

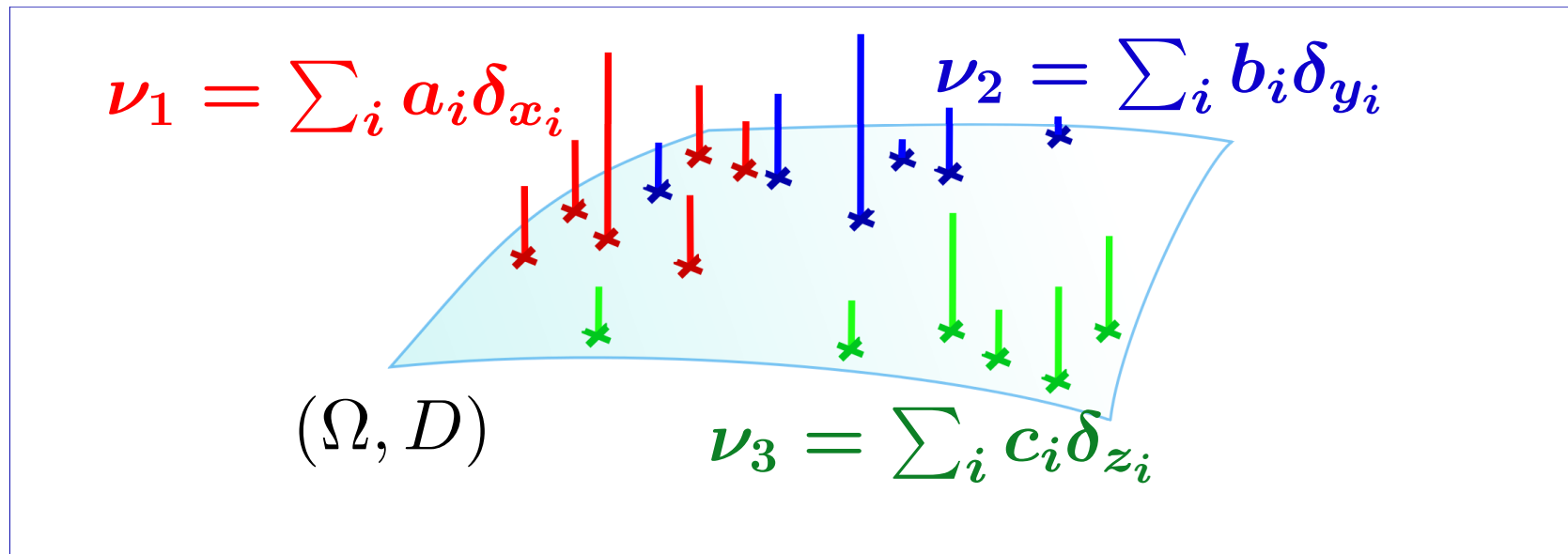
**ICML 2014**

**Fast Computation of  
Wasserstein Barycenters**

**M. Cuturi<sup>1</sup> A. Doucet<sup>2</sup>**

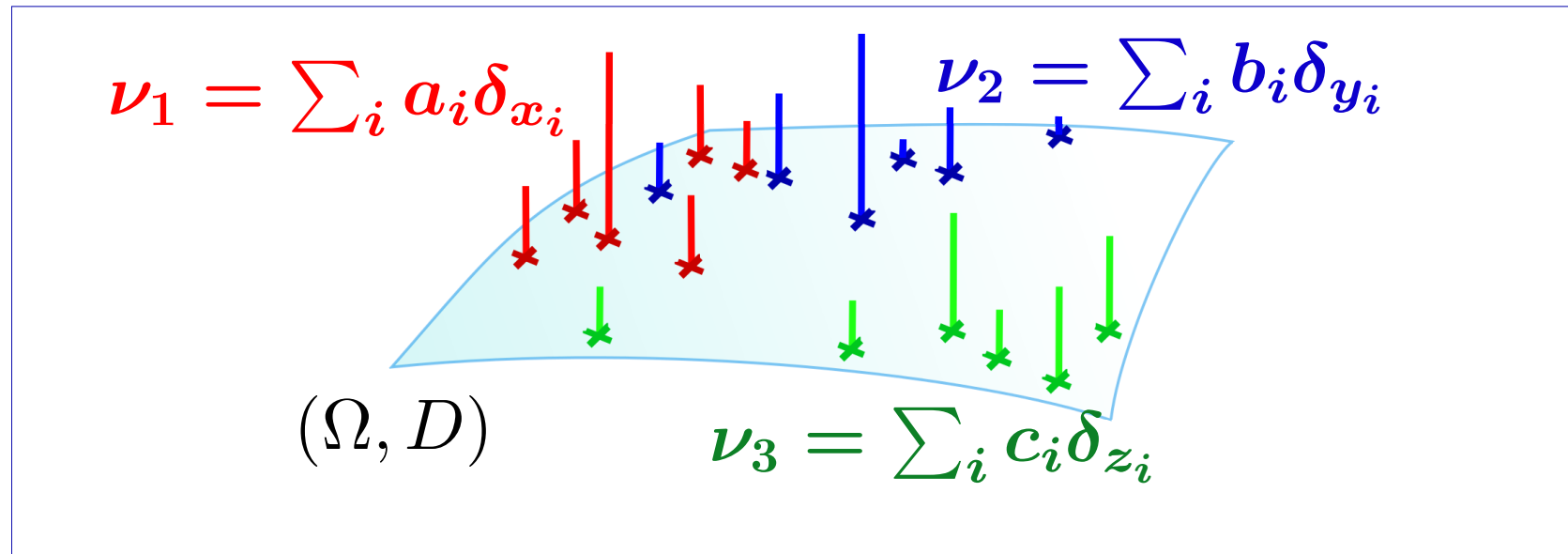
**<sup>1</sup>Kyoto University <sup>2</sup>University of Oxford**

# Problem: Average $N$ Probability Measures



- $\{\Omega, D\}$  a metric space
- $\{\nu_1, \dots, \nu_N\}$  family of empirical probability measures.

# Problem: Average $N$ Probability Measures



Can we **summarize** the  $\{\nu_i\}$  as an “**average**” or a “**barycentric**” single empirical probability measure?

*interest in ML: empirical measure = dataset,*

*histogram/bags-of-features, single observation with uncertainty*

# Euclidean Means for Vectors

- For **vectors**  $\{x_1, \dots, x_N\}$  in a Hilbert, their average is

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

# Euclidean Means for Vectors

- For **vectors**  $\{x_1, \dots, x_N\}$  in a Hilbert, their average is

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

- behind this formula lies a **variational** problem

$$\bar{x} = \operatorname{argmin}_{u \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N \|u - x_i\|_2^2$$

# Euclidean Means for Measures

- For **probability measures**  $\{\nu_i\}_{i=1..N}$ , we can also use:

$$\mu = \frac{1}{N} \sum_{i=1}^N \nu_i,$$

- as well as, using a **smoothing kernel**  $k = e^{-D^2/\sigma}$ ,

$$\mu = \frac{1}{N} \sum_{i=1}^N (k * \nu_i)$$

(a.k.a *RKHS* mean map [**Gretton'07**])

# Other Means for Probabilities

- Other means can be defined using other **metrics** or **divergences**:

$$\operatorname{argmin}_{\mu \in P(\Omega)} \sum_{i=1}^N \Delta(\mu, \nu_i).$$

- **KL, Symmetrized KL** [Nielsen'12]
- **Bregman Divergence** [Bhanerjee'05]
- **Wasserstein Distance** (a.k.a *EMD*) [Agueh'11]

# Wasserstein Barycenter Problem

- **[Agueh'11]** defined

$$\operatorname{argmin}_{\mu \in P(\Omega)} \sum_{i=1}^N W_p^p(\mu, \nu_i),$$

provided theoretical analysis, unicity of solution.

- Simple cases ( $N = 2$ , multivariate Gaussians) covered.
- **very challenging computational problem.**



# Our Contribution

- **First computational approach** to solve efficiently **variational Wasserstein problems**,
- including the Wasserstein barycenter problem,

$$\operatorname{argmin}_{\mu \in P(\Omega)} \sum_{i=1}^N \mathbf{W}_p^p(\mu, \nu_i),$$

- that is applicable for arbitrary  $(\Omega, D)$  and  $p > 0$ , using **entropy-smoothed optimal transport [Cuturi'13]**.

