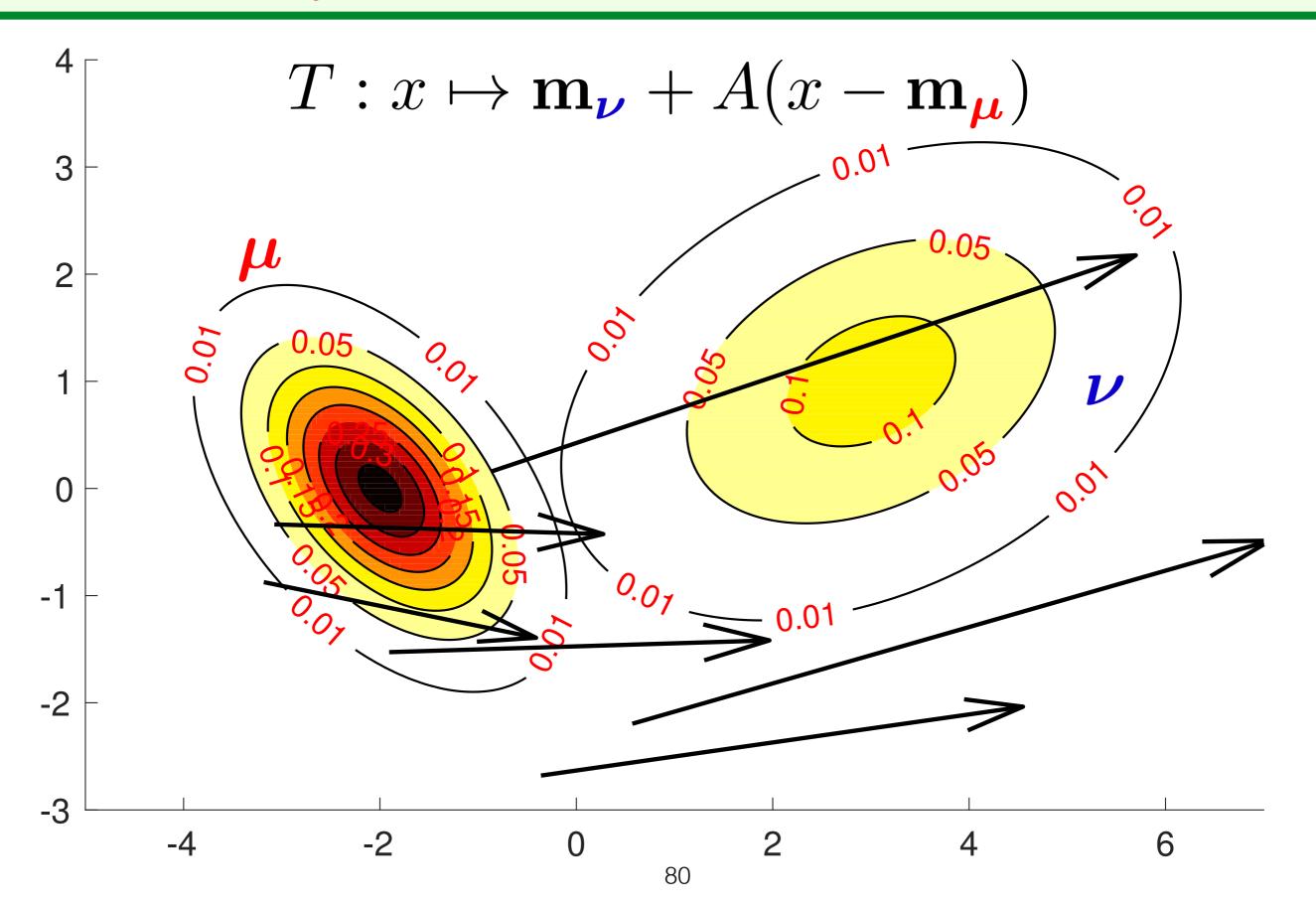
where
$$B$$
 is the Bures $Y := CX + b \sim \mathcal{N}(Cm + b, C\Sigma C^T)$.

If $X \sim \mathcal{N}(m, \Sigma)$ then $Y := CX + b \sim \mathcal{N}(\Sigma^{1/2}\Sigma_{\nu}\Sigma_{\mu}^{1/2})^{1/2}$.

The map $T: x \mapsto \mathbf{m}_{\nu} + A(x - \mathbf{m}_{\mu})$ is optimal,

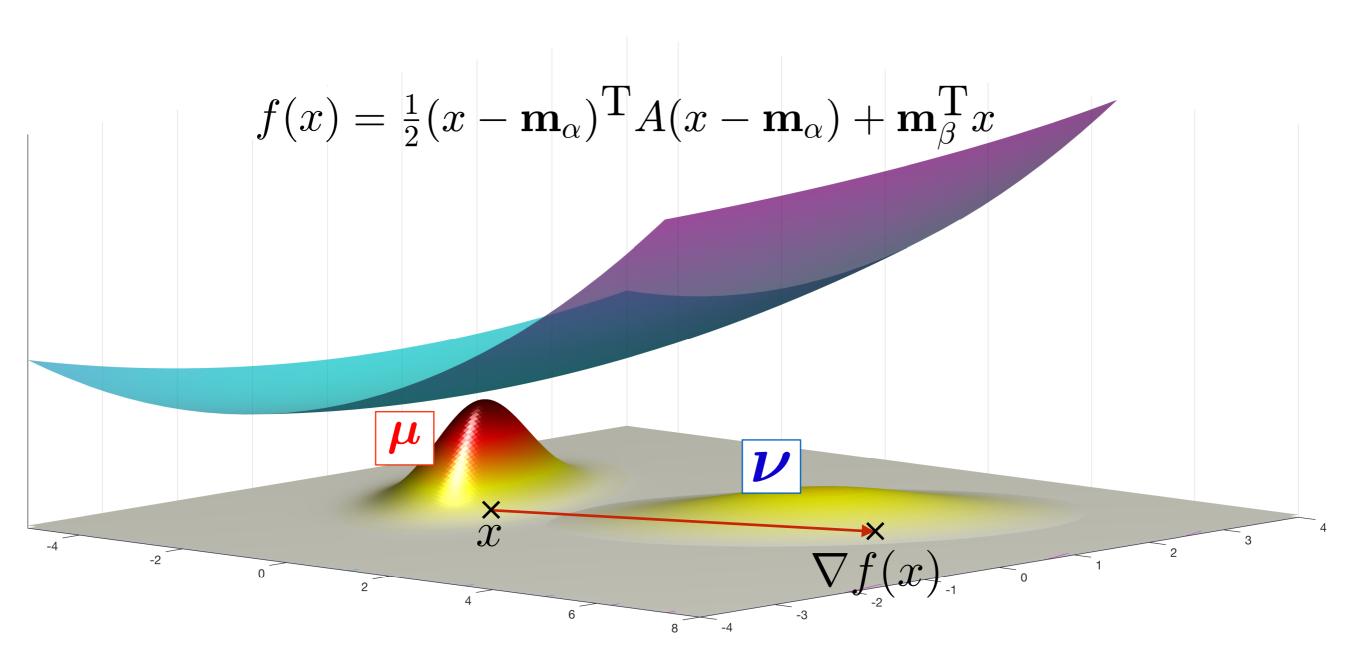
where
$$A = \Sigma_{\mu}^{-\frac{1}{2}} \left(\Sigma_{\mu}^{\frac{1}{2}} \Sigma_{\nu} \Sigma_{\mu}^{\frac{1}{2}} \right)^{\frac{1}{2}} \Sigma_{\mu}^{-\frac{1}{2}}.$$

