

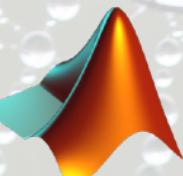
# Numerical Optimal Transport

<http://optimaltransport.github.io>

## *Entropic Regularization*

Gabriel Peyré

[www.numerical-tours.com](http://www.numerical-tours.com)



**ENS**

ÉCOLE NORMALE  
SUPÉRIEURE

# Overview

---

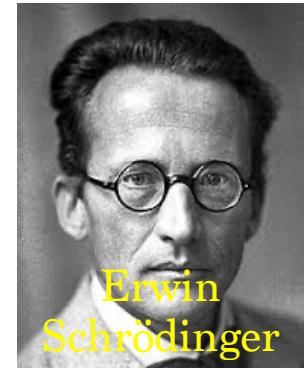
- **Entropic Regularization and Sinkhorn**
- Convergence Analysis
- Sinkhorn Divergences
- Generative Model Fitting

# Entropic Regularization

*Schrödinger's problem:*

[1931]

$$\min_{\mathbf{P} \mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b}} \sum_{i,j} d(\mathbf{x}_i, \mathbf{y}_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j})$$

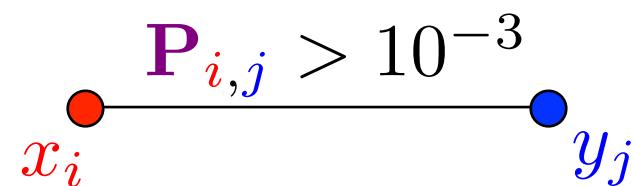
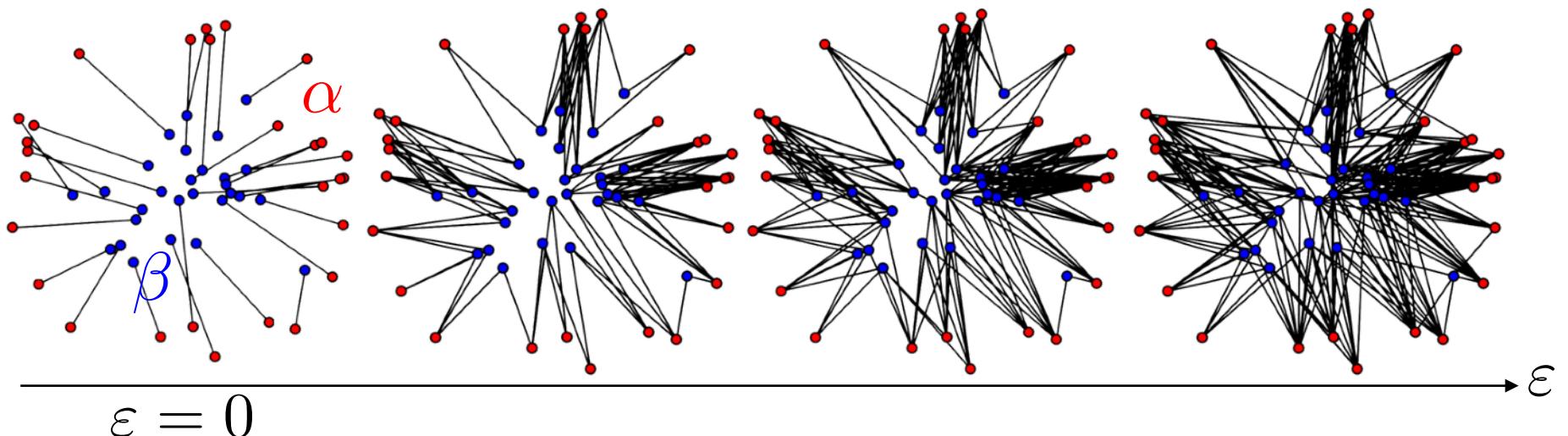
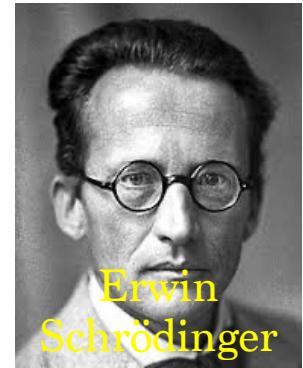


# Entropic Regularization

Schrödinger's problem:

[1931]

$$\min_{\mathbf{P} \mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b}} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j})$$



# Entropic Regularization: General Case

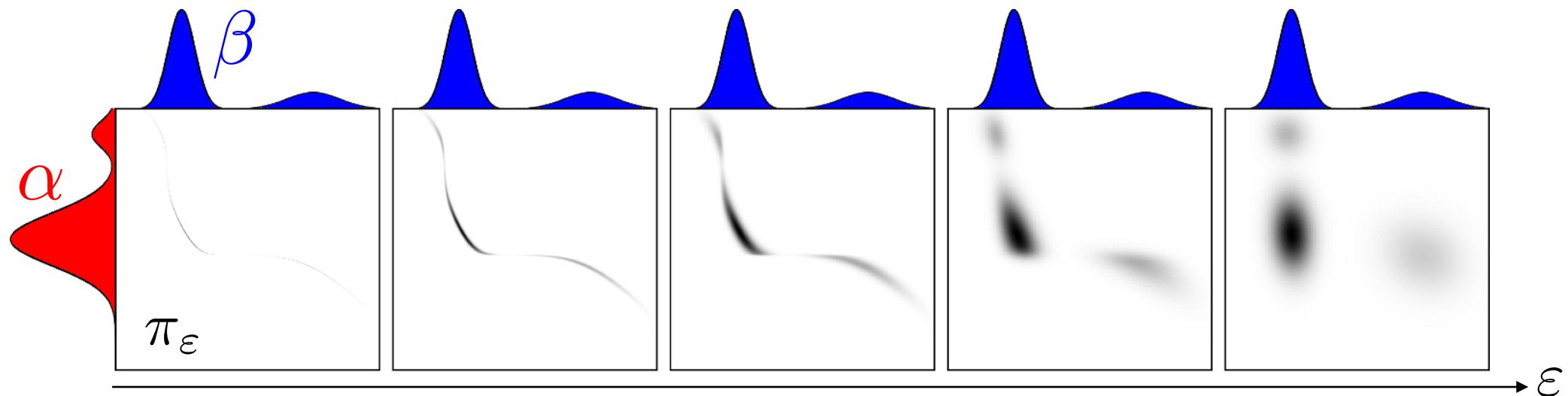
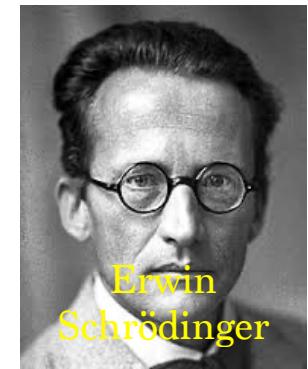
$$\pi = \sum_{i,j} \mathbf{P}_{i,j} \delta_{\mathbf{x}_i, \mathbf{y}_j}$$

Relative-entropy:  $\text{KL}(\pi | \alpha \otimes \beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}^2} \log \left( \frac{d\pi}{d\alpha d\beta}(x, y) \right) d\pi(x, y)$

Schrödinger's problem:

[1931]

$$W_{\varepsilon, p}^p(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d^p(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta)$$

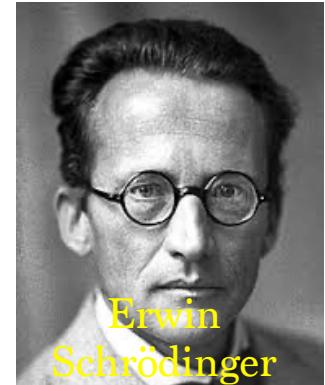


# Probabilistic Interpretation

Relative-entropy:  $\text{KL}(\pi|\alpha \otimes \beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}^2} \log \left( \frac{d\pi}{d\alpha d\beta}(x, y) \right) d\pi(x, y)$

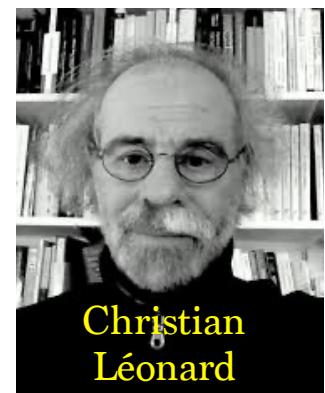
Schrödinger's problem: [1931]

$$W_{\varepsilon, p}^p(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d^p(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi|\alpha \otimes \beta)$$



$$\min_{(X, Y)} \{ \mathbb{E}(c(X, Y)) + \varepsilon I(X, Y) ; X \sim \alpha, Y \sim \beta \}$$

Mutual information

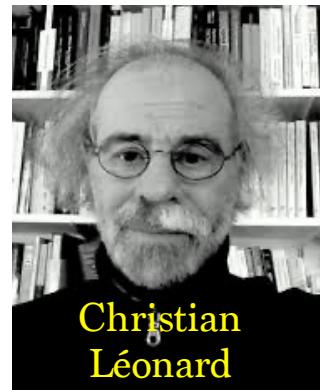
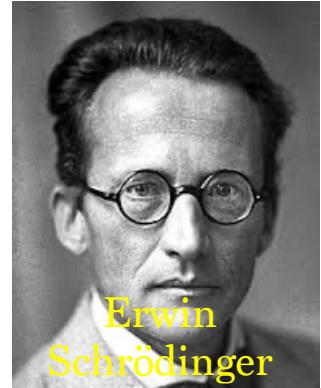


# Probabilistic Interpretation

Relative-entropy:  $\text{KL}(\pi|\alpha \otimes \beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}^2} \log \left( \frac{d\pi}{d\alpha d\beta}(x, y) \right) d\pi(x, y)$

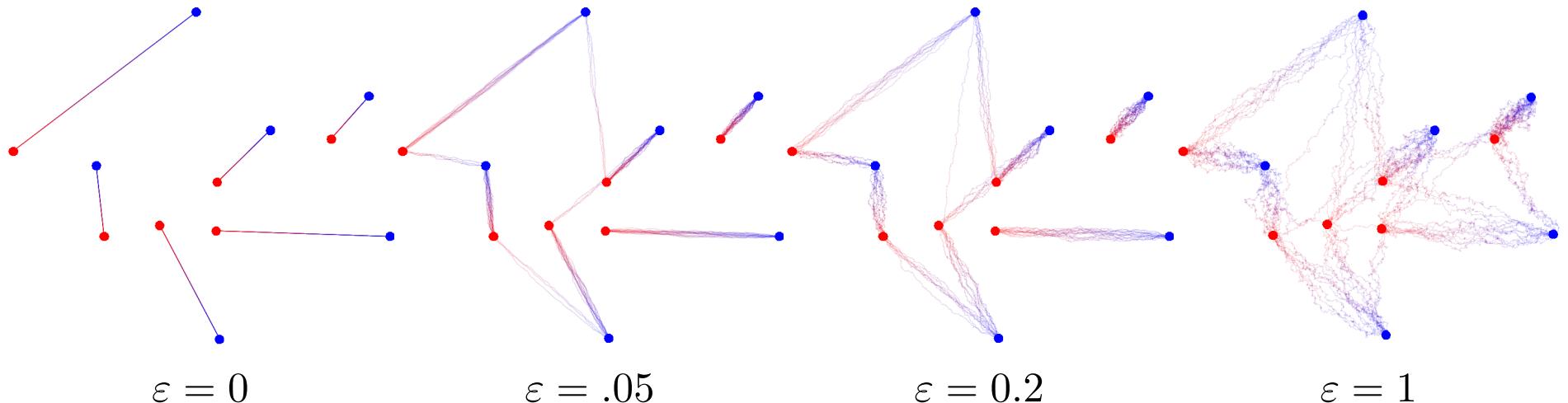
Schrödinger's problem: [1931]

$$W_{\varepsilon, p}^p(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d^p(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi|\alpha \otimes \beta)$$

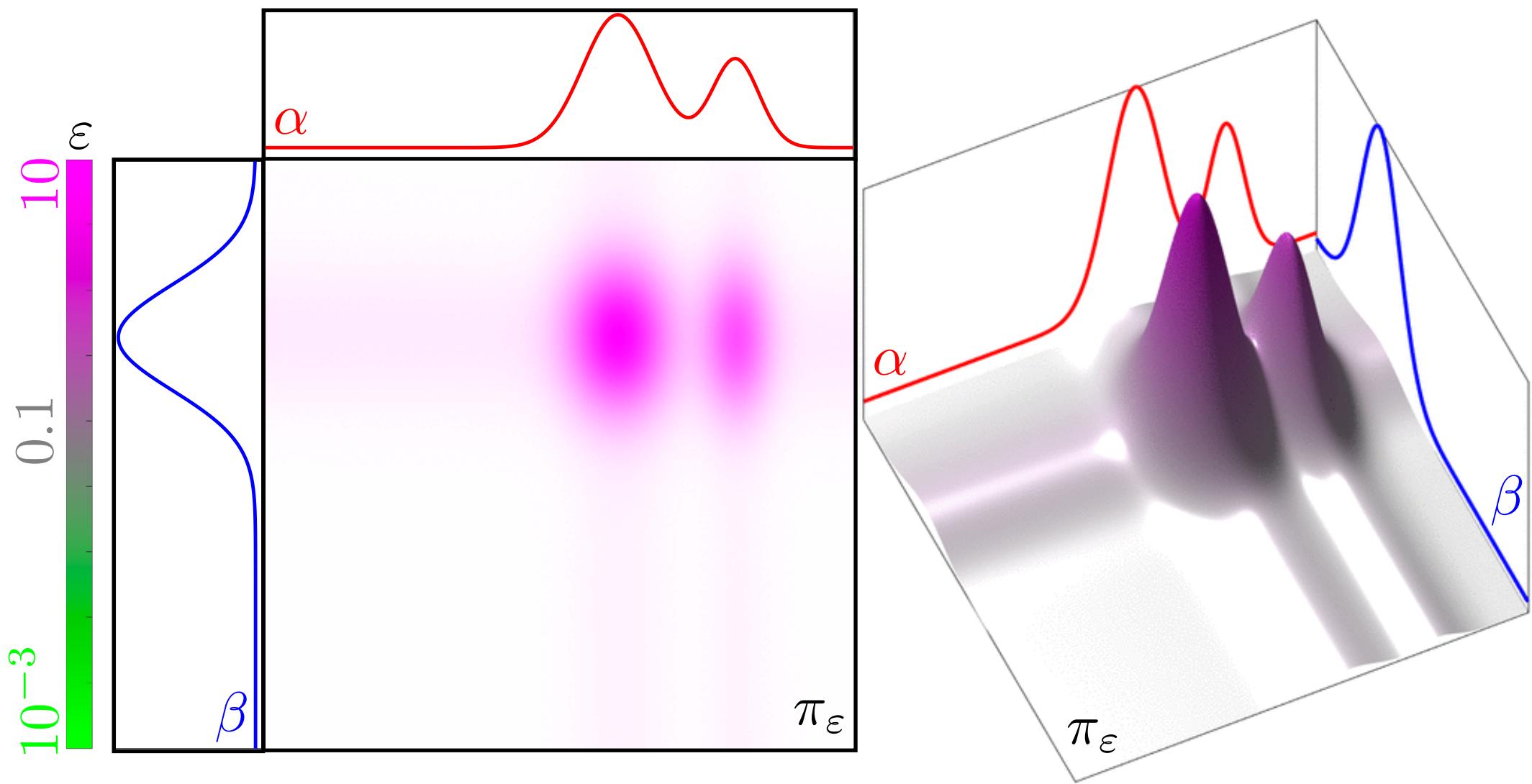


$$\min_{(X, Y)} \{ \mathbb{E}(c(X, Y)) + \varepsilon I(X, Y) ; X \sim \alpha, Y \sim \beta \}$$

Mutual information



# Impact of Regularization



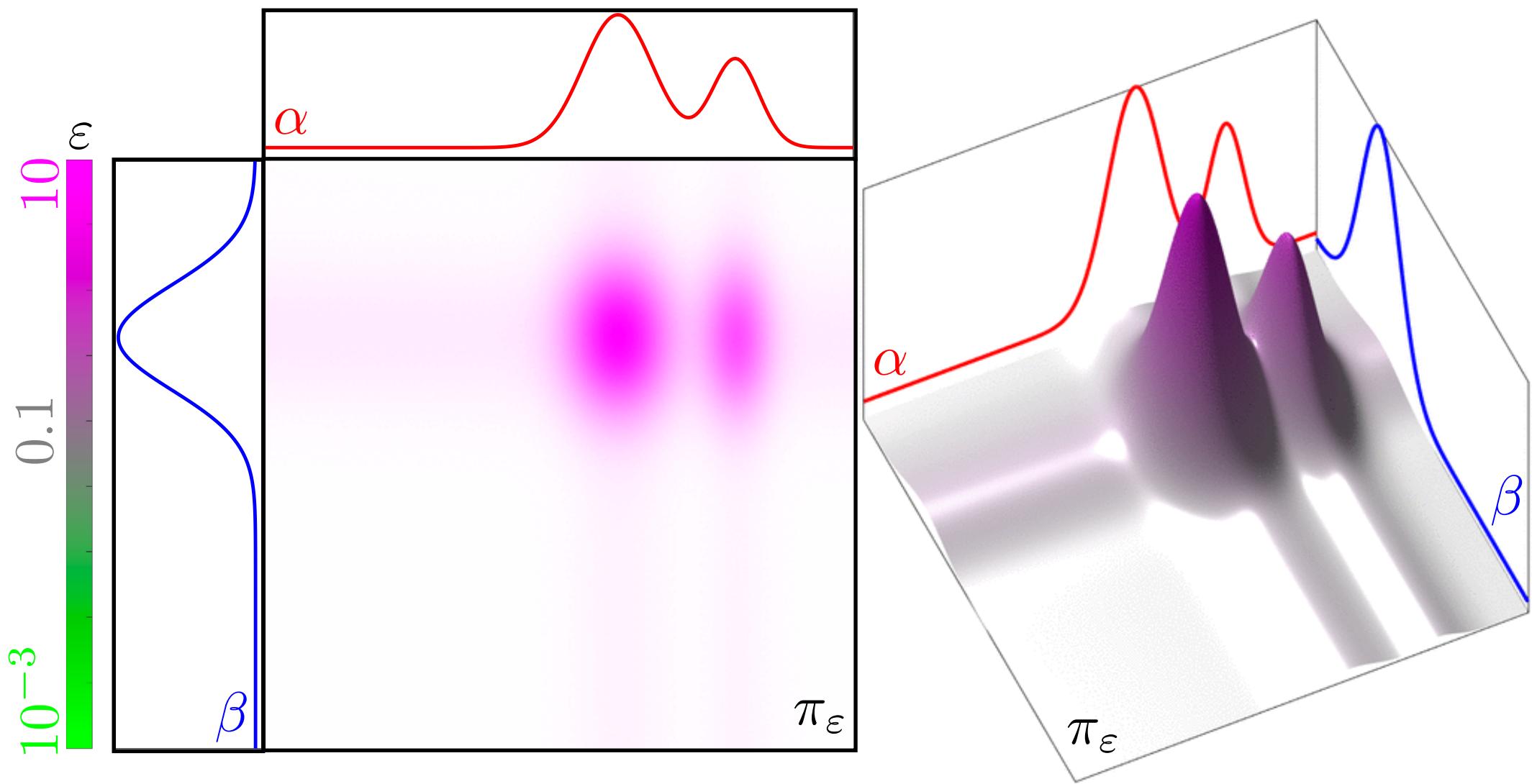
$$\pi_\varepsilon = \operatorname{argmin}_\pi \left\{ \int_{\mathbb{R}^2} \left( \|x - y\|^2 + \varepsilon \log \left( \frac{d\pi}{d\alpha d\beta}(x, y) \right) \right) d\pi(x, y) ; \pi_1 = \alpha, \pi_2 = \beta \right\}$$

*Theorem:*

$$\pi_\varepsilon \xrightarrow{\varepsilon \rightarrow +\infty} \alpha \otimes \beta$$

$$\pi_\varepsilon \xrightarrow{\varepsilon \rightarrow 0} \pi_{\text{OT}}$$

# Impact of Regularization



$$\pi_\varepsilon = \operatorname{argmin}_\pi \left\{ \int_{\mathbb{R}^2} \left( \|x - y\|^2 + \varepsilon \log \left( \frac{d\pi}{d\alpha d\beta}(x, y) \right) \right) d\pi(x, y) ; \pi_1 = \alpha, \pi_2 = \beta \right\}$$

*Theorem:*

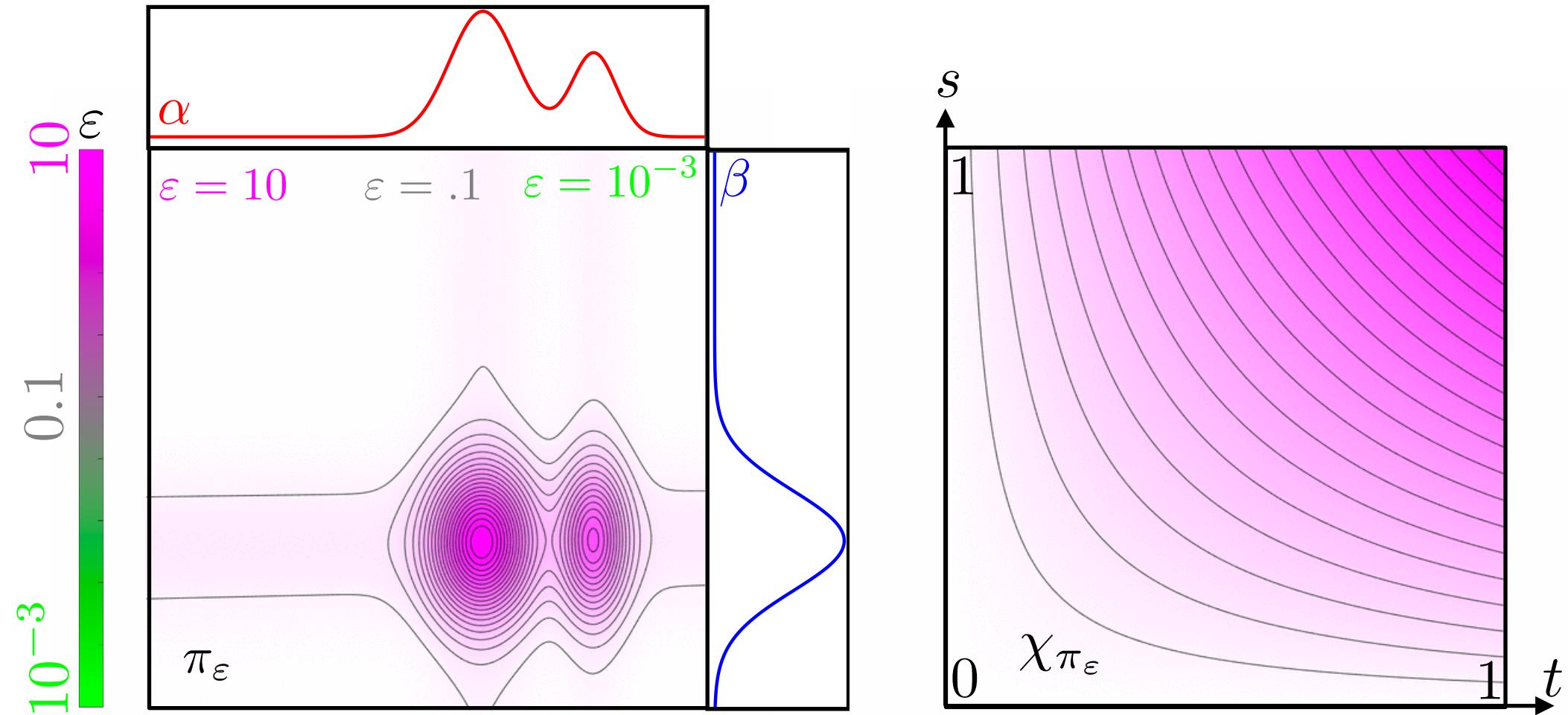
$$\pi_\varepsilon \xrightarrow{\varepsilon \rightarrow +\infty} \alpha \otimes \beta$$

$$\pi_\varepsilon \xrightarrow{\varepsilon \rightarrow 0} \pi_{\text{OT}}$$

# Impact of Regularization

Cumulative:  $C_\pi(x, y) \stackrel{\text{def.}}{=} \int_{-\infty}^x \int_{-\infty}^y d\pi(x, y)$

Copula:  $\chi_\pi(s, t) \stackrel{\text{def.}}{=} C_\pi(C_\alpha^{-1}(s), C_\beta^{-1}(t))$

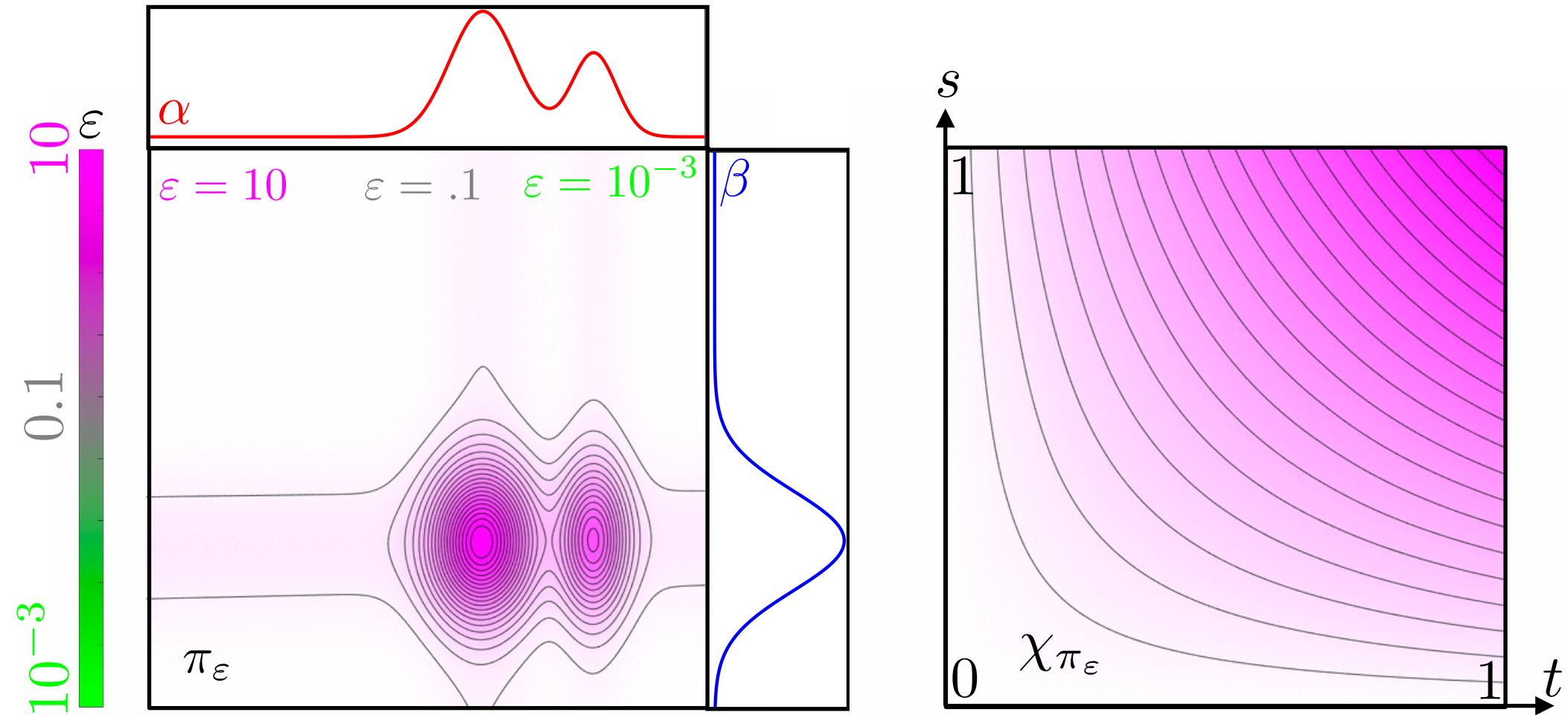


Theorem:  $\chi_{\pi_\varepsilon}(s, t)$   $\begin{cases} \xrightarrow{\varepsilon \rightarrow 0} \min(s, t) & \text{(dependence)} \\ \xrightarrow{\varepsilon \rightarrow +\infty} st & \text{(independence)} \end{cases}$

# Impact of Regularization

Cumulative:  $C_\pi(x, y) \stackrel{\text{def.}}{=} \int_{-\infty}^x \int_{-\infty}^y d\pi(x, y)$

Copula:  $\chi_\pi(s, t) \stackrel{\text{def.}}{=} C_\pi(C_\alpha^{-1}(s), C_\beta^{-1}(t))$



Theorem:  $\chi_{\pi_\varepsilon}(s, t)$   $\begin{cases} \xrightarrow{\varepsilon \rightarrow 0} \min(s, t) & \text{(dependence)} \\ \xrightarrow{\varepsilon \rightarrow +\infty} st & \text{(independence)} \end{cases}$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(\mathbf{x}_i, \mathbf{y}_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P} \mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*  $\begin{cases} \mathbf{P} \mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \mathbf{P} \text{ solution} \Leftrightarrow \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases}$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(\mathbf{x}_i, \mathbf{y}_j)^p}{\varepsilon}}$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(\mathbf{x}_i, \mathbf{y}_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*  $\begin{cases} \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \mathbf{P} \text{ solution} \Leftrightarrow \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases}$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(\mathbf{x}_i, \mathbf{y}_j)^p}{\varepsilon}}$

$$\mathbf{P} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \implies \mathbf{a} = \mathbf{P}\mathbf{1} = \text{diag}(\mathbf{u})(\mathbf{K}\mathbf{v}) = \mathbf{u} \odot (\mathbf{K}\mathbf{v})$$

Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$

Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(\mathbf{x}_i, \mathbf{y}_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*  $\begin{cases} \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \mathbf{P} \text{ solution} \Leftrightarrow \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases}$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(\mathbf{x}_i, \mathbf{y}_j)^p}{\varepsilon}}$

$$\mathbf{P} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \implies \mathbf{a} = \mathbf{P}\mathbf{1} = \text{diag}(\mathbf{u})(\mathbf{K}\mathbf{v}) = \mathbf{u} \odot (\mathbf{K}\mathbf{v})$$

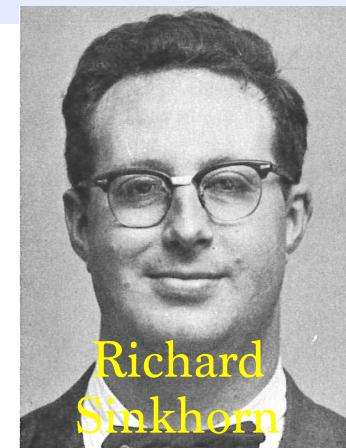
Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$

Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

Sinkhorn iterations:

$$\mathbf{u} \leftarrow \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}}$$

$$\mathbf{v} \leftarrow \frac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}}$$



*Theorem:* [Sinkhorn 1964]  $(\mathbf{u}, \mathbf{v})$  converges.

# Sinkhorn's Algorithm

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(\mathbf{x}_i, \mathbf{y}_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*  $\begin{cases} \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \mathbf{P} \text{ solution} \Leftrightarrow \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases}$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(\mathbf{x}_i, \mathbf{y}_j)^p}{\varepsilon}}$

$$\mathbf{P} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \implies \mathbf{a} = \mathbf{P}\mathbf{1} = \text{diag}(\mathbf{u})(\mathbf{K}\mathbf{v}) = \mathbf{u} \odot (\mathbf{K}\mathbf{v})$$

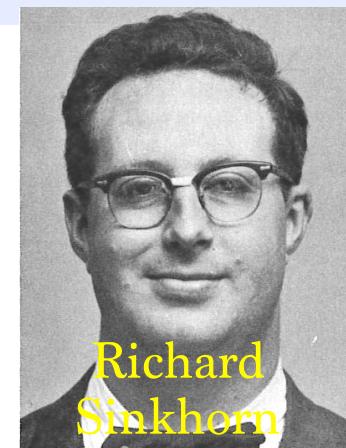
Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$

Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

Sinkhorn iterations:

$$\mathbf{u} \leftarrow \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}}$$

$$\mathbf{v} \leftarrow \frac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}}$$



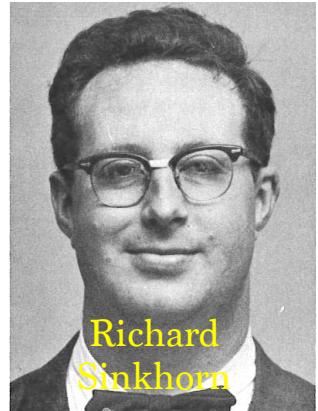
*Theorem:* [Sinkhorn 1964]  $(\mathbf{u}, \mathbf{v})$  converges.

Matrix/vector multiplications:  $\rightarrow O(n^2/\varepsilon^2)$  complexity.

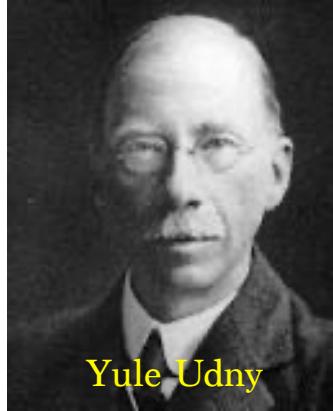
$\rightarrow$  Parallelizable on GPUs.

$\rightarrow$  Convolution on regular grids, separable kernels.

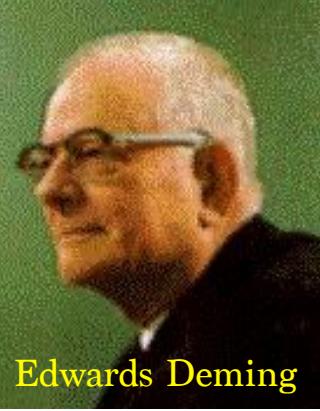
# Sinkhorn, IPFP, RAS, ...



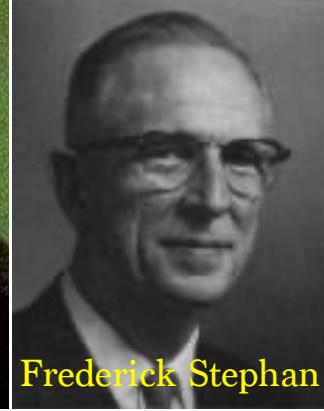
Richard  
Sinkhorn



Yule Udney



Edwards Deming



Frederick Stephan

## Many names:

Sinkhorn algorithm

Udny 1912

DAD scaling

Kruithof, 1937

Iterative proportional fitting

Deming and Stephan, 1940

Biproportional fitting

Sinkhorn 1964

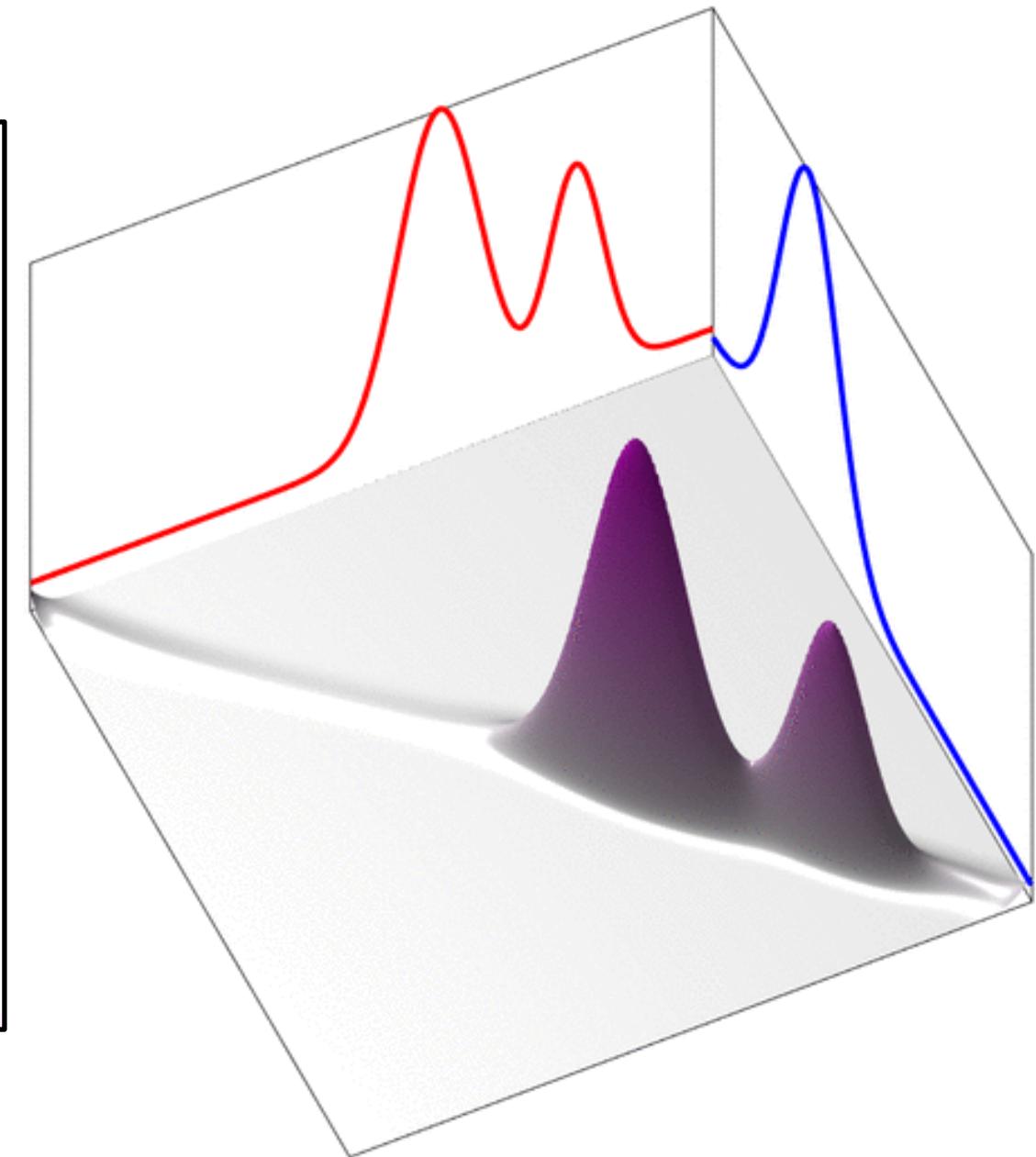
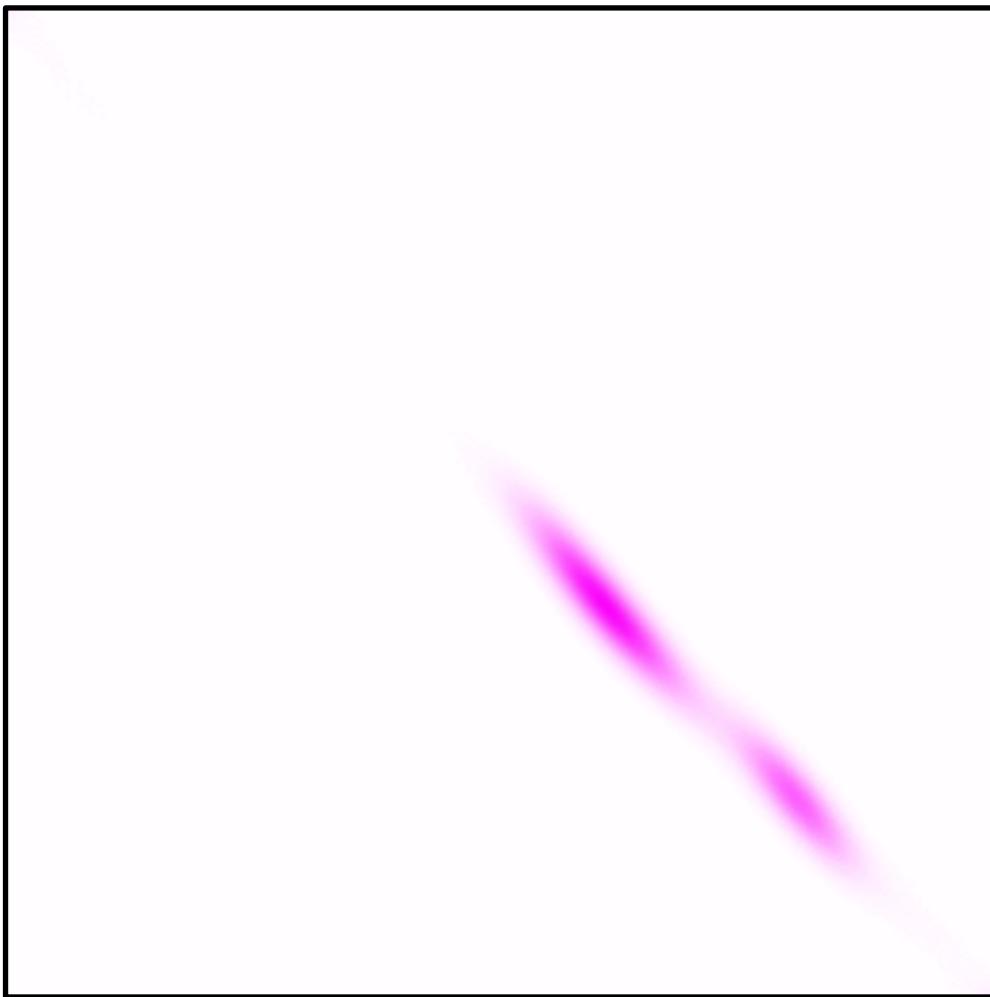
RAS algorithm

Matrix scaling

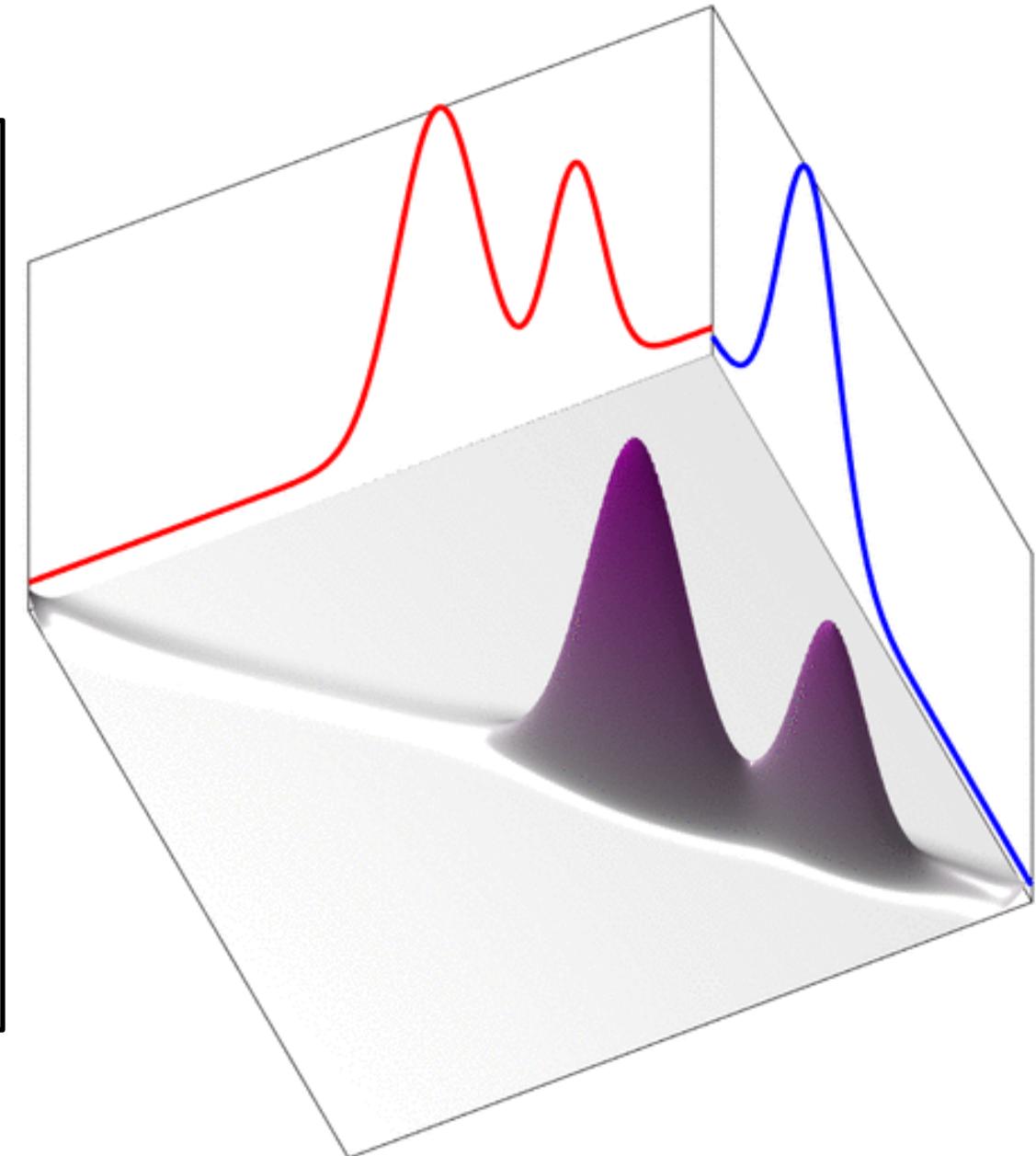
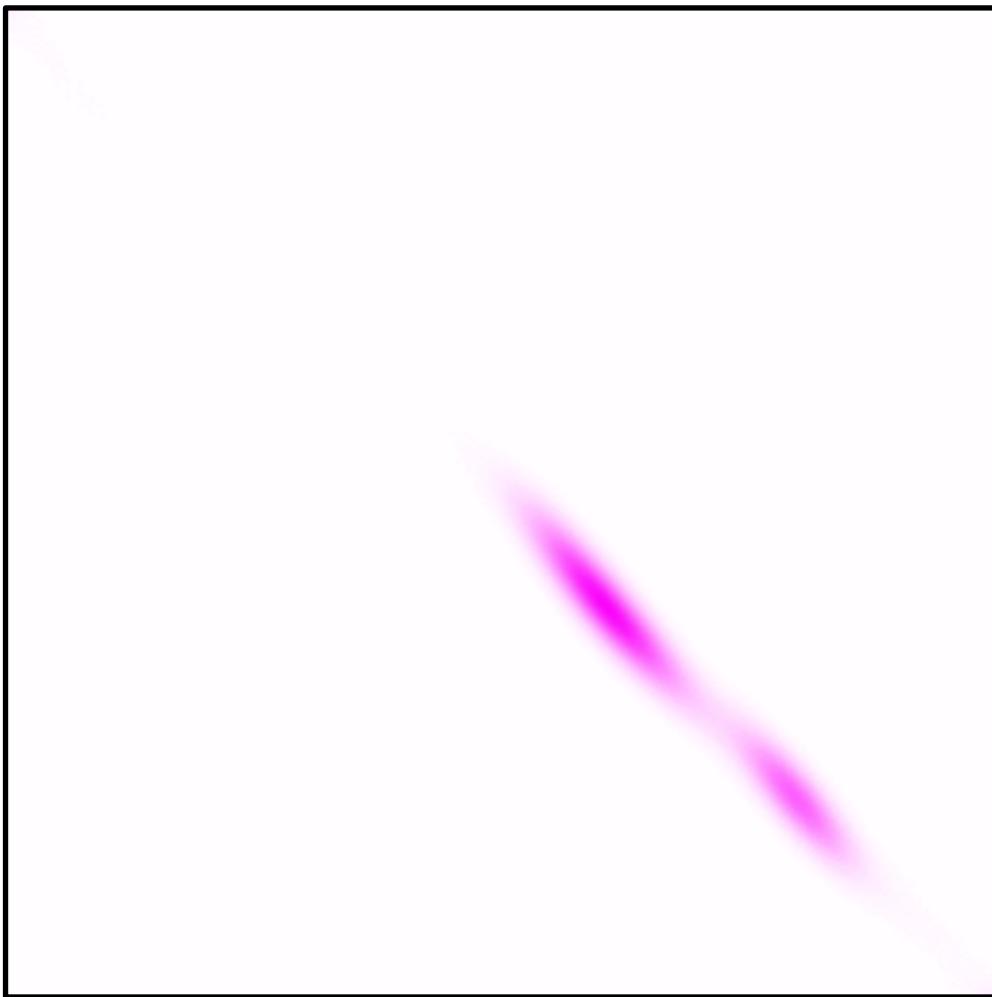
Gravity model

*Special thanks to Tony Silvetti-Falls for the picture of Sinkhorn*

# Sinkhorn Evolution

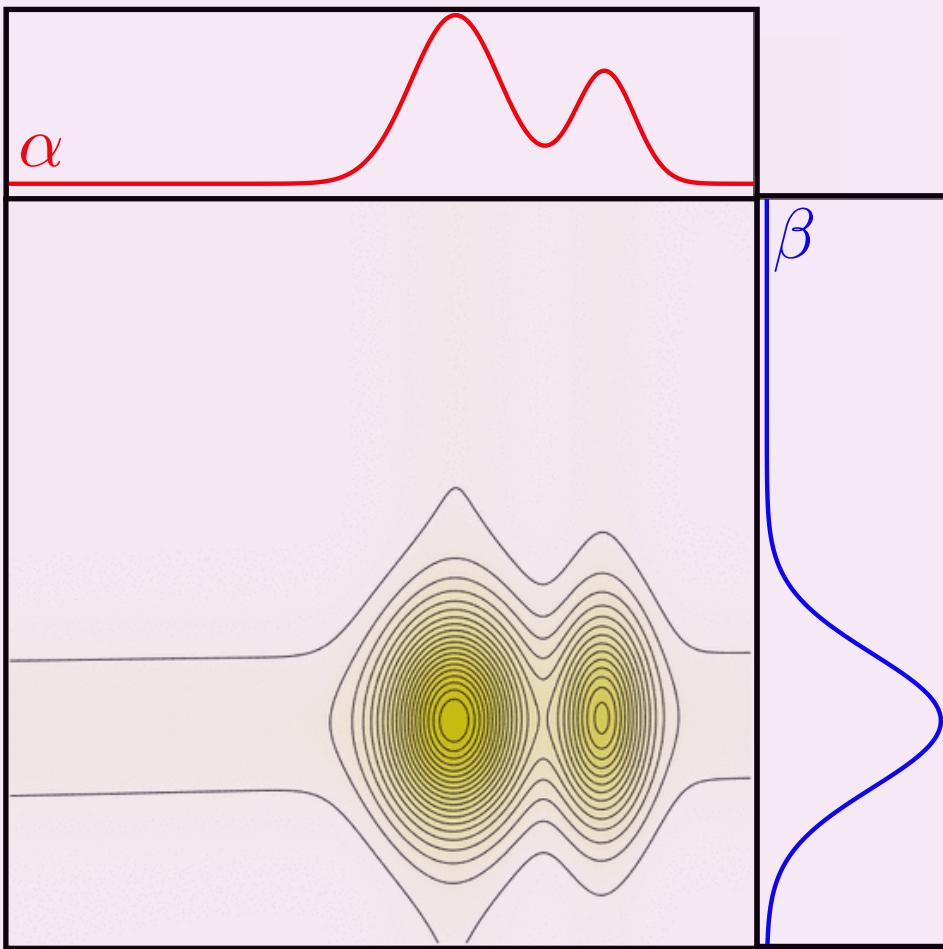


# Sinkhorn Evolution



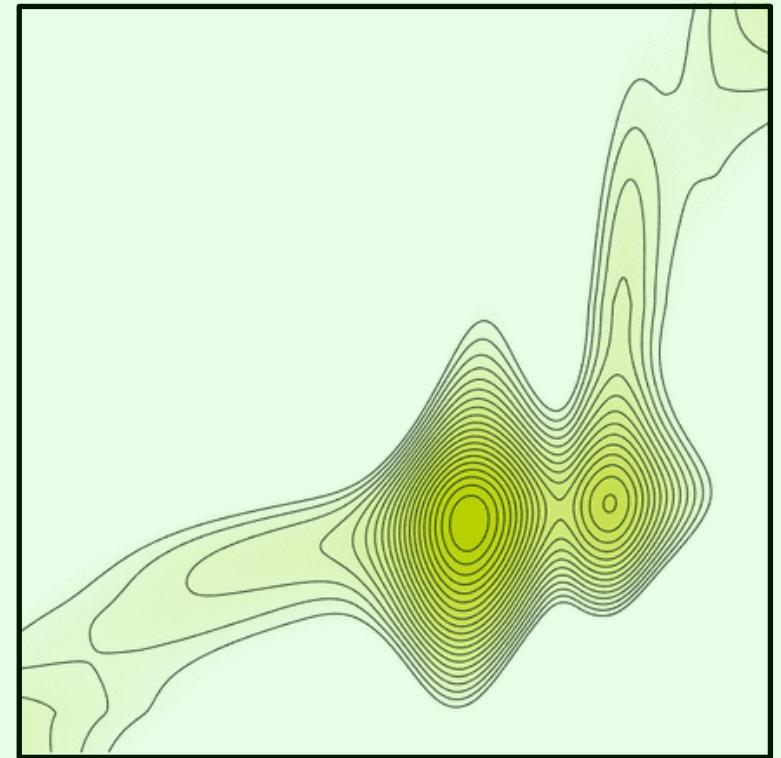
# Other Regularizations

$$\min_{\pi} \left\{ \int_{\mathbb{R}^2} \|x - y\|^2 d\pi(x, y) + \varepsilon R(\pi) ; \pi_1 = \alpha, \pi_2 = \beta \right\}$$



$$R(\pi) = \int \log \left( \frac{d\pi}{dxdy} \right) d\pi(x, y)$$

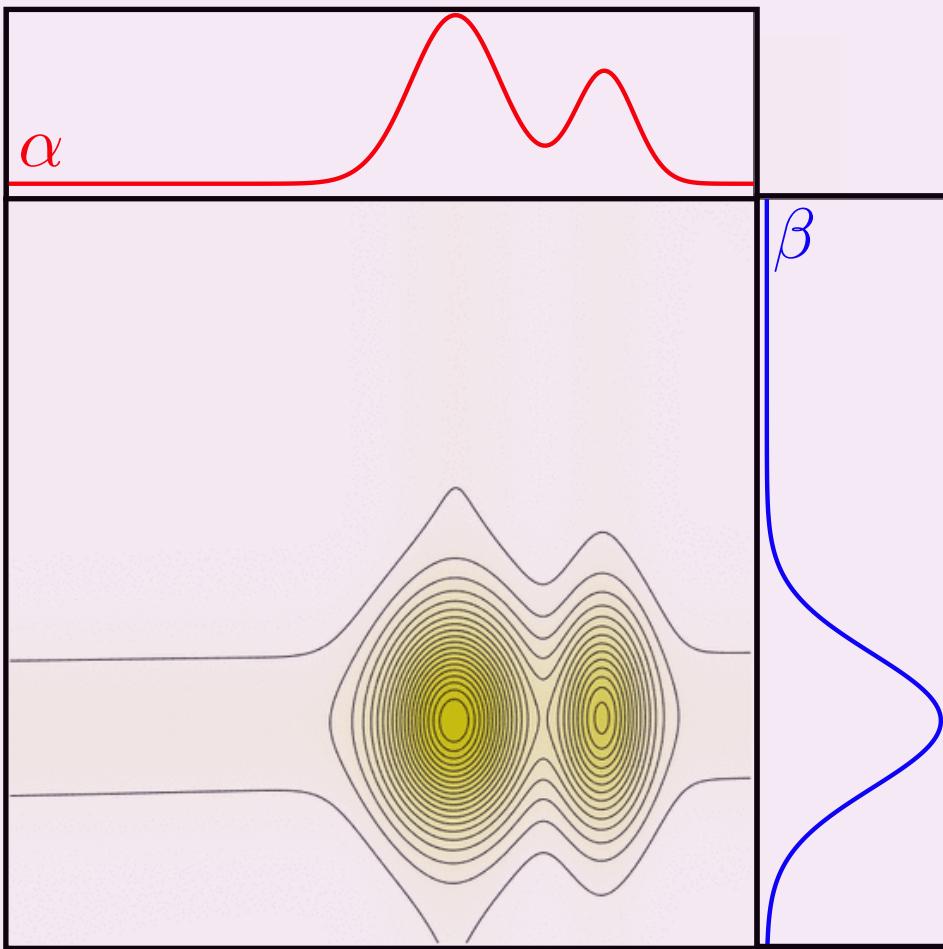
Dykstra's algorithm



$$R(\pi) = \int \left( \frac{d\pi}{dxdy} \right)^2 dx dy$$

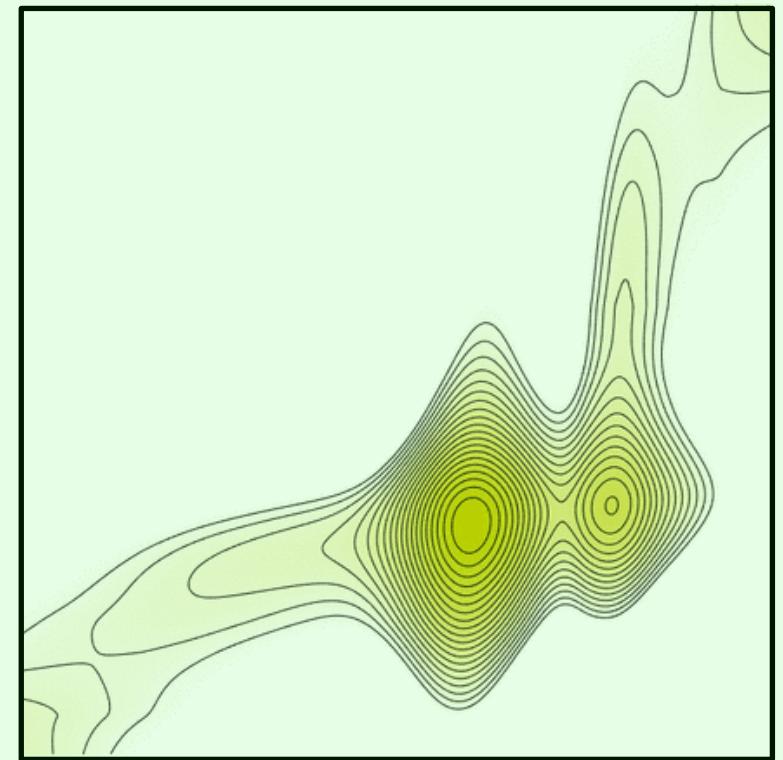
# Other Regularizations

$$\min_{\pi} \left\{ \int_{\mathbb{R}^2} \|x - y\|^2 d\pi(x, y) + \varepsilon R(\pi) ; \pi_1 = \alpha, \pi_2 = \beta \right\}$$



$$R(\pi) = \int \log \left( \frac{d\pi}{dxdy} \right) d\pi(x, y)$$

Dykstra's algorithm



$$R(\pi) = \int \left( \frac{d\pi}{dxdy} \right)^2 dx dy$$

# Extension: Unbalanced OT

$$W_p^{\tau,p}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi} \int d^p d\pi + \tau \text{KL}(\pi_1 | \alpha) + \tau \text{KL}(\pi_2 | \beta)$$

[Liero, Mielke, Savaré 2015]

See also:

[Chizat, Schmitzer, Peyré, Vialard 2015]  
[Kondratyev, Monsaingeon, Vorotnikov 2015]

# Extension: Unbalanced OT

$$W_p^{\tau,p}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi} \int d^p d\pi + \tau \text{KL}(\pi_1 | \alpha) + \tau \text{KL}(\pi_2 | \beta)$$

[Liero, Mielke, Savaré 2015]

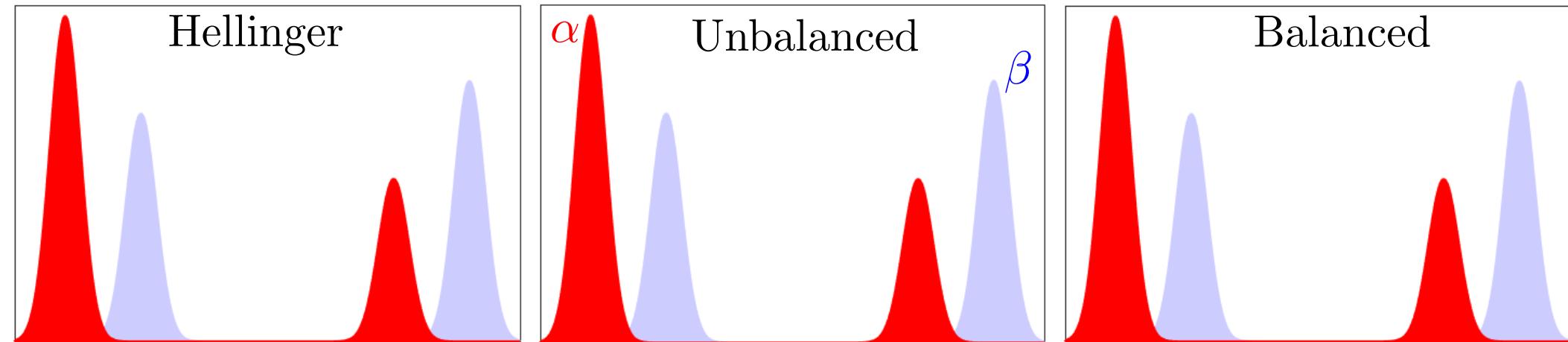
See also: [Chizat, Schmitzer, Peyré, Vialard 2015]  
[Kondratyev, Monsaingeon, Vorotnikov 2015]

$$\int (\sqrt{\alpha} - \sqrt{\beta})^2 \xleftarrow{\tau \rightarrow 0} W_p^{\tau,p}(\alpha, \beta) \xrightarrow{\tau \rightarrow +\infty} W_p^p(\alpha, \beta)$$

Hellinger

Unbalanced

Balanced



# Extension: Unbalanced OT

$$W_p^{\tau,p}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi} \int d^p d\pi + \tau \text{KL}(\pi_1 | \alpha) + \tau \text{KL}(\pi_2 | \beta)$$

[Liero, Mielke, Savaré 2015]

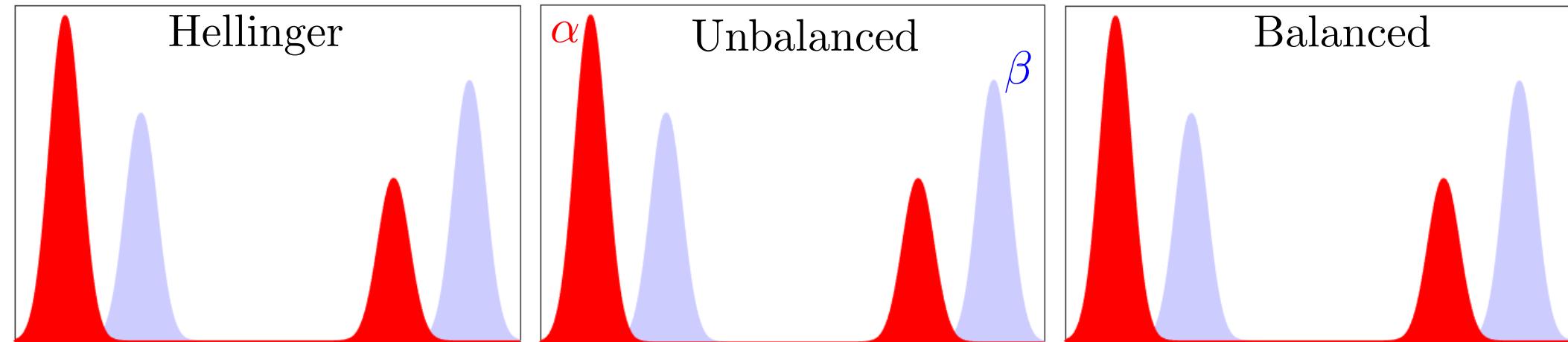
See also: [Chizat, Schmitzer, Peyré, Vialard 2015]  
 [Kondratyev, Monsaingeon, Vorotnikov 2015]

$$\int (\sqrt{\alpha} - \sqrt{\beta})^2 \xleftarrow{\tau \rightarrow 0} W_p^{\tau,p}(\alpha, \beta) \xrightarrow{\tau \rightarrow +\infty} W_p^p(\alpha, \beta)$$

Hellinger

Unbalanced

Balanced



$$W_{\varepsilon,p}^{\tau,p}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi} \int d^p d\pi + \tau \text{KL}(\pi_1 | \alpha) + \tau \text{KL}(\pi_2 | \beta) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta)$$

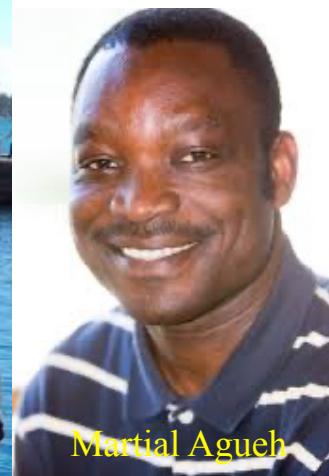
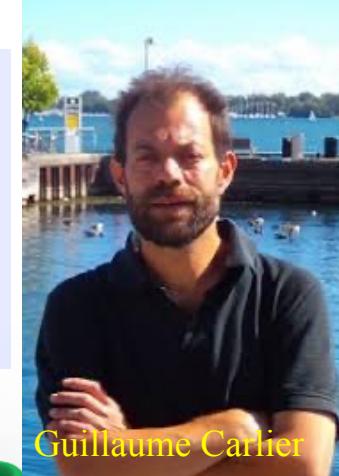
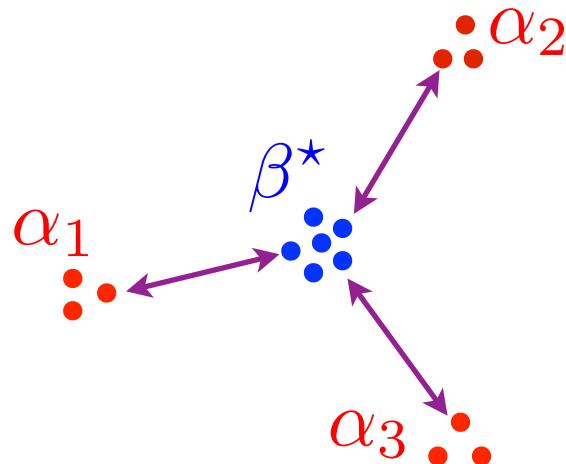
Sinkhorn's algorithm:

$$\mathbf{u} \leftarrow \left( \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}} \right)^{1+\frac{\varepsilon}{\tau}} \quad \longleftrightarrow \quad \mathbf{v} \leftarrow \left( \frac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}} \right)^{1+\frac{\varepsilon}{\tau}}$$

# Extension: Wasserstein Barycenters

Barycenters of measures  $(\alpha_s)_s$ :  $\sum_s \lambda_s = 1$

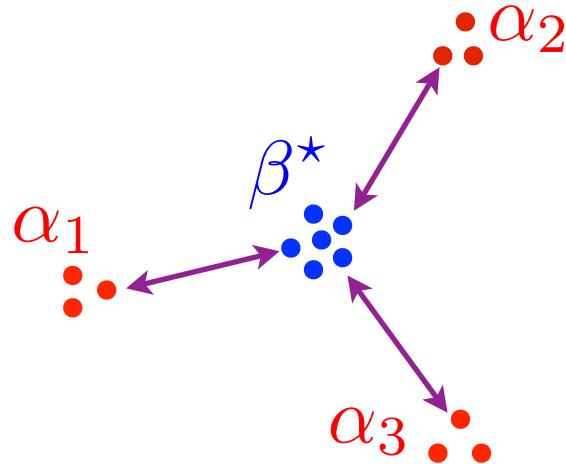
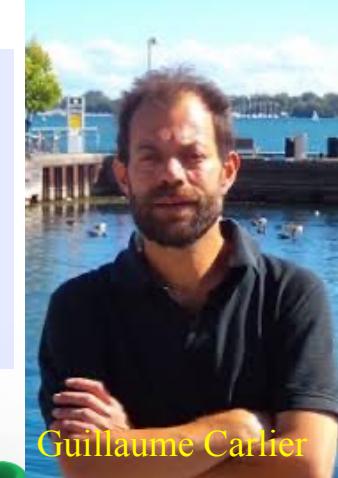
$$\beta^* \in \operatorname{argmin}_{\beta} \sum_s \lambda_s W_p^p(\alpha_s, \beta)$$



# Extension: Wasserstein Barycenters

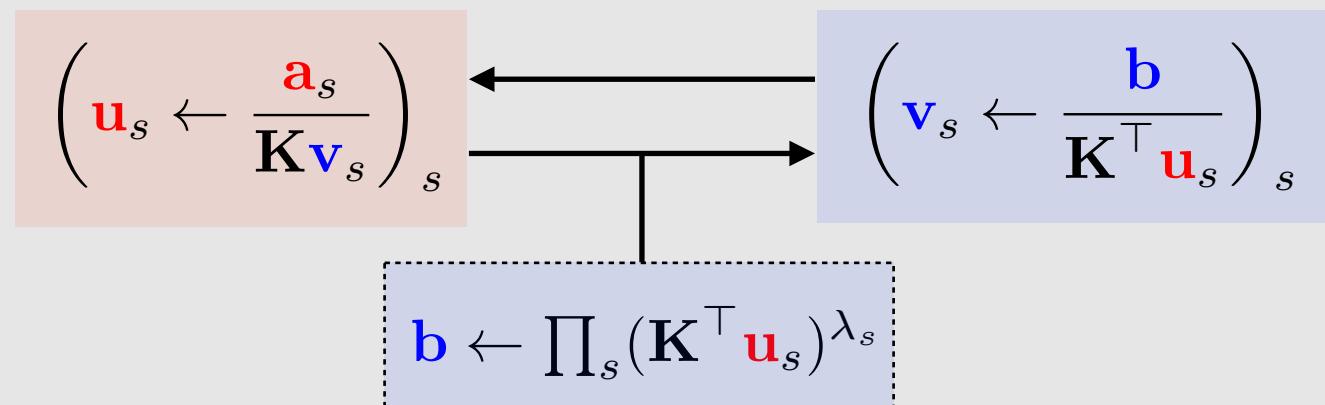
Barycenters of measures  $(\alpha_s)_s$ :  $\sum_s \lambda_s = 1$

$$\beta^* \in \operatorname{argmin}_{\beta} \sum_s \lambda_s W_p^p(\alpha_s, \beta)$$



[Solomon et al, SIGGRAPH 2015]

Sinkhorn's algorithm:



# Overview

---

- Entropic Regularization and Sinkhorn
- **Convergence Analysis**
- Sinkhorn Divergences
- Generative Model Fitting

# Bregman Iterative Projections

$$\langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \text{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) = \varepsilon \text{KL}(\mathbf{P}|\mathbf{K}) + \text{cst} \quad \text{where} \quad \mathbf{K}_{i,j} = e^{-\frac{\mathbf{C}_{i,j}}{\varepsilon}} \mathbf{a}_i \mathbf{b}_j$$

$$\text{Shrödinger problem: } \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \text{KL}(\mathbf{P}|\mathbf{K})$$

Constraints :

$$\mathcal{C}_{\mathbf{a}}^1 \cup \mathcal{C}_{\mathbf{b}}^2$$

$$\mathcal{C}_{\mathbf{a}}^1 \stackrel{\text{def.}}{=} \{\mathbf{P} : \mathbf{P}\mathbb{1}_m = \mathbf{a}\}$$

$$\mathcal{C}_{\mathbf{b}}^2 \stackrel{\text{def.}}{=} \left\{ \mathbf{P} : \mathbf{P}^T \mathbb{1}_m = \mathbf{b} \right\}$$

# Bregman Iterative Projections

$$\langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \text{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) = \varepsilon \text{KL}(\mathbf{P}|\mathbf{K}) + \text{cst} \quad \text{where} \quad \mathbf{K}_{i,j} = e^{-\frac{\mathbf{C}_{i,j}}{\varepsilon}} \mathbf{a}_i \mathbf{b}_j$$

*Shrödinger problem:*  $\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \text{KL}(\mathbf{P}|\mathbf{K})$

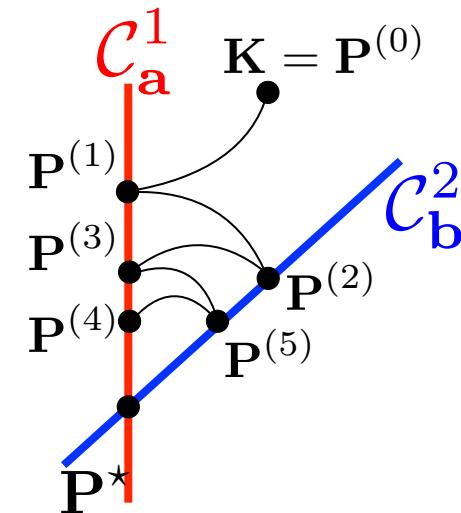
Constraints :  
 $\mathcal{C}_{\mathbf{a}}^1 \cup \mathcal{C}_{\mathbf{b}}^2$

$$\mathcal{C}_{\mathbf{a}}^1 \stackrel{\text{def.}}{=} \{\mathbf{P} : \mathbf{P}\mathbb{1}_m = \mathbf{a}\}$$

$$\mathcal{C}_{\mathbf{b}}^2 \stackrel{\text{def.}}{=} \left\{ \mathbf{P} : \mathbf{P}^T \mathbb{1}_m = \mathbf{b} \right\}$$

Iterative projections:  $\mathbf{P}^{(\ell+1)} \stackrel{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{a}}^1}^{\text{KL}}(\mathbf{P}^{(\ell)})$  and  $\mathbf{P}^{(\ell+2)} \stackrel{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{b}}^2}^{\text{KL}}(\mathbf{P}^{(\ell+1)})$

Theorem:  $\mathbf{P}^{(\ell)} \rightarrow \mathbf{P}^* = \underset{\mathbf{P} \in \mathcal{C}_{\mathbf{a}}^1 \cap \mathcal{C}_{\mathbf{b}}^2}{\operatorname{argmin}} \text{KL}(\mathbf{P}|\mathbf{K})$   
 For affine  $(\mathcal{C}_{\mathbf{a}}^1, \mathcal{C}_{\mathbf{b}}^2)$ ,



# Bregman Iterative Projections

$$\langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \text{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) = \varepsilon \text{KL}(\mathbf{P}|\mathbf{K}) + \text{cst} \quad \text{where} \quad \mathbf{K}_{i,j} = e^{-\frac{\mathbf{C}_{i,j}}{\varepsilon}} \mathbf{a}_i \mathbf{b}_j$$

*Shrödinger problem:*  $\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \text{KL}(\mathbf{P}|\mathbf{K})$

Constraints :  
 $\mathcal{C}_{\mathbf{a}}^1 \cup \mathcal{C}_{\mathbf{b}}^2$

$$\mathcal{C}_{\mathbf{a}}^1 \stackrel{\text{def.}}{=} \{\mathbf{P} : \mathbf{P}\mathbb{1}_m = \mathbf{a}\}$$

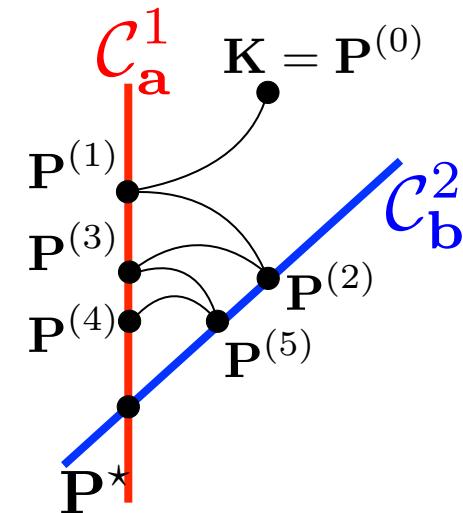
$$\mathcal{C}_{\mathbf{b}}^2 \stackrel{\text{def.}}{=} \left\{ \mathbf{P} : \mathbf{P}^T \mathbb{1}_m = \mathbf{b} \right\}$$

[Bregman, 1967] Iterative projections:  $\mathbf{P}^{(\ell+1)} \stackrel{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{a}}^1}^{\mathbf{KL}}(\mathbf{P}^{(\ell)})$  and  $\mathbf{P}^{(\ell+2)} \stackrel{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{b}}^2}^{\mathbf{KL}}(\mathbf{P}^{(\ell+1)})$

Theorem:  $\mathbf{P}^{(\ell)} \rightarrow \mathbf{P}^* = \underset{\mathbf{P} \in \mathcal{C}_{\mathbf{a}}^1 \cap \mathcal{C}_{\mathbf{b}}^2}{\operatorname{argmin}} \text{KL}(\mathbf{P}|\mathbf{K})$   
 For affine  $(\mathcal{C}_{\mathbf{a}}^1, \mathcal{C}_{\mathbf{b}}^2)$ ,

Sinkhorn  $\iff$  iterative projections.

$$\mathbf{P}^{(2\ell)} \stackrel{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell)}) \mathbf{K} \text{diag}(\mathbf{v}^{(\ell)}), \quad \mathbf{P}^{(2\ell+1)} \stackrel{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell+1)}) \mathbf{K} \text{diag}(\mathbf{v}^{(\ell)})$$



# Dual Analysis: Alternate Maximization

Dual problem:

$$W_p^\varepsilon(\alpha, \beta)^p \stackrel{\text{def.}}{=} \sup_{(\mathbf{f}, \mathbf{g}) \in \mathcal{C}(\mathcal{X})^2} \int \mathbf{f} d\alpha + \int \mathbf{g} d\beta + \varepsilon \int_{\mathcal{X}^2} (1 - e^{\frac{-d^p + \mathbf{f} \oplus \mathbf{g}}{\varepsilon}}) d\alpha \otimes d\beta$$

Primal-dual relations:  $d\pi^*(x, y) = e^{\frac{\mathbf{f}^*(x) + \mathbf{g}^*(y) - c(x, y)}{\varepsilon}} d\alpha(x) d\beta(y)$

# Dual Analysis: Alternate Maximization

Dual problem:

$$W_p^\varepsilon(\alpha, \beta)^p \stackrel{\text{def.}}{=} \sup_{(\mathbf{f}, \mathbf{g}) \in \mathcal{C}(\mathcal{X})^2} \int \mathbf{f} d\alpha + \int \mathbf{g} d\beta + \varepsilon \int_{\mathcal{X}^2} (1 - e^{\frac{-d^p + \mathbf{f} \oplus \mathbf{g}}{\varepsilon}}) d\alpha \otimes d\beta$$

Primal-dual relations:  $d\pi^*(x, y) = e^{\frac{\mathbf{f}^*(x) + \mathbf{g}^*(y) - c(x, y)}{\varepsilon}} d\alpha(x) d\beta(y)$

Soft min:

$$\min_{\alpha}^\varepsilon(h) \stackrel{\text{def.}}{=} -\varepsilon \log \int_{\mathcal{X}} e^{-h/\varepsilon} d\alpha \xrightarrow{\varepsilon \rightarrow 0} \min_{\text{supp}(\alpha)} h$$

Soft  $c$ -transforms:

$$\begin{aligned} \mathbf{f}^{c, \varepsilon}(y) &\stackrel{\text{def.}}{=} \min_{\alpha}^\varepsilon(d^p(\cdot, y) - \mathbf{f}) \\ \mathbf{g}^{c, \varepsilon}(x) &\stackrel{\text{def.}}{=} \min_{\beta}^\varepsilon(d^p(x, \cdot) - \mathbf{g}) \end{aligned}$$

# Dual Analysis: Alternate Maximization

Dual problem:

$$W_p^\varepsilon(\alpha, \beta)^p \stackrel{\text{def.}}{=} \sup_{(\mathbf{f}, \mathbf{g}) \in \mathcal{C}(\mathcal{X})^2} \int \mathbf{f} d\alpha + \int \mathbf{g} d\beta + \varepsilon \int_{\mathcal{X}^2} (1 - e^{\frac{-d^p + \mathbf{f} \oplus \mathbf{g}}{\varepsilon}}) d\alpha \otimes d\beta$$

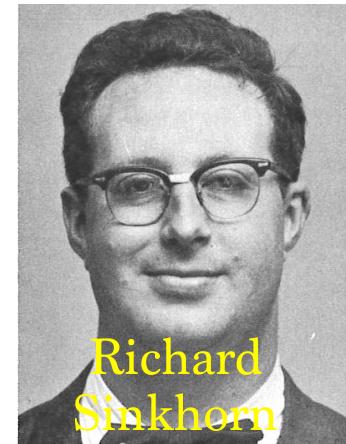
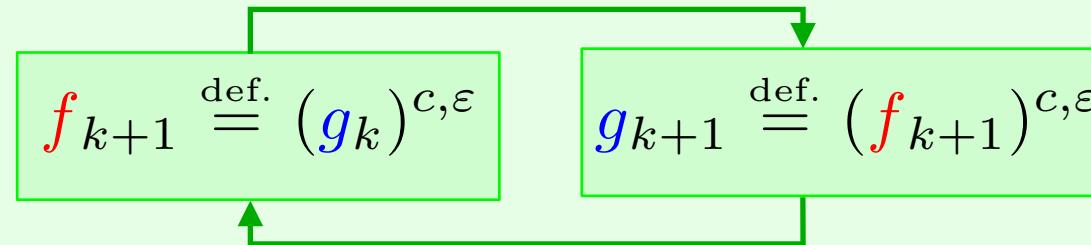
Primal-dual relations:  $d\pi^*(x, y) = e^{\frac{\mathbf{f}^*(x) + \mathbf{g}^*(y) - c(x, y)}{\varepsilon}} d\alpha(x) d\beta(y)$

Soft min:  $\min_{\alpha}^\varepsilon(h) \stackrel{\text{def.}}{=} -\varepsilon \log \int_{\mathcal{X}} e^{-h/\varepsilon} d\alpha \xrightarrow{\varepsilon \rightarrow 0} \min_{\text{supp}(\alpha)} h$

Soft  $c$ -transforms:

$$\begin{aligned} \mathbf{f}^{c, \varepsilon}(y) &\stackrel{\text{def.}}{=} \min_{\alpha}^\varepsilon(d^p(\cdot, y) - \mathbf{f}) \\ \mathbf{g}^{c, \varepsilon}(x) &\stackrel{\text{def.}}{=} \min_{\beta}^\varepsilon(d^p(x, \cdot) - \mathbf{g}) \end{aligned}$$

Sinkhorn's algorithm:

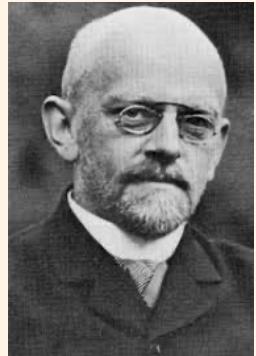


*Proposition:*  $f_0 = 0 \leq f_1 \leq f_2 \leq \dots$  If  $c$  is bounded,  $f_k \rightarrow f^*$ .

[Robert Fortet 1938]

# Hilbert Projective Metric

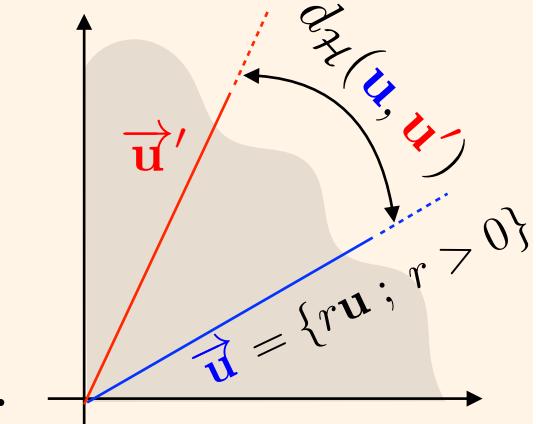
Hilbert's projective metric:  $\forall (\mathbf{u}, \mathbf{u}') \in (\mathbb{R}_{+,*}^n)^2$



$$d_{\mathcal{H}}(\mathbf{u}, \mathbf{u}') \stackrel{\text{def.}}{=} \|\log(\mathbf{u}) - \log(\mathbf{u}')\|_V$$

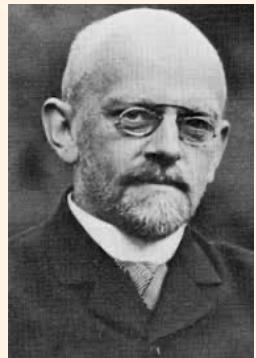
$$\|f\|_V \stackrel{\text{def.}}{=} \max(f) - \min(f)$$

$d_{\mathcal{H}}$  is a distance on the set of rays  $\overrightarrow{\mathbf{u}}$ .



# Hilbert Projective Metric

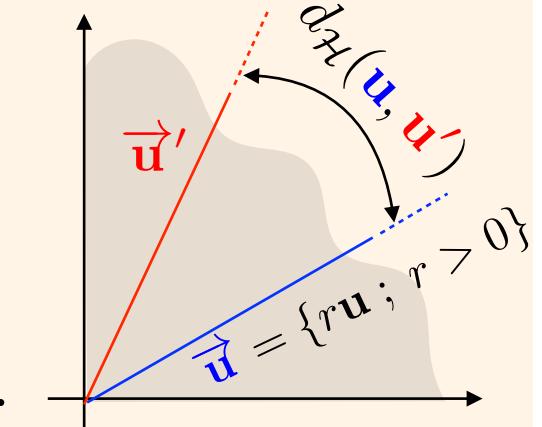
*Hilbert's projective metric:*  $\forall (\mathbf{u}, \mathbf{u}') \in (\mathbb{R}_{+,*}^n)^2$



$$d_{\mathcal{H}}(\mathbf{u}, \mathbf{u}') \stackrel{\text{def.}}{=} \|\log(\mathbf{u}) - \log(\mathbf{u}')\|_V$$

$$\|f\|_V \stackrel{\text{def.}}{=} \max(f) - \min(f)$$

$d_{\mathcal{H}}$  is a distance on the set of rays  $\overrightarrow{\mathbf{u}}$ .



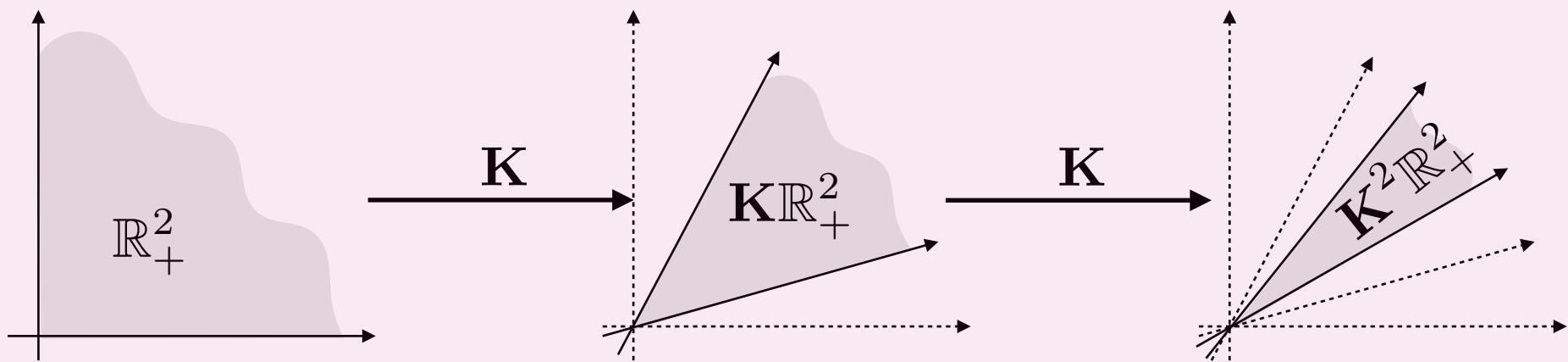
*Birkhoff's contraction theorem:*



**Theorem 1.1.** Let  $\mathbf{K} \in \mathbb{R}_{+,*}^{n \times m}$ , then for  $(\mathbf{v}, \mathbf{v}') \in (\mathbb{R}_{+,*}^m)^2$

$$d_{\mathcal{H}}(\mathbf{K}\mathbf{v}, \mathbf{K}\mathbf{v}') \leq \lambda(\mathbf{K})d_{\mathcal{H}}(\mathbf{v}, \mathbf{v}')$$

where  $\begin{cases} \lambda(\mathbf{K}) \stackrel{\text{def.}}{=} \frac{\sqrt{\eta(\mathbf{K})}-1}{\sqrt{\eta(\mathbf{K})}+1} < 1 \\ \eta(\mathbf{K}) \stackrel{\text{def.}}{=} \max_{i,j,k,\ell} \frac{\mathbf{K}_{i,k}\mathbf{K}_{j,\ell}}{\mathbf{K}_{j,k}\mathbf{K}_{i,\ell}}. \end{cases}$



# Perron Frobenius

Simplex:  $\Sigma_k = \{p \in \mathbb{R}_+^k ; \sum_i p_i = 1\}$

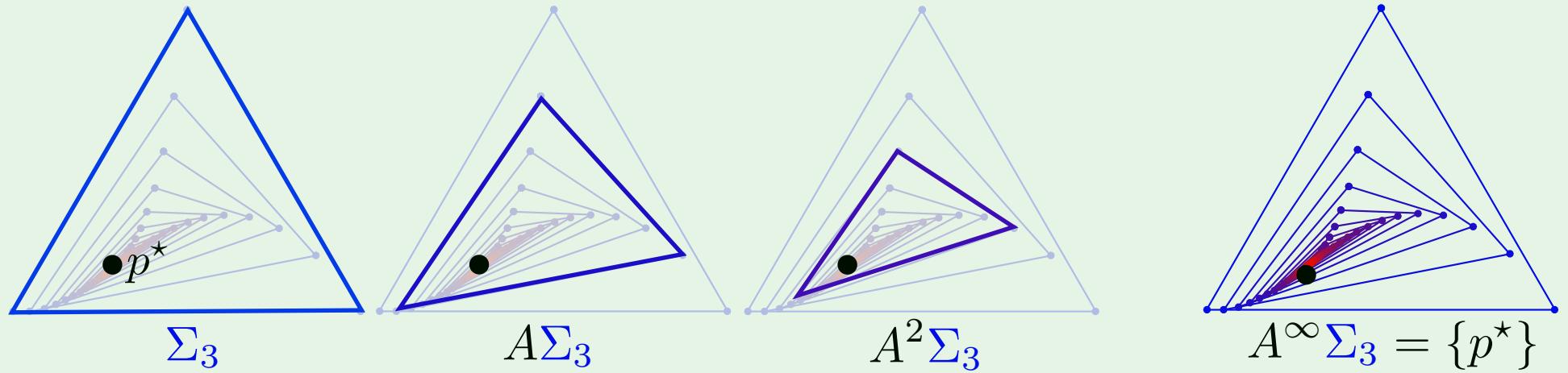
$$A : \Sigma_k \rightarrow \Sigma_k$$

Stochastic matrix:  $A \in \mathbb{R}_+^n, A^\top \mathbf{1}_k = \mathbf{1}_k$

*Theorem:* [Perron-Frobenius]

If  $A > 0$ ,  $\exists! p^*$ ,  $Ap^* = p^*$ .

$\exists \rho \in [0, 1[, \|A^k p - p^*\| \leq \rho^k$



# Sinkhorn under Hilbert's Metric

Sinkhorn iterations:

$$\mathbf{u}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}^{(\ell)}}$$

$$\text{and} \quad \mathbf{v}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{b}}{\mathbf{K}^T \mathbf{u}^{(\ell+1)}}$$

[Franklin and Lorenz, 1989]

*Theorem:* One has  $(\mathbf{u}^{(\ell)}, \mathbf{v}^{(\ell)}) \rightarrow (\mathbf{u}^*, \mathbf{v}^*)$

$$d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) = O(\lambda(\mathbf{K})^{2\ell}), \quad d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^*) = O(\lambda(\mathbf{K})^{2\ell}).$$

$$d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) \leq \frac{d_{\mathcal{H}}(\mathbf{P}^{(\ell)} \mathbb{1}_m, \mathbf{a})}{1 - \lambda(\mathbf{K})^2}$$

$$d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^*) \leq \frac{d_{\mathcal{H}}(\mathbf{P}^{(\ell), T} \mathbb{1}_n, \mathbf{b})}{1 - \lambda(\mathbf{K})^2}$$

$$\|\log(\mathbf{P}^{(\ell)}) - \log(\mathbf{P}^*)\|_\infty \leq d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^*) + d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^*)$$

# Hilbert Metric Analysis

Dual cost:  $W^{(k)}(\alpha, \beta) \stackrel{\text{def.}}{=} \int f_k d\alpha + \int g_k d\beta$

*Theorem:*  $|W^{(k)}(\alpha, \beta) - W_p^\varepsilon(\alpha, \beta)^p| \leq C(1 - e^{-\frac{\|d\|_\infty^p}{\varepsilon}})^k$

Fast in term of  $k$  . . . slow in term of  $\varepsilon$   
→ useless to approximate  $W_p$

# Hilbert Metric Analysis

Dual cost:  $W^{(k)}(\alpha, \beta) \stackrel{\text{def.}}{=} \int f_k d\alpha + \int g_k d\beta$

*Theorem:*  $|W^{(k)}(\alpha, \beta) - W_p^\varepsilon(\alpha, \beta)^p| \leq C(1 - e^{-\frac{\|d\|_\infty^p}{\varepsilon}})^k$

Fast in term of  $k$  . . . slow in term of  $\varepsilