

# OSL 2015

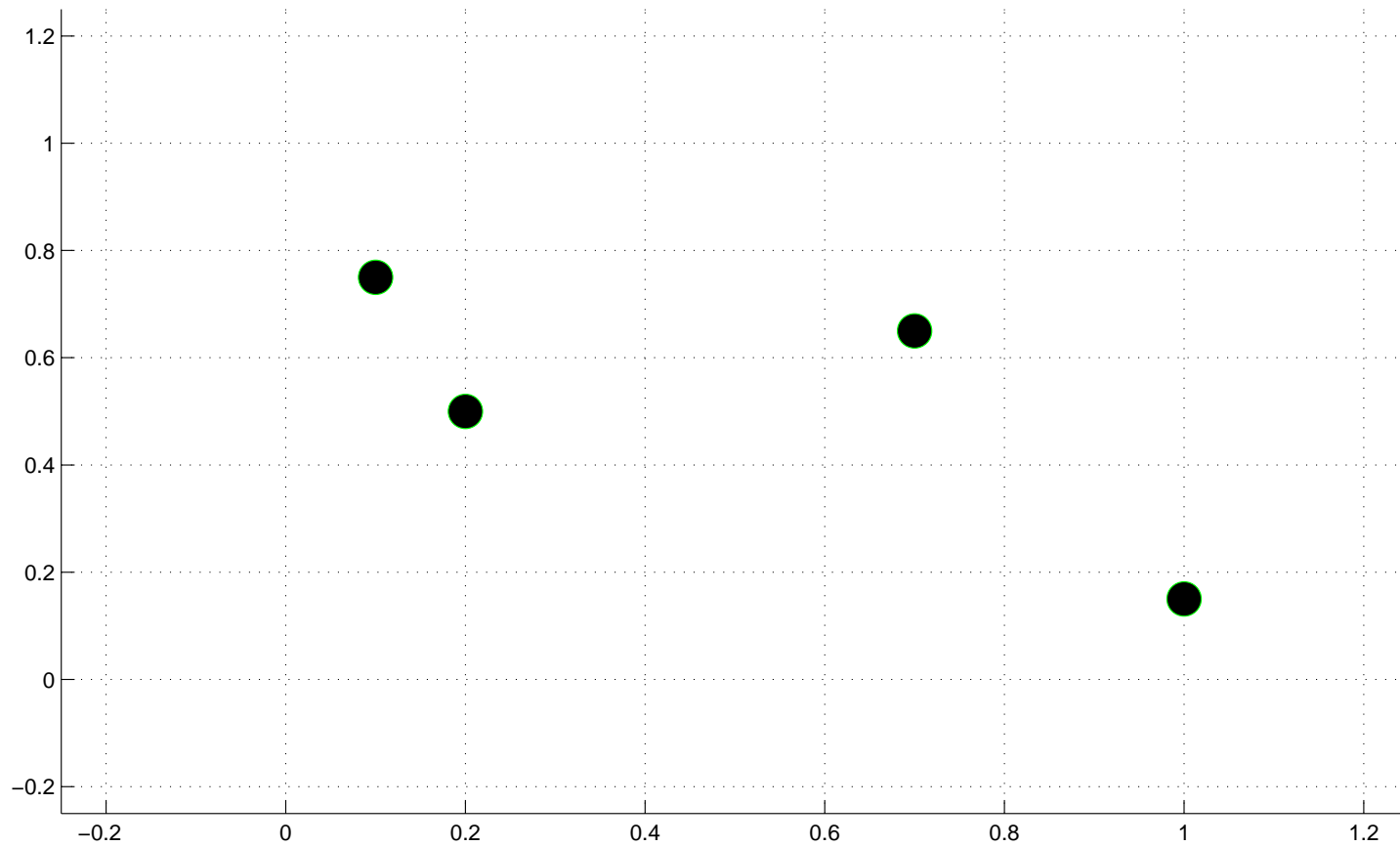
## The Wasserstein Barycenter Problem

**Marco Cuturi**

`mcuturi@i.kyoto-u.ac.jp`

Joint work with G. Peyré, G. Carlier, J.D. Benamou, L. Nenna,  
A. Gramfort, J. Solomon, ...

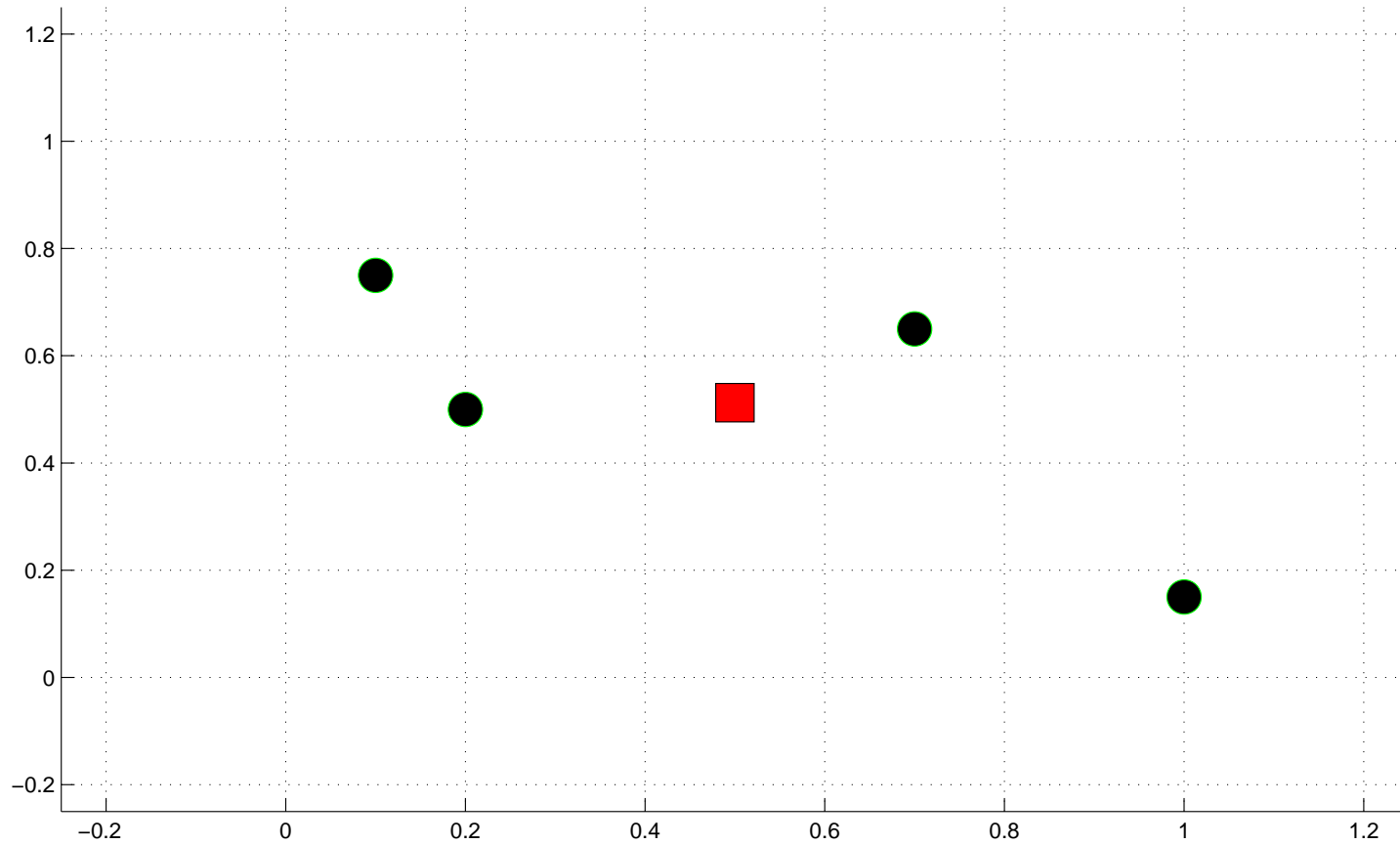
# Motivation



4 points in  $\mathbb{R}^2$

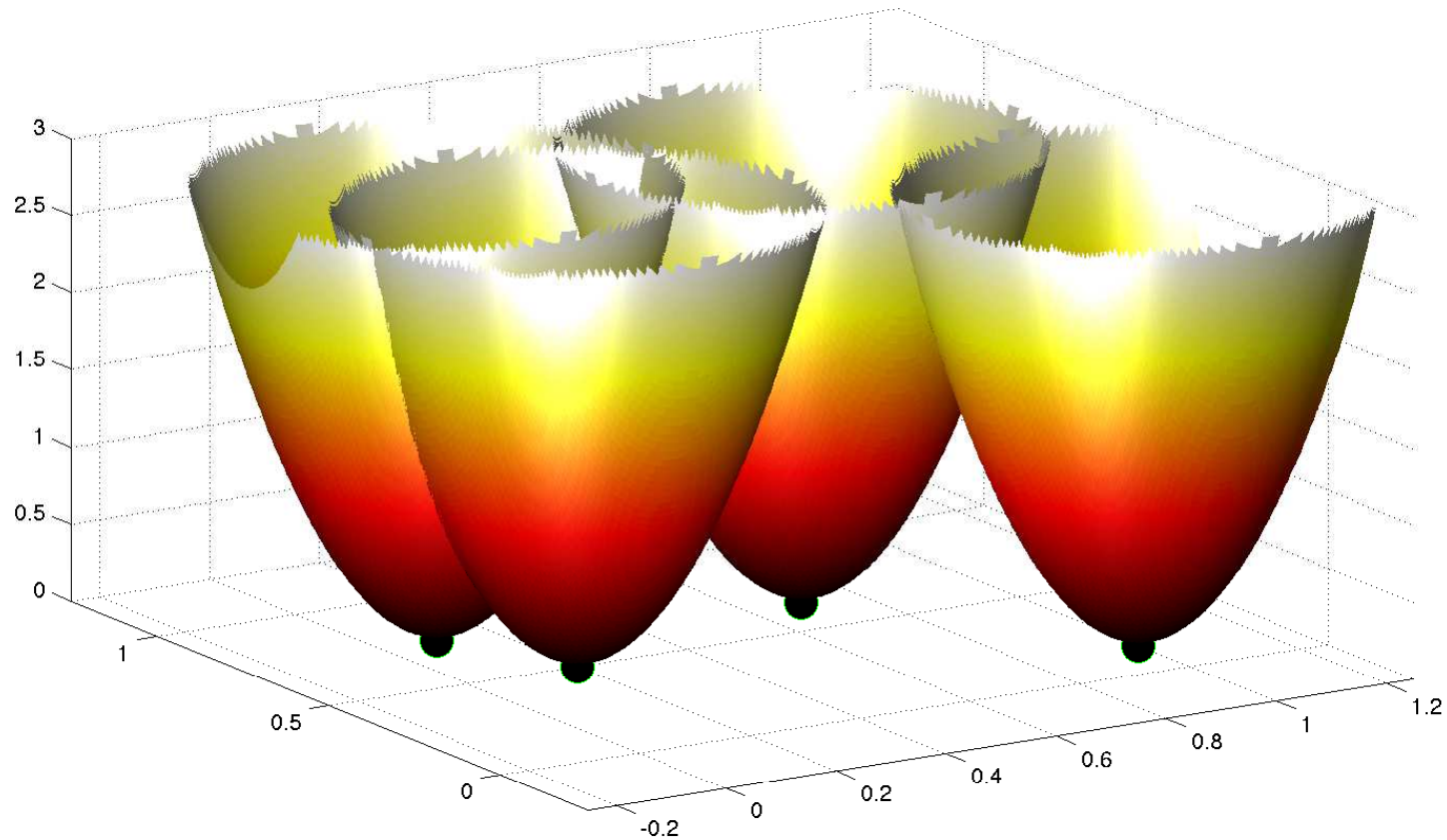
$x_1, x_2, x_3, x_4$

# Mean



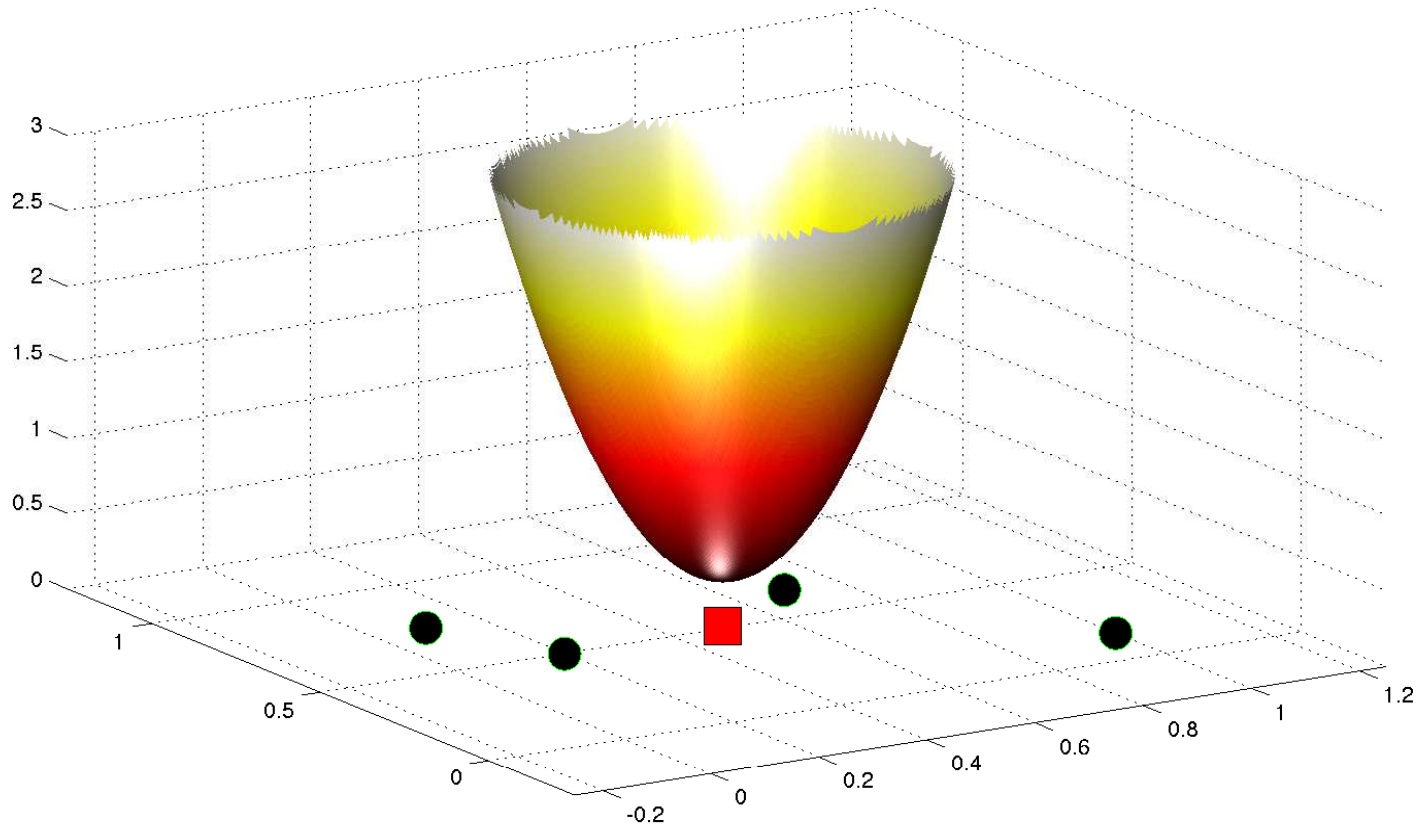
Their **mean** is  $(x_1 + x_2 + x_3 + x_4) / 4$ .

# Computing Means



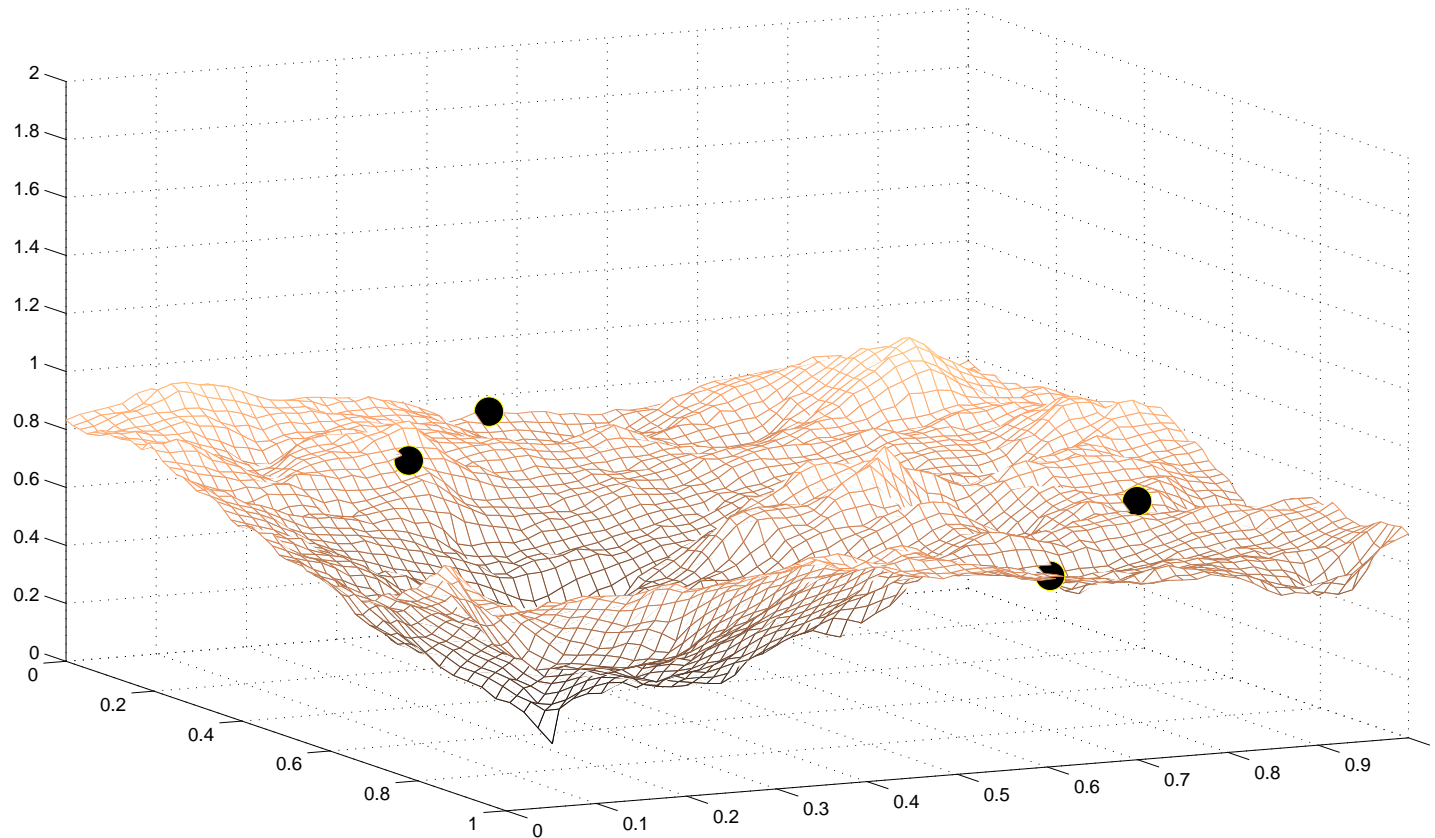
Consider for each point the function  $\|\cdot - x_i\|_2^2$

# Computing Means



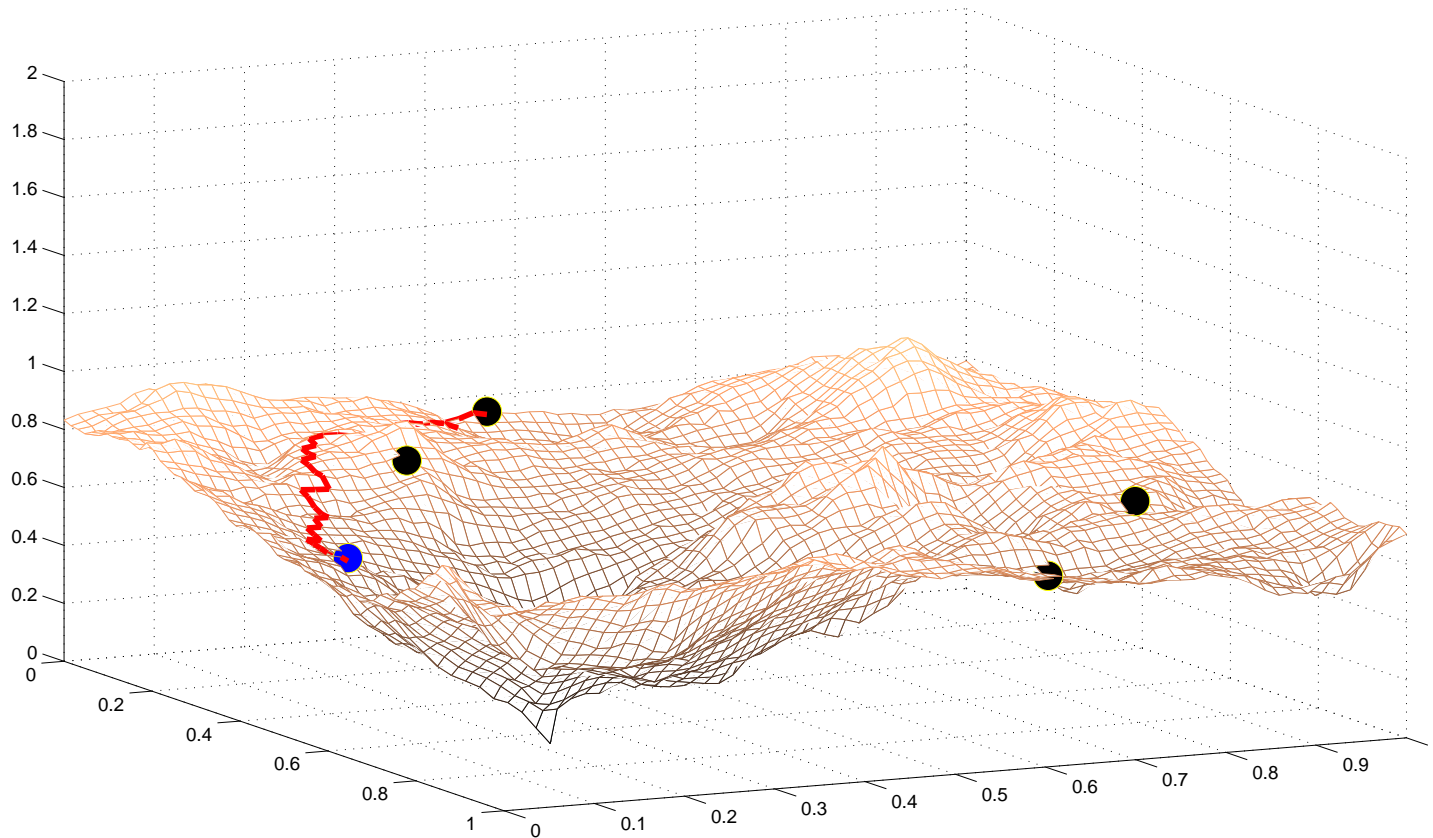
The **mean** is the  $\operatorname{argmin} \frac{1}{4} \sum_{i=1}^4 \|\cdot - x_i\|_2^2$ .