Easy (2): Gaussian Measures

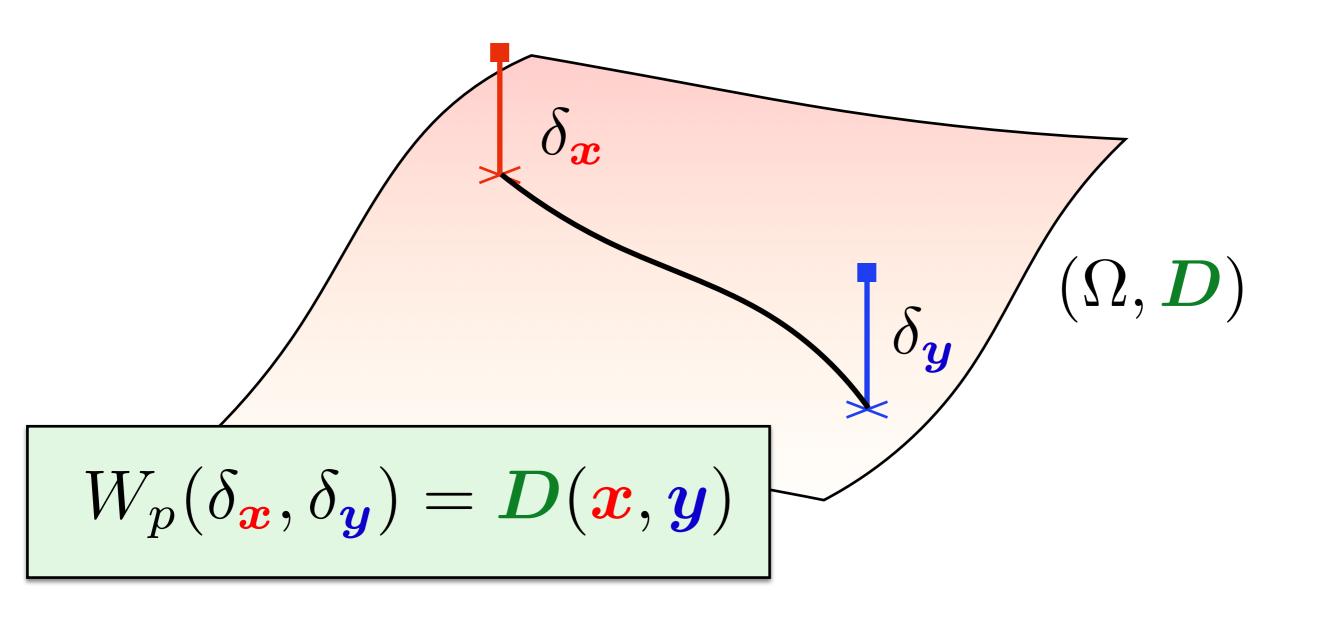


Easy (2): Gaussian Measures

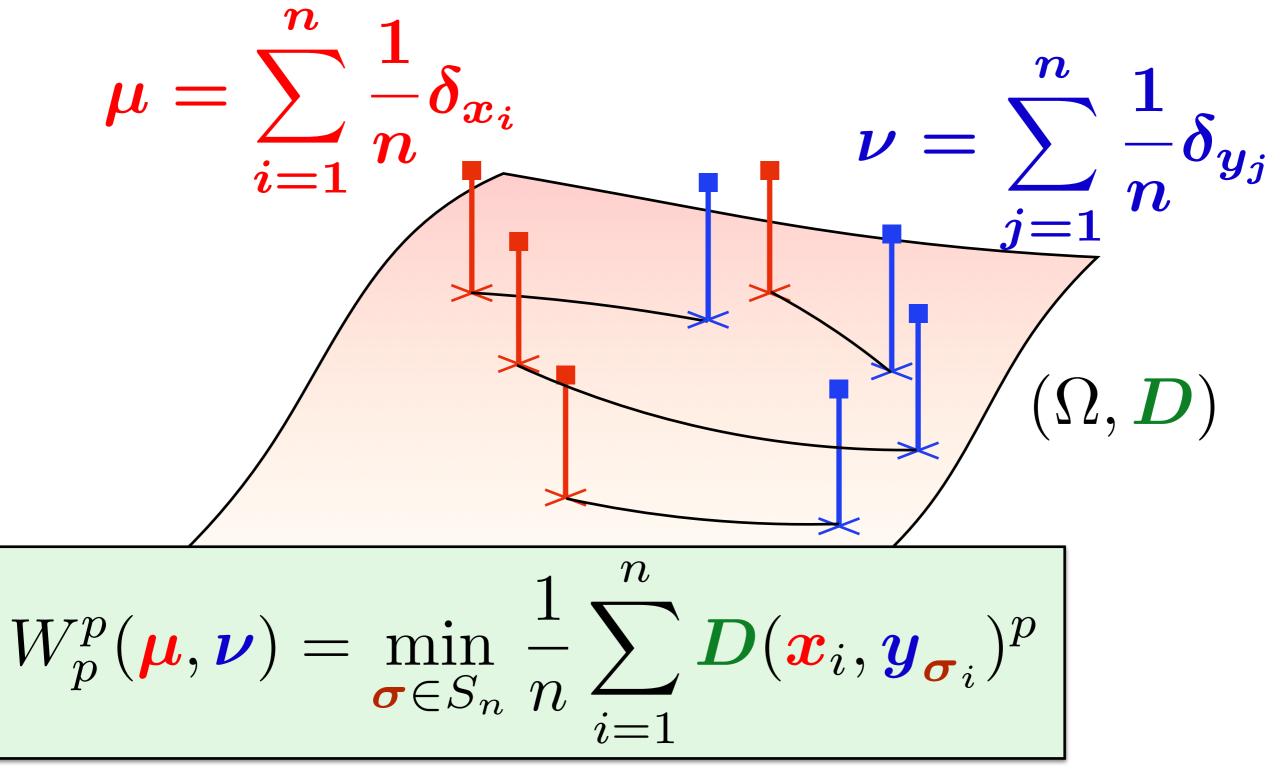
$$T = \nabla \psi : x \mapsto \mathbf{m}_{\nu} + A(x - \mathbf{m}_{\mu})$$



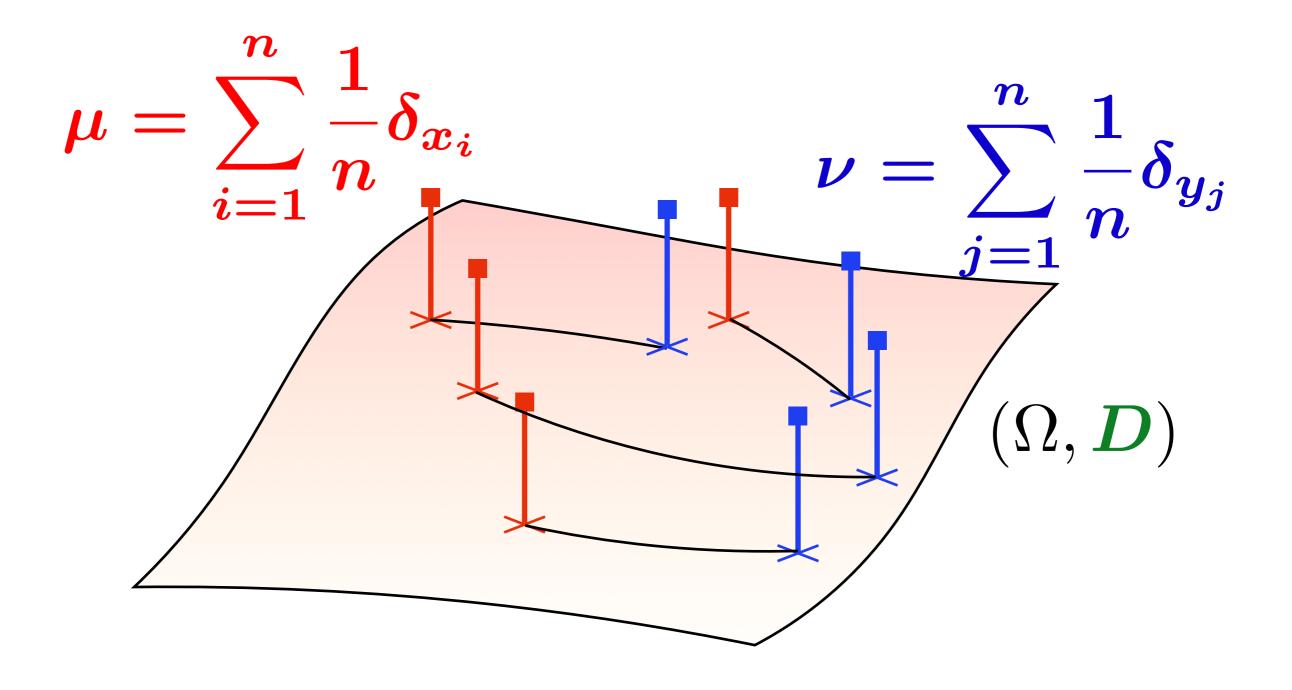
Wasserstein Between Two Diracs



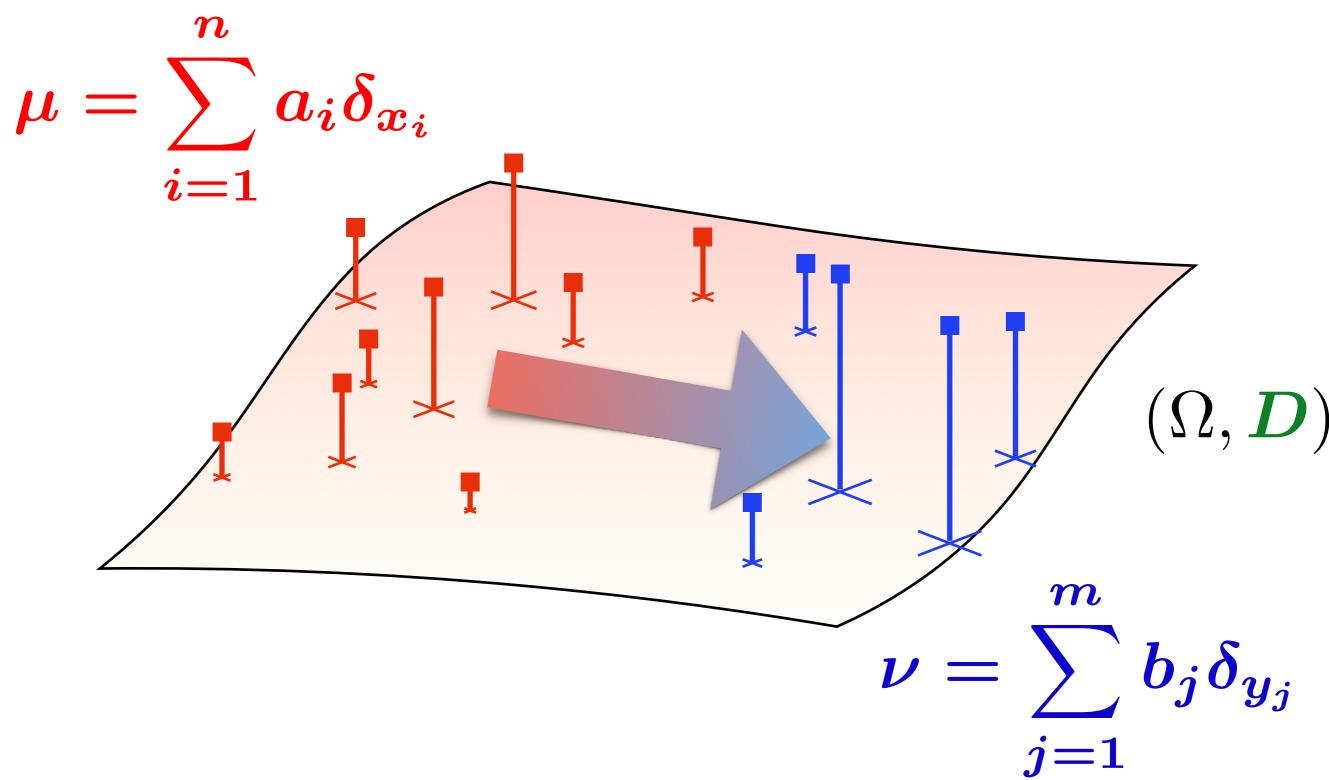
Linear Assignment Wasserstein



Linear Assignment Wasserstein



OT on Two Empirical Measures



Wasserstein on Empirical Measures

Consider
$$\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 and $\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$.
$$M_{XY} \stackrel{\text{def}}{=} [D(x_i, y_j)^p]_{ij}$$

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

Def. Optimal Transport Problem

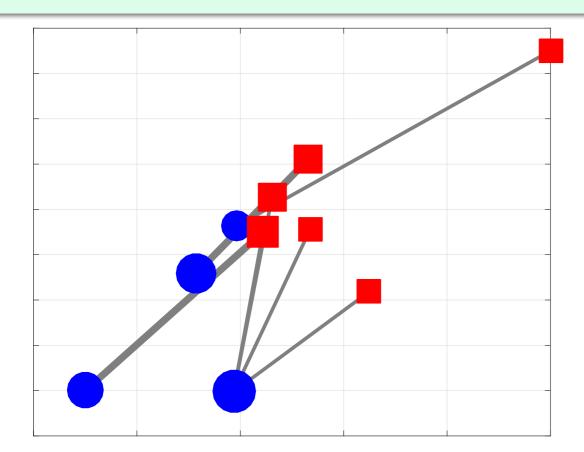
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

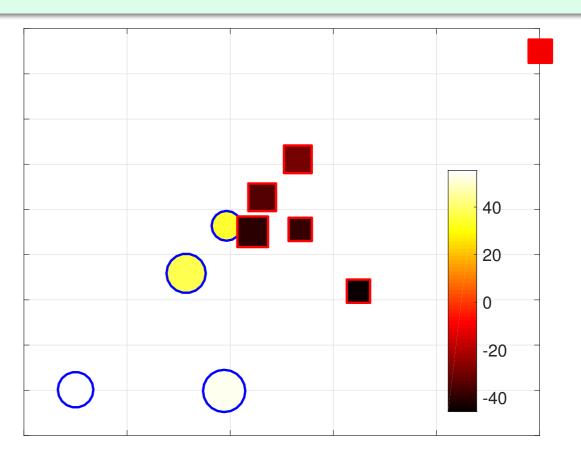
Dual Kantorovich Problem

$$W_p^p(oldsymbol{\mu},oldsymbol{
u}) = \min_{oldsymbol{P}\in\mathbb{R}_+^{n imes m}} \left\langle oldsymbol{P}, M_{oldsymbol{XY}}
ight
angle \ oldsymbol{P} oldsymbol{1}_m = oldsymbol{a}, oldsymbol{P}^T oldsymbol{1}_n = oldsymbol{b}$$