# Our Contribution

- **First computational approach** to solve efficiently **variational Wasserstein problems**,
- including the Wasserstein barycenter problem,

$$\operatorname{argmin}_{\mu \in P(\Omega)} \sum_{i=1}^N W_p^p(\mu, \nu_i),$$

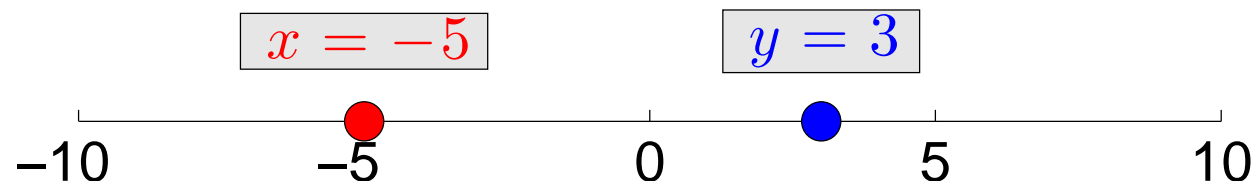
- that is applicable for arbitrary  $(\Omega, D)$  and  $p > 0$ , using **entropy-smoothed optimal transport [Cuturi'13]**.

([Rabin'12, Bonneel'14] studied case  $\Omega = \mathbb{R}^2$ )

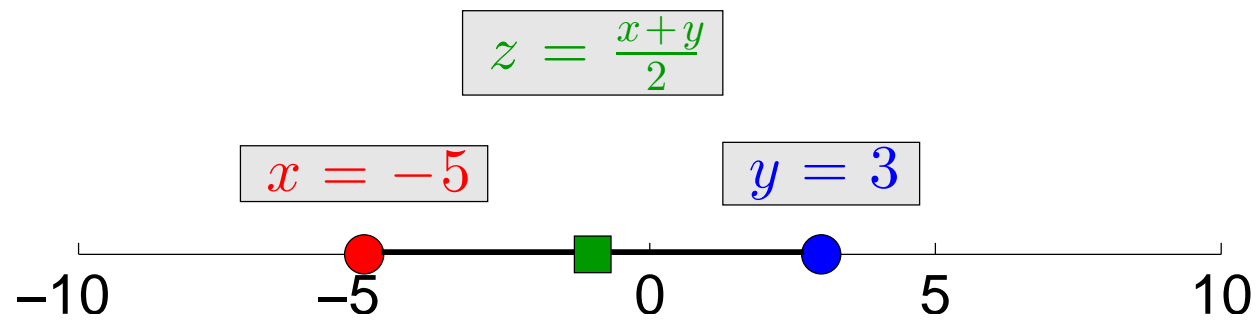
---

# Motivating Examples

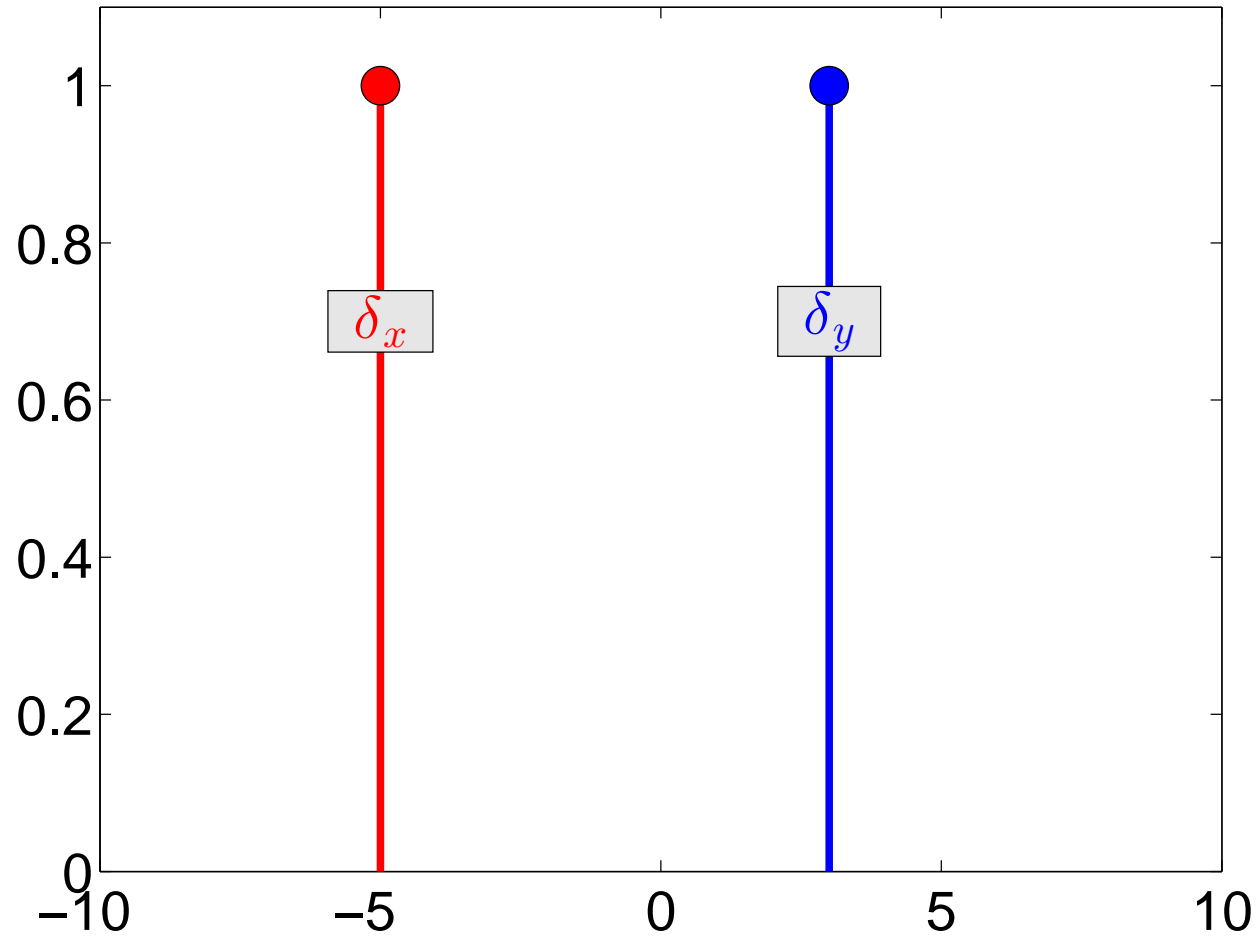
# 2 Points on the Real Line



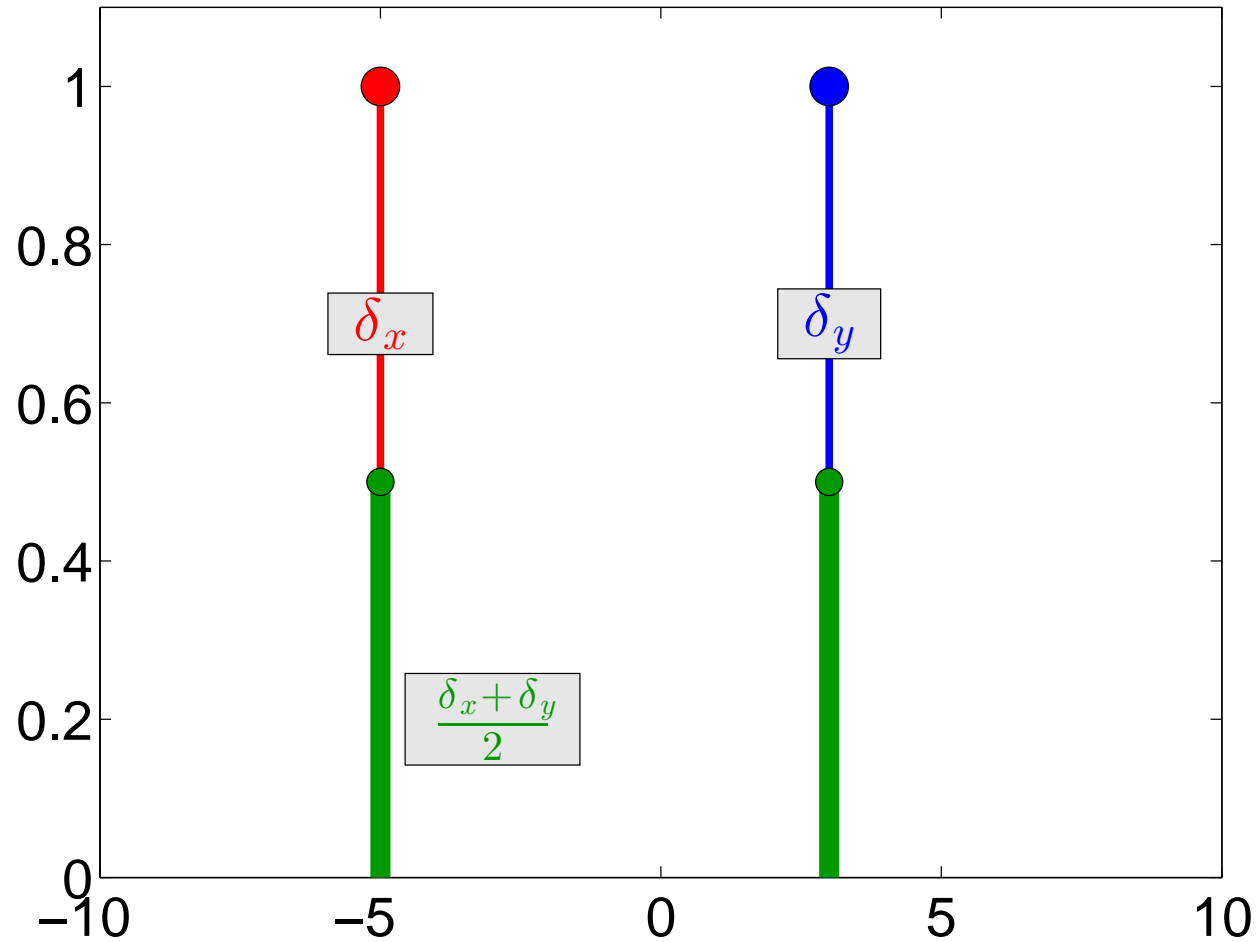
# Their Average



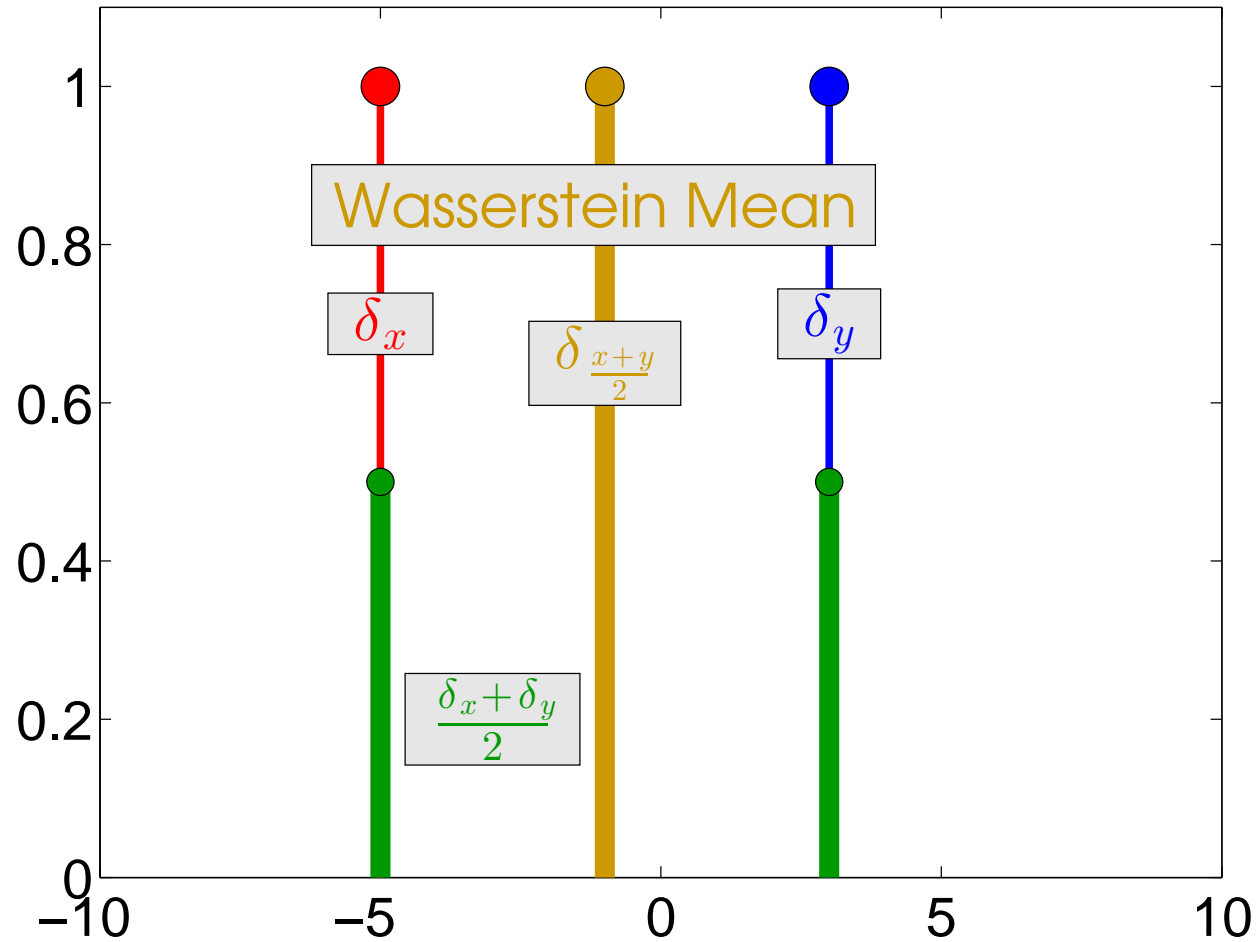
# 2 Points as Diracs



# Euclidean Mean of Diracs

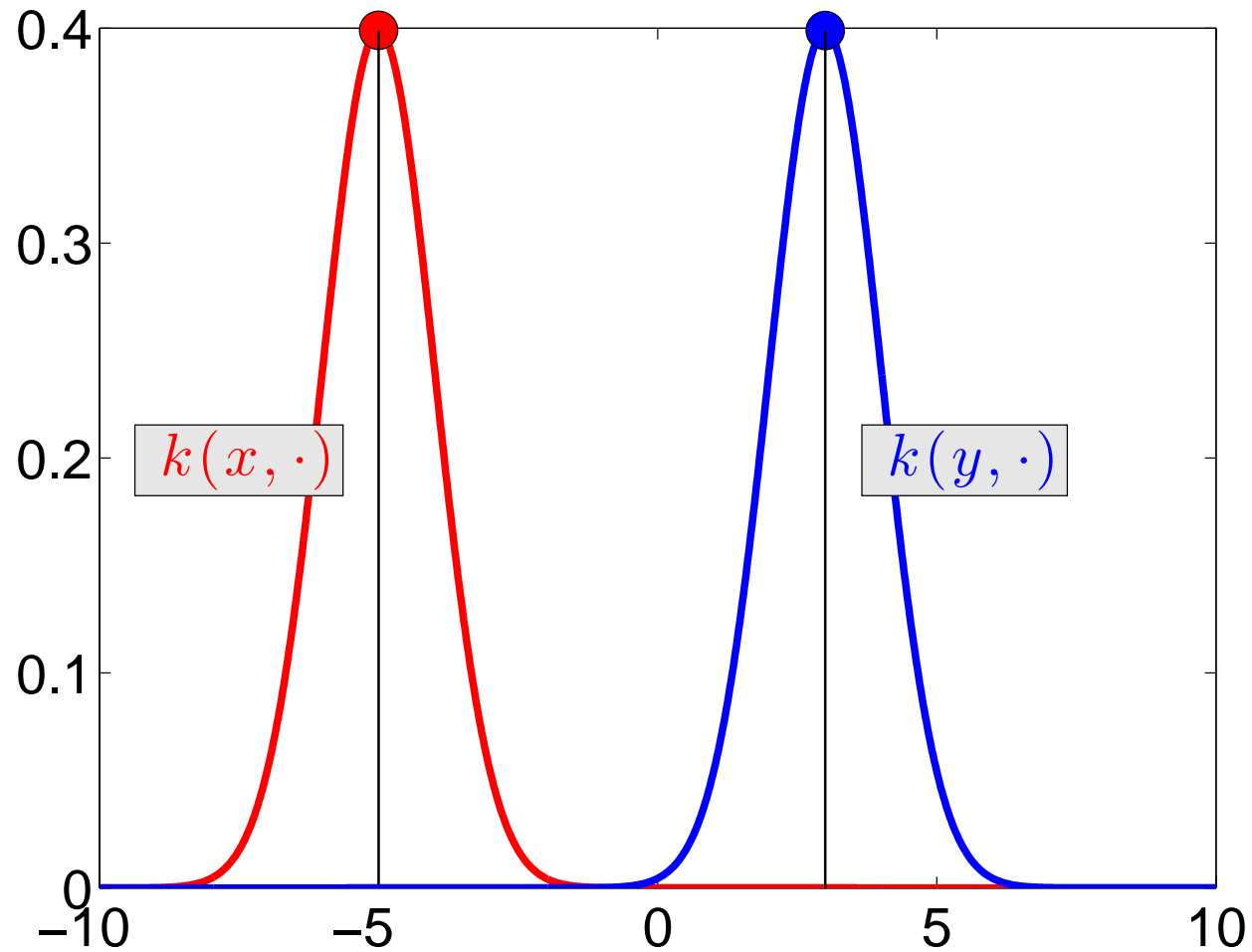


# Wasserstein Mean of Diracs

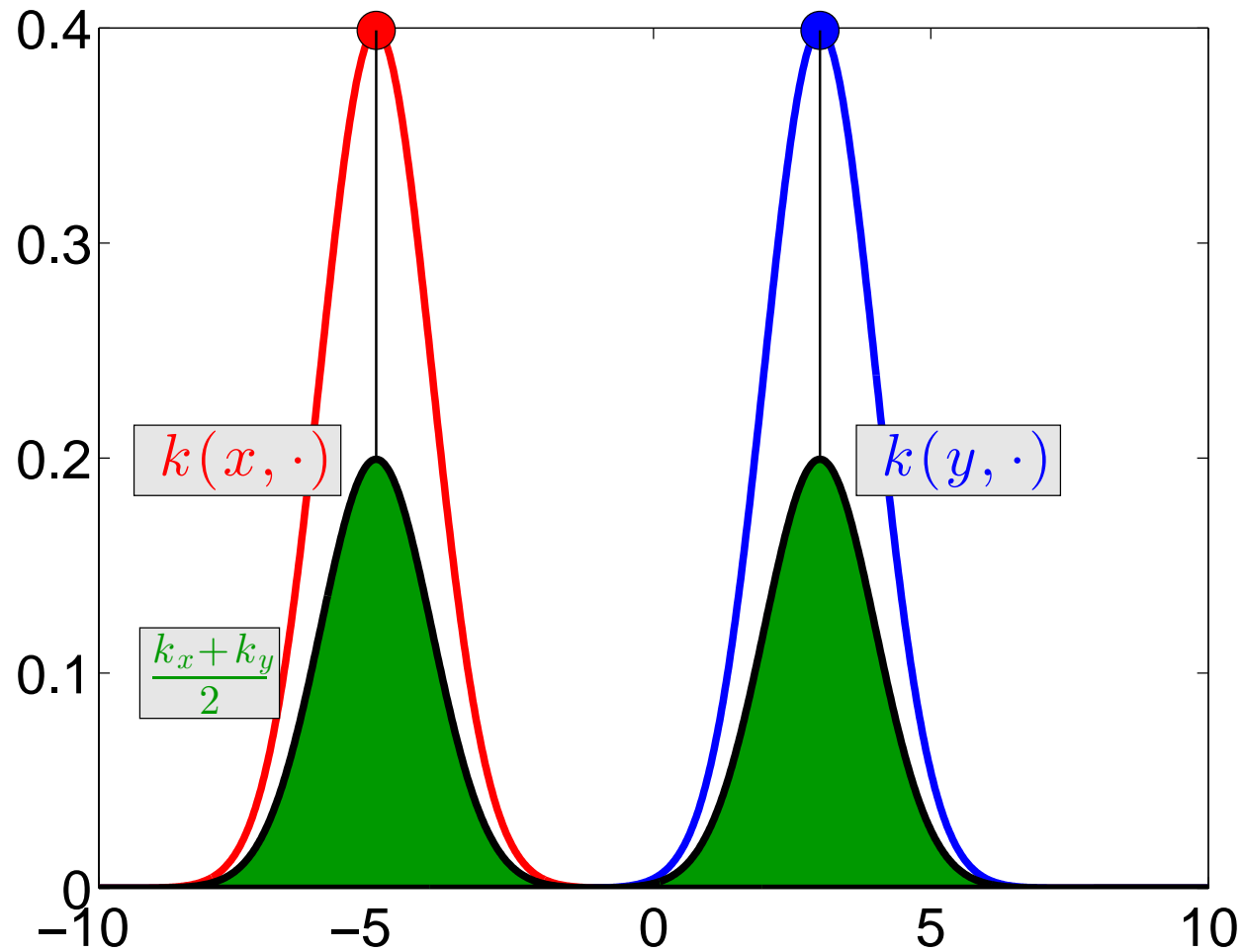




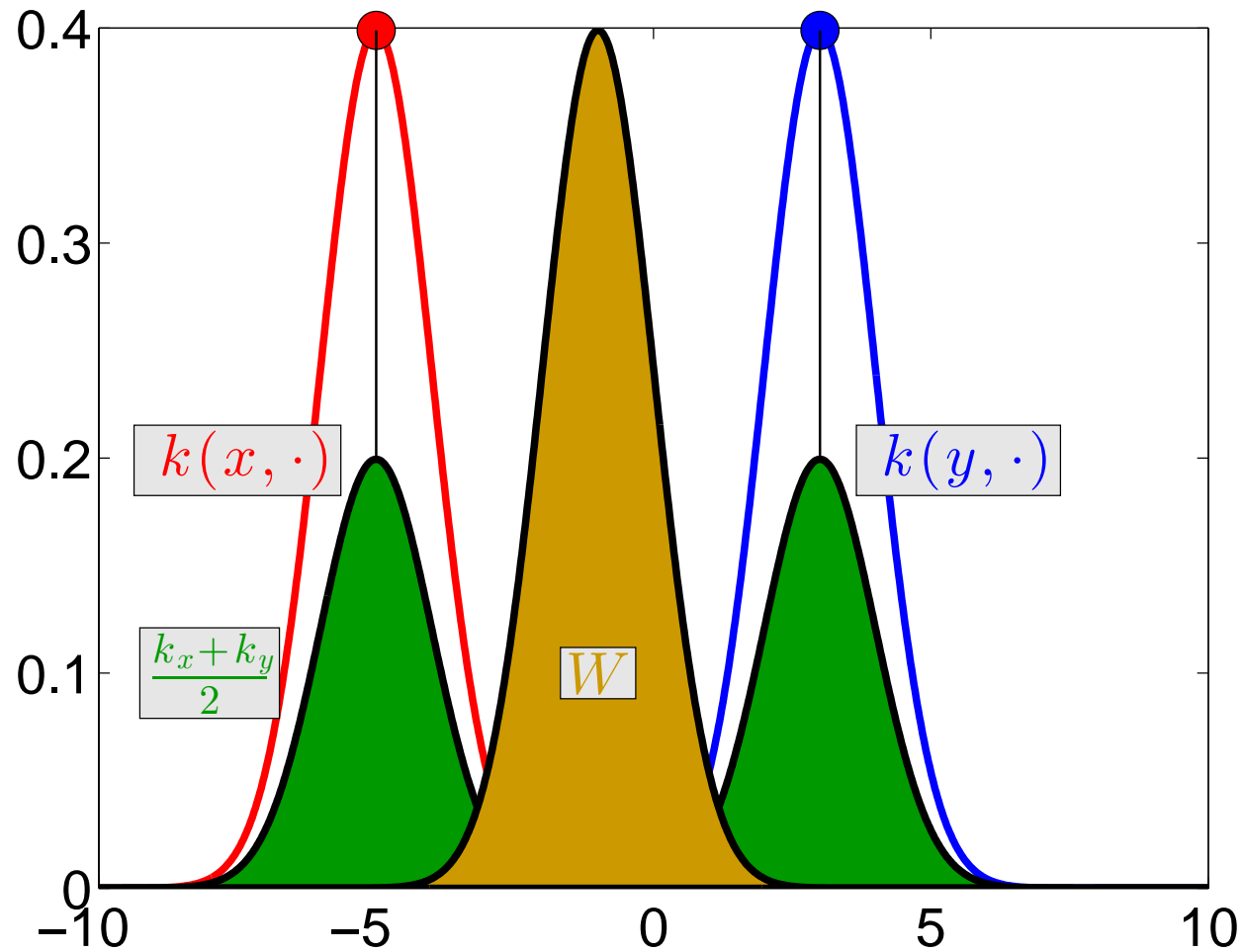
# Smoothed Measures (RKHS mean map)



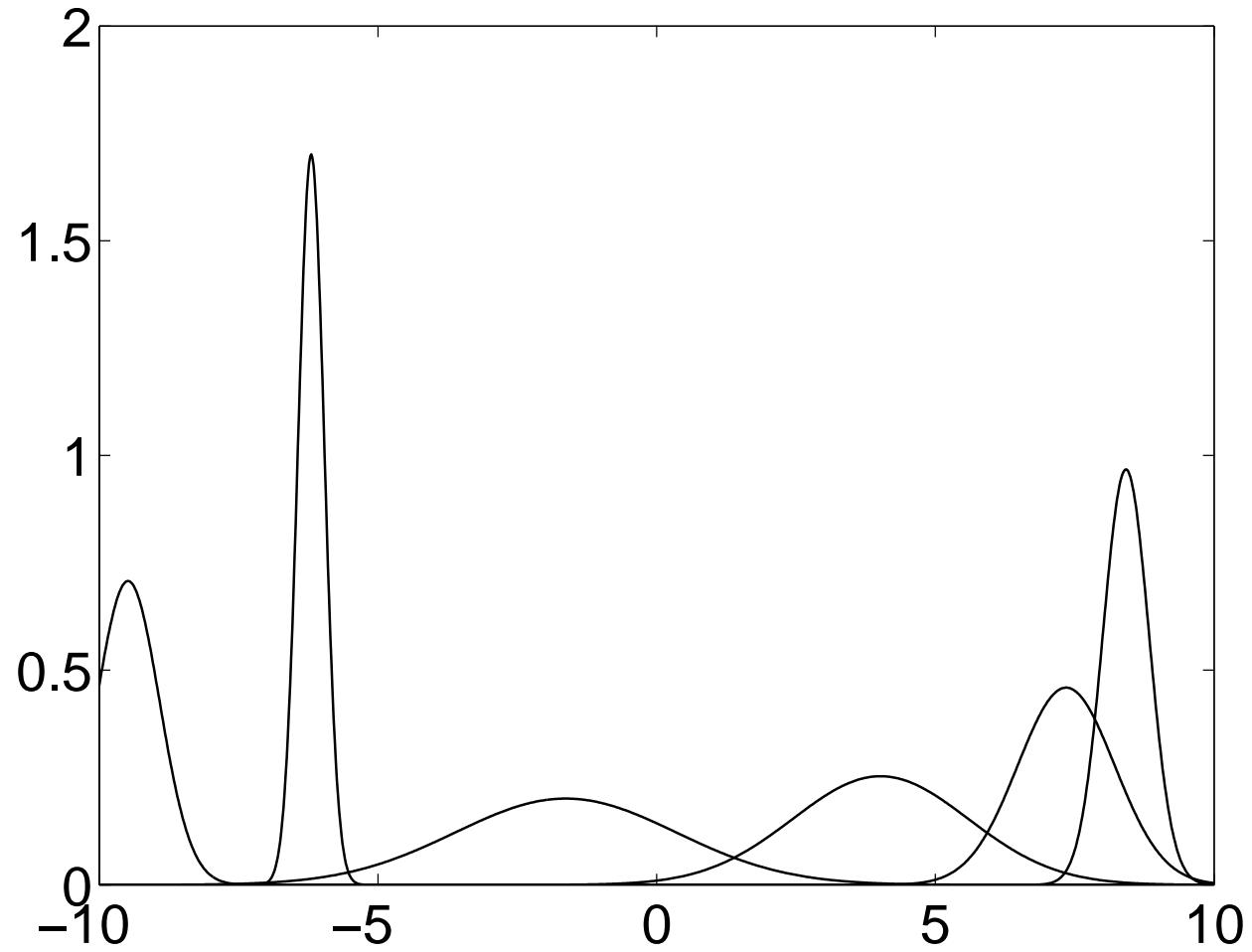
# Euclidean Mean of 2 Gaussians



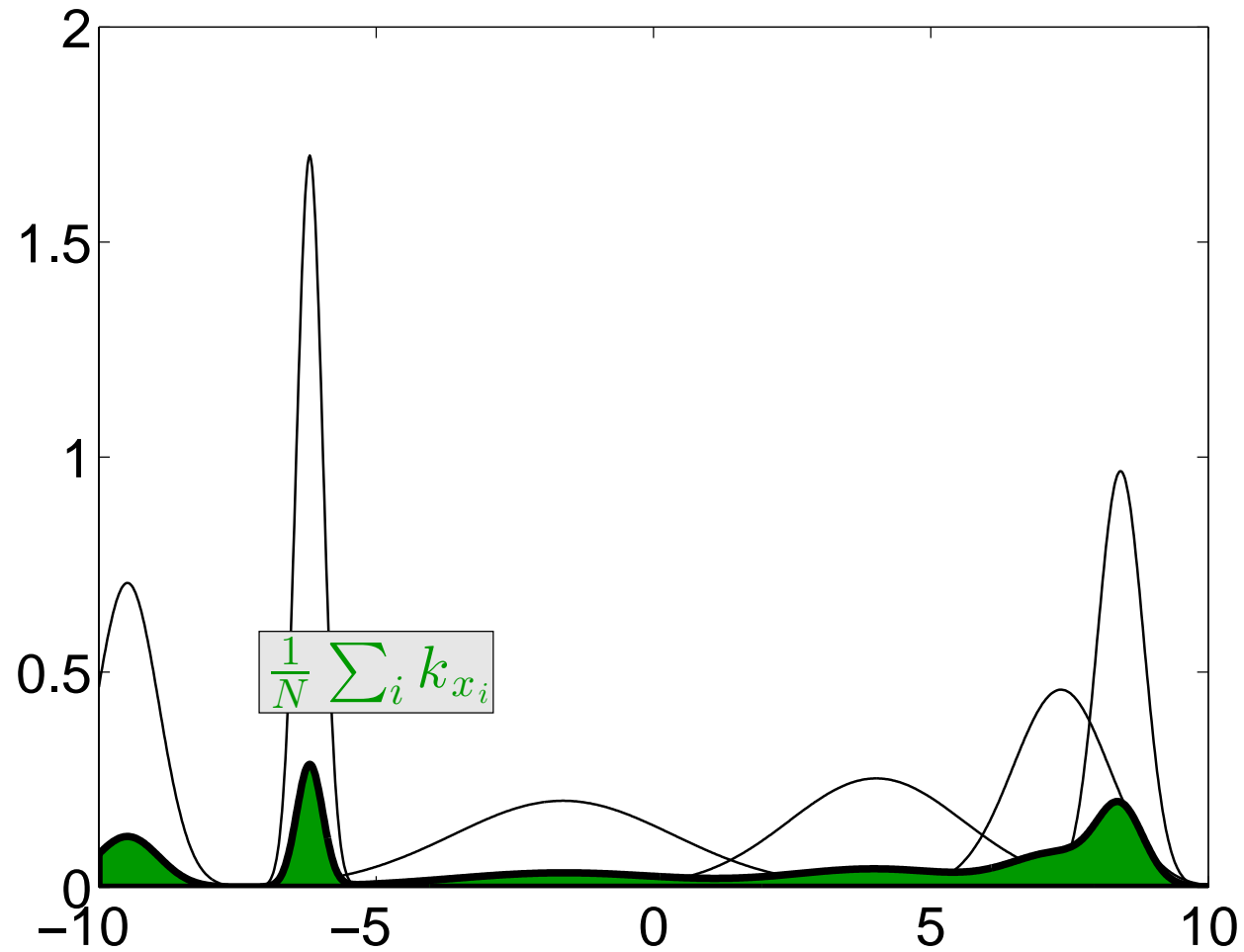
# Wasserstein Mean of 2 Gaussians



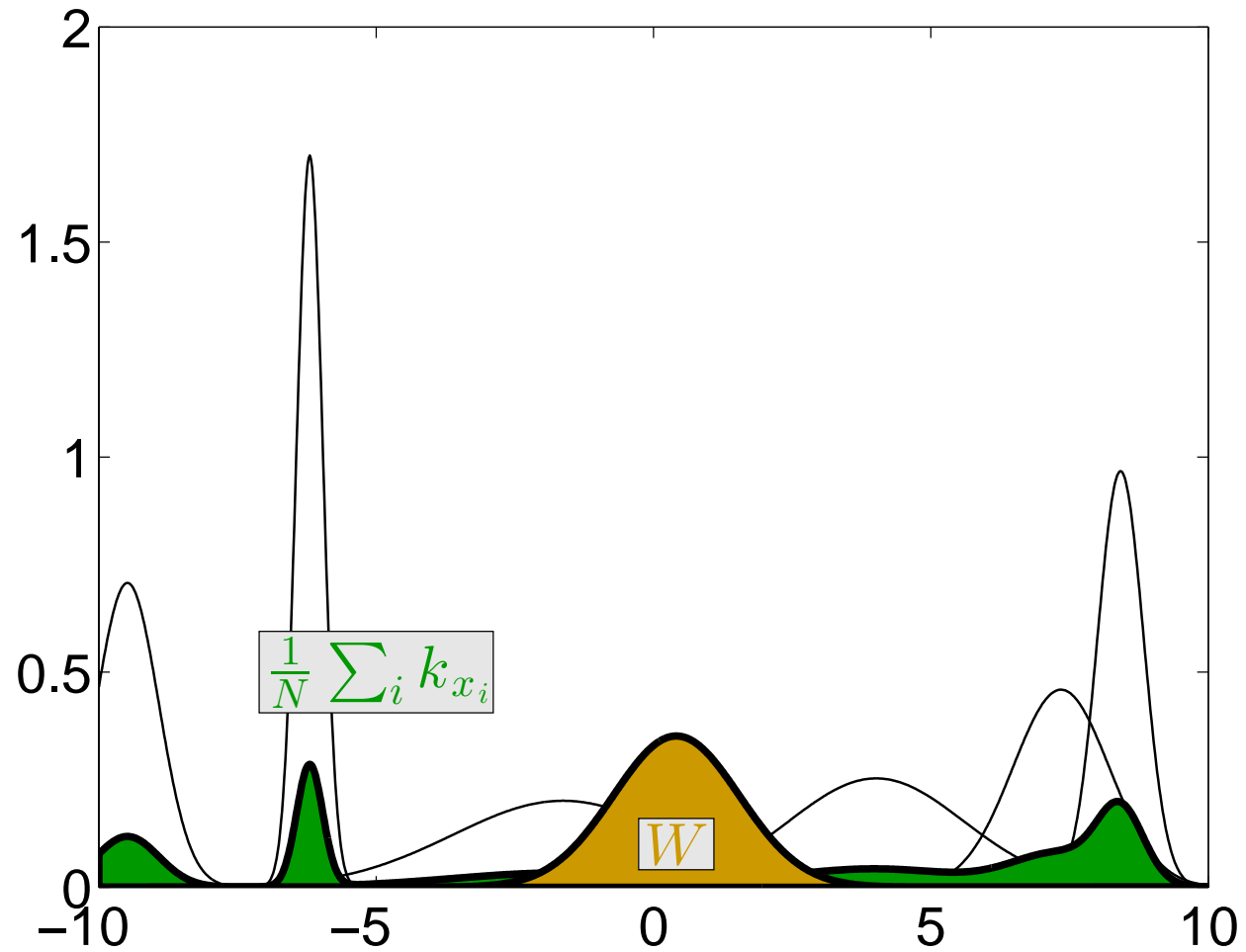
# 6 Gaussians



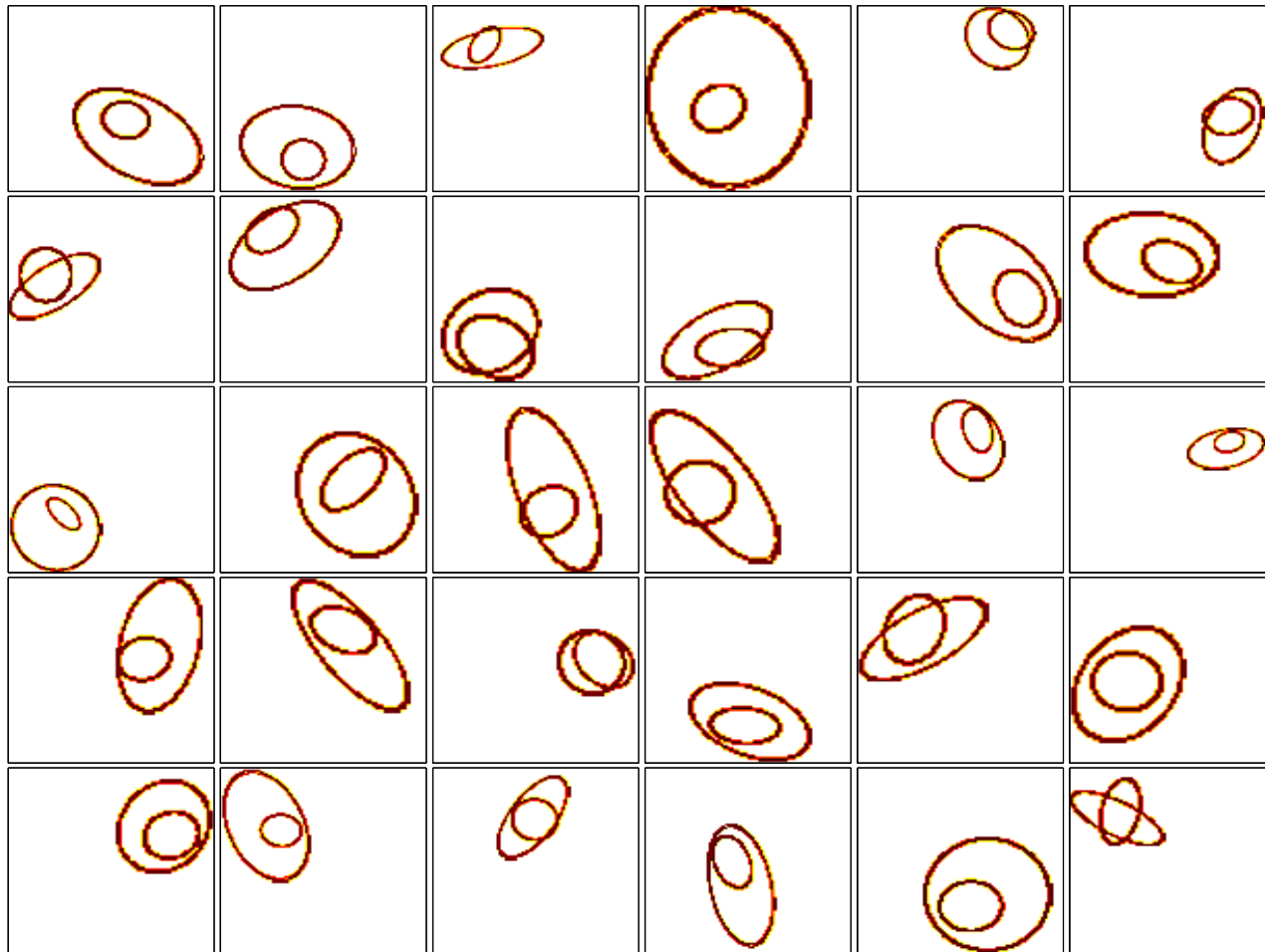
# Euclidean Mean



# Wasserstein Mean

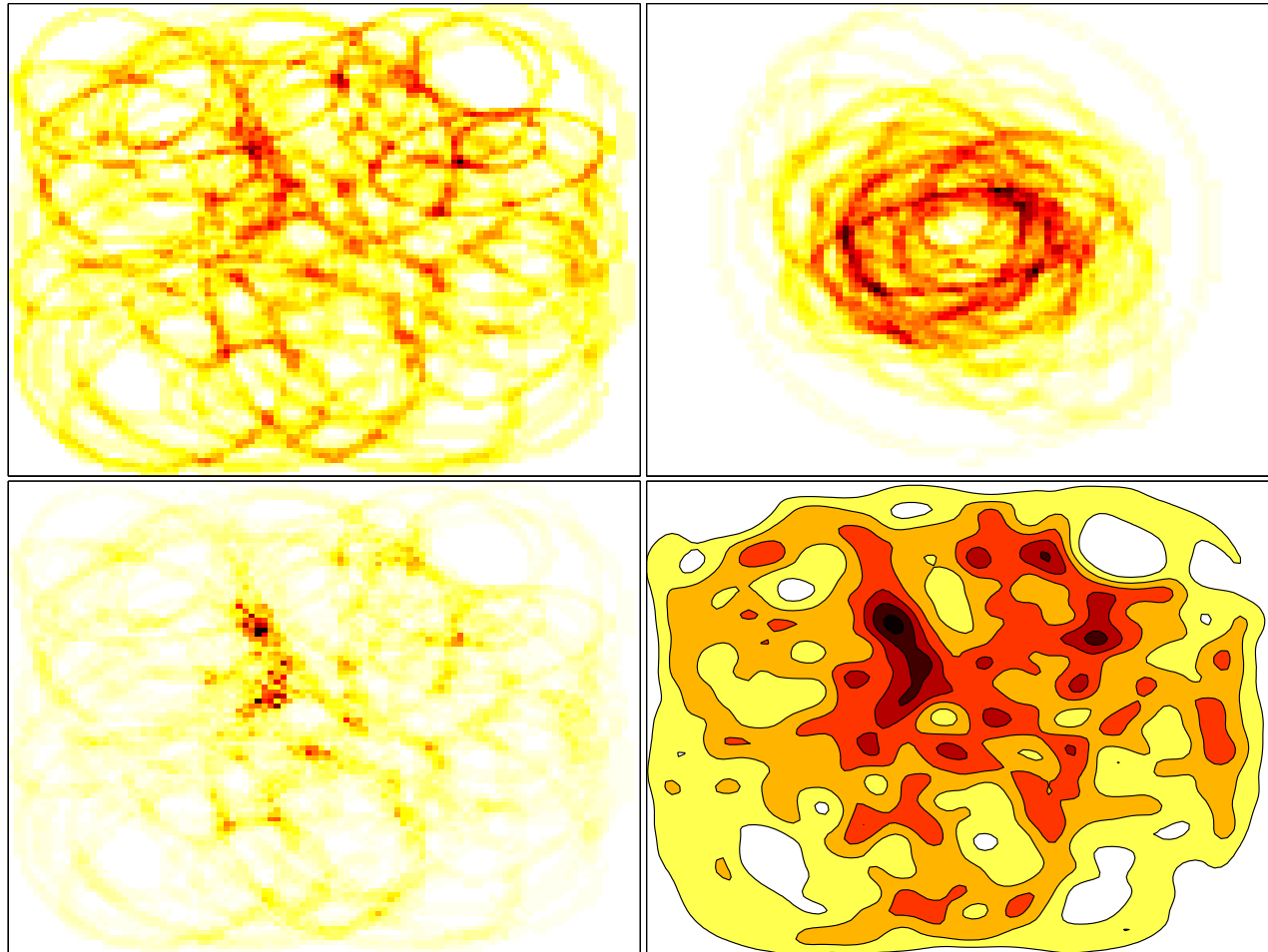


# Motivation in 2D



$$\{\nu_1, \dots, \nu_{30}\} \in P([0, 1]^2).$$

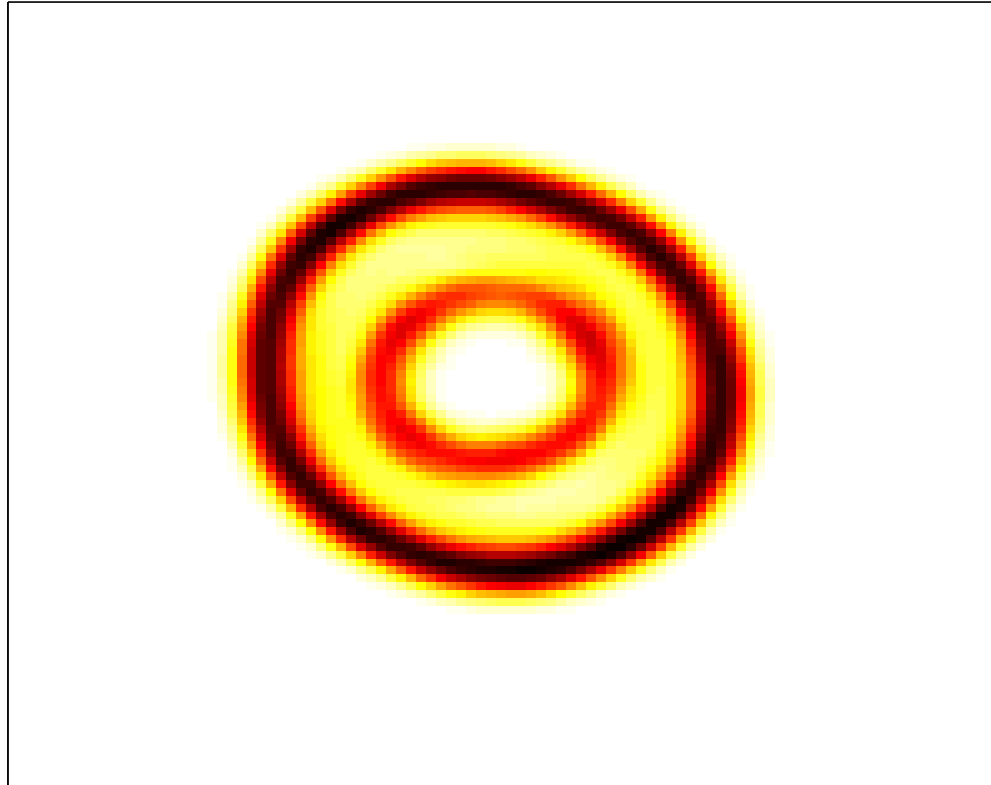
# Euclidean / Centered / Jeffrey / RKHS



Euclidean distance / recentered,  
Sym. Kullback / RKHS Mean Map



# Wasserstein Barycenter

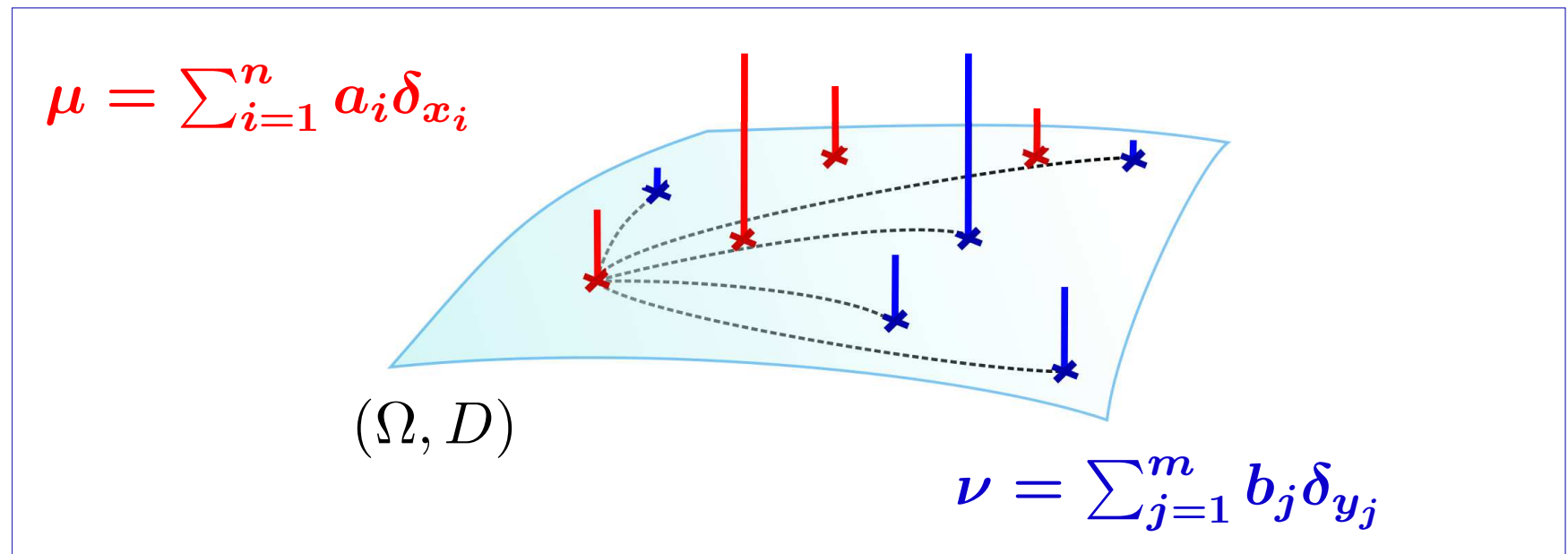


2-Wasserstein barycenter  
**(computed with our method)**

---

# Variational Perspective on the Wasserstein Distance

# Wasserstein for Empirical Measures



- $(\Omega, D)$  metric.  $p \geq 1$ .
- Two empirical measures  $\mu, \nu$ .

$p$ -Wasserstein distance  $W_p(\mu, \nu)$ ?

# Computing $p$ Wasserstein Distances

$$\mu = \sum_{i=1}^n \mathbf{a}_i \delta_{\mathbf{x}_i}, \quad \nu = \sum_{j=1}^m \mathbf{b}_j \delta_{\mathbf{y}_j},$$

$W_p(\mu, \nu)$  is the solution of a linear program involving:

1.  $M_{\mathbf{X}\mathbf{Y}} \stackrel{\text{def}}{=} [D(\mathbf{x}_i, \mathbf{y}_j)^p]_{ij} \in \mathbb{R}^{n \times m}$
2.  $U(\mathbf{a}, \mathbf{b}) \stackrel{\text{def}}{=} \{T \in \mathbb{R}_+^{n \times m} \mid T\mathbf{1}_m = \mathbf{a}, T^T\mathbf{1}_n = \mathbf{b}\}.$

# Computing the OT Distance

- $p$ -Wasserstein is the solution (primal or dual LP):

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

# Computing the OT Distance

- $p$ -Wasserstein is the solution (primal or dual LP):

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

Changes in  $f(\mathbf{a}, \mathbf{X}) \stackrel{\text{def}}{=} W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu})$  as  $\mathbf{a}$  &  $\mathbf{X}$  change?

# Wasserstein (Sub)differentiability

$$f(\mathbf{a}, X) = \max_{(\alpha, \beta) \in C_{M_{\mathbf{X}\mathbf{Y}}}} \alpha^T \mathbf{a} + \beta^T \mathbf{b}$$

- $\partial f|_{\mathbf{a}} = \boldsymbol{\alpha}^*$ : the *dual optimum*  $\boldsymbol{\alpha}^*$  is a subgradient.

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

- $\partial f|_X = Y \mathbf{T}^{*T} \text{diag}(\mathbf{a}^{-1})$ : *primal optimum*  $\mathbf{T}^{*T}$  yields a subgradient (when  $D = \text{Euclidean}$ ,  $p = 2$ ).

# Average of Wasserstein Distances

$$\begin{aligned} g(\mathbf{a}, \mathbf{X}) &\stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}_i) \\ &= \frac{1}{N} \sum_{i=1}^N \text{primal}(\mathbf{a}, \mathbf{b}_i, M_{\mathbf{X}\mathbf{Y}_i}) \end{aligned}$$

- $a \rightarrow g(a, X)$  is **convex**, non-smooth
- $X \rightarrow g(a, X)$  is **not convex**, non-smooth



# Wasserstein Barycenter Problem

$$\min_{\mathbf{a}} g(\mathbf{a}, X) = \frac{1}{N} \sum_{i=1}^N \text{primal}(\mathbf{a}, \mathbf{b}_i, M_{X\mathbf{Y}_i})$$

- $\mathbf{a} \rightarrow g(\mathbf{a}, X)$  is **convex**
  - subgradient method works (in theory).
  - Great if  $X$  is fixed (b-o-w or discretized  $\Omega$ )!
  - **Need to solve  $\{\alpha_i^*\}$  at each subgradient step.**

# Wasserstein Barycenter Problem

$$\min_{\mathbf{X}} g(a, \mathbf{X}) = \frac{1}{N} \sum_{i=1}^N \text{primal}(a, \mathbf{b}_i, M_{\mathbf{X}\mathbf{Y}_i})$$

- $\mathbf{X} \rightarrow g(a, \mathbf{X})$  is **not convex**
  - (and so far only applicable when  $\Omega$  is  $\mathbb{R}^d$ .)
  - local minimum with subgradient method
  - **Need to compute  $\{T_i^*\}$  at each subgradient step**

# To recapitulate...

$$\min_{a, X} g(a, X)$$

- **convex** w.r.t weights  $a$ , not locations  $X$ .
- **only subgradients** ( $g$  is usually very degenerate).
- **computationally intractable** (cost of OT  $\approx n^3 \log n$ )
- **computationally inefficient** (hard to parallelize)

# Solution: Entropic Smoothing

**Original** primal problem gives us  $T^*$ :

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

**Original** dual problem gives us  $\alpha^*$ :

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

# Solution: Entropic Smoothing

**Smoothed** ( $\lambda > 0$ ) primal problem gives us  $T_\lambda^*$ :

$$\text{primal}_\lambda(a, b, M_{XY}) = \min_{T \in U(a,b)} \langle T, M_{XY} \rangle - \frac{1}{\lambda} h(T)$$

**Smoothed** dual problem gives us  $\alpha_\lambda^*$ :

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

# Benefits of Smoothing [Cuturi'13]

- Objective now **strongly convex** vs. **piecewise linear**: infinitely more efficient in practice [**Nesterov'05**].
- Primal/dual smoothed optima  $\alpha_\lambda^*$ ,  $T_\lambda^*$  can be solved
  - **In  $O(n^2)$**  with **Sinkhorn's algorithm**,
  - in **parallel on GPGPUs** for **any metric** on finite  $\Omega$ ,
  - **millions of time faster** than simplex,
  - can deal with **large dimensions** ( $\approx 20.000$ ).

# To conclude...

- our approach also **generalizes  $k$ -means**
  - can consider weight constraints (see paper),
  - can quantize simultaneously different datasets
- Versatile and scalable approach for **other variational Wasserstein problems** (*e.g.* Wasserstein propagation [**Solomon'14**])
- Future applications to visualization of measures on Riemannian manifolds, data-fusion, inference...