# Means in Metric Spaces



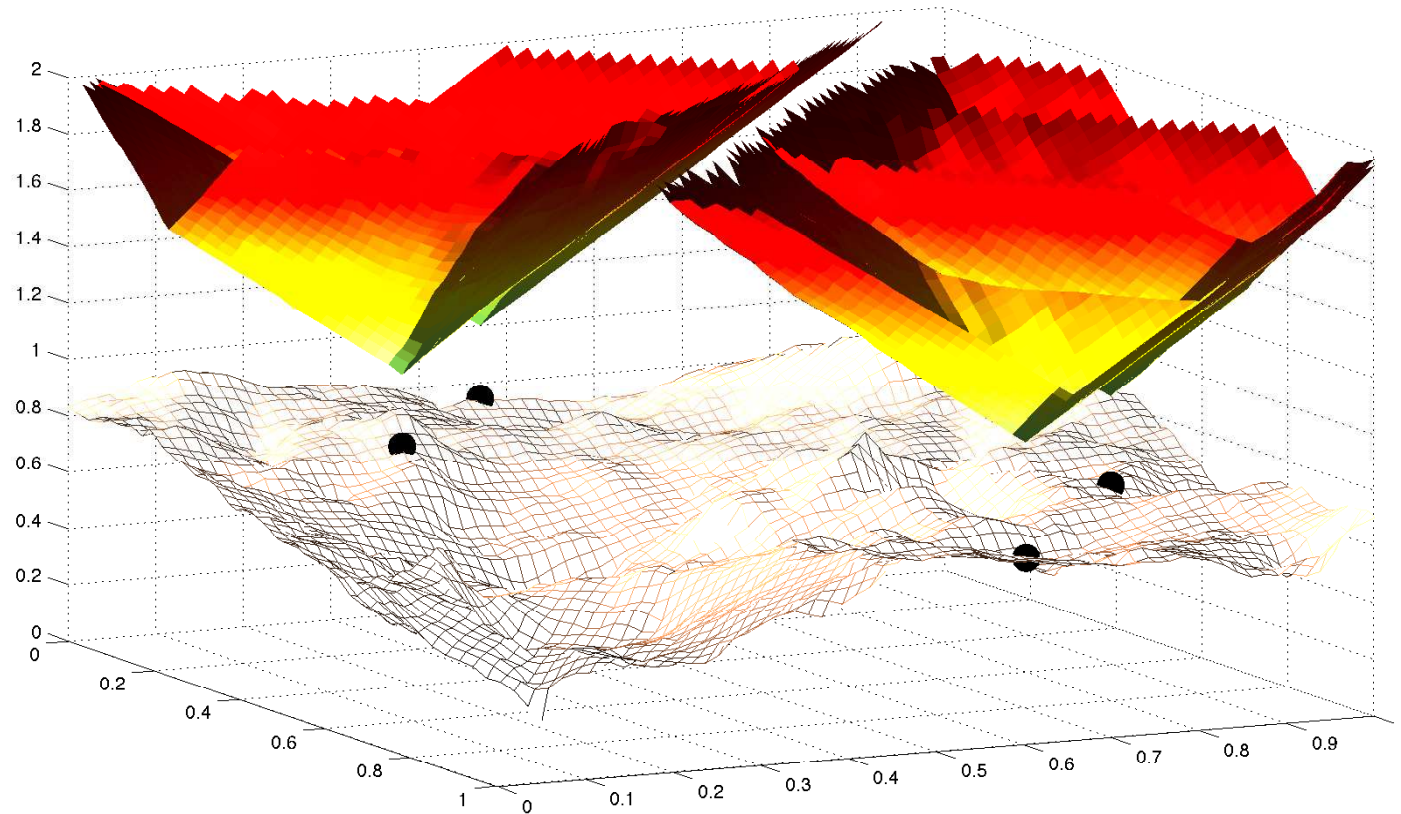
Means can be defined using any distance/divergence/discrepancy.

# Means in Metric Spaces



Using *e.g.* geodesic distances. Here  $\Delta(\bullet, \bullet) = 0.994$

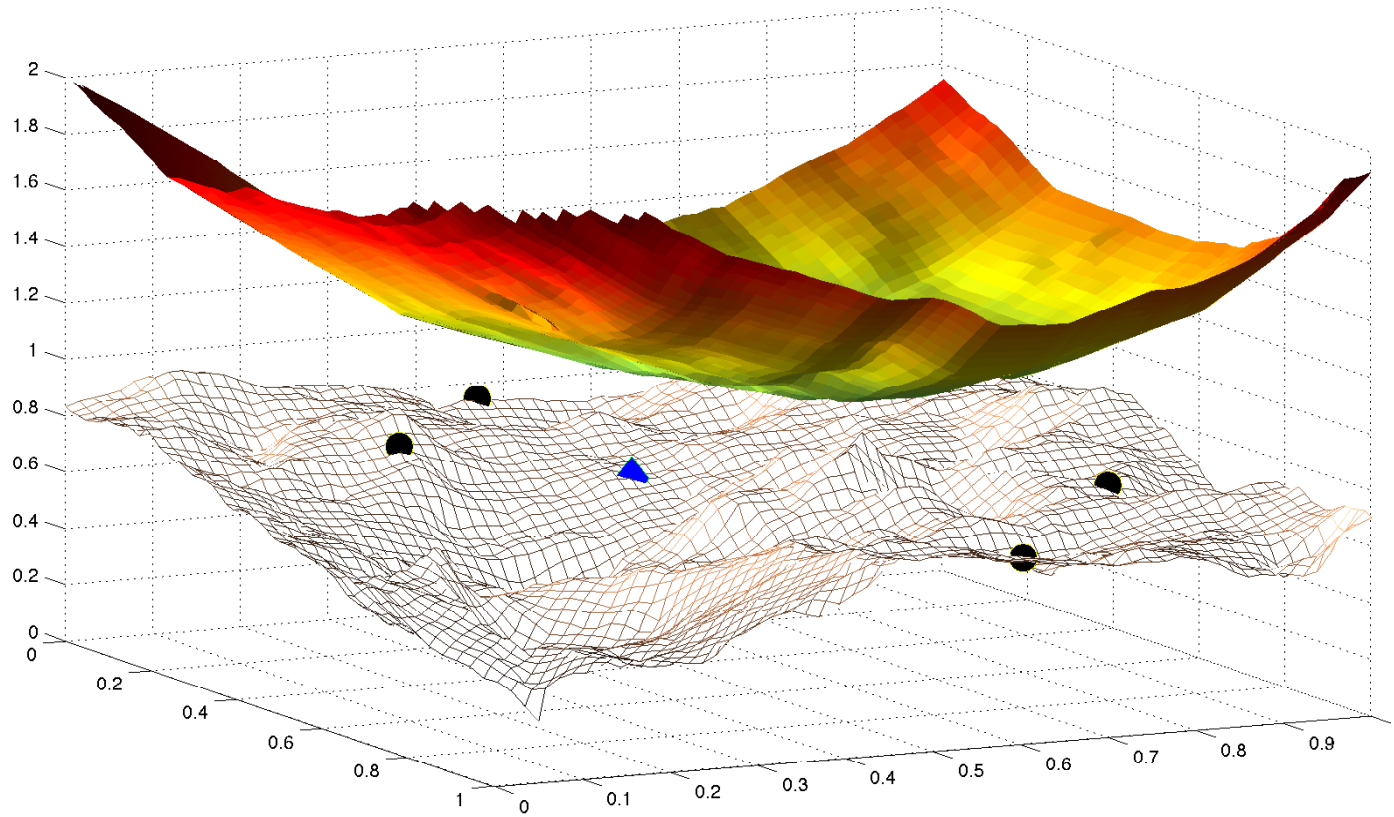
# Means in Metric Spaces



Consider the distance functions  $\Delta(\cdot, x_i), i = 1, 2, 3, 4$ .

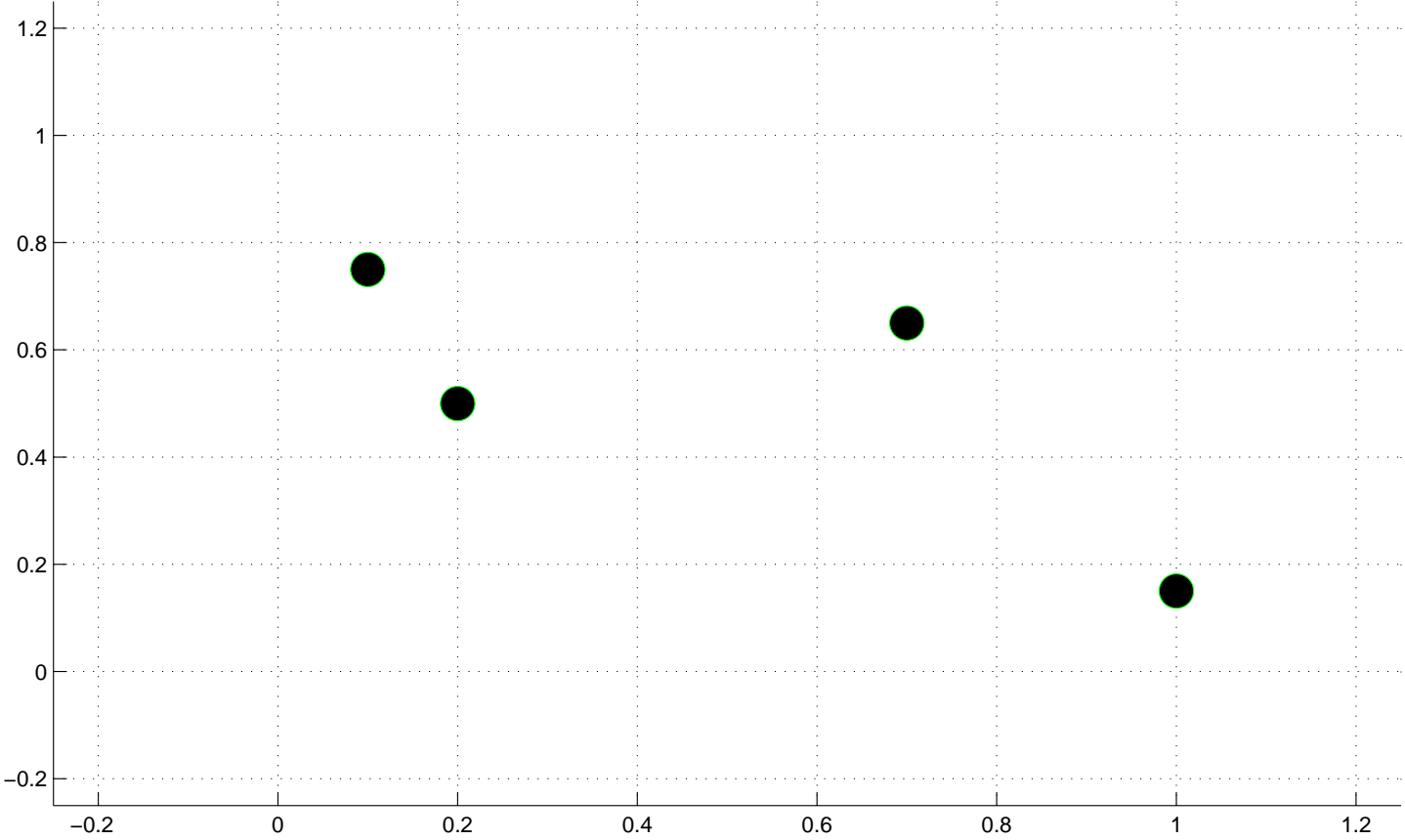


# Means in Metric Spaces

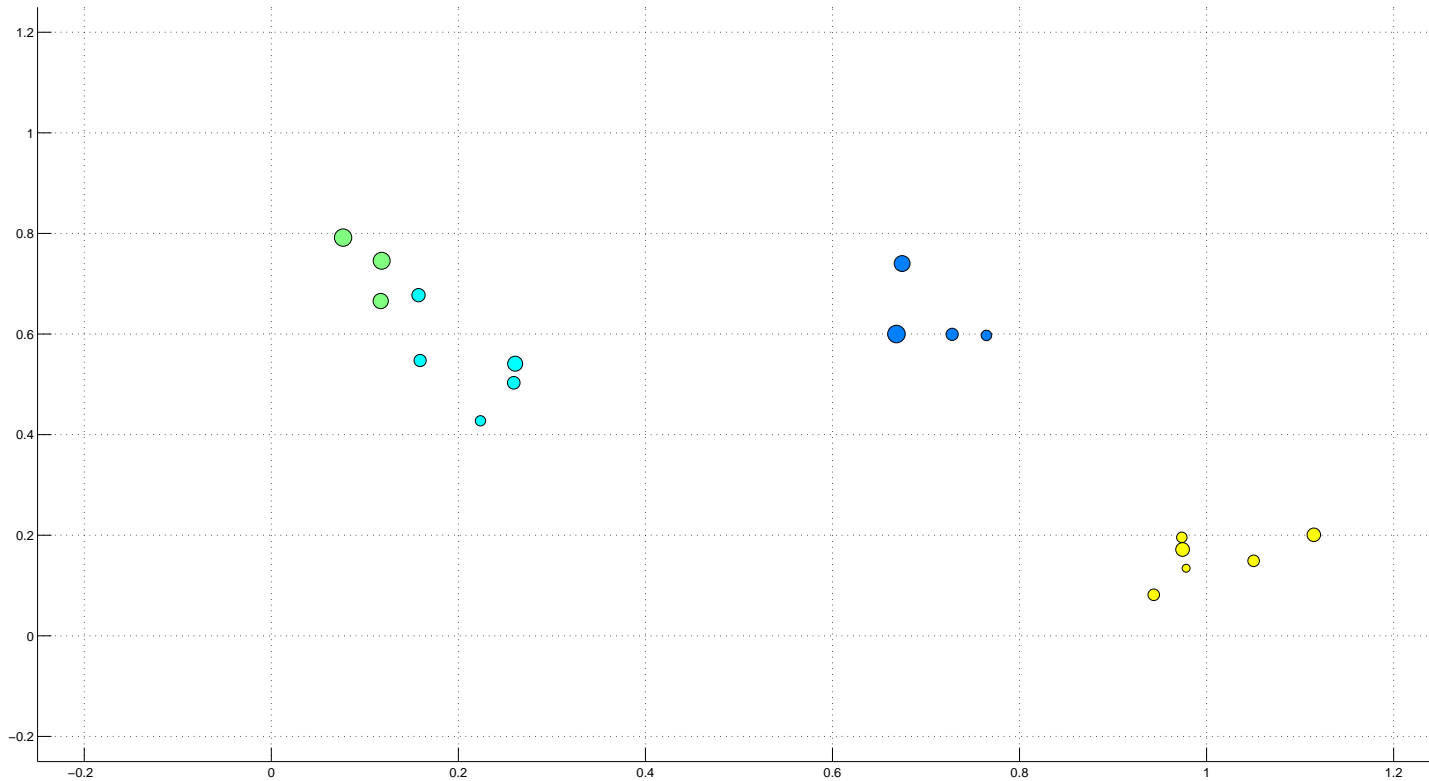


$$\blacklozenge = \operatorname{argmin} \frac{1}{N} \sum_{i=1}^N \Delta(\cdot, x_i).$$

# From points

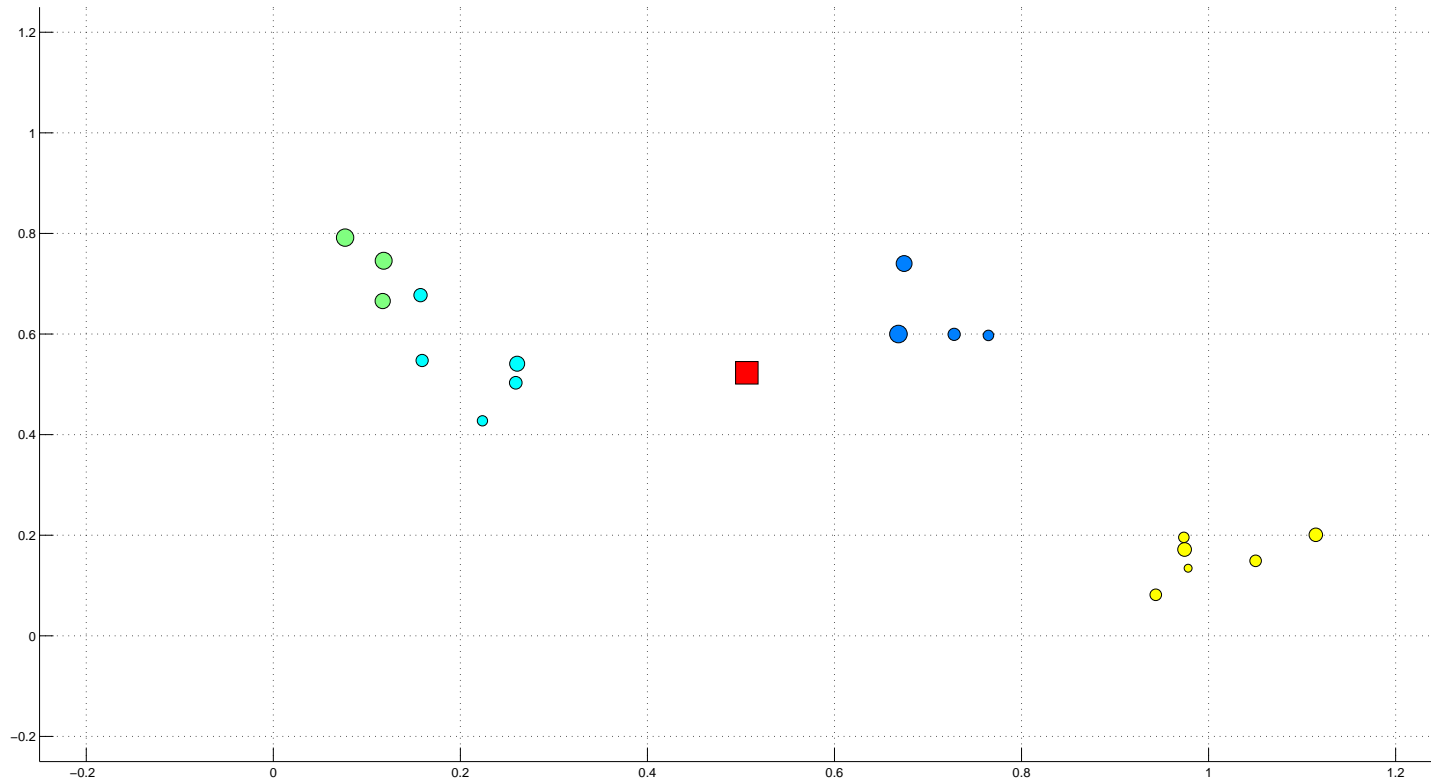


# to Probability Measures



Assume that each datum is now an **empirical measure**.  
What could be the mean of these 4 measures?

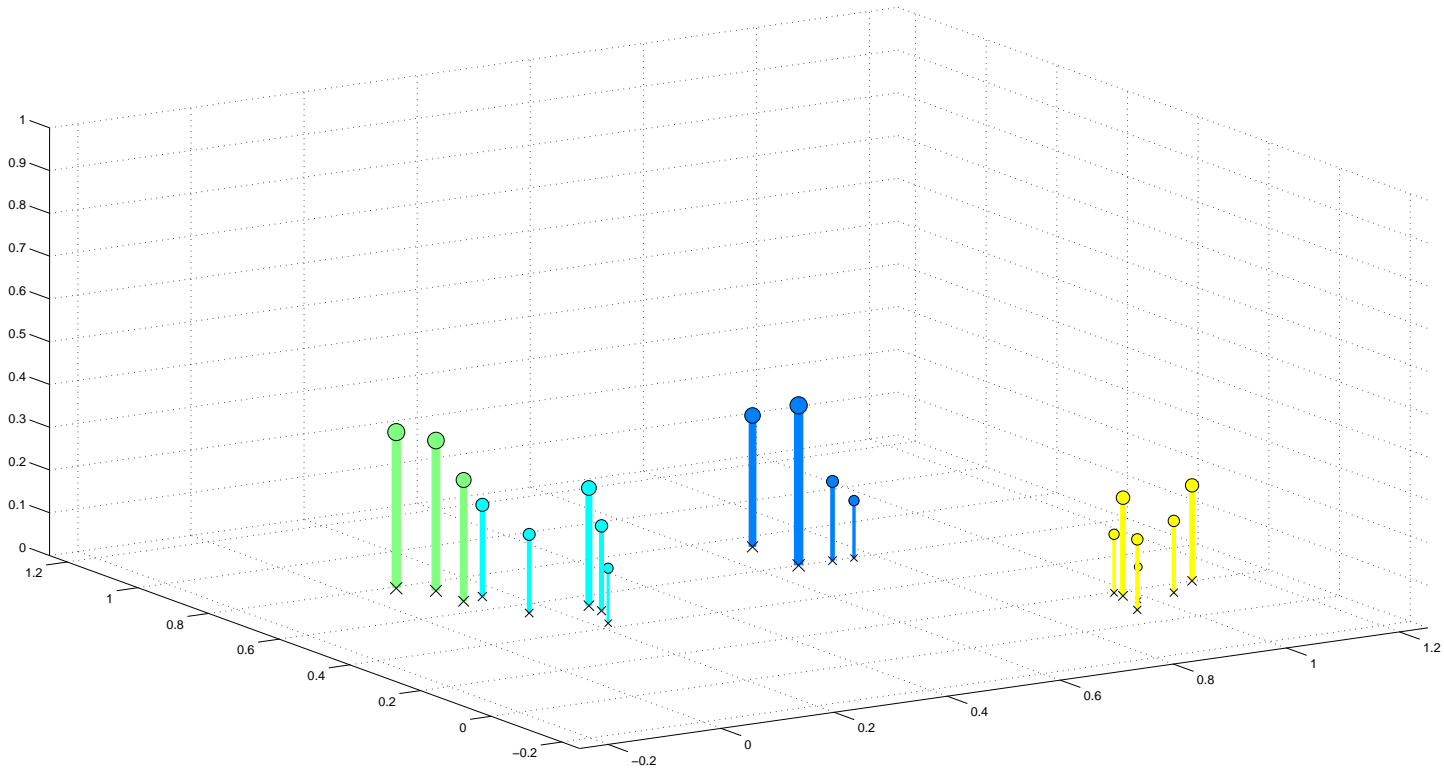
# 1. Naive Averaging



■ = naive mean of *all* observations.

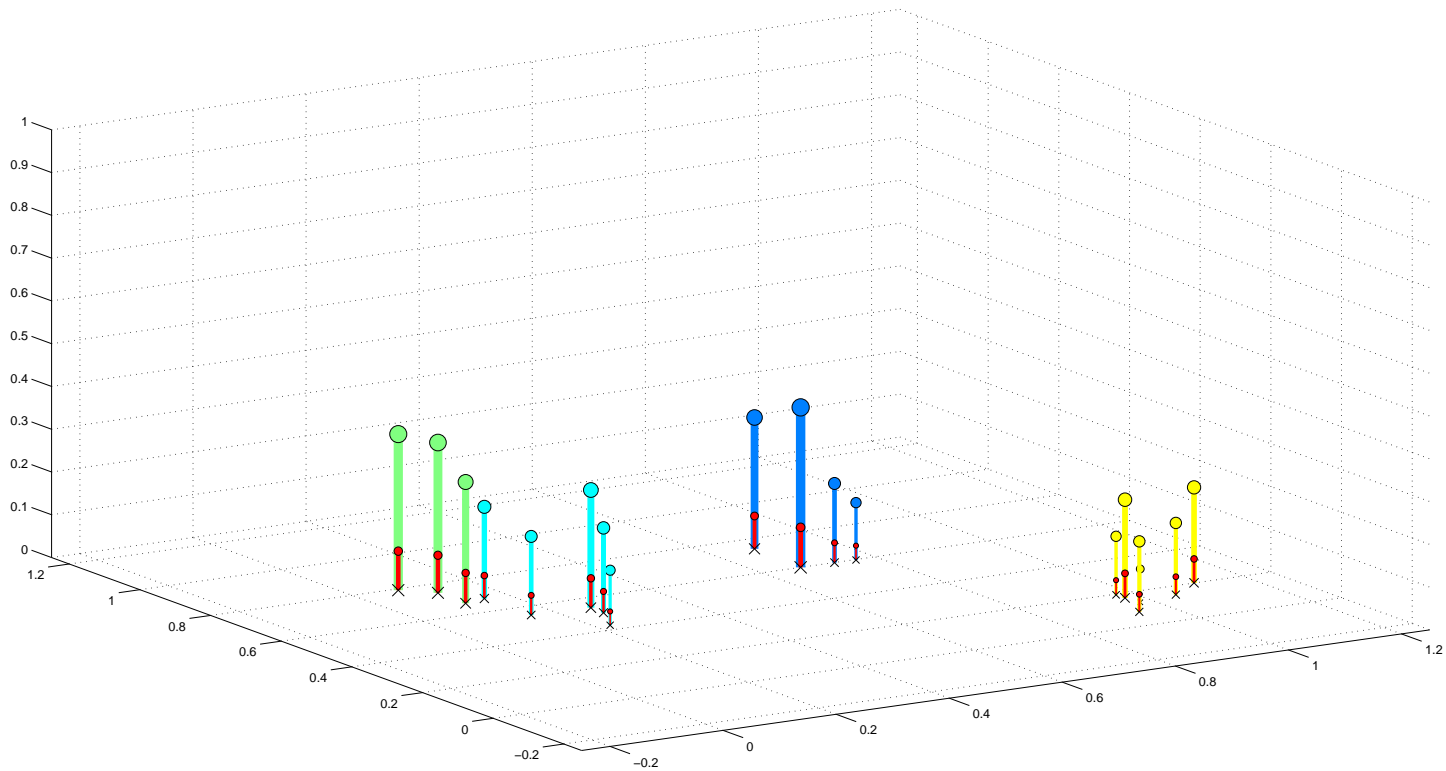
**Mean of 4 measures = a point?**

# Averaging Probability Measures



Same measures, in a 3D perspective.

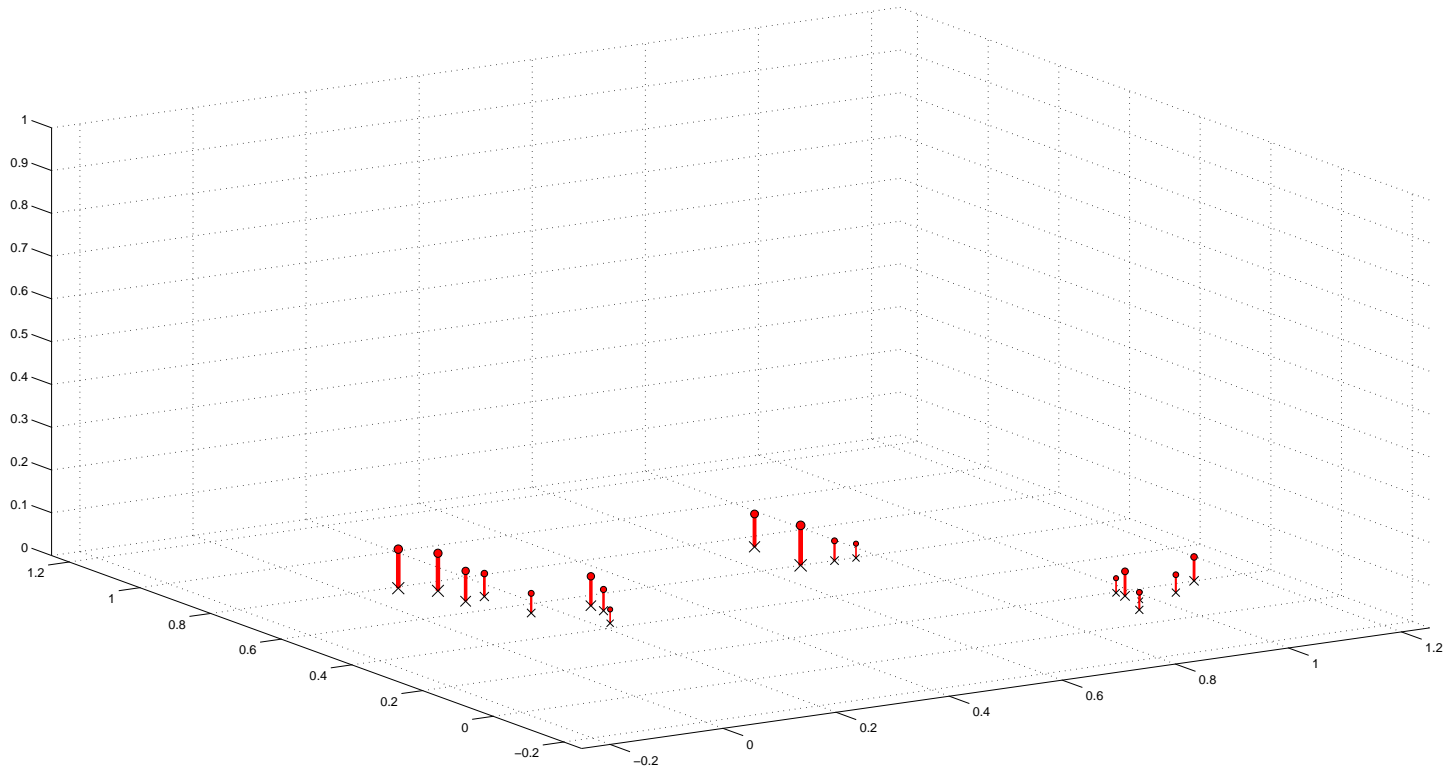
## 2. Naive Averaging



**Euclidean mean** of measures is their sum /  $N$ .

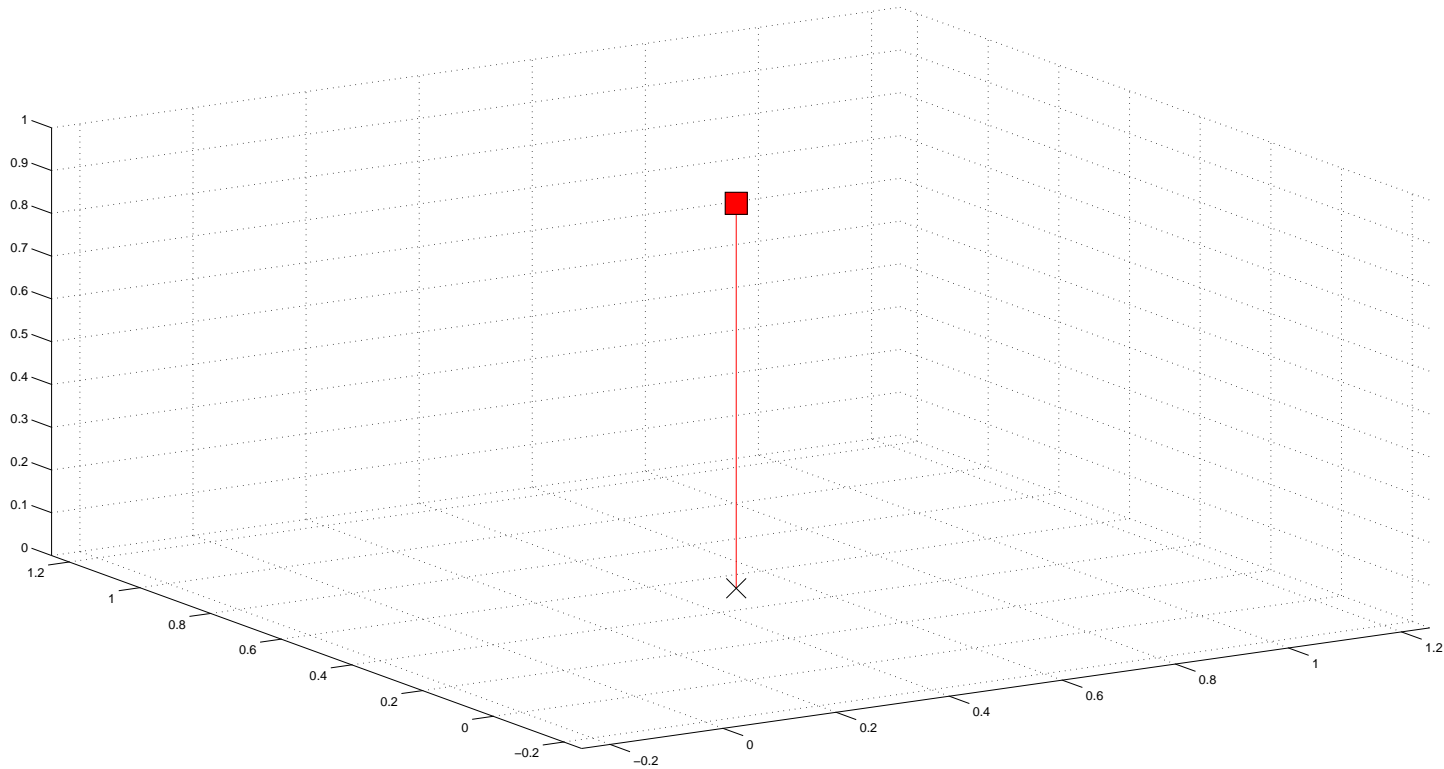
$$\text{Here, } \Delta(\mu, \nu) = \int_{\mathbb{R}^2} [d\mu - d\nu]^2.$$

# Focus on uncertainty



...but geometric knowledge ignored.

# Focus on geometry



...but uncertainty is lost.



# Problem of interest

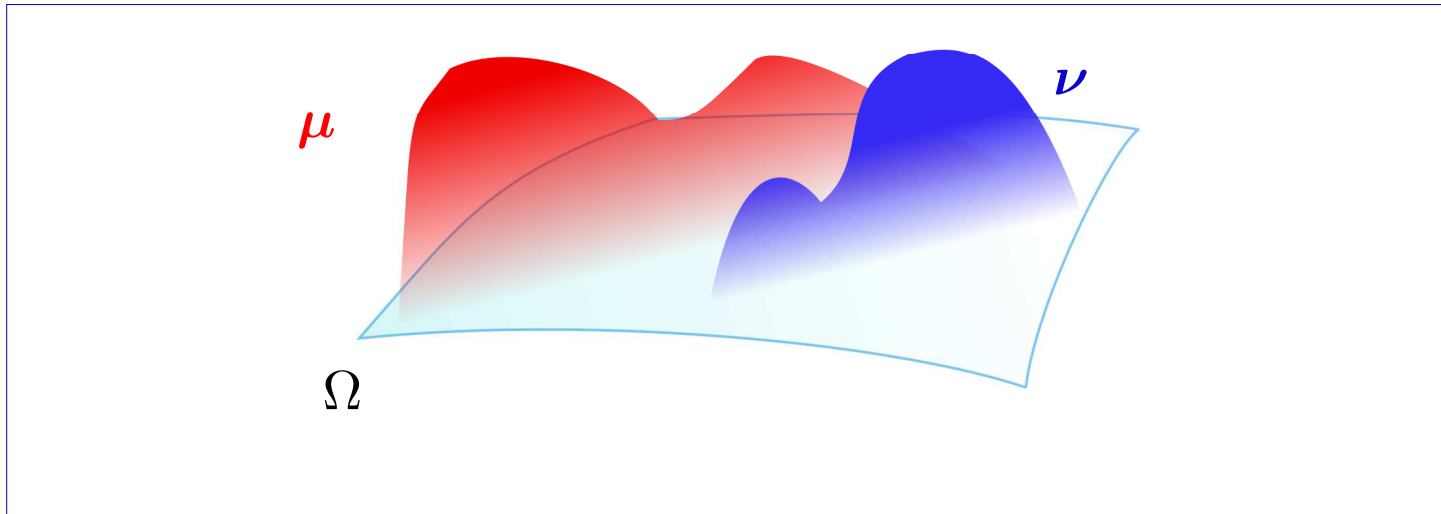
Given a discrepancy function  $\Delta$  between probabilities, compute **their mean**:  $\operatorname{argmin} \sum_i \Delta(\cdot, \nu_i)$

- The idea is useful, sometimes tractable & appears in
  - Bregman clustering for histograms [**Banerjee'05**].
  - Topic modeling [**Blei & al.'03**].
  - Clustering problems ( $k$ -means).
- Our goal in this talk: study the case  $\Delta = \text{Wasserstein}$

---

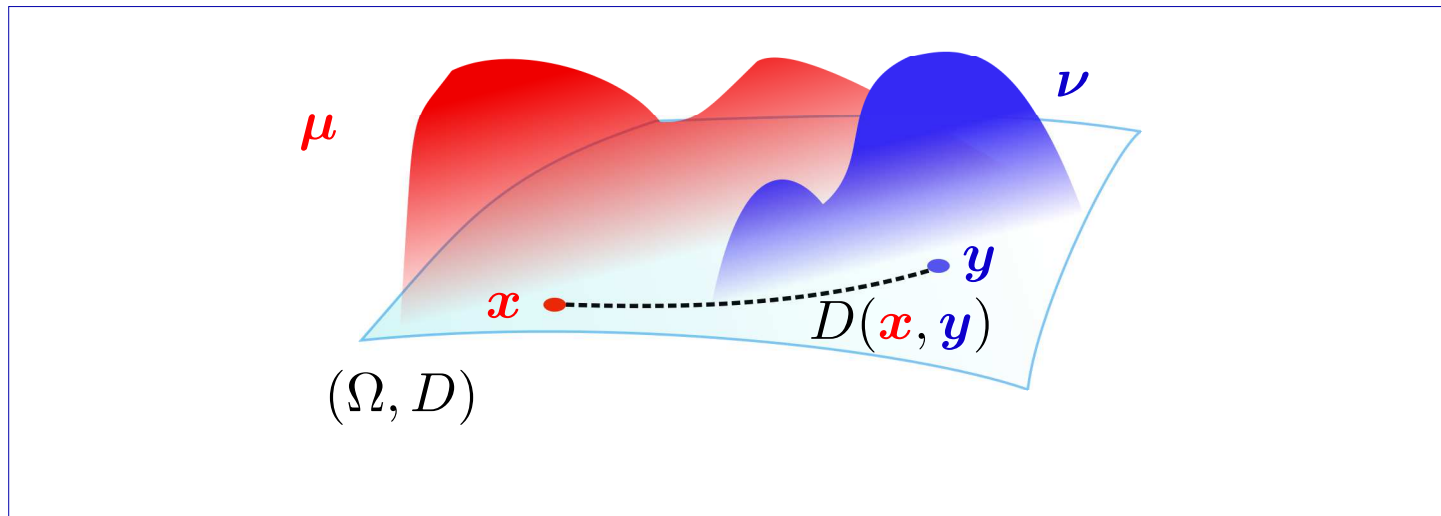
# Wasserstein Distances

# Comparing Two Measures



Two measures  $\mu, \nu \in P(\Omega)$ .

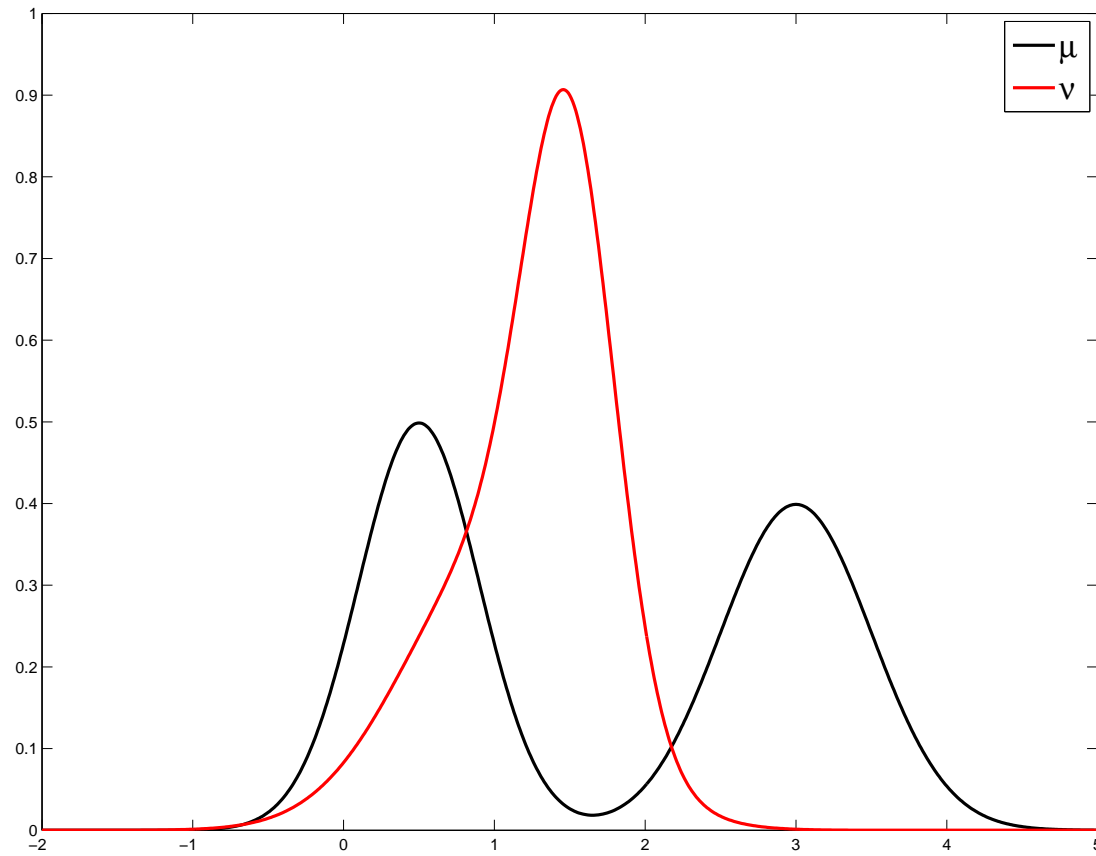
# The Optimal Transport Approach



Optimal Transport distances rely **on 2 key concepts**:

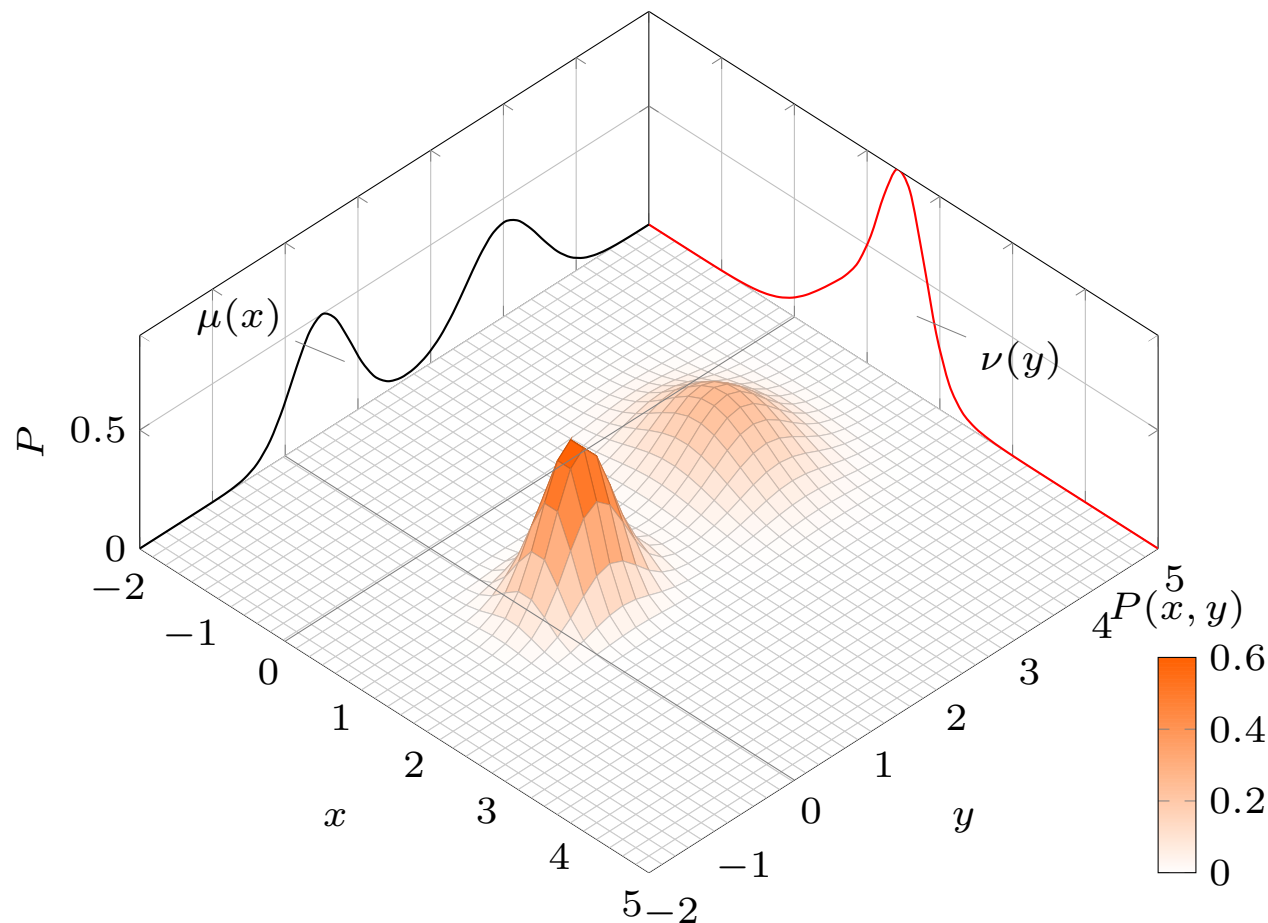
- A **metric**  $D : \Omega \times \Omega \rightarrow \mathbb{R}_+$  ;
- $\Pi(\mu, \nu)$ : **joint probabilities** with marginals  $\mu, \nu$ .

# Joint Probabilities of $(\mu, \nu)$



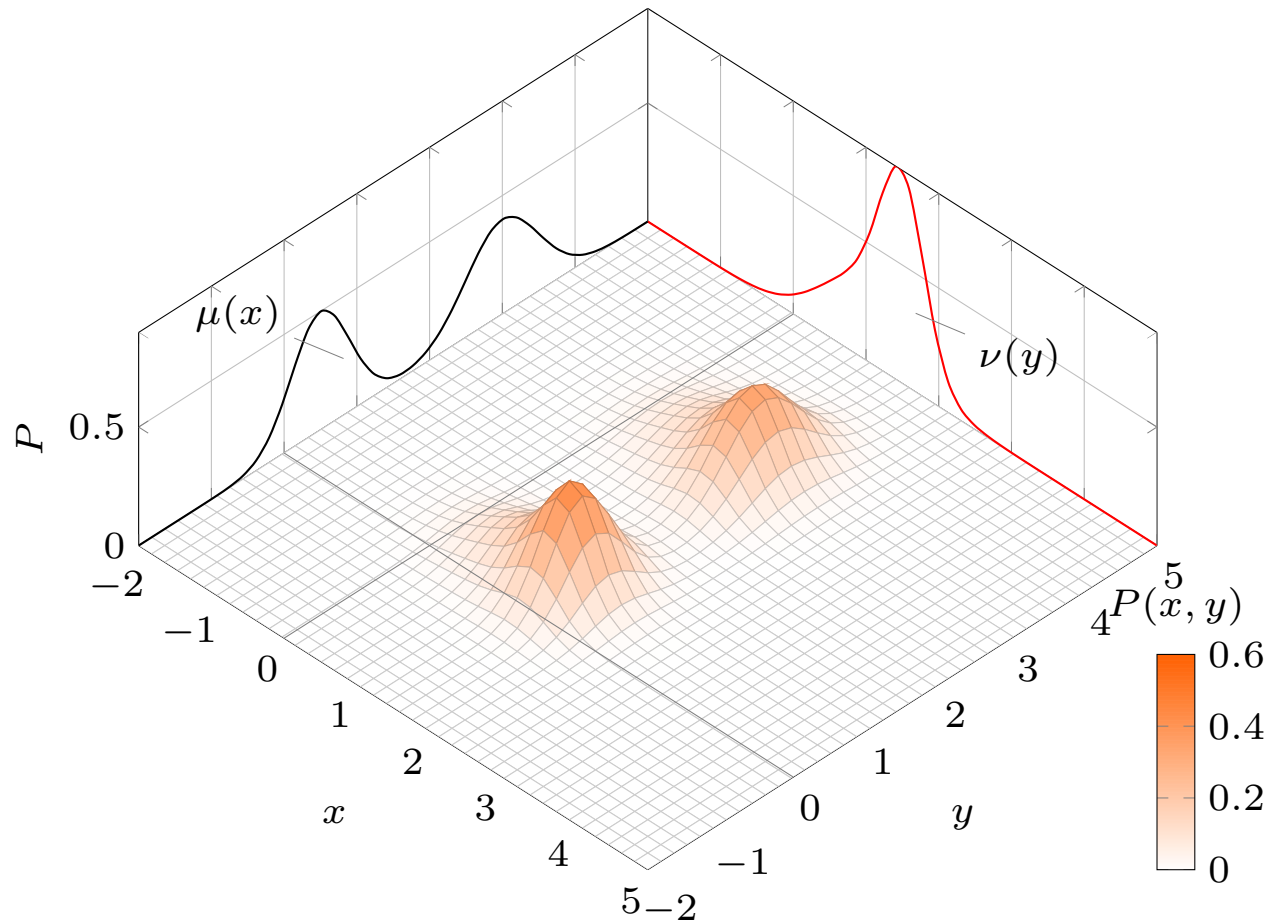
Consider  $\mu, \nu$  two measures on the real line.

# Joint Probabilities of $(\mu, \nu)$



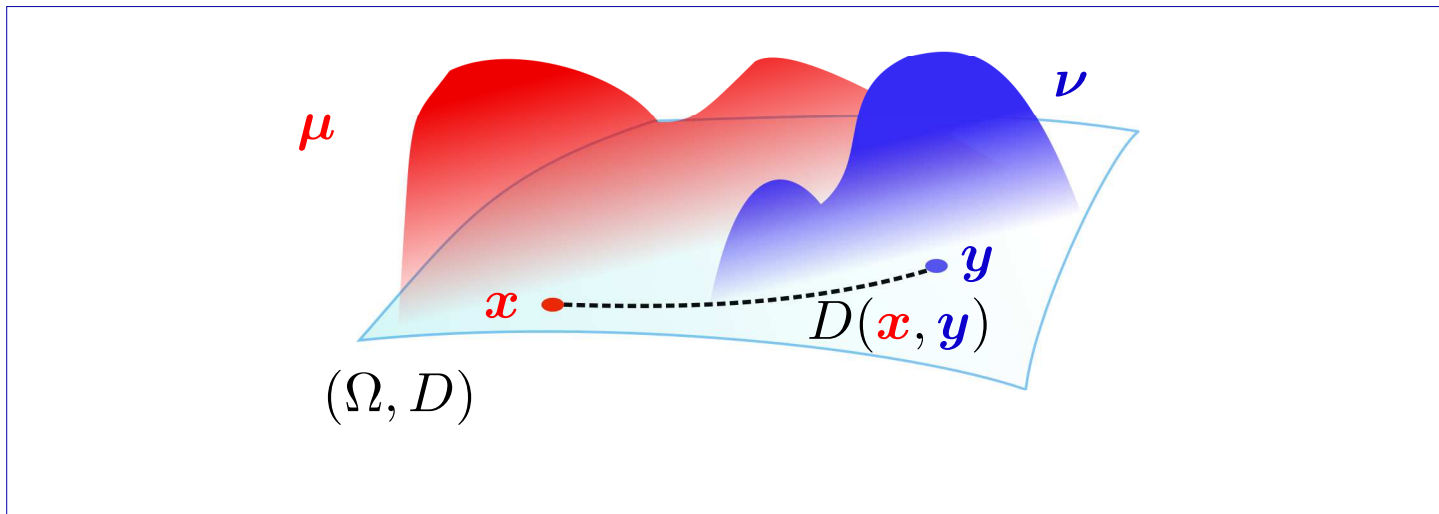
$\Pi(\mu, \nu)$  = probability measures on  $\Omega^2$   
with marginals  $\mu$  and  $\nu$ .

# Joint Probabilities of $(\mu, \nu)$



$\Pi(\mu, \nu)$  = probability measures on  $\Omega^2$   
with marginals  $\mu$  and  $\nu$ .

# Optimal Transport Distance



$p$ -Wasserstein (or OT) distance, assuming  $p \geq 1$ , is:

$$W_p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \left( \inf_{P \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \mathbb{E}_P[D(X, Y)^p] \right)^{1/p}.$$

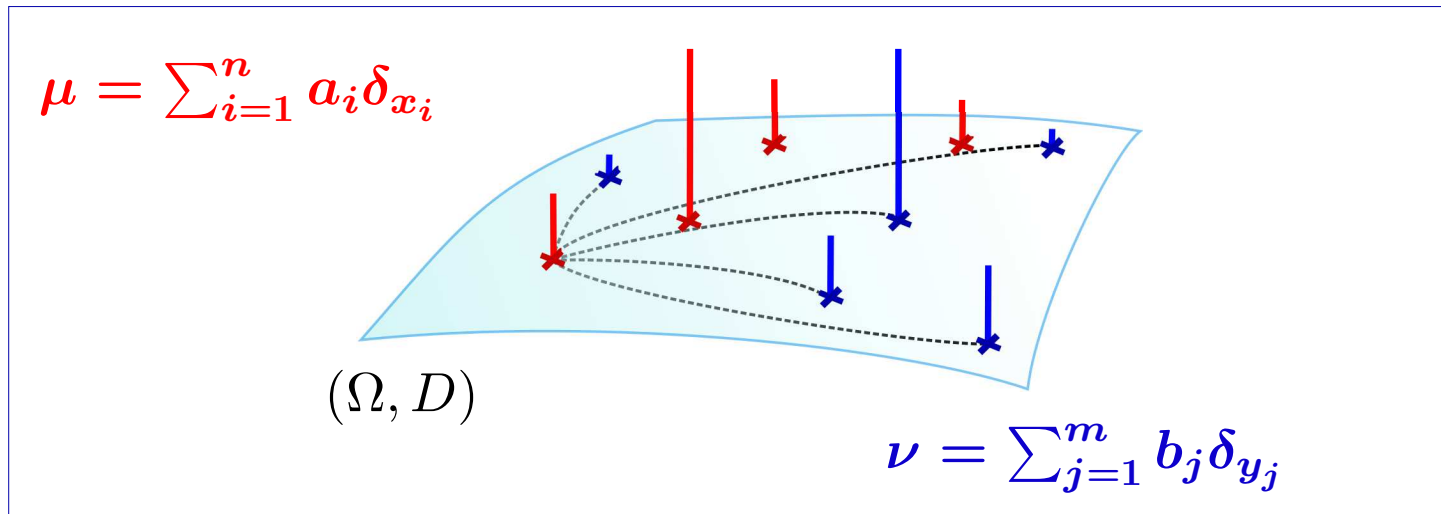


# (Historical Parenthesis)

Monge-Kantorovich, Kantorovich-Rubinstein, Wasserstein, Earth Mover's Distance, Mallows

- Monge 1781 *Mémoire sur la théorie des déblais et des remblais*
- **Optimization & Operations Research**
  - Kantorovich'42, Dantzig'47, Ford Fulkerson'55, *etc.*
- **Probability & Statistical Physics**
  - Rachev'92, Talagrand'96, Villani'09
- **Computer Vision**: Rubner et al'98

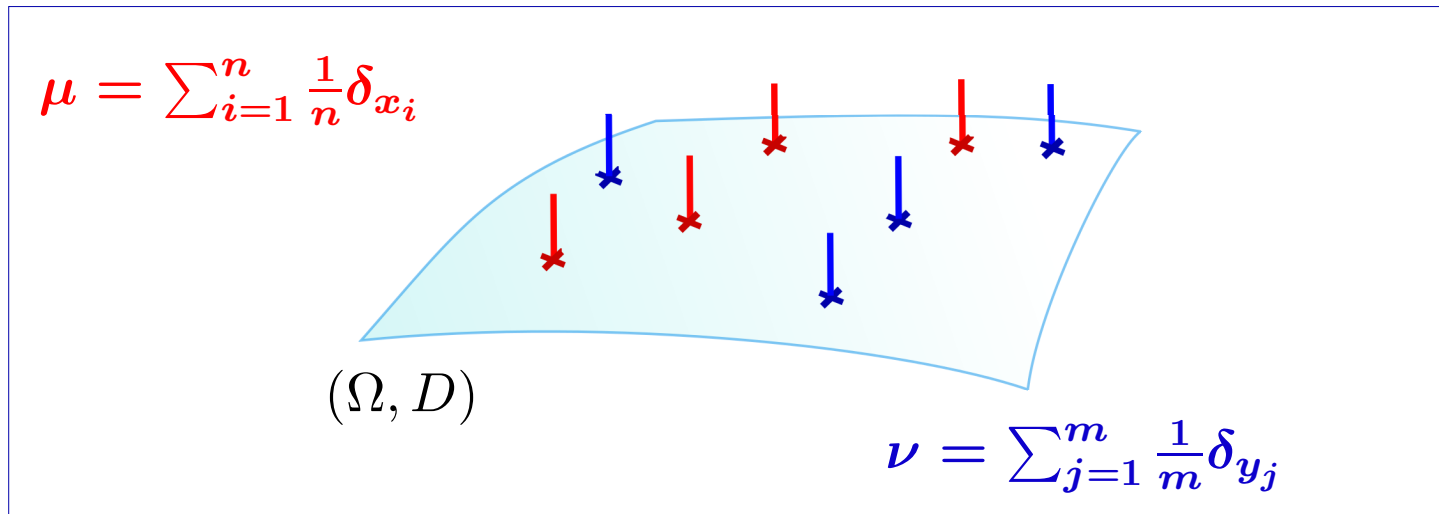
# OT Distance for Empirical Measures



$$W_p(\mu, \nu) = \left( \inf_{P \in \Pi(\mu, \nu)} \mathbb{E}_P[D(X, Y)^p] \right)^{1/p}.$$

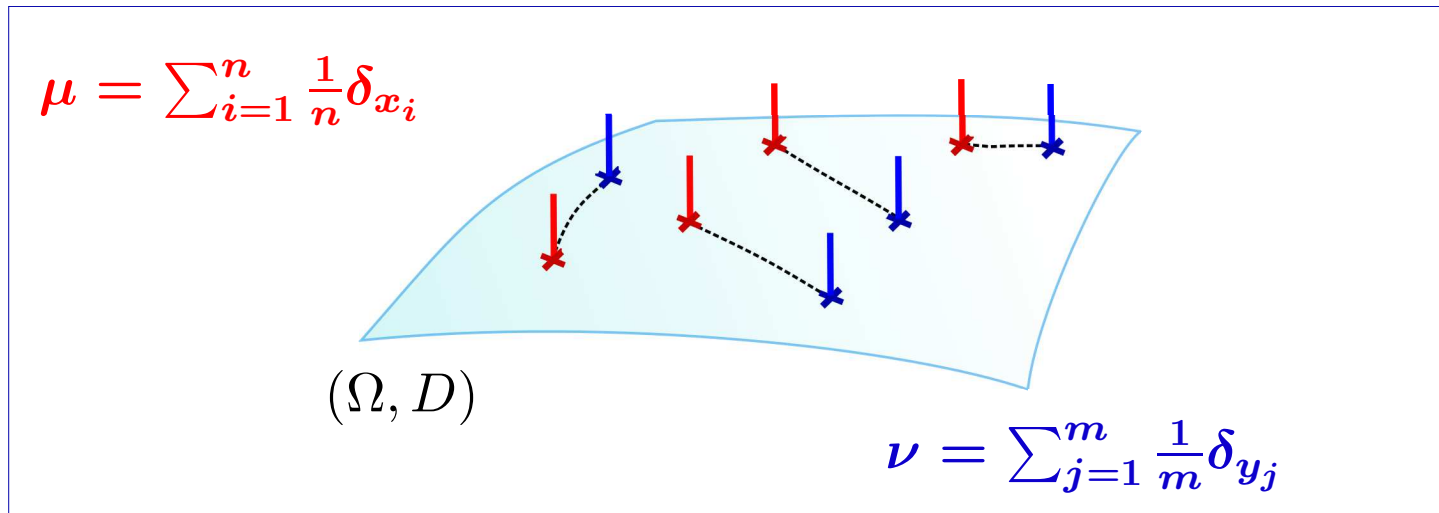
Algorithmically?

# OT Distance for Empirical Measures



Suppose  $n = m$  and all weights are uniform

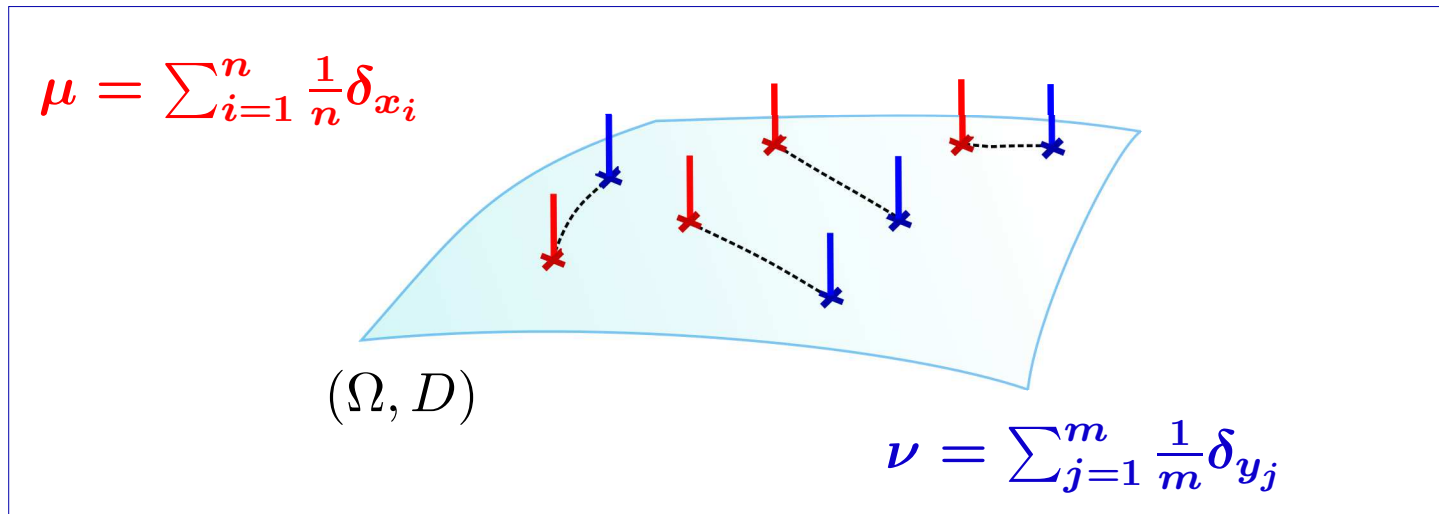
# OT Distance for Empirical Measures



Then  $W_p^p = \text{optimal matching cost}$   
(solved for instance with Hungarian algorithm)

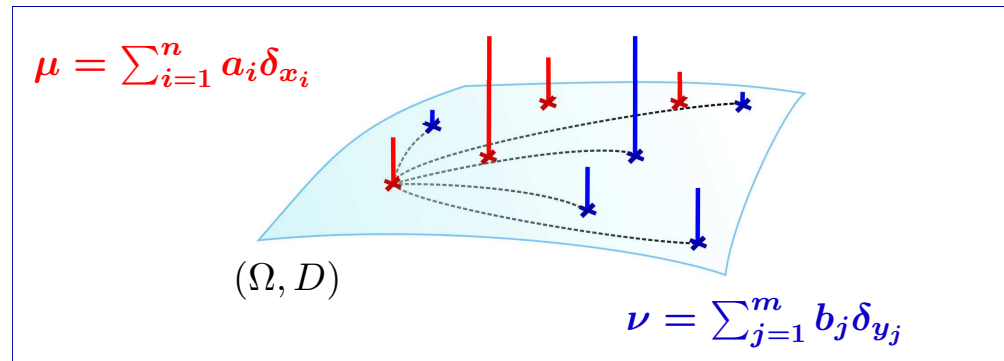
$$\left( \min_{\sigma \in S_n} \frac{1}{n} \sum_{i=1}^n D(x_i, y_{\sigma_i})^p \right)^{1/p}$$

# OT Distance for Empirical Measures



As soon as  $n \neq m$  or weights are non uniform, optimal matching does not make sense.

# Computing the OT Distance



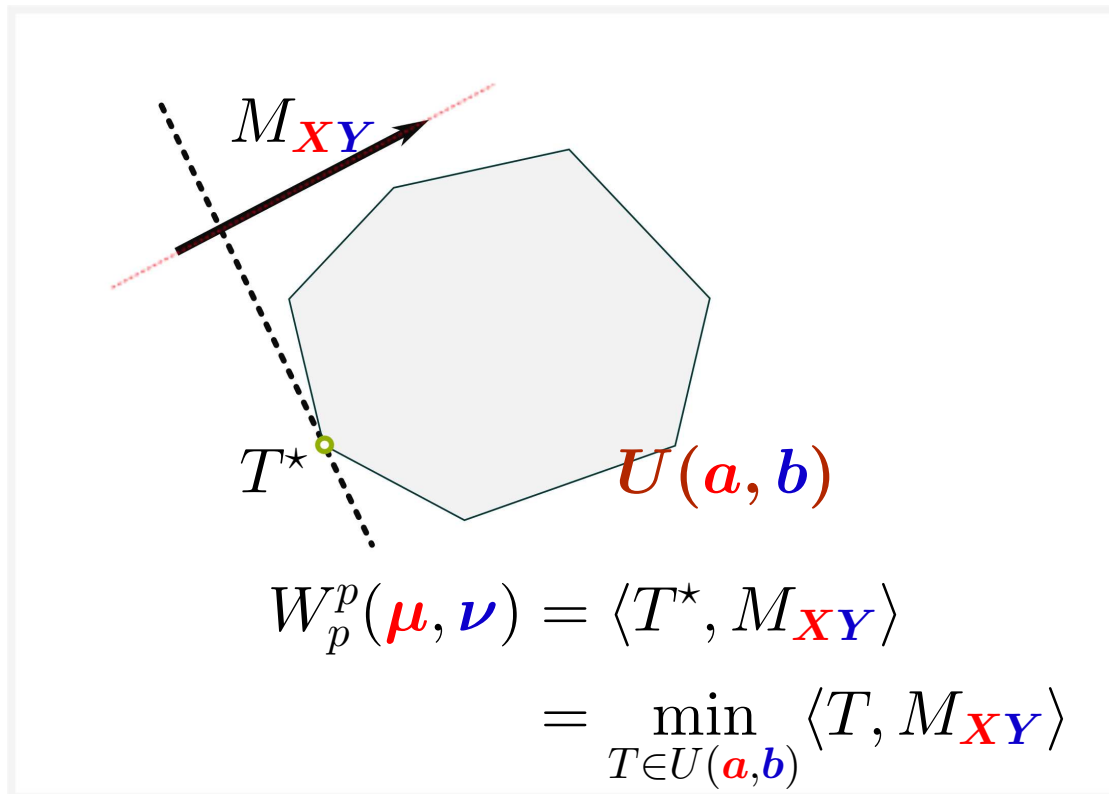
$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu})$  can be cast as a linear program in  $\mathbb{R}^{n \times m}$ :

1.  $M_{\mathbf{X}\mathbf{Y}} \stackrel{\text{def}}{=} [D(\mathbf{x}_i, \mathbf{y}_j)^p]_{ij} \in \mathbb{R}^{n \times m}$  (*metric information*)
2. Transportation Polytope (*joint probabilities*)

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

# Computing $p$ -Wasserstein Distances

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \text{primal}(\mathbf{a}, \mathbf{b}, M_{\mathbf{XY}}) \stackrel{\text{def}}{=} \min_{T \in U(\mathbf{a}, \mathbf{b})} \langle T, M_{\mathbf{XY}} \rangle$$



# [Kantorovich'42] Duality

- This primal problem has an equivalent, dual LP:

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \left\{ \begin{array}{l} \text{primal}(\mathbf{a}, \mathbf{b}, M_{\mathbf{XY}}) \stackrel{\text{def}}{=} \min_{T \in U(\mathbf{a}, \mathbf{b})} \langle T, M_{\mathbf{XY}} \rangle \\ \text{or} \\ \text{dual}(\mathbf{a}, \mathbf{b}, M_{\mathbf{XY}}) \stackrel{\text{def}}{=} \max_{(\alpha, \beta) \in C_M} \alpha^T \mathbf{a} + \beta^T \mathbf{b}, \\ \text{where } C_M = \{(\alpha, \beta) \in \mathbb{R}^{n+m} \mid \alpha_i + \beta_j \leq M_{ij}\} \end{array} \right.$$



# [Kantorovich'42] Duality

- This primal problem has an equivalent, dual LP:

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \left\{ \begin{array}{l} \text{primal}(\mathbf{a}, \mathbf{b}, M_{\mathbf{XY}}) \stackrel{\text{def}}{=} \min_{T \in U(\mathbf{a}, \mathbf{b})} \langle T, M_{\mathbf{XY}} \rangle \\ \text{or} \\ \text{dual}(\mathbf{a}, \mathbf{b}, M_{\mathbf{XY}}) \stackrel{\text{def}}{=} \max_{(\alpha, \beta) \in C_{M_{\mathbf{XY}}}} \alpha^T \mathbf{a} + \beta^T \mathbf{b}, \\ \text{where } C_M = \{(\alpha, \beta) \in \mathbb{R}^{n+m} \mid \alpha_i + \beta_j \leq M_{ij}\} \end{array} \right.$$

**⚠ Both problems require  $O(n^3 \log(n))$  operations.  
Typically solved using the network simplex.**

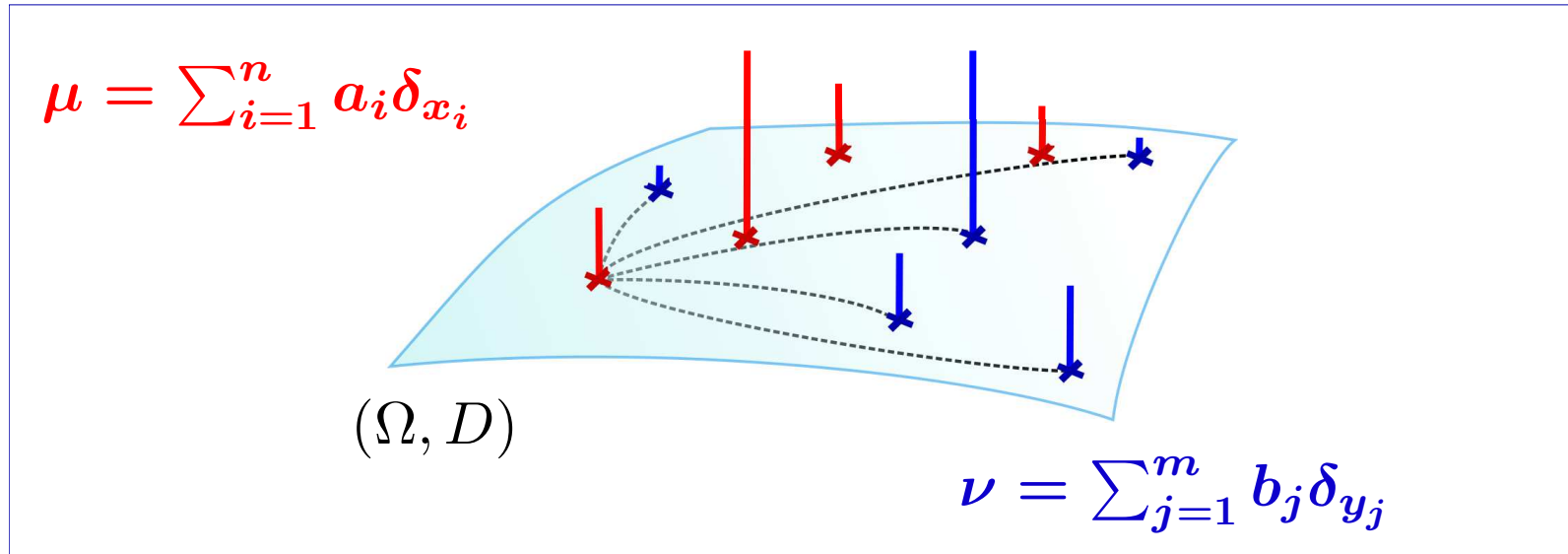
# Wasserstein Barycenter Problem (WBP)

- [Agueh'11] introduced the WBP:

$$\operatorname{argmin}_{\mu \in P(\Omega)} C(\boldsymbol{\mu}) \stackrel{\text{def}}{=} \sum_{i=1}^N W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}_i),$$

- Can be solved with a **multi-marginal** OT problem.
- **Intractable**: LP of  $\prod_i \text{card}(\text{supp}(\nu_i))$  variables.

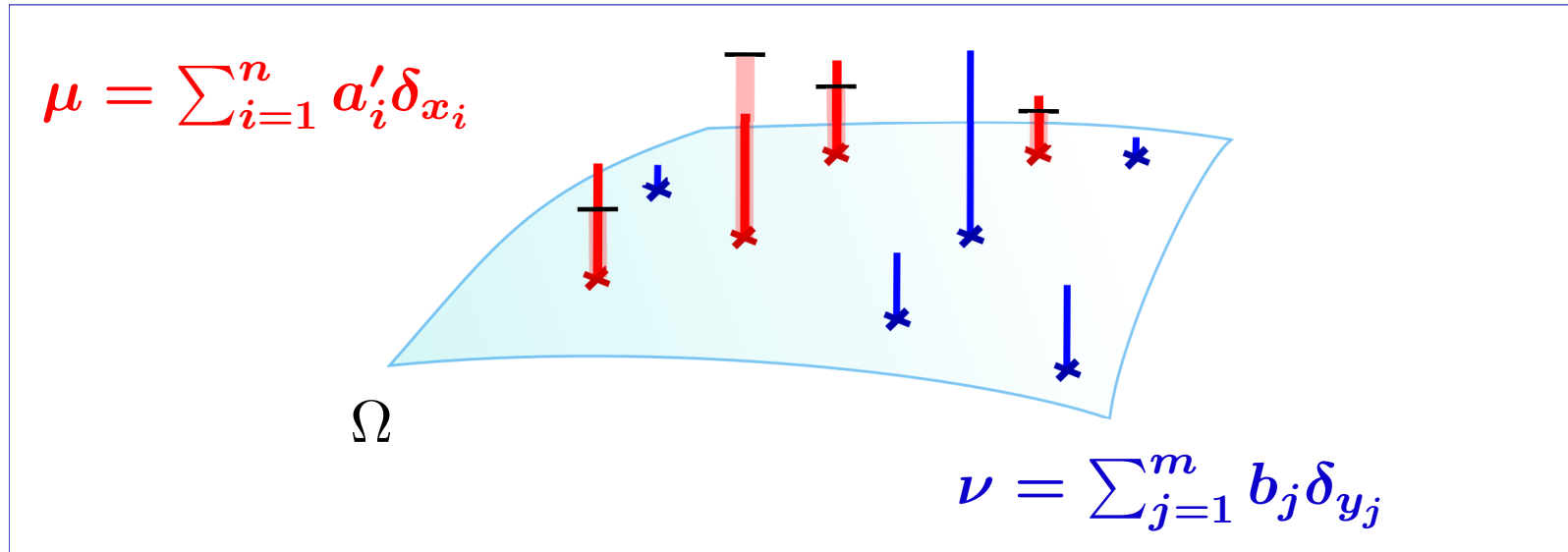
# Differentiability w.r.t. $X$ or $a$



To solve it **numerically**, we must understand how

$f_{\nu}(\mathbf{a}, \mathbf{X}) \stackrel{\text{def}}{=} W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu})$  varies when  $\mathbf{a}$  &  $\mathbf{X}$  varies.

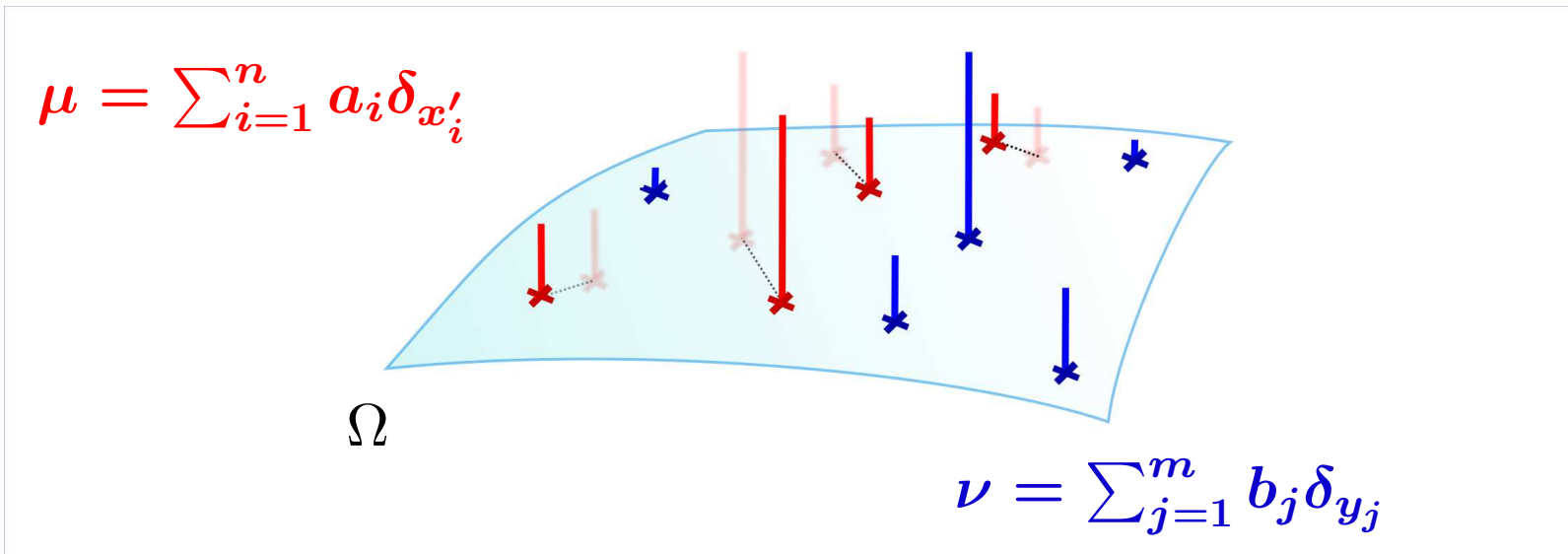
# Differentiability w.r.t. $X$ or $a$



## 1. Infinitesimal Variation in Weights

$$f_{\nu}(a', X)?, \quad \text{if } a' \approx a$$

# Differentiability w.r.t. $X$ or $a$

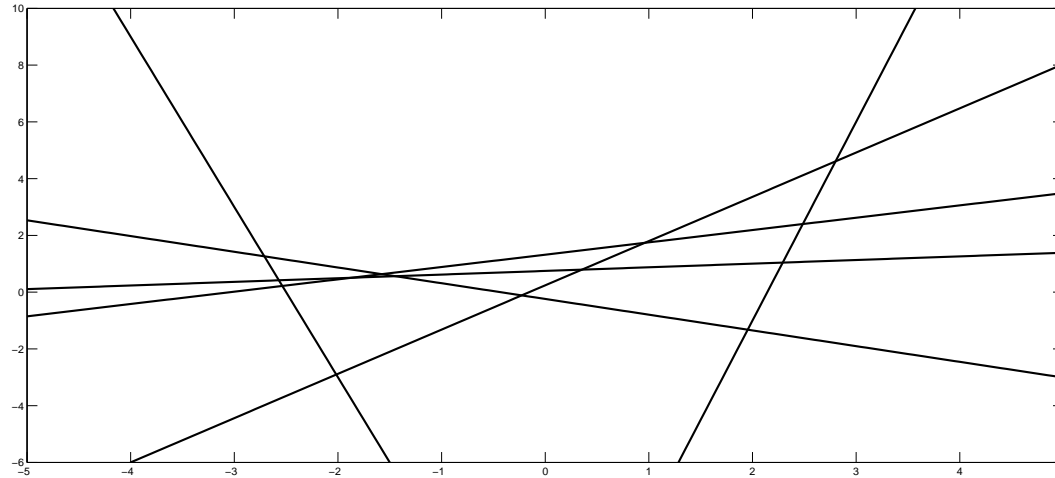


## 2. Infinitesimal Variation in Locations

$$f_{\nu}(a, X')?, \quad \text{if } X' \approx X$$

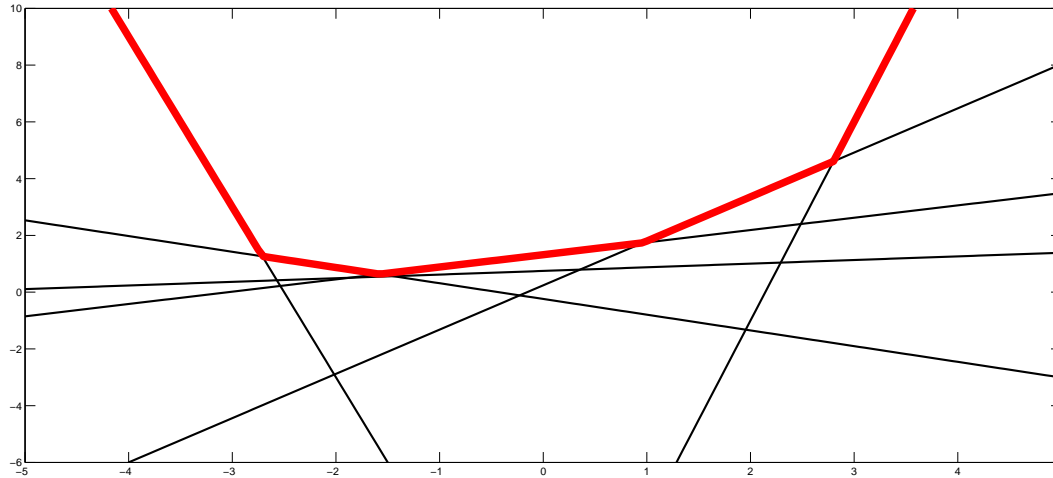
# Using the dual, $\partial|_a$

$$f_{\nu}(\mathbf{a}, X) = \max_{(\alpha, \beta) \in C_{MXY}} \alpha^T \mathbf{a} + \beta^T \mathbf{b}$$



# Using the dual, $\partial|_a$

$$f_\nu(\mathbf{a}, X) = \max_{(\alpha, \beta) \in C_{M \times Y}} \alpha^T \mathbf{a} + \beta^T \mathbf{b}$$



$\mathbf{a} \mapsto f_\nu(\mathbf{a}, X)$  is a **convex non-smooth** map.  
The *dual optimum*  $\alpha^*$  is a subgradient  $f_\nu(\mathbf{a}, X)$ .

## Using the primal $\partial|_X$

$$f_\nu(a, \mathbf{X}) = \min_{T \in U(\mathbf{a}, \mathbf{b})} \langle T, M_{\mathbf{X}\mathbf{Y}} \rangle$$

- More involved computations. Tractable when  $D = \text{Euclidean}$ ,  $p = 2$ .
- Convex quadratic + piecewise linear concave of  $\mathbf{X}$
- $\partial f_\nu|_X = \mathbf{Y}\mathbf{T}^{*T} \text{diag}(a^{-1})$ : *optimal transport*  $\mathbf{T}^{*T}$  yields a subgradient.



To sum up: (1) the WBP is challenging

$$C(\mathbf{a}, \mathbf{X}) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}_i) = \frac{1}{N} \sum_{i=1}^N f_{\boldsymbol{\nu}_i}(\mathbf{a}, \mathbf{X})$$

- $\mathbf{a} \rightarrow C(\mathbf{a}, \mathbf{X})$  is **convex**, **non-smooth**, **computing one subgradient requires solving  $N$  OT problems!**
- $\mathbf{X} \rightarrow C(\mathbf{a}, \mathbf{X})$  is **not convex**, **non-smooth**

## (2) the WBP is unstable

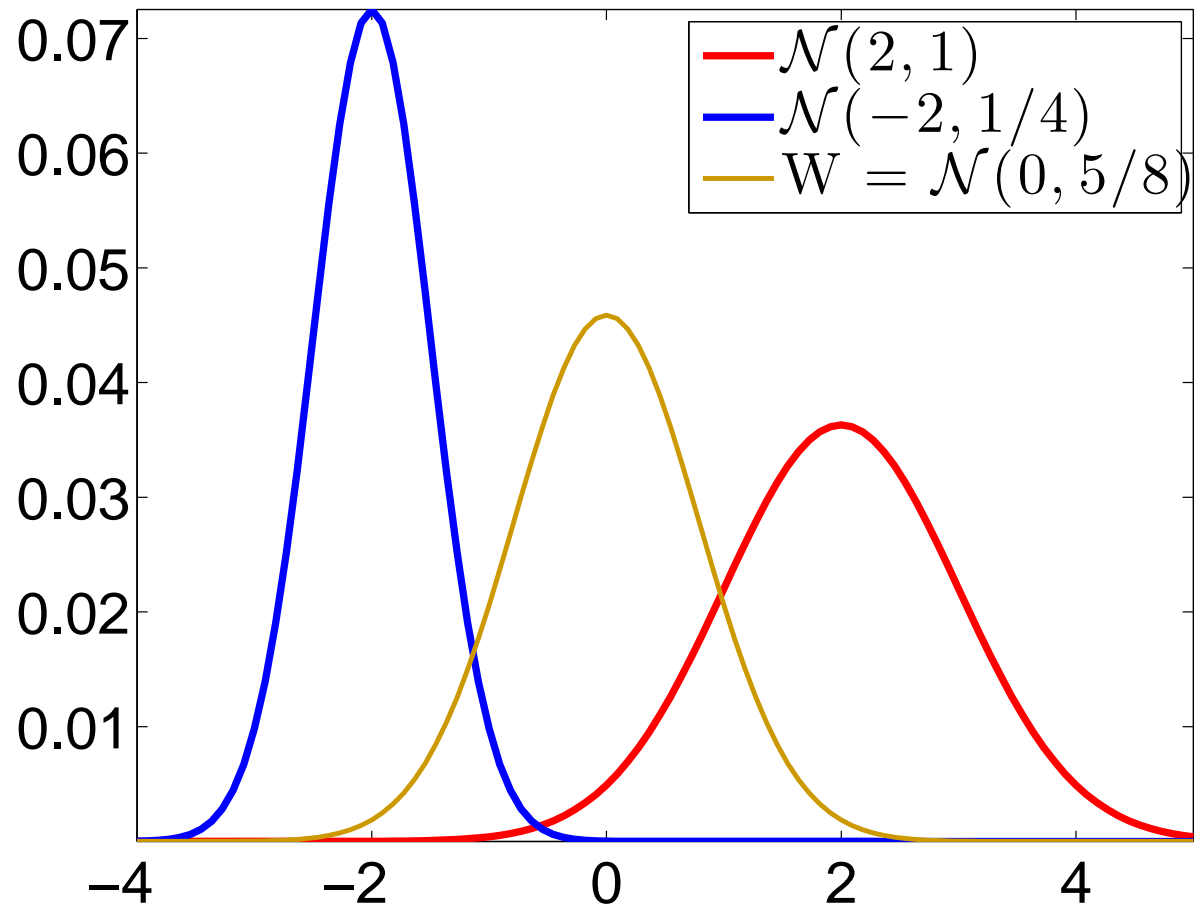
- Assume  $X = Y_1 = \dots = Y_N$  (fixed grid).

$$C(\mathbf{a}) = \frac{1}{N} \sum_{i=1}^N \text{primal}(\mathbf{a}, \mathbf{b}_i, M) = \frac{1}{N} \sum_{i=1}^N \min_{T_i \in U(\mathbf{a}, \mathbf{b}_i)} \langle T_i, M \rangle$$

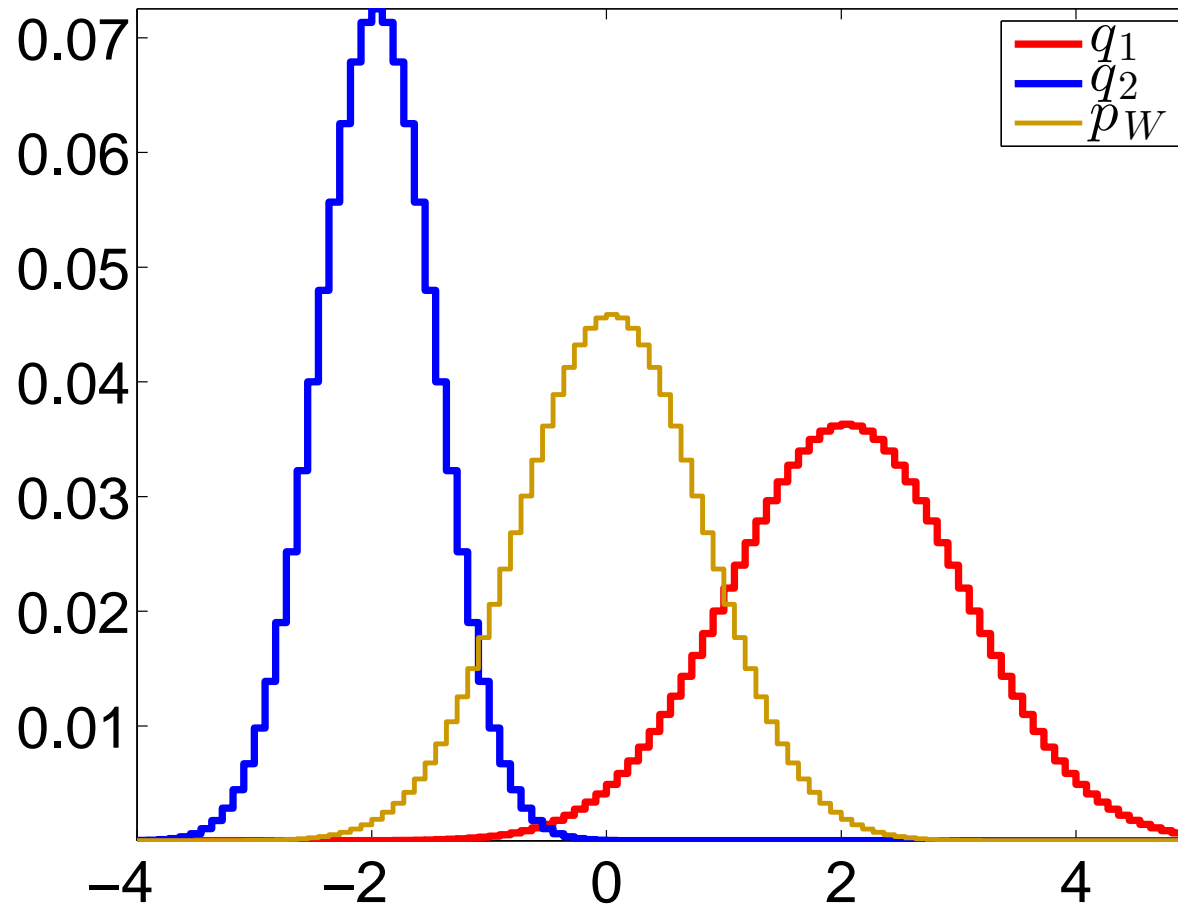
- In that case, the WBP can be solved as a large LP:

$$\begin{aligned} & \min_{T_1, \dots, T_N, \mathbf{a}} \sum_{i=1}^N \langle T_i, M \rangle \\ & \text{s.t. } T_i^T \mathbf{1}_d = \mathbf{b}_i, \forall i \leq N, \\ & T_1 \mathbf{1}_d = \dots = T_N \mathbf{1}_d = \mathbf{a}. \end{aligned}$$

# Averaging Two Gaussians

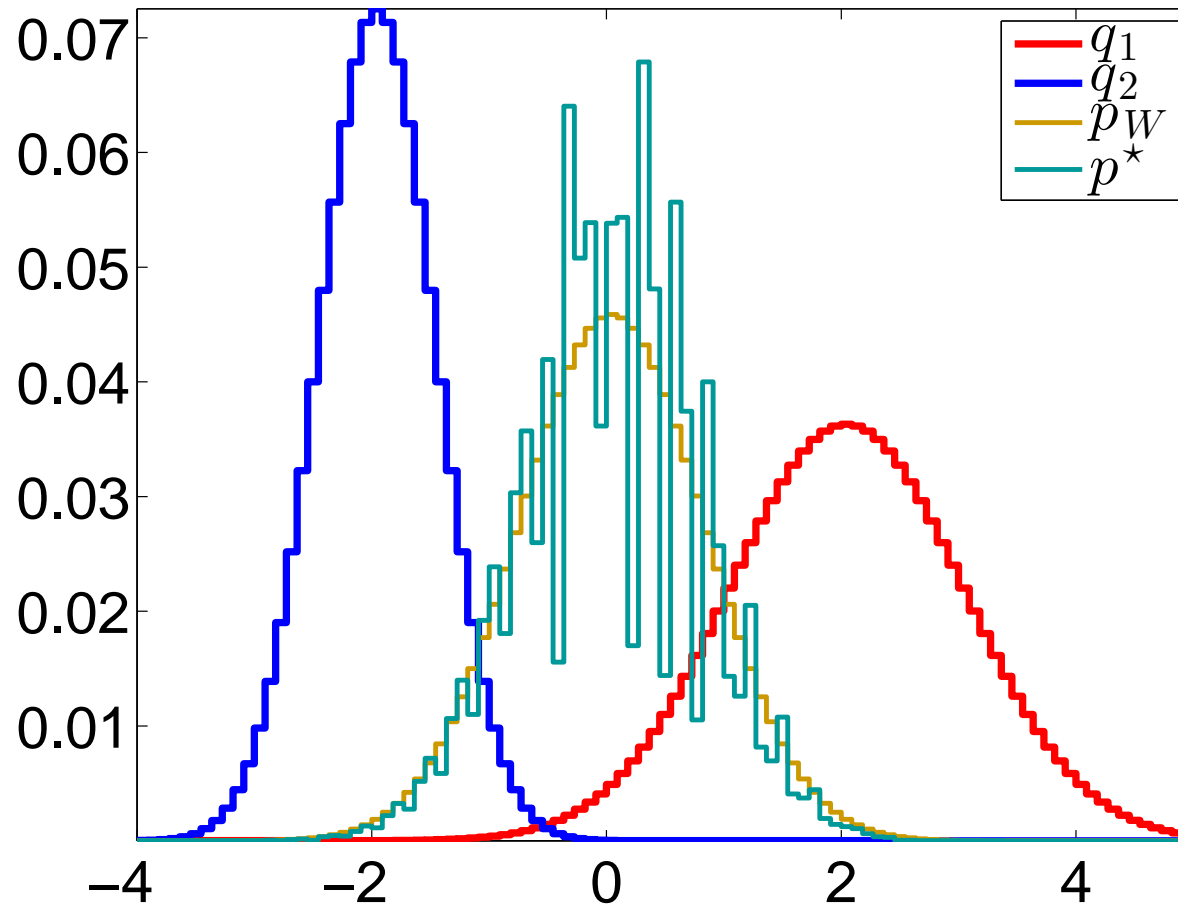


# Discretized



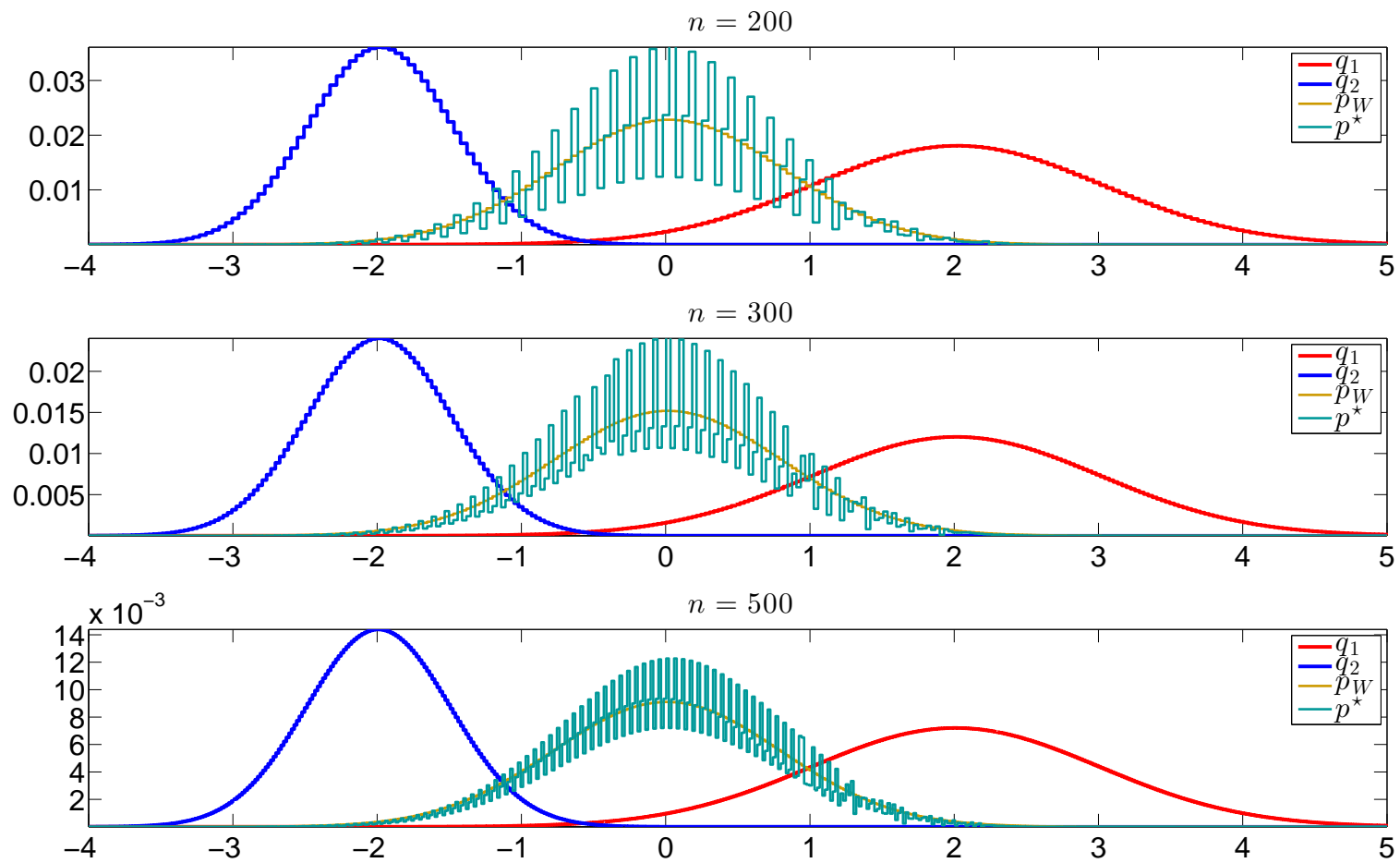
$p_W$  is the discrete equivalent of the true barycenter.

# Exact Solution



$p^*$  is the solution to that LP

# Does not get much better with large $n$ ...



---

# Entropic Smoothing of OT

# Smoothing solves (almost) everything

**Original** OT primal:

$$\text{primal}(a, b, M_{XY}) = \min_{T \in U(a,b)} \langle T, M_{XY} \rangle$$

**Original** OT Kantorovich dual:

$$\text{dual}(a, b, M_{XY}) = \max_{(\alpha, \beta), \alpha_i + \beta_j \leq M_{ij}} \alpha^T a + \beta^T b$$



# Smoothing solves (almost) everything

**Entropy-smoothed** ( $\gamma > 0$ ) primal problem:

$$\text{primal}_\gamma(a, b, M_{XY}) = \min_{\mathbf{T} \in U(a, b)} \langle \mathbf{T}, M_{XY} \rangle - \gamma H(\mathbf{T})$$

**Smoothed** dual problem:

$$\text{dual}_\gamma(a, b, M_{XY}) = \max_{(\boldsymbol{\alpha}, \boldsymbol{\beta})} \boldsymbol{\alpha}^T a + \boldsymbol{\beta}^T b - \gamma \sum_{i \leq n, j \leq m} e^{-(M_{ij} - \alpha_i - \beta_j) / \gamma}$$

# Smoothing solves (almost) everything

**Entropy-smoothed** ( $\gamma > 0$ ) primal problem:

$$\text{primal}_{\gamma}(a, b, M_{XY}) = \min_{T \in U(a,b)} \mathbf{KL}(T \| e^{-M_{XY}/\gamma})$$

**Smoothed** dual problem:

$$\text{dual}_{\gamma}(a, b, M_{XY}) = \max_{(\alpha, \beta)} \alpha^T a + \beta^T b - \gamma \sum_{i \leq n, j \leq m} e^{-(M_{ij} - \alpha_i - \beta_j)/\gamma}$$

# Why is entropy a good regularizer for OT?

The penalized problem

$$T_\gamma = \operatorname{argmin}_{T \in U(r,c)} \langle P, M \rangle - \gamma H(T)$$

implies that  $T_\gamma$  has the form: (first order cond.)

$$\exists \mathbf{u} \in \mathbb{R}_n^+, \mathbf{v} \in \mathbb{R}_m^+ \mid T_\gamma = \mathbf{diag}(\mathbf{u}) e^{-M/\gamma} \mathbf{diag}(\mathbf{v}).$$

**Gravity Model in Transportation [Wilson'69]**  
**Schrödinger Problem ['32]**

# Sinkhorn - Matrix Scaling

**Theorem 1** (Sinkhorn'62). *For any  $n \times m$  matrix  $A$  with positive entries, any  $\mathbf{r}$  and  $\mathbf{c}$  in the simplex,  $\exists! \mathbf{u} \in \mathbb{R}_n^+, \mathbf{v} \in \mathbb{R}_m^+$  such that*

$$\begin{bmatrix} \mathbf{u}_1 & 0 & \dots & 0 \\ 0 & \mathbf{u}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{u}_n \end{bmatrix} \begin{bmatrix} A \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 & 0 & \dots & 0 \\ 0 & \mathbf{v}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{v}_m \end{bmatrix} \in U(\mathbf{r}, \mathbf{c})$$

$\mathbf{u}, \mathbf{v}$  can be computed in  $O(nm)$  time using the Sinkhorn fixed-point iteration.

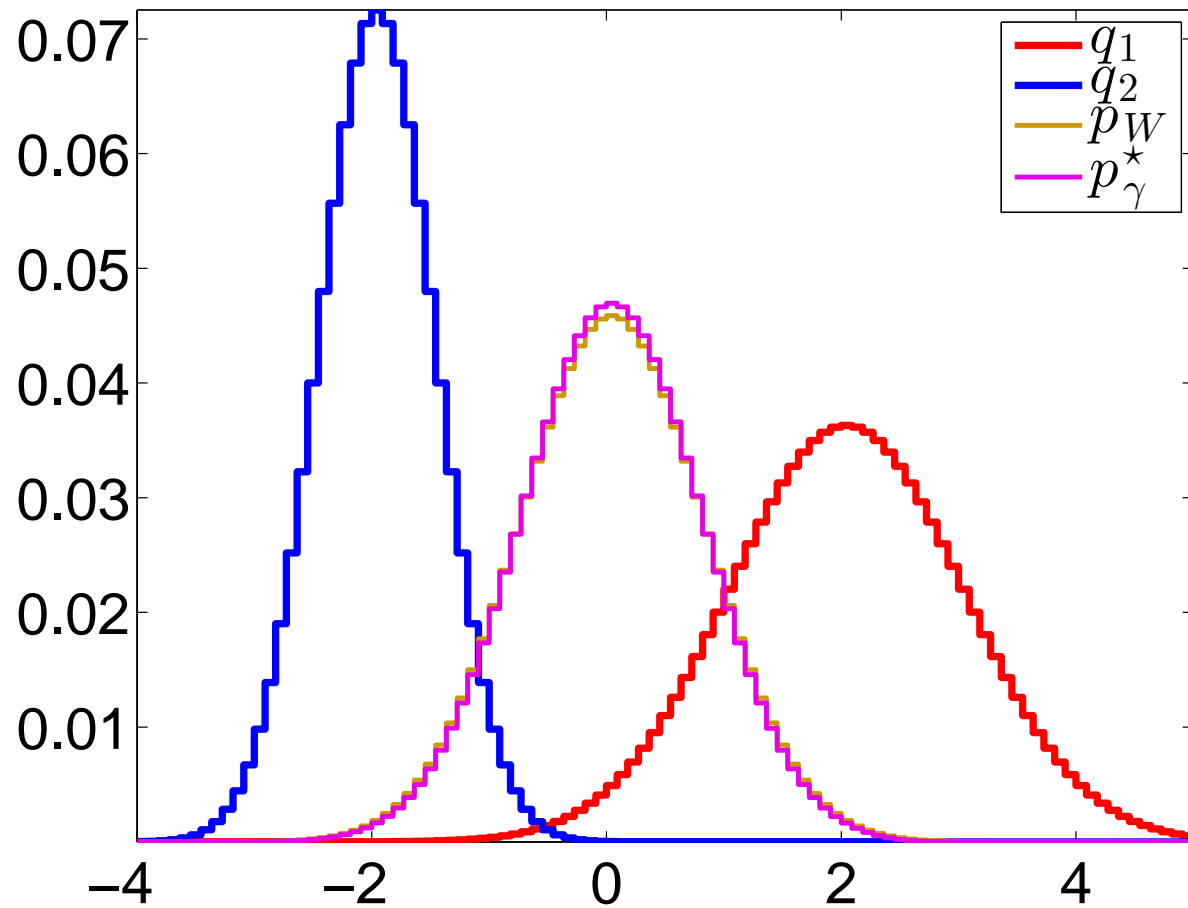
# Sinkhorn Algorithm

1. Set  $K = \exp(-M/\gamma)$  (note: if  $M$  is Euclidean metric, this is a Gaussian convolution...)
2. Seed initial random values for  $\mathbf{u}$ .
3. Loop until convergence
  - (a) Set  $\mathbf{v} \leftarrow (K' \mathbf{u}^{-1}) ./ c$
  - (b) Set  $\mathbf{u} \leftarrow (K \mathbf{v}^{-1}) ./ r$
  - $T_\gamma^* = \text{diag}(\mathbf{u}^*) K \text{diag}(\mathbf{v}^*), \alpha_\gamma^* = \log(u^*)/\gamma.$
  - $W_\gamma(\mu, \nu) = \langle T_\gamma, M \rangle = \mathbf{u}^{*T} (K. * M) \mathbf{v}^*$

# Benefits of Smoothing [C.'13]


- These OT problems are **strongly convex** vs. **LPs**.  
**Unicity** of solutions, **differentiable**.
- Considerably more efficient in practice [**Nesterov'05**].
- Primal/dual smoothed optima  $\alpha_\gamma^*$ ,  $T_\gamma^*$  can be solved
  - **In  $O(n^2)$**  with **Sinkhorn's (IPFP) algorithm**,
  - in **parallel on GPGPUs** for **any metric** on finite  $\Omega$ ,
  - **millions of time faster** than simplex,
  - can deal with **large dimensions** ( $\approx 50.000$  so far).

# Our Solution (using regularization)



# 1. Smoothed primal [C.Doucet'14]

$$C(\mathbf{a}) = \frac{1}{N} \sum_{i=1}^N \text{primal}_{\gamma}(\mathbf{a}, \mathbf{b}_i, M)$$

- (Projected) gradient descent:
  - Solve  $N$  smoothed (dual) OT problems  $\alpha_{i,\gamma}^*$
  - Update  $\mathbf{a}$  using gradient  $\frac{1}{N} \sum_i \alpha_{i,\gamma}^*$
-  each step requires **computing**  $\alpha_{i,\gamma}^*$ .



## 2. Dual approach [C. Peyré'14]

- The Fenchel-Legendre conjugate of

$$f_{\mathbf{b}}(\mathbf{a}) = \text{primal}_{\gamma}(\mathbf{a}, \mathbf{b}, M),$$

namely

$$f_{\mathbf{b}}^*(\mathbf{g}) = \max_{\mathbf{p} \in \Sigma_n} \langle \mathbf{g}, \mathbf{p} \rangle - f_{\mathbf{b}}(\mathbf{p}).$$

has a **closed form**

$$f_{\mathbf{b}}^*(\mathbf{g}) = \gamma \left( H(\mathbf{b}) + \langle \mathbf{b}, \log e^{-M/\gamma} e^{\mathbf{g}/\gamma} \rangle \right)$$

## 2. Dual approach [C. Peyré'14]

- The original problem in splitted form:

$$\min_{\mathbf{a}_1, \dots, \mathbf{a}_N \in \Sigma_n} \sum_{i=1}^N f_{b_i}(\mathbf{a}_i) \text{ subj. to } \mathbf{a}_1 = \dots = \mathbf{a}_N$$

- can be replaced with an easier problem:

$$\min_{\mathbf{g}_1, \dots, \mathbf{g}_N \in \mathbb{R}^n} \sum_{i=1}^N f_{b_i}^*(\mathbf{g}_i) \text{ subj. to } \sum_{i=1}^N \mathbf{g}_i = 0.$$

gradient/Hessian explicit, equality constraint  $\rightarrow$  **truncated Newton**.

at convergence, all  $\nabla f_{b_i}^*(\mathbf{g}_i)$  are equal to solution  $\mathbf{a}^*$ .

### 3. Generalized KL Projections [NBCCP'14]

- Idea: generalize **KL** projection for two marginals

$$\operatorname{argmin}_{T \in U(\mathbf{a}, \mathbf{b}_i)} \mathbf{KL}(T | e^{-M/\gamma})$$

- to alternated **KL** projections for  $N + 1$  common (unknown) one.

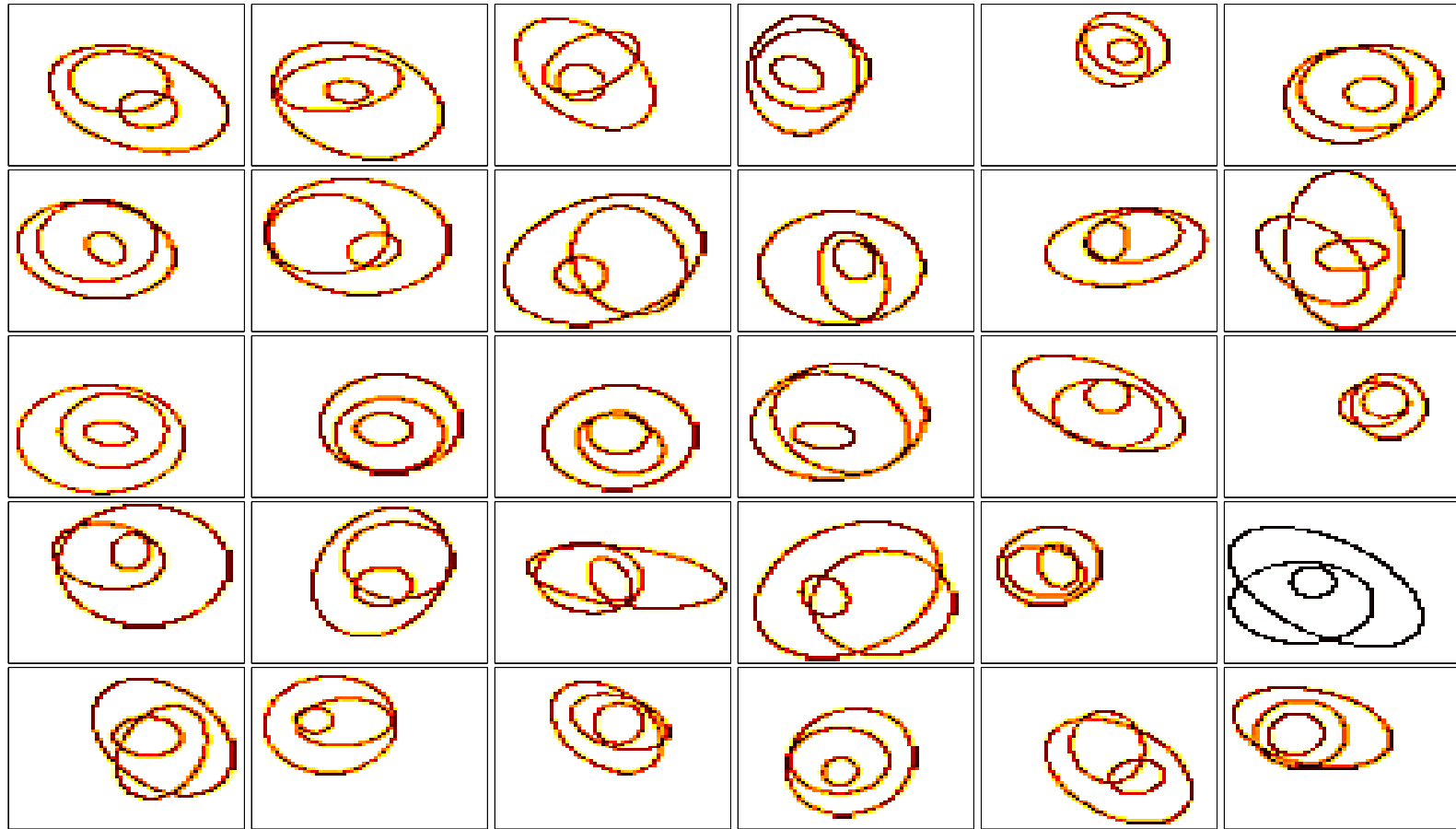
$$\operatorname{argmin}_{T_i^T \mathbf{1} = \mathbf{b}_i, T_i \mathbf{1} = T_{i+1} \mathbf{1}} \sum_i \mathbf{KL}(T_i | e^{-M/\gamma})$$

- 2 lines of matlab code.

---

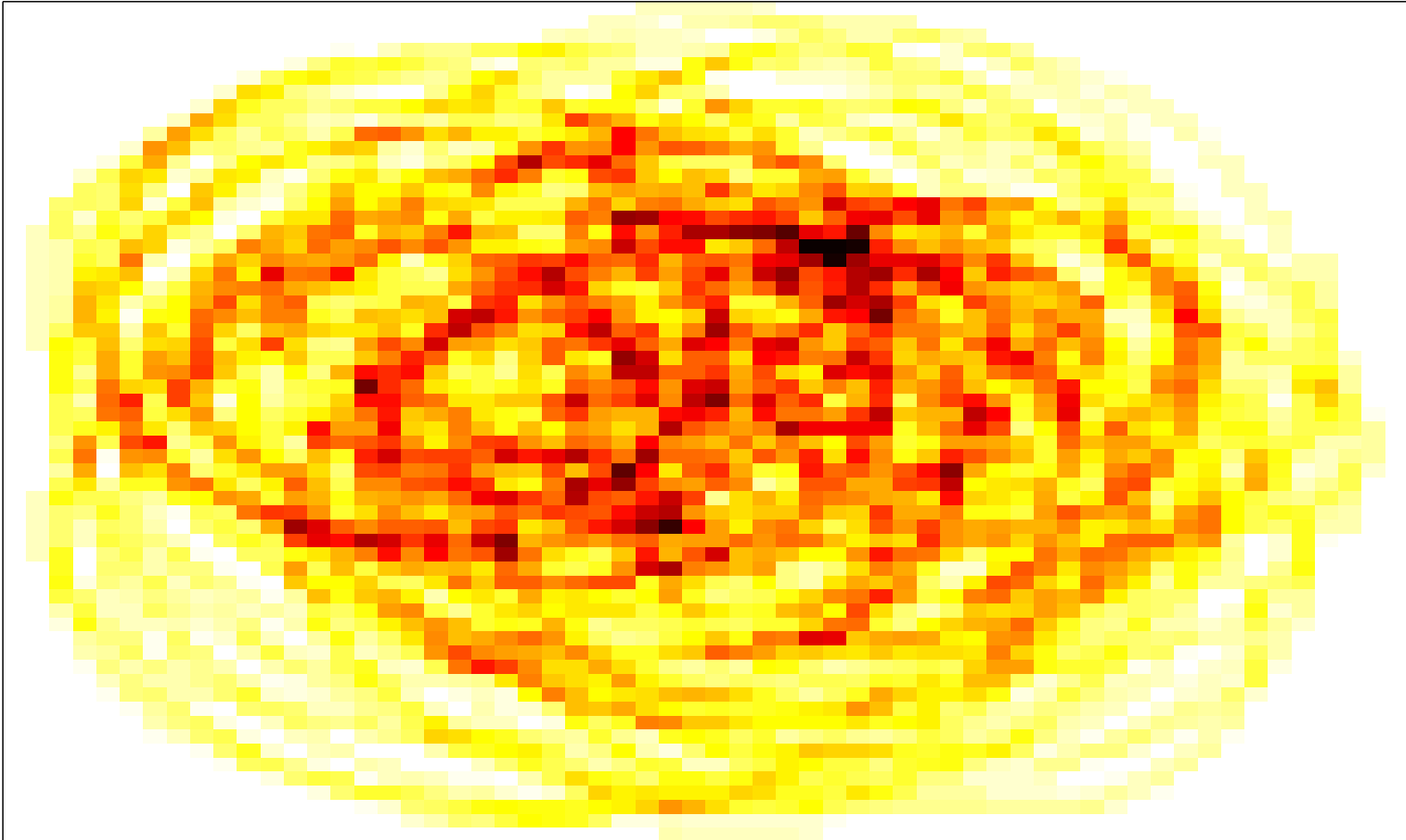
# Applications

# Averaging 30 Measures

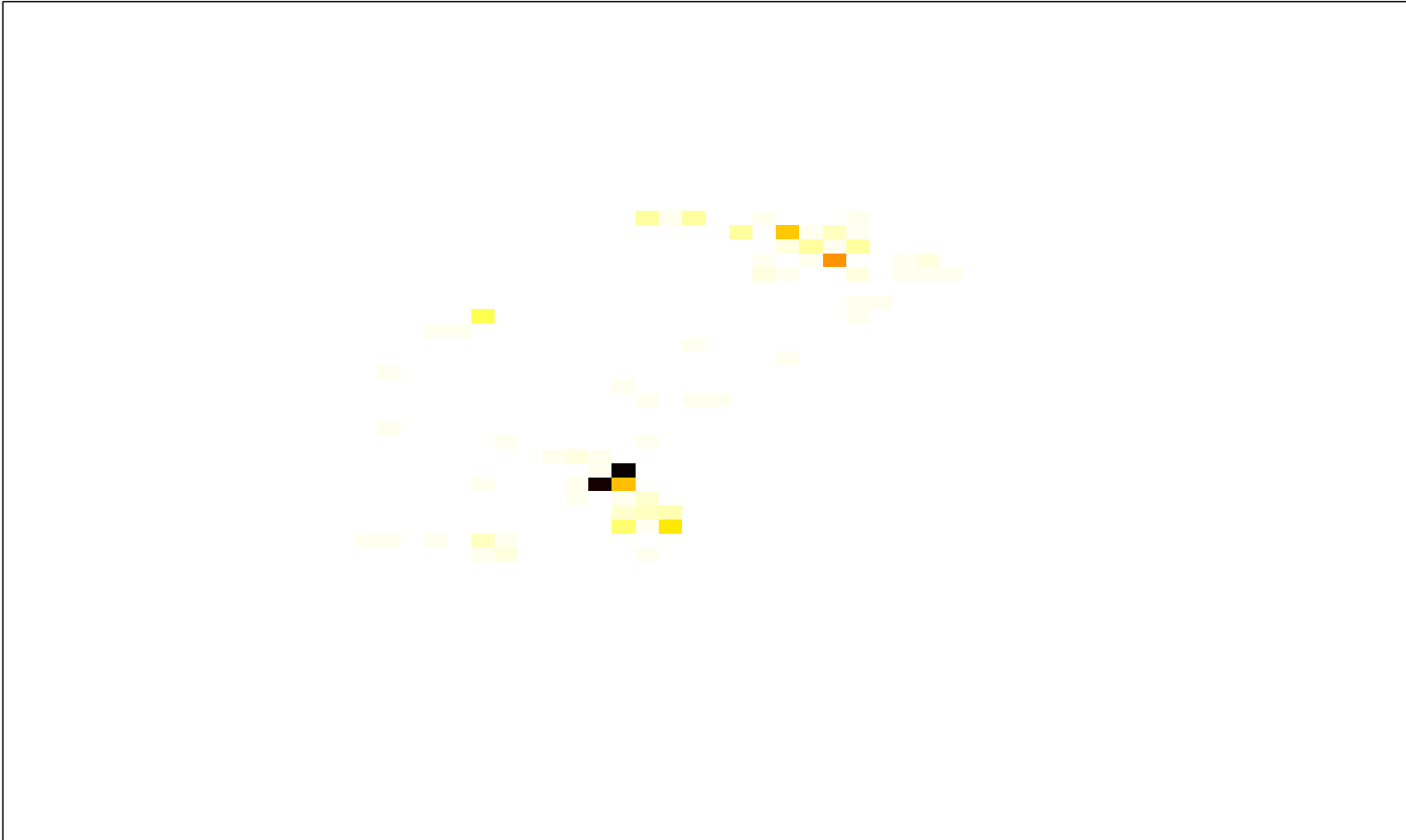


30 measures on  $\mathbb{R}^2$ .

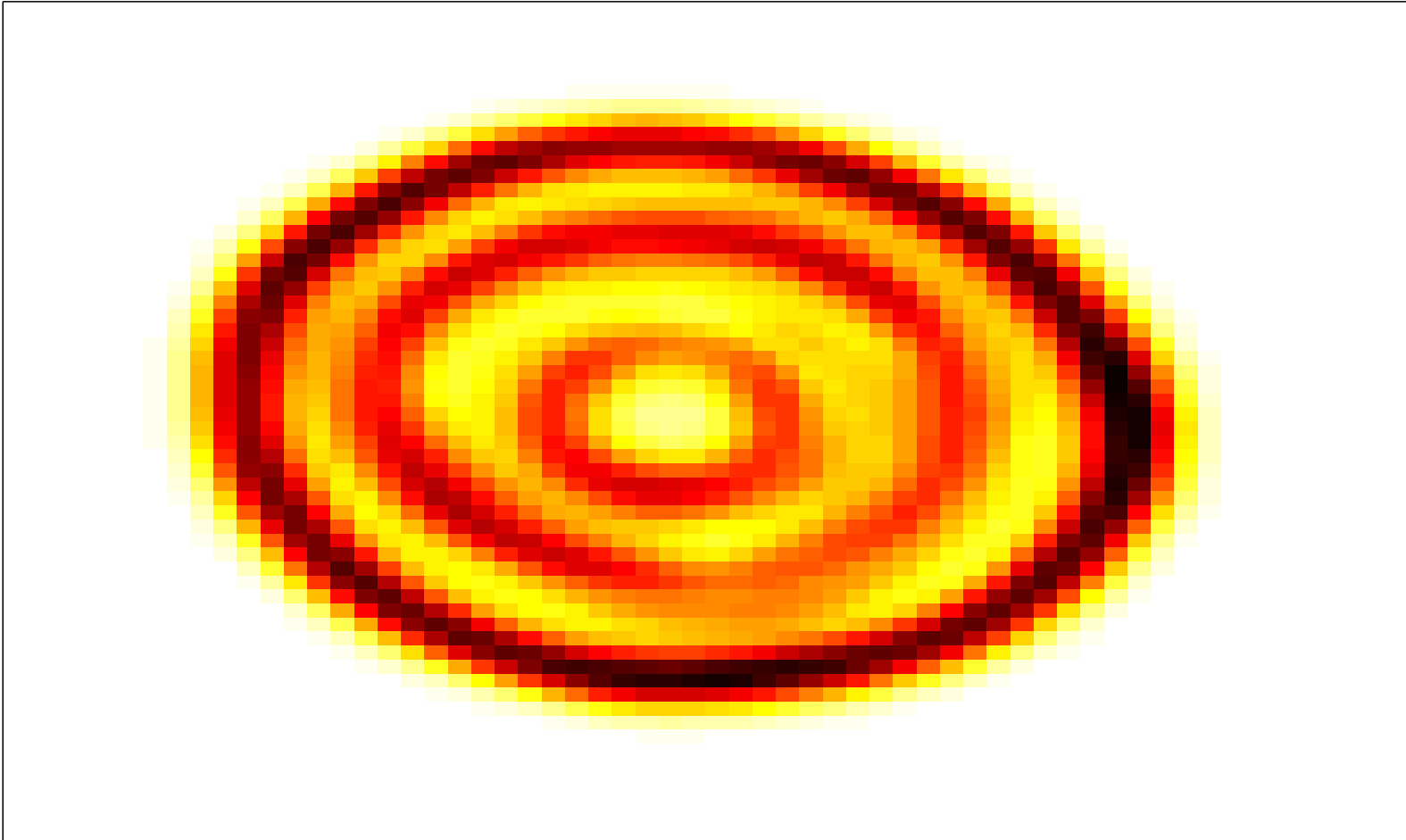
# Euclidean Mean



# Symmetric KL Mean

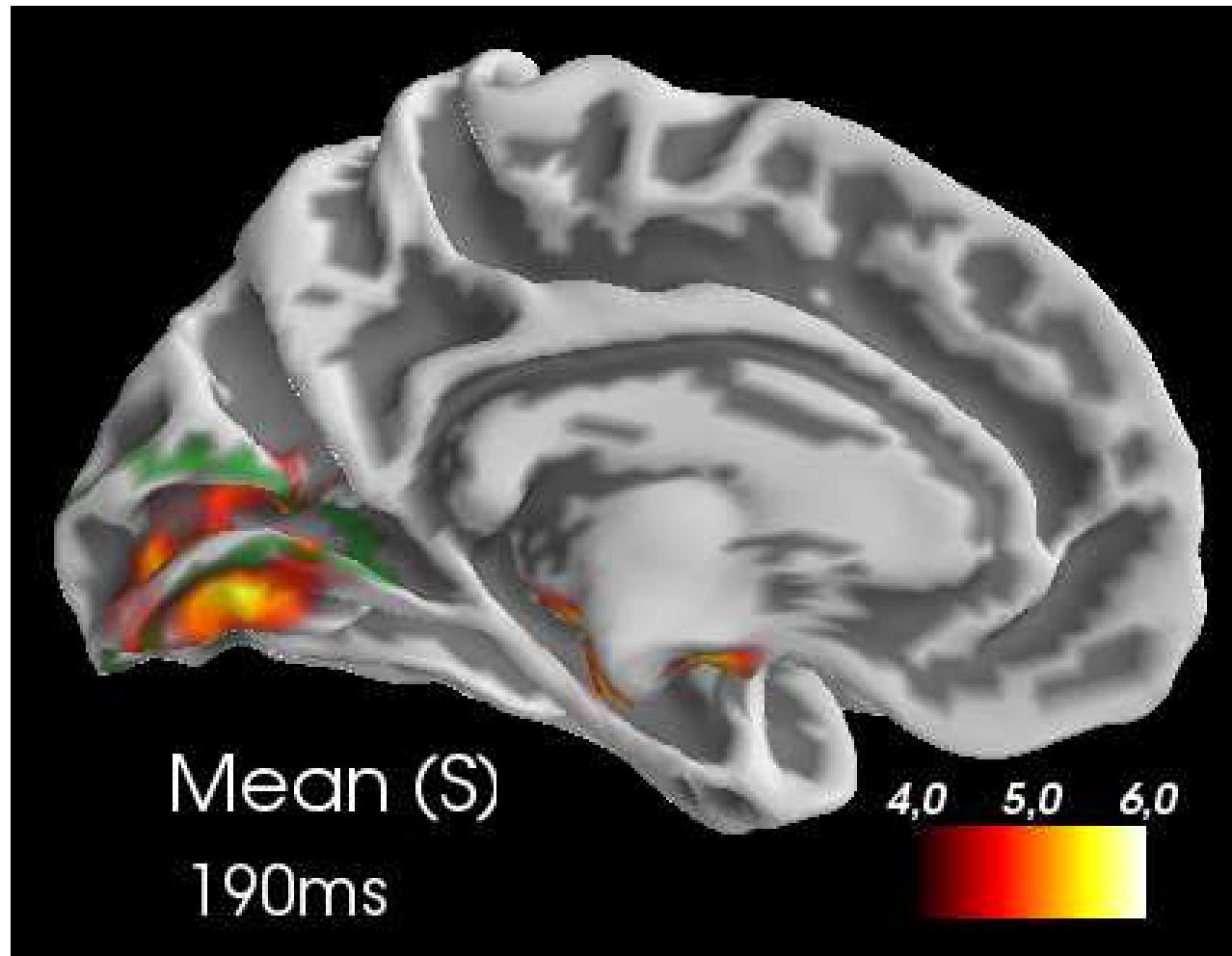


# 2-Wasserstein



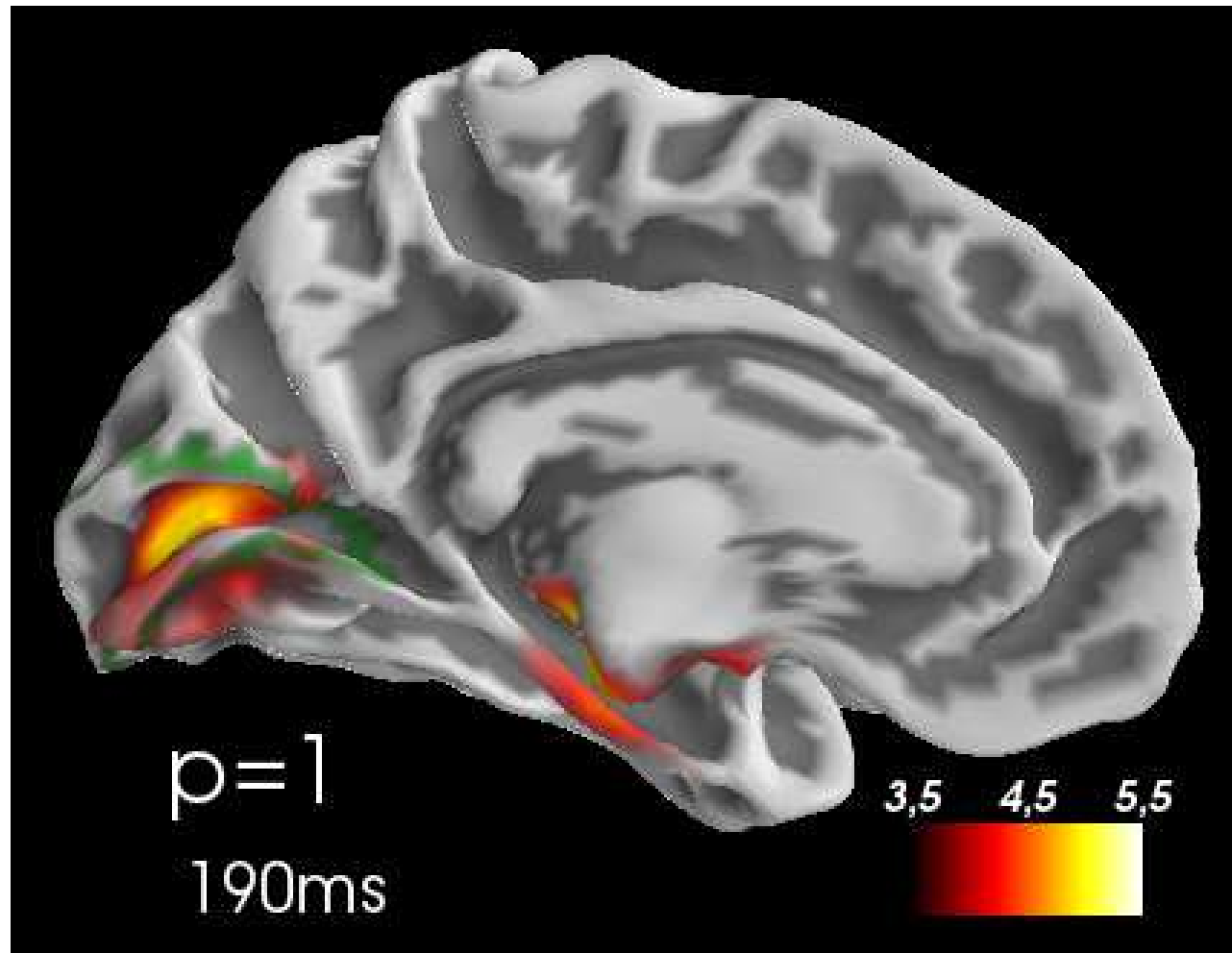


# Averaging Brain Activations real

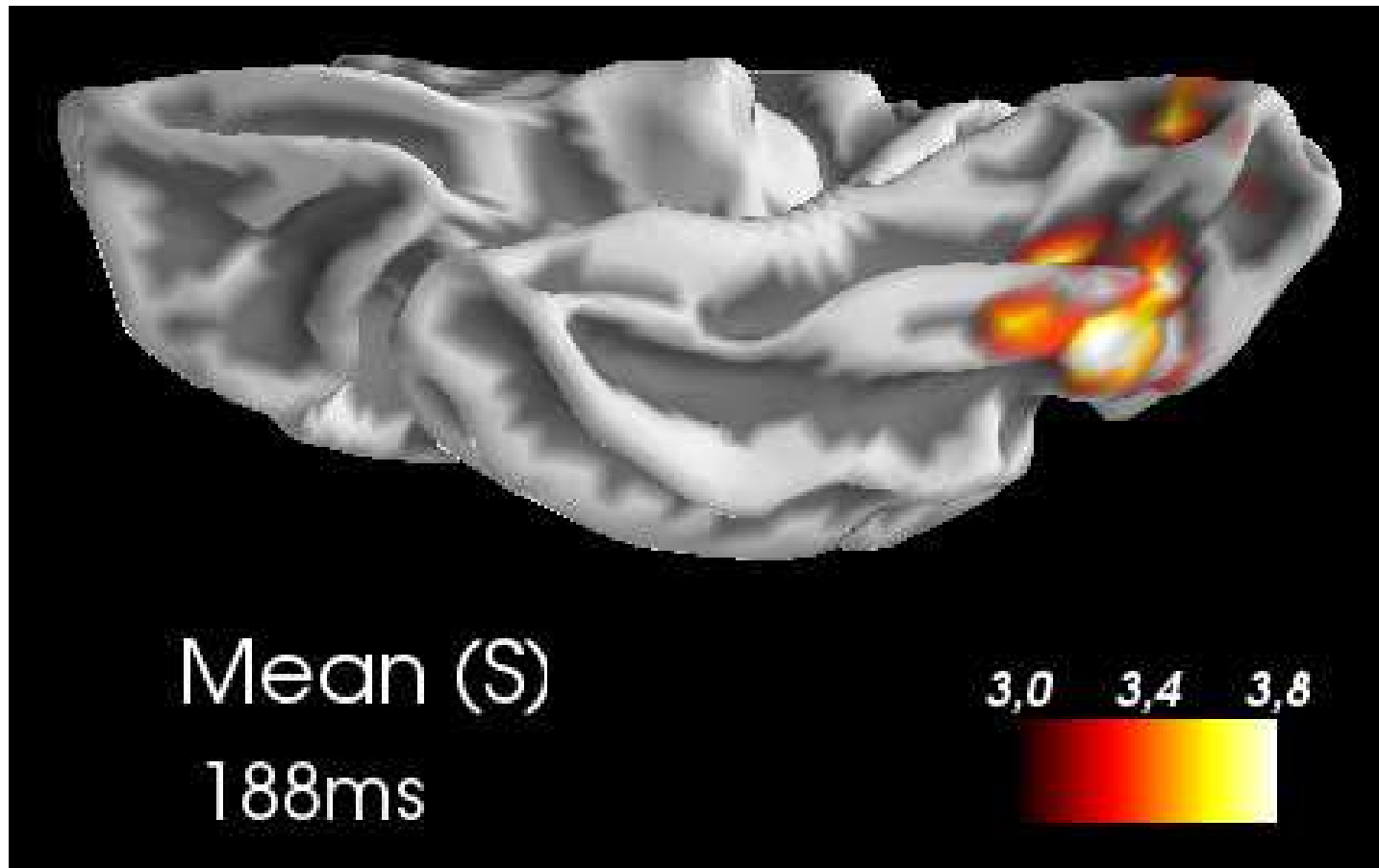


MEG ERF data, N=16. Left, medial view. **border** of (V1)

# 1-Wasserstein

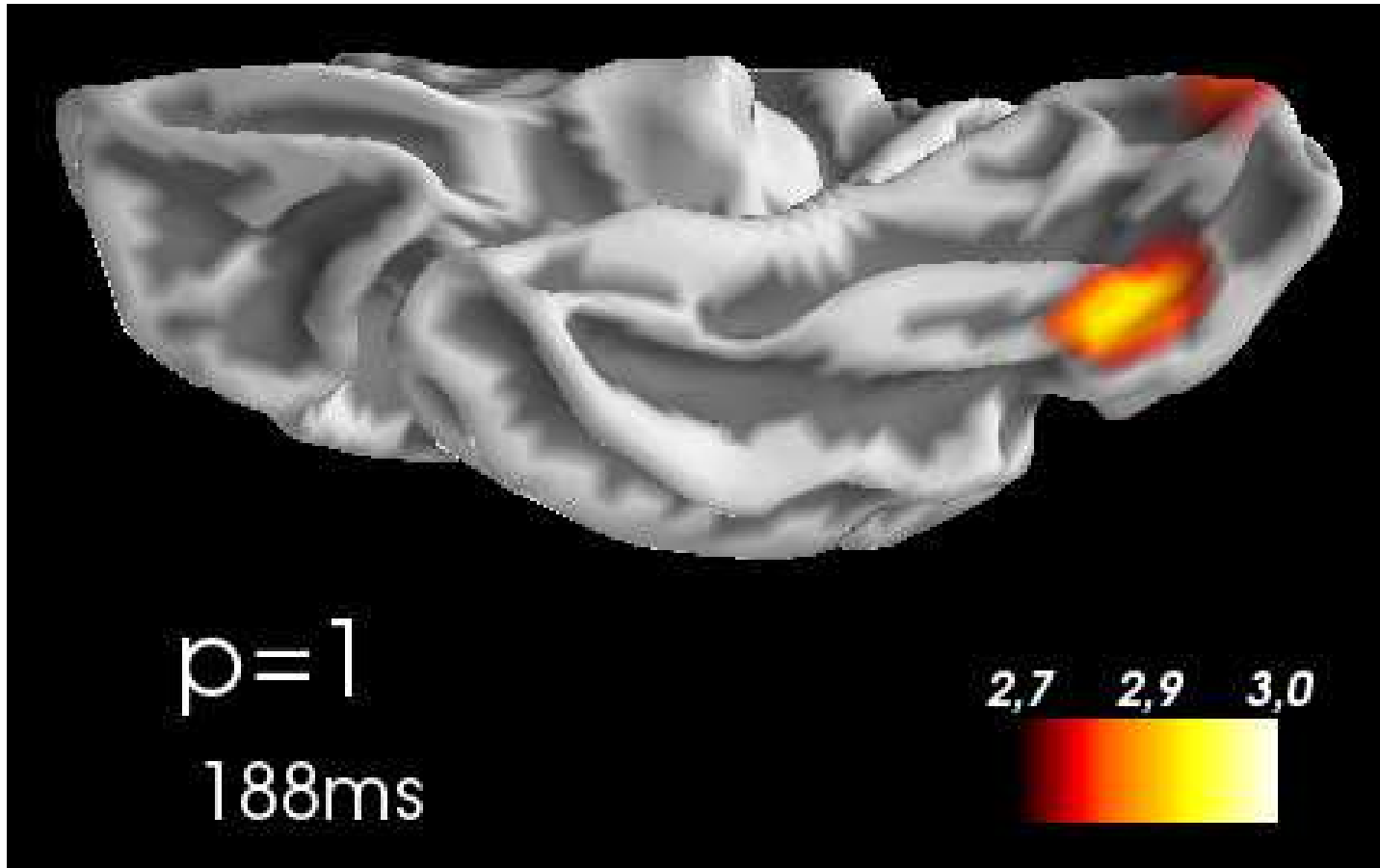


# Averaging Brain Activations real



Right, ventral view

# Averaging Brain Activations real



Centered on the Fusiform gyrus

# Averaging Text Histograms

- Using GLOVE embeddings for words, 2-Wasserstein.

# Graphics

# Graphics

# Graphics



# Graphics

# Graphics

# Graphics

**End**