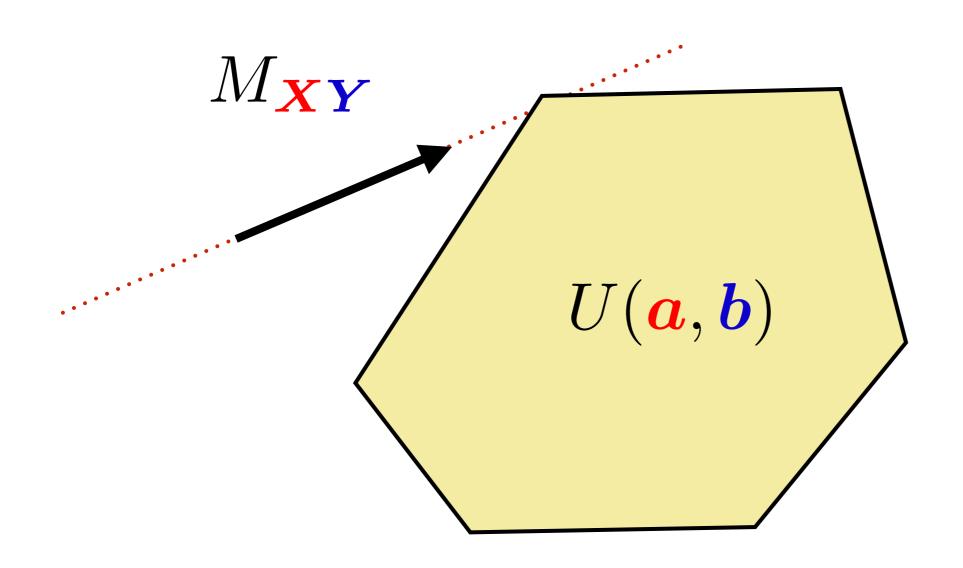
Def. Dual OT problem

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \boldsymbol{\alpha_i} + \boldsymbol{\beta_j} \le D(\boldsymbol{x_i}, \boldsymbol{y_j})^p}} \alpha^T \boldsymbol{a} + \beta^T \boldsymbol{b}$$

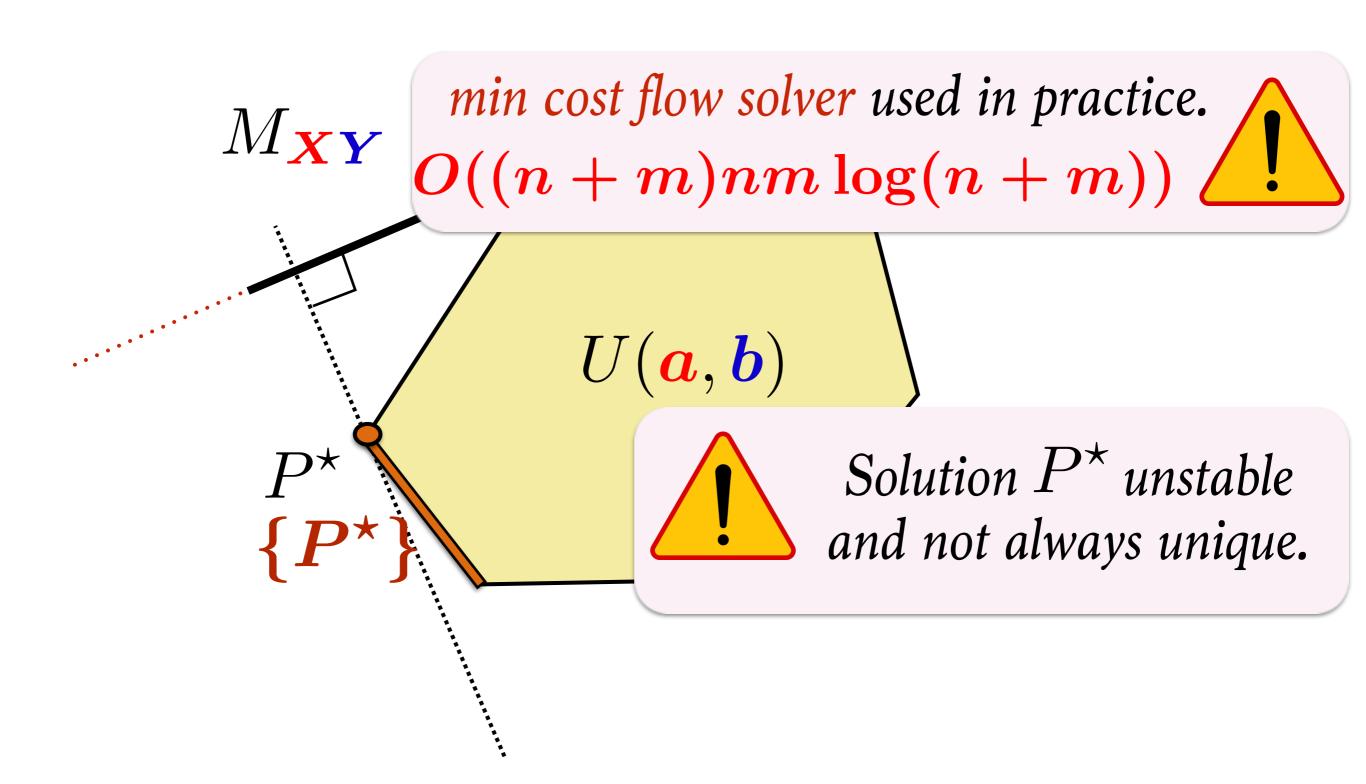




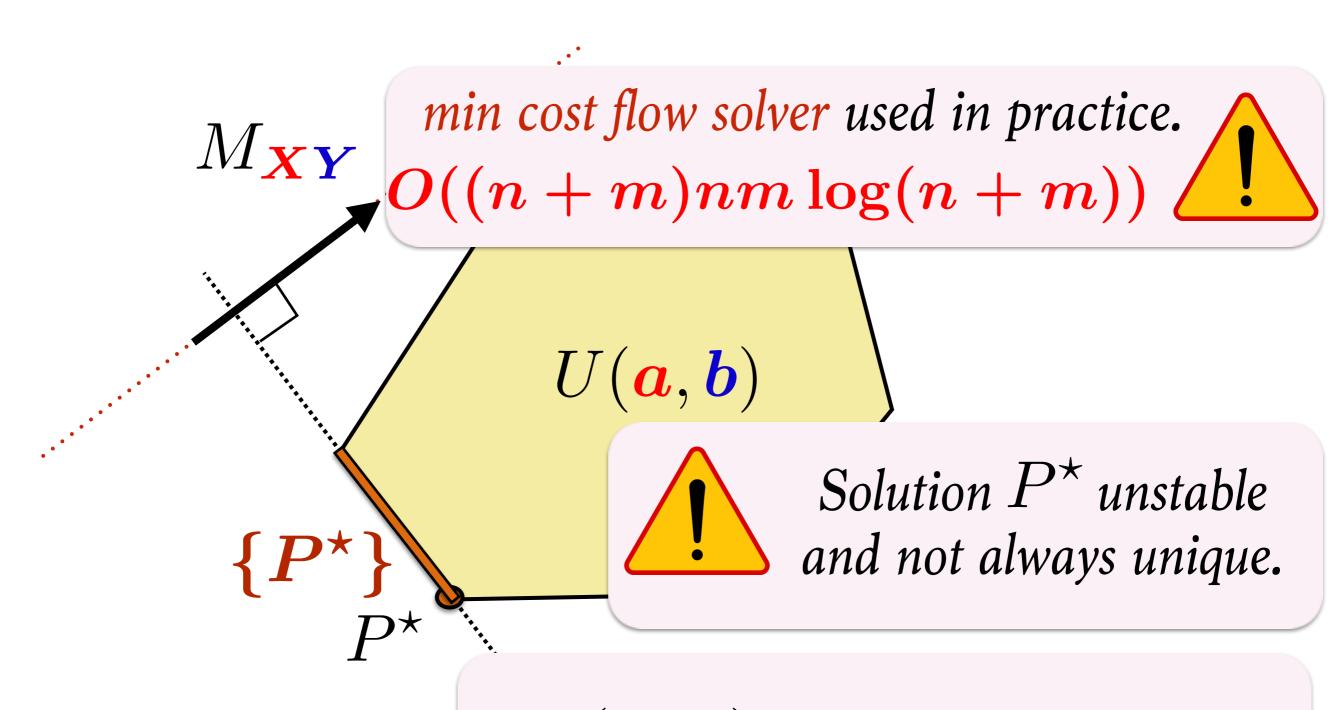
Solving the OT Problem



Solving the OT Problem



Solving the OT Problem



 $W_p^p(\mu, \nu)$ not differentiable.

Discrete OT Problem

```
c emd.c
   U- - C- #- 0 2
        emd.c
        Last update: 3/14/98
        An implementation of the Earth Movers Distance.
        Based of the solution for the Transportation problem as described in
        "Introduction to Mathematical Programming" by F. S. Hillier and
        G. J. Lieberman, McGraw-Hill, 1990.
10
11
        Copyright (C) 1998 Yossi Rubner
12
        Computer Science Department, Stanford University
13
        E-Mail: rubner@cs.stanford.edu URL: http://vision.stanford.edu/~rubner
14
15
    /*#include <stdio.h>
17
    #include <stdlib.h>*/
    #include <math.h>
18
    #include "emd.h"
    #define DEBUG_LEVEL 0
23
     DEBUG LEVEL:
25
       0 = NO MESSAGES
26
       1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
       2 = PRINT THE RESULT AFTER EVERY ITERATION
       3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
29
       4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
30
31
32
    #define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSIBLE DUMMY FEATURE */
34
35
    /* NEW TYPES DEFINITION */
36
37
     /* node1_t IS USED FOR SINGLE-LINKED LISTS */
    typedef struct node1_t {
      int i;
      double val;
     struct node1_t *Next;
    } node1_t;
43
     /* node1_t IS USED FOR DOUBLE-LINKED LISTS */
    typedef struct node2_t {
      int i, j;
47
      double val;
      struct node2_t *NextC;
                                         /* NEXT COLUMN */
48
49
      struct node2_t *NextR;
                                          /* NEXT ROW */
    } node2_t;
52
53
    /* GLOBAL VARIABLE DECLARATION */
54
55
    static int _n1, _n2;
                                                 /* SIGNATURES SIZES */
    static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1];/* THE COST MATRIX */
    static node2_t _X[MAX_SIG_SIZE1*2];
                                          /* THE BASIC VARIABLES VECTOR */
5.8
```

3. Computing OT for data sciences

- The need for regularization
- Entropic regularization
- Differentiation

What matters for practitioners?

i.i.d samples
$$x_1, \ldots, x_n \sim \mu, y_1, \ldots, y_m \sim \nu$$
,

$$\hat{\boldsymbol{\mu}}_{\boldsymbol{n}} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i} \delta_{\boldsymbol{x}_{\boldsymbol{i}}}, \hat{\boldsymbol{\nu}}_{\boldsymbol{m}} \stackrel{\text{def}}{=} \frac{1}{m} \sum_{j} \delta_{\boldsymbol{y}_{\boldsymbol{j}}}$$

Computational properties

Compute/approximate $W_p(\hat{\mu}_n, \hat{\nu}_m)$?

Statistical properties

$$\mathbb{E}\left[|\boldsymbol{W_p}(\boldsymbol{\mu},\boldsymbol{\nu}) - \boldsymbol{W_p}(\boldsymbol{\hat{\mu}_n},\boldsymbol{\hat{\nu}_m})|\right] \leq f(n,m)?$$

What matters for practitioners?

i.i.d samples
$$x_1, \ldots, x_n \sim \mu, y_1, \ldots, y_m \sim \nu$$
,

$$\hat{\boldsymbol{\mu}}_{\boldsymbol{n}} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i} \delta_{\boldsymbol{x}_{\boldsymbol{i}}}, \hat{\boldsymbol{\nu}}_{\boldsymbol{m}} \stackrel{\text{def}}{=} \frac{1}{m} \sum_{j} \delta_{\boldsymbol{y}_{\boldsymbol{j}}}$$

Computational properties



$$O((n+m)nm\log(n+m))$$

Statistical properties

$$\mathbb{E}\left[|\boldsymbol{W_p}(\boldsymbol{\mu},\boldsymbol{\nu}) - \boldsymbol{W_p}(\boldsymbol{\hat{\mu}_n},\boldsymbol{\hat{\nu}_m})|\right] \leq f(n,m)?$$

Sample Complexity



If
$$\Omega = \mathbb{R}^d, d > 3$$

$$\mathbb{E}\left[|\boldsymbol{W_p}(\boldsymbol{\mu},\boldsymbol{\nu}) - \boldsymbol{W_p}(\boldsymbol{\hat{\mu}_n},\boldsymbol{\hat{\nu}_n})|\right] = O(n^{-1/d})$$

- •[Dudley'69][Dereich+'11][Fournier+'13] & others..
- [Weed/Bach'17]: sharper results when measures' support has "low effective d" in metric spaces
- [Weed/Berthet'19] for smooth densities
- Lower bounds: optimal quantization error.

From theory to practice?

$$\hat{\boldsymbol{\mu}}_{\boldsymbol{n}} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i} \delta_{\boldsymbol{x_i}}, \hat{\boldsymbol{\nu}}_{\boldsymbol{m}} \stackrel{\text{def}}{=} \frac{1}{m} \sum_{j} \delta_{\boldsymbol{y_j}}$$

Computational properties



$$O((n+m)nm\log(n+m))$$

Statistical properties



$$\mathbb{E}\left[|\boldsymbol{W_p}(\boldsymbol{\mu},\boldsymbol{\nu}) - \boldsymbol{W_p}(\hat{\boldsymbol{\mu_n}},\hat{\boldsymbol{\nu_n}})|\right] = O(n^{-1/d})$$

For data sciences, we <u>must</u> regularize the problem to improve on either/both aspects.

Many ways to regularize (dual) OT

$$\sup_{\boldsymbol{\varphi}(x)+\boldsymbol{\psi}(y)\leq\boldsymbol{c}(x,y)}\int \boldsymbol{\varphi}d\boldsymbol{\mu}+\int \boldsymbol{\psi}d\boldsymbol{\nu}.$$

• RKHS for potentials [GCBP'16], dualize/smooth indicator constraint, [VMVRB'21]

$$W_1(\mu, \nu) = \sup_{\varphi \text{ 1-Lipschitz}} \int_{\varphi} \varphi(d\mu - d\nu).$$

• Parameterize functions using ReLU Deep net with bounded weights [Arjovsky+'17] or use Wavelet decompositions [Shirdonkhar+'08] for low *d*.

Many ways to regularize (primal) OT

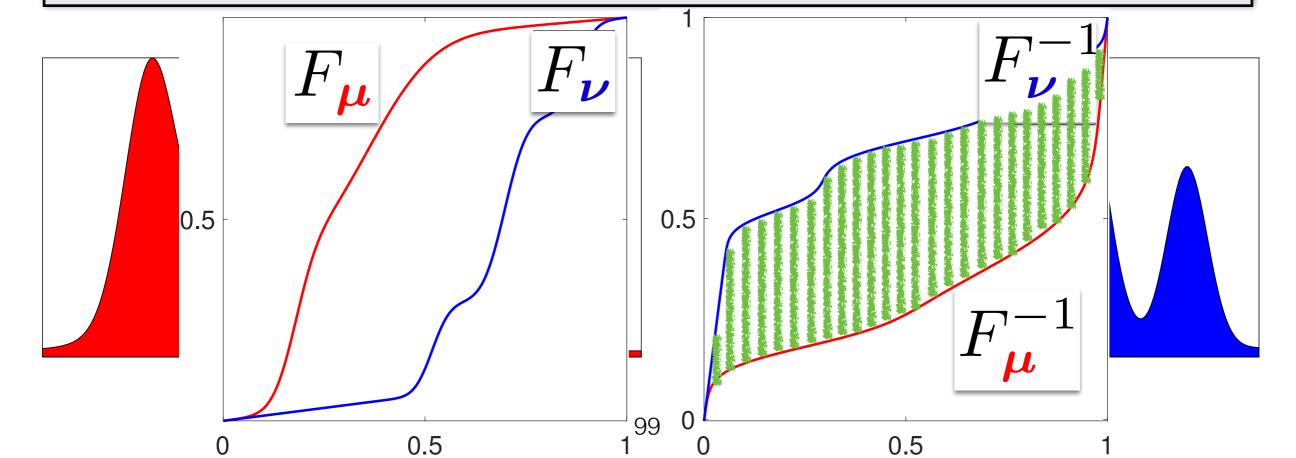
$$\inf_{\boldsymbol{P}\in\Pi(\boldsymbol{\mu},\boldsymbol{\nu})}\int\int \boldsymbol{c}(x,y)\boldsymbol{P}(dx,dy).$$

- •Change cost function: threshold metric [Pele+'09], use geodesic distance on graphs [Beckman'52] [Lin+'07], [Solomon+'14], simplifies the LP.
- Quantize measures first [Canas+'12]; use Gaussians [Gelbrich'92][Chen+17]; projections [Rabin+'11] & k-dimensional subspaces [Paty+'19][Weed+'19].
- •Add regularization on coupling [C'13][GP'16] [GCBP'16] [GCBCP'19] [DPR'16] [BSR'18]

A Different Route: Projection

Remark. If $\Omega = \mathbb{R}$, c(x,y) = c(|x-y|), c convex, F_{μ}^{-1} , F_{ν}^{-1} quantile functions,

$$W(\boldsymbol{\mu}, \boldsymbol{\nu}) = \int_0^1 \boldsymbol{c}(|F_{\boldsymbol{\mu}}^{-1}(x) - F_{\boldsymbol{\nu}}^{-1}(x)|) dx$$



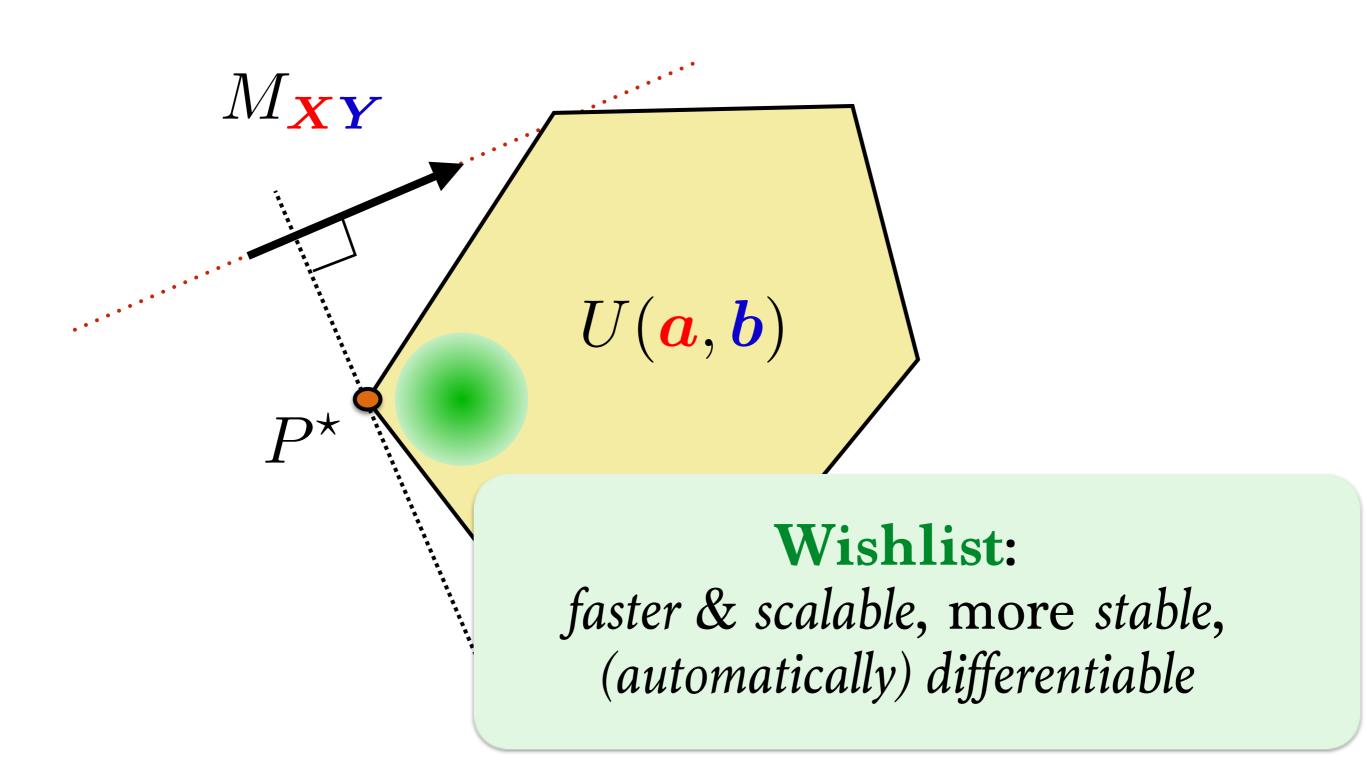
A simple baseline

Sliced Wasserstein Distance [Rabin+'11]

$$SW(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathbb{E}_{\theta \sim \mathcal{S}^{d-1}} \left[\int_0^1 \boldsymbol{c}(|F_{\theta_{\sharp}^T \boldsymbol{\mu}}^{-1}(x) - F_{\theta_{\sharp}^T \boldsymbol{\nu}}^{-1}(x)|) dx \right]$$

- Dodges the high-dimensionality curse, by simplifying considerably the measures.
- Effective in practice, fast and easy, but far from OT. Induced matchings are extremely blurry.

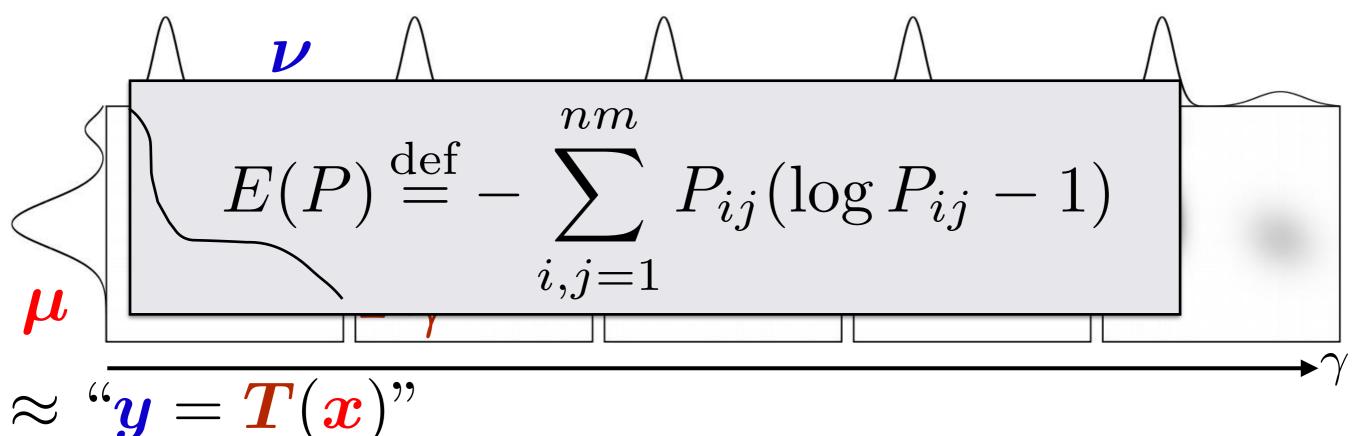
Regularization on the Primal



Entropic Regularization [Wilson'62]

Def. Regularized Wasserstein, $\gamma \geq 0$

$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$



Note: Unique optimal solution because of strong concavity of entropy

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \underset{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})}{\operatorname{argmin}} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$
 then $\exists ! \boldsymbol{u} \in \mathbb{R}^n_+, \boldsymbol{v} \in \mathbb{R}^m_+, \text{ such that}$

$$P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$$

$$L(P, \alpha, \beta) = \sum_{ij} P_{ij} M_{ij} + \gamma P_{ij} (\log P_{ij} - 1) + \alpha^T (P \mathbf{1} - \mathbf{a}) + \beta^T (P^T \mathbf{1} - \mathbf{b})$$
$$\partial L/\partial P_{ij} = M_{ij} + \gamma \log P_{ij} + \alpha_i + \beta_j$$

$$(\partial L/\partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma}} e^{-\frac{M_{ij}}{\gamma}} e^{\frac{\beta_j}{\gamma}} = u_i K_{ij}v_j$$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^n_+, \boldsymbol{v} \in \mathbb{R}^m_+$, such that
$$P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$$

Sinkhorn's Algorithm: Repeat v

1.
$$u = a/Kv$$

2.
$$\mathbf{v} = \mathbf{b}/K^T\mathbf{u}$$

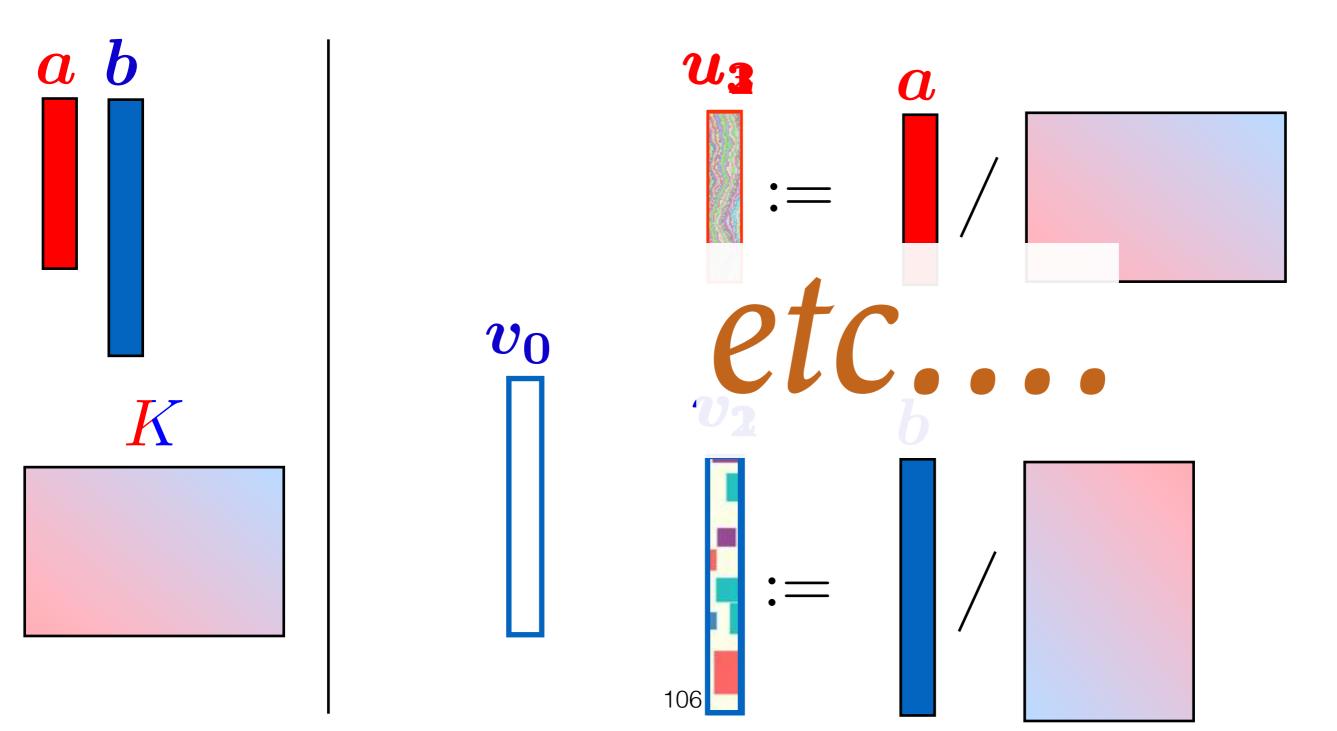
Sinkhorn's Algorithm: Repeat

1.
$$\boldsymbol{u} = \boldsymbol{a}/K\boldsymbol{v}$$
2. $\boldsymbol{v} = \boldsymbol{b}/K^T\boldsymbol{u}$

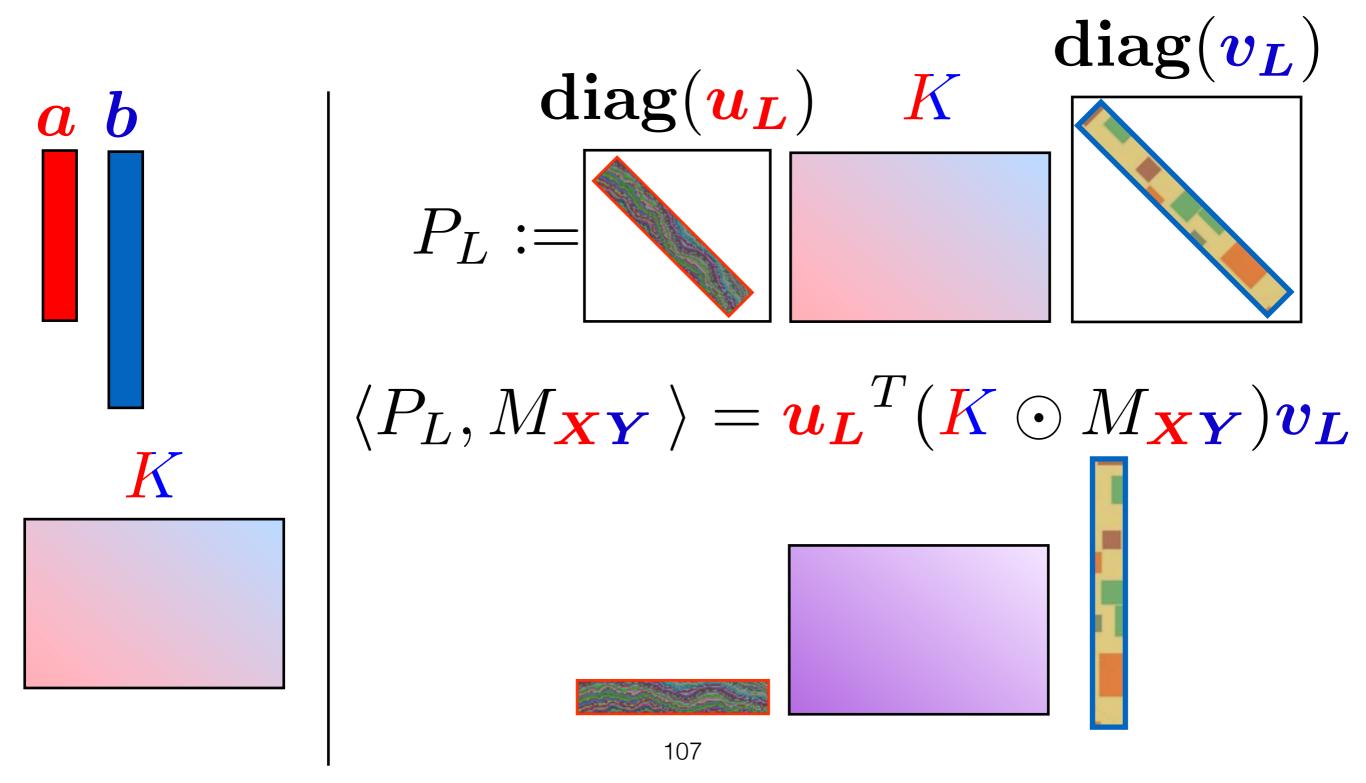
- [Sinkhorn'64] proved first convergence result [Lorenz+'89] characterised linear convergence
- Recent wave of great results by [Altschuler+'17]
 [Dvurechensky+18][Lin+19]
- O(nm) complexity, GPGPU parallel [C'13].
- $O(n \log n)$ on gridded spaces using convolutions. [Solomon'+15]

• [Sinkhorn'64] fixed-point iterations for (u, v)

$$oldsymbol{u}\leftarrow oldsymbol{a}/Koldsymbol{v}, \quad oldsymbol{v}\leftarrow oldsymbol{b}/K^Toldsymbol{u}$$

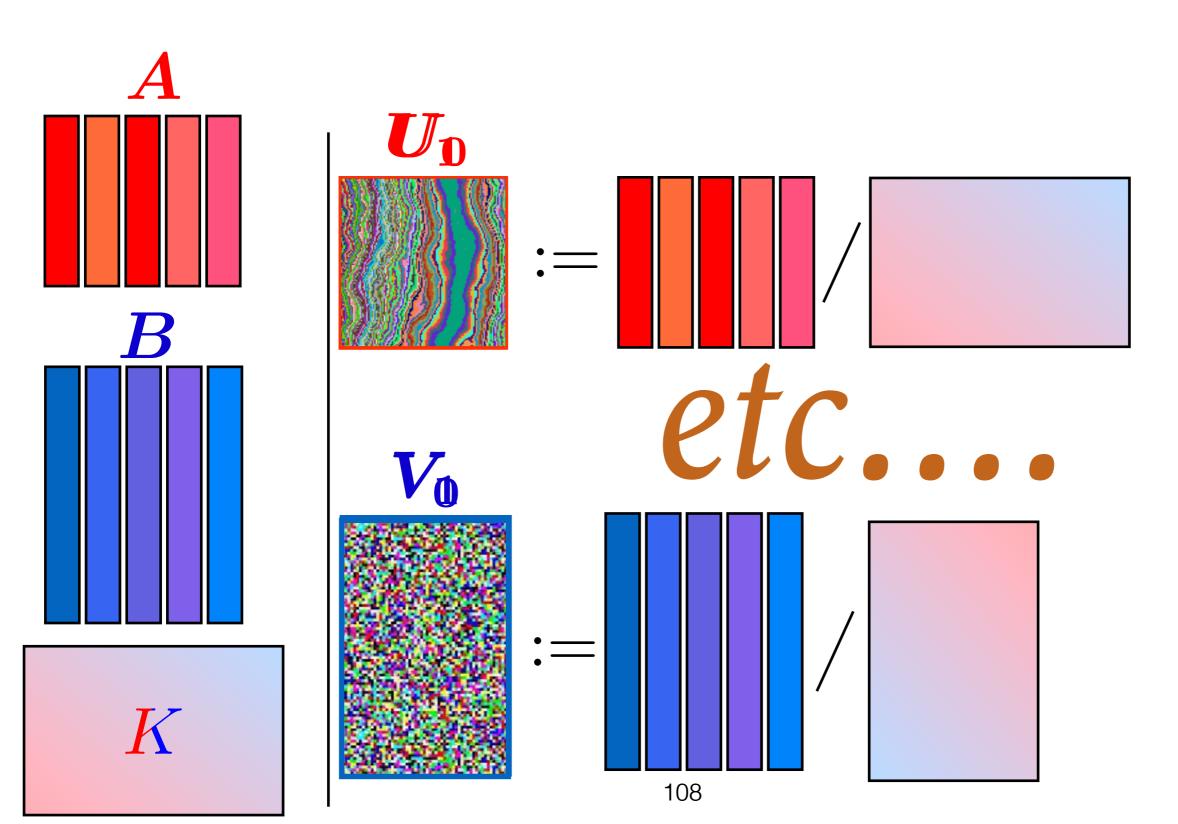


• [Sinkhorn'64] fixed-point iterations.

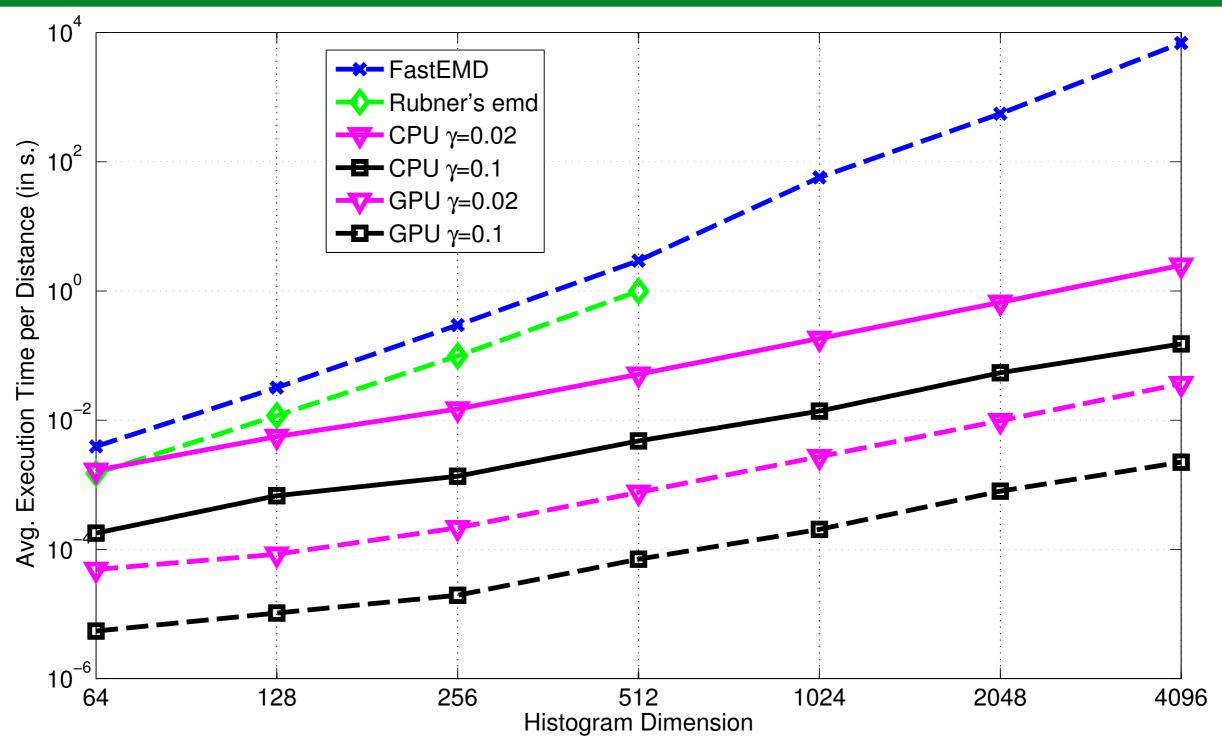


Also embarrassingly parallel

• [Sinkhorn'64] with matrix fixed-point iterations

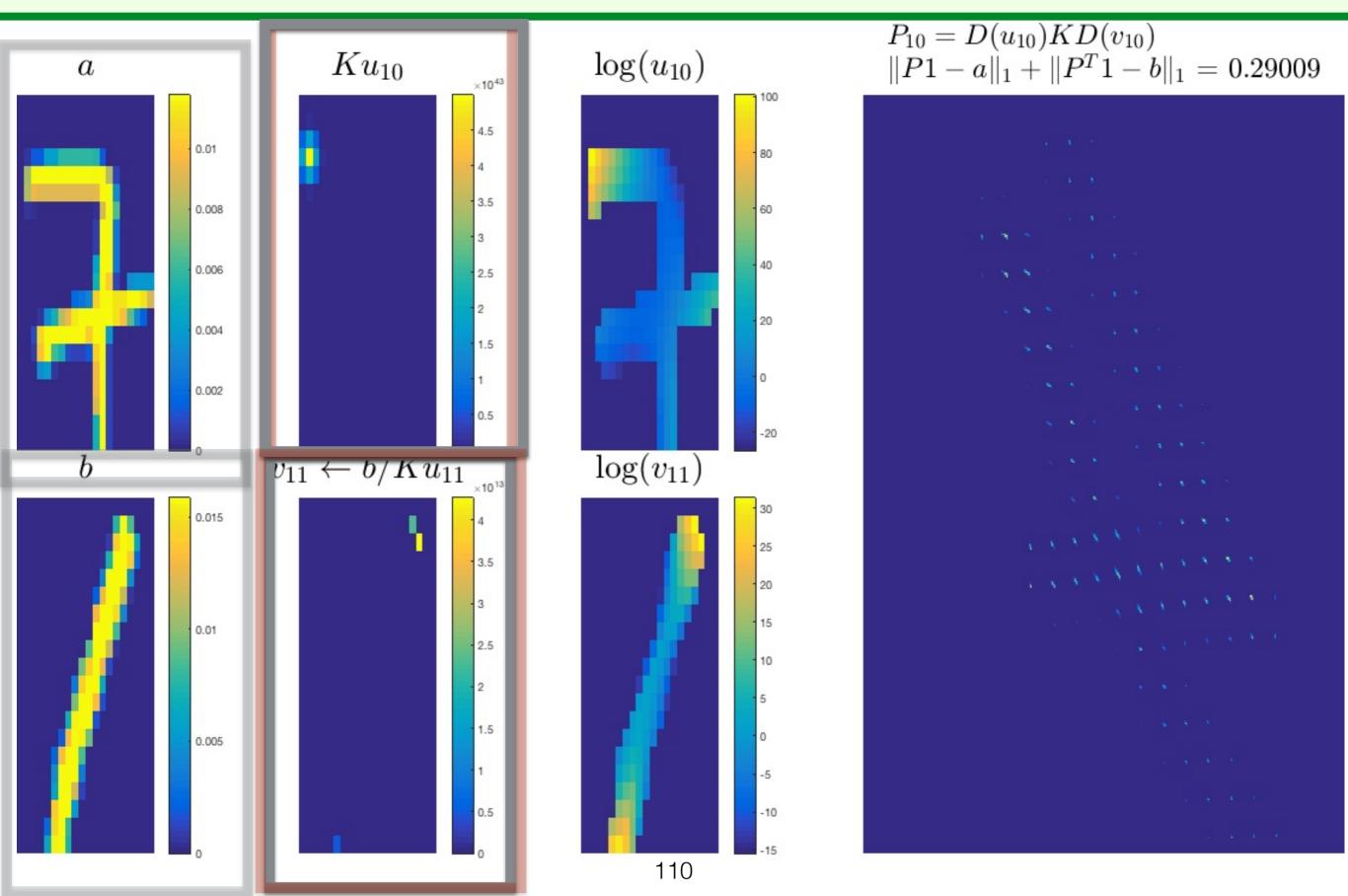


Very Fast EMD Approx. Solver



Setup. (Ω, \mathbf{D}) is a random graph with shortest path metric, histograms sampled uniformly on simplex, Sinkhorn tolerance 10^{-2} .

Very Fast EMD Approx. Solver



Sinkhorn as a Dual Algorithm

Def. Regularized Wasserstein, $\gamma \geq 0$

$$W_{\gamma}(oldsymbol{\mu}, oldsymbol{
u}) \stackrel{\mathrm{def}}{=} \min_{oldsymbol{P} \in U(oldsymbol{a}, oldsymbol{b})} \langle oldsymbol{P}, M_{oldsymbol{XY}}
angle - \gamma E(oldsymbol{P})$$
 regularized discrete primal

$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \boldsymbol{\alpha}^T \boldsymbol{a} + \boldsymbol{\beta}^T \boldsymbol{b} - \gamma (e^{\boldsymbol{\alpha}/\gamma})^T K(e^{\boldsymbol{\beta}/\gamma})$$
 where $K = \left[e^{-\frac{D^p(\boldsymbol{x_i}, \boldsymbol{y_j})}{\gamma}}\right]_{ij}$

Sinkhorn = *Block Coordinate Ascent* on Dual

Block Coordinate Ascent, a.k.a Sinkhorn

$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \boldsymbol{\alpha}^{T} \boldsymbol{a} + \boldsymbol{\beta}^{T} \boldsymbol{b} - \gamma (e^{\boldsymbol{\alpha}/\gamma})^{T} K(e^{\boldsymbol{\beta}/\gamma})$$

REGULARIZED DISCRETE DUAL

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \boldsymbol{\alpha}^T \boldsymbol{a} + \boldsymbol{\beta}^T \boldsymbol{b} - \gamma (e^{\boldsymbol{\alpha}/\gamma})^T K e^{\boldsymbol{\beta}/\gamma}$$

$$\nabla_{\alpha} \mathcal{E} = \mathbf{a} - e^{\alpha/\gamma} \odot K e^{\beta/\gamma}$$

$$\alpha \leftarrow \gamma \left(\log \alpha - \log K(e^{\beta/\gamma}) \right)$$

$$\nabla_{\beta} \mathcal{E} = \mathbf{b} - e^{\beta/\gamma} \odot K^T e^{\alpha/\gamma}$$

$$|\beta \leftarrow \gamma \left(\log b - \log K^T (e^{\alpha/\gamma}) \right)|$$

Block Coordinate Ascent, a.k.a Sinkhorn

$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \boldsymbol{\alpha}^{T} \boldsymbol{a} + \boldsymbol{\beta}^{T} \boldsymbol{b} - \gamma (e^{\boldsymbol{\alpha}/\gamma})^{T} K(e^{\boldsymbol{\beta}/\gamma})$$

$$(\boldsymbol{u}, \boldsymbol{v}) \stackrel{\text{def}}{=} (e^{\boldsymbol{\alpha}/\gamma}, e^{\boldsymbol{\beta}/\gamma})$$

$$\boldsymbol{\alpha} \leftarrow \gamma \left(\log \boldsymbol{a} - \log K(e^{\boldsymbol{\beta}/\gamma}) \right)$$
 $\boldsymbol{u} \leftarrow \frac{\boldsymbol{a}}{K\boldsymbol{v}}$

$$oldsymbol{u} \leftarrow rac{oldsymbol{a}}{Koldsymbol{v}}$$

$$\beta \leftarrow \gamma \left(\log \boldsymbol{b} - \log K^T(e^{\boldsymbol{\alpha}/\gamma}) \right) \quad \boldsymbol{v} \leftarrow \frac{\boldsymbol{b}}{K^T \boldsymbol{u}}$$

$$oldsymbol{v} \leftarrow rac{oldsymbol{b}}{K^T oldsymbol{u}}$$

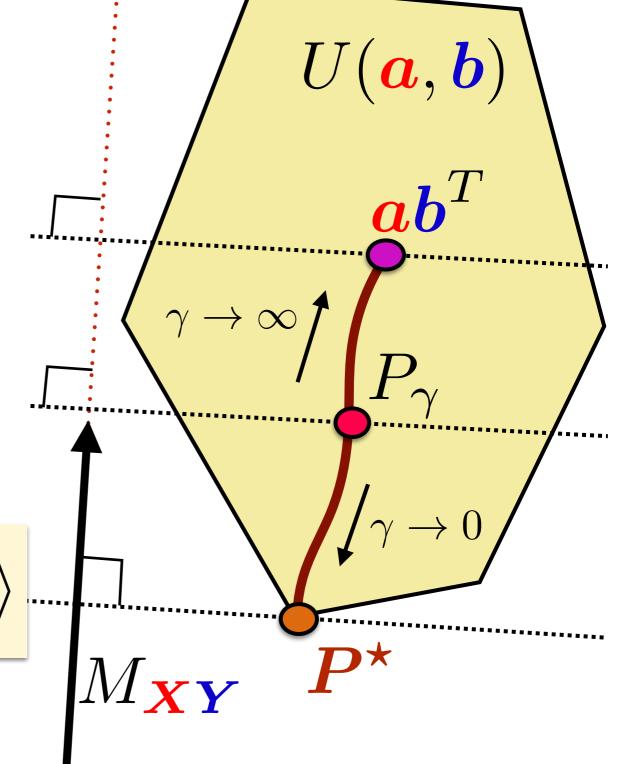
Sinkhorn, Wand Energy Distance

$$\mu = \sum_{i=1}^n a_i \delta_{x_i} \quad \nu = \sum_{j=1}^m b_j \delta_{y_j}$$

$$\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \boldsymbol{a} \boldsymbol{b}^T, M_{\boldsymbol{X} \boldsymbol{Y}} \rangle$$

$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle P_{\gamma}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

$$W^p(\mu, \nu) = \langle P^{\star}, M_{XY} \rangle$$



Sinkhorn, Wand Energy Distance

$$\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \boldsymbol{a} \boldsymbol{b}^T, M_{\boldsymbol{X} \boldsymbol{Y}} \rangle$$

$$ED(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2} (\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle P_{\gamma}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(P_{\gamma})$$

$$\overline{W}_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\mu}) + W_{\gamma}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$W^p(\mu, \nu) = \langle P^{\star}, M_{XY} \rangle$$

Sinkhorn, Wand Energy Distance

$$\mathcal{MMD}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2} (\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$\gamma \to \infty$$

$$\overline{W}_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\mu}) + W_{\gamma}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$\gamma \rightarrow 0$$

$$W^p(\mu, \nu) = \langle P^{\star}, M_{XY} \rangle$$

How to compare them?

i.i.d samples $x_1, \ldots, x_n \sim \mu, y_1, \ldots, y_m \sim \nu$,

$$\hat{\boldsymbol{\mu}}_{\boldsymbol{n}} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i} \delta_{\boldsymbol{x}_{\boldsymbol{i}}}, \hat{\boldsymbol{\nu}}_{\boldsymbol{m}} \stackrel{\text{def}}{=} \frac{1}{m} \sum_{j} \delta_{\boldsymbol{y}_{\boldsymbol{j}}}$$

Computational properties

Effort to compute/approximate $\Delta(\hat{\mu}_n, \hat{\nu}_m)$?

Statistical properties

$$|\Delta(\mu, \nu) - \Delta(\hat{\mu}_n, \hat{\nu}_n)| \le f(n)$$
?

Sinkhorn in between W and MMD

$$\mathcal{MMD}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2} (\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$(n+m)^{2}$$

$$(n+m)^2$$
 $O(1/\sqrt{n})$ [see Arthur]

$$\bar{W}_{\gamma}(\mu, \nu) = W_{\gamma}(\mu, \nu) - \frac{1}{2}(W_{\gamma}(\mu, \mu) + W_{\gamma}(\nu, \nu))$$

$$O((n+m)^2)$$

$$O((n+m)^2)$$

$$O\left(\frac{1}{\gamma^{d/2}\sqrt{n}}\right)$$
[GCBCP'18]
[FSVATP'18]

$$W^p(\mu, \nu) = \langle P^{\star}, M_{XY} \rangle$$

$$O((n+m)nm\log(n+m)) \qquad O(1/n^{1/d})$$

Differentiability of W

$$W((a + \Delta a, X), (b, Y)) = W(a, X), (b, Y) + ??$$

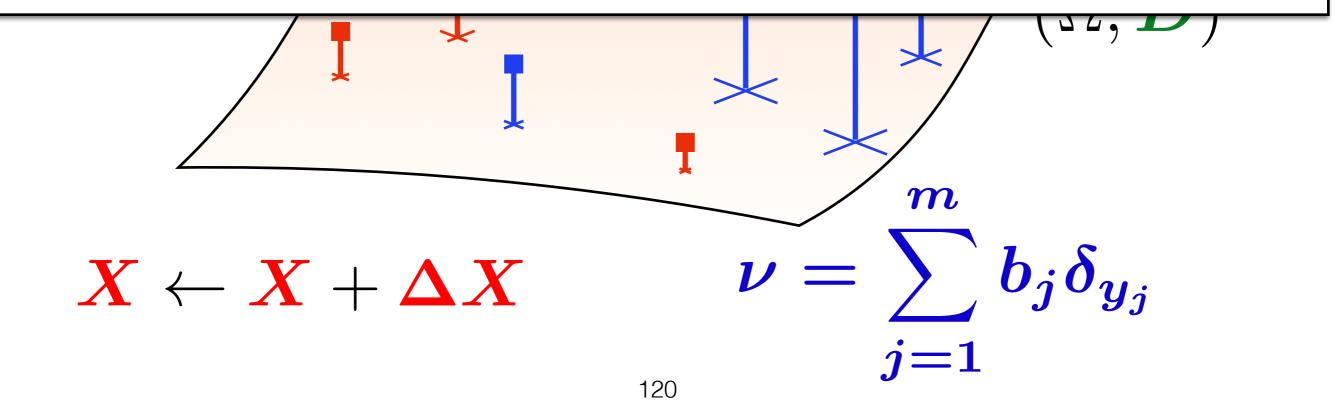
$$\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 (Ω, D) $a \leftarrow a + \Delta a$ $\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$

Sinkhorn ----> Differentiability

$$W((a, X + \Delta X), (b, Y)) = W((a, X), (b, Y)) + ??$$

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

Changes in objective can be handled with Danskin's theorem.



How to decrease W? change weights

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \boldsymbol{\alpha} \oplus \boldsymbol{\beta} \le M_{\boldsymbol{X}\boldsymbol{Y}}}} \boldsymbol{\alpha}^T \boldsymbol{a} + \boldsymbol{\beta}^T \boldsymbol{b}.$$

Prop.
$$W(\mu, \nu)$$
 is convex w.r.t. \boldsymbol{a} ,
$$\partial_{\boldsymbol{a}} W = \arg_{\boldsymbol{\alpha}} \max_{\boldsymbol{\alpha} \in \boldsymbol{\beta} \leq M_{\boldsymbol{X}\boldsymbol{Y}}} \boldsymbol{\alpha}^T \boldsymbol{a} + \boldsymbol{\beta}^T \boldsymbol{b}.$$

Prop. $W_{\gamma}(\mu, \nu)$ is convex and differentiable w.r.t. \boldsymbol{a} , $\nabla_{\boldsymbol{a}} W_{\gamma} = \boldsymbol{\alpha}_{\gamma}^{\star} = \gamma \log \boldsymbol{u}$

How to decrease W? change locations

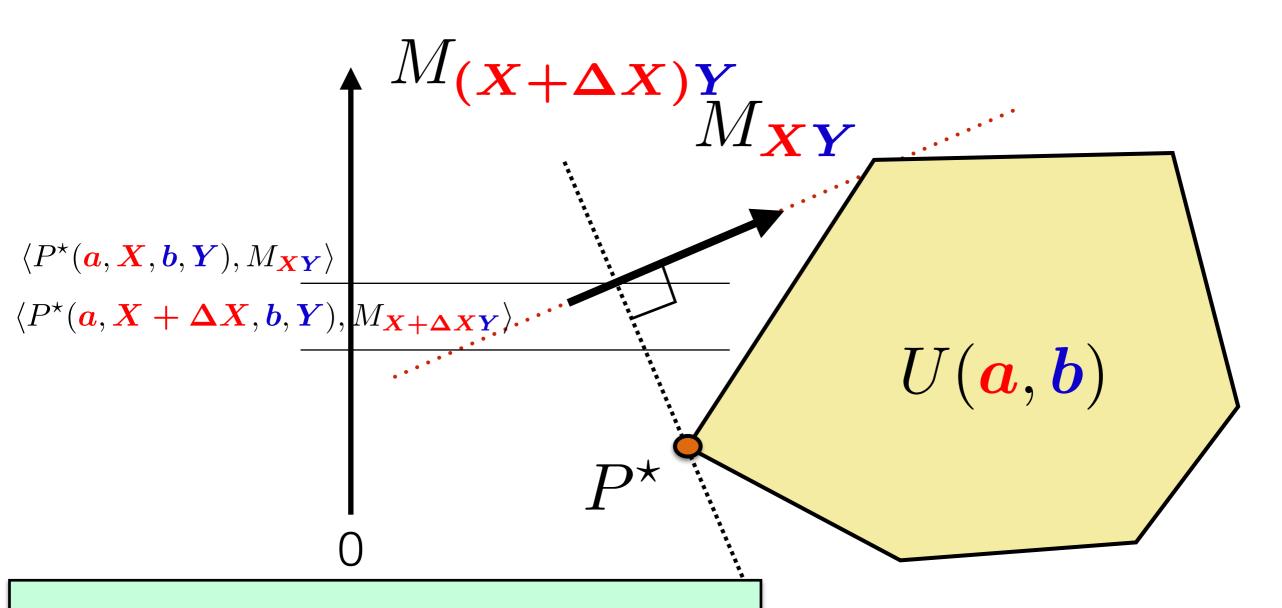
$$W_{\mathbf{2}}^{\mathbf{2}}(\boldsymbol{\mu}, \boldsymbol{
u}) = \min_{\substack{\boldsymbol{P} \in \mathbb{R}_{+}^{n imes m} \\ \boldsymbol{P} \mathbf{1}_{m} = \boldsymbol{a}, \boldsymbol{P}^{T} \mathbf{1}_{n} = \boldsymbol{b}}} \langle \boldsymbol{P}, \mathbf{1}_{n} \mathbf{1}_{d}^{T} \boldsymbol{X}^{2} + \boldsymbol{Y}^{2T} \mathbf{1}_{d} \mathbf{1}_{m} - 2 \boldsymbol{X}^{T} \boldsymbol{Y} \rangle$$

Prop. $p = 2, \Omega = \mathbb{R}^d$. $W(\mu, \nu)$ decreases if $X \leftarrow Y P^{\star T} \mathbf{D}(a^{-1})$.

Prop.
$$p=2, \Omega=\mathbb{R}^d$$
. $W_{\gamma}(\mu,\nu)$ is differentiable w.r.t. X , with

$$\nabla_{\mathbf{X}} W_{\gamma} = \mathbf{X} - \mathbf{Y} P_{\gamma}^T \mathbf{D}(\mathbf{a}^{-1}).$$

Solving the OT Problem



Computing $\nabla \blacksquare W$ is easy.

here \blacksquare can be anything,

e.g. \boldsymbol{a} or \boldsymbol{X} , parameter $\boldsymbol{\theta}$ for cost $\boldsymbol{c}_{\boldsymbol{\theta}}$

Computing $\frac{\partial P^*}{\partial \blacksquare}$ is harder.

Sinkhorn: A Programmer View

Def. For $L \geq 1$, define

$$P_L \stackrel{\text{def}}{=} \operatorname{diag}(u_L) K \operatorname{diag}(v_L),$$

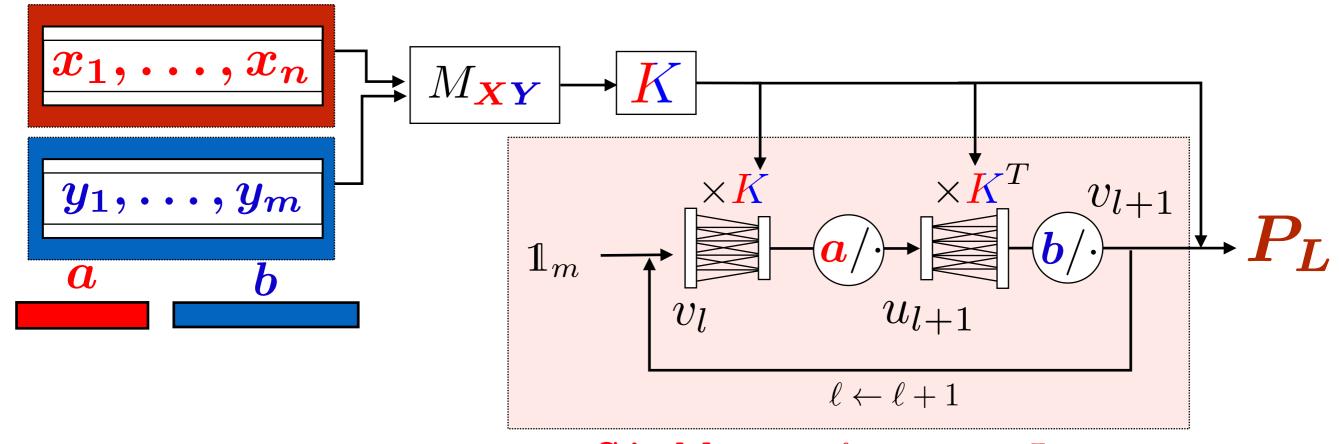
where

$$\mathbf{v_0} = \mathbf{1}_m; l \geq 0, \mathbf{u_l} \stackrel{\text{def}}{=} \mathbf{a} / K \mathbf{v_l}, \mathbf{v_{l+1}} \stackrel{\text{def}}{=} \mathbf{b} / K^T \mathbf{u_l}.$$

Prop. $\frac{\partial P_L}{\partial X}$, $\frac{\partial P_L}{\partial a}$ can be computed recursively, in O(L) kernel $K \times$ vector products.

Sinkhorn: A Programmer View

Def. For
$$L \geq 1$$
, define
$$W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P_L}, M_{\boldsymbol{XY}} \rangle,$$



Sinkhorn $\ell = 1, \dots, L-1$

[Adams'11] [Hashimoto'16] [Bonneel'16] [Shalit'16]

Sinkhorn: A Mathematician View

$$F: \boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{c}, \varepsilon \to (\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$$

$$H(\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{c}, \varepsilon; \boldsymbol{\alpha}, \boldsymbol{\beta}) := \begin{bmatrix} e^{\frac{\boldsymbol{\alpha} \oplus \boldsymbol{\beta} - \boldsymbol{C}}{\varepsilon}} \mathbf{1}_m - \boldsymbol{a} \\ e^{\frac{\boldsymbol{\beta} \oplus \boldsymbol{\alpha} - \boldsymbol{C}^T}{\varepsilon}} \mathbf{1}_n - \boldsymbol{b} \end{bmatrix} = 0$$

At the optimum,

$$H(\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{c}, \boldsymbol{\varepsilon}; F(\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{c}, \boldsymbol{\varepsilon})) = \mathbf{0}$$

Sinkhorn: A Mathematician View

$$F: \boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{c}, \varepsilon \to (\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$$

$$H(\boldsymbol{\mu}, \boldsymbol{\nu}, \boldsymbol{c}, \varepsilon; \boldsymbol{\alpha}, \boldsymbol{\beta}) := \begin{bmatrix} e^{\frac{\boldsymbol{\alpha} \oplus \boldsymbol{\beta} - \boldsymbol{C}}{\varepsilon}} \mathbf{1}_m = \boldsymbol{a} \\ e^{\frac{\boldsymbol{\beta} \oplus \boldsymbol{\alpha} - \boldsymbol{C}^T}{\varepsilon}} \mathbf{1}_n = \boldsymbol{b} \end{bmatrix} = 0$$

Using the implicit function theorem

$$J_{F,\blacksquare} = -\left(J_{H,(\boldsymbol{\alpha},\boldsymbol{\beta})}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare}(\blacksquare,(\boldsymbol{\alpha}^*,\boldsymbol{\beta}^*))^{-1}J_{H,\blacksquare$$

$$J_{F,\blacksquare}^T = -J_{H,\blacksquare}(\blacksquare, (\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)^T \left(J_{H,(\boldsymbol{\alpha},\boldsymbol{\beta})}(\blacksquare, (\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)\right)^{-T})$$

The best of both worlds

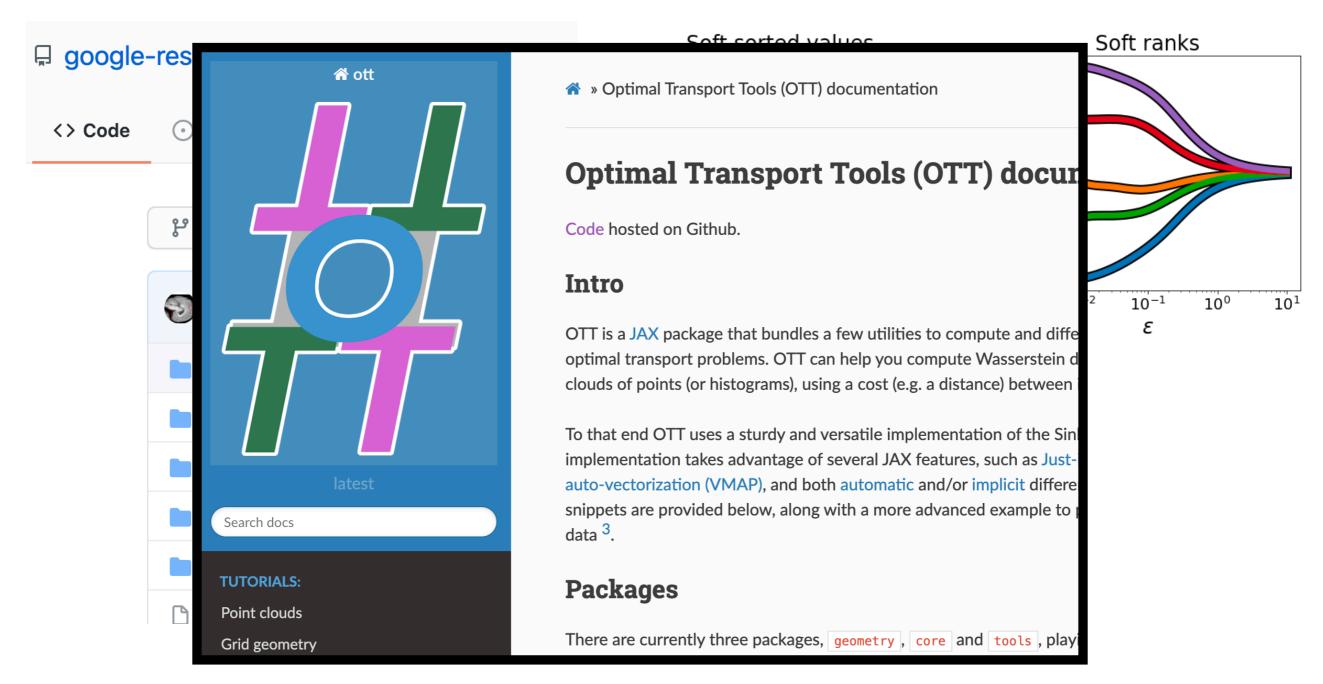
When computing the gradient of a loss whose evaluation depends on the optimal transport solution, backward mode differentiation is more efficient. This involves evaluating:

$$J_{F,lacksquare}^T\mathbf{z}$$

$$-J_{H,\blacksquare}(\blacksquare,(\pmb{\alpha}^*,\pmb{\beta}^*)^T \left(\left(J_{H,(\pmb{\alpha},\pmb{\beta})}(\blacksquare,(\pmb{\alpha}^*,\pmb{\beta}^*)\right)^T \right)^{-1}\mathbf{z}$$
jax.vjp
jax.linalg.sparse.cg

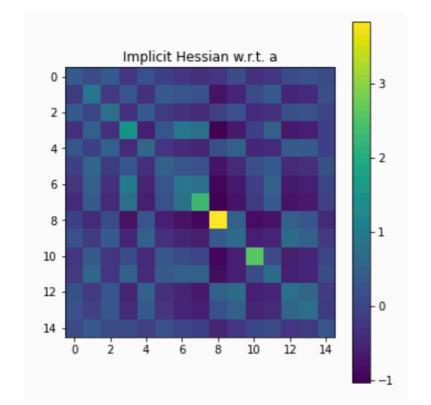
For more details

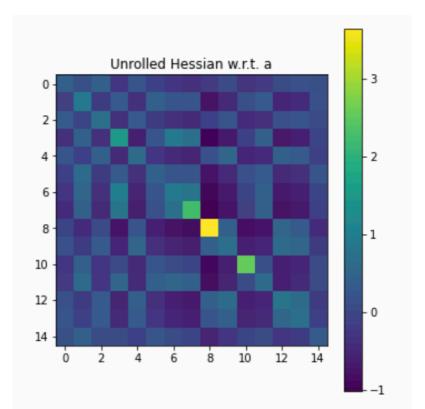
https://github.com/google-research/ott



For instance...computing Hessians.

https://ott-jax.readthedocs.io/en/stable/notebooks/Hessians.html





To conclude

To sum up

- Optimal matchings used everywhere in ML, to
 - Compute a distance (W) between measures.
 - register/map/match points, to disambiguate.
- Regularization is needed to:
 - Scale up/robustify these tools;
 - Ensure gradients/jacobians exist and are not 0.
- To differentiate through regularised OT, either
 - Automatic differentiation (unrolling)
 - Implicit differentiation

A few research topics around OT

- Sinkhorn fixed-point iterations speed
 - Using convolutions [Solomon+15]
 - Faster kernel multiplications (many refs!)
 - Accelerations (Anderson [Chizat+20])
 - Improved convergence proofs

- Generalize Sinkhorn to continuous spaces?
 (Schrödinger bridges) [Chen+14] [Bortoli+21]
- Other regularizers (L2) / constraints (low-rank)

A few research topics around OT

Unbalanced OT formulations

- Penalized [FZMAP15][Chizat+15], approaches can handle (with Sinkhorn) different marginals
- Brenier / Monge topics
 - Estimate Monge maps using ICNN [VM+19]
 - Links with Normalizing Flows [Huang+20]
- Proximal optimization in spaces of measures.
 - [JKO'98] starting to make an impact in practice, see for instance [Bunne+21]

A few research topics around OT

OT for Heterogeneous spaces

- Use quadratic assignment problem to compute machines (GW, [Memoli'11])
- Computational challenges, stats unknown.

Exploit matching differentiability

- soft-sorting [CTV19] and -quantiles [CTNV'20]
- to impute missing values [Muzellec+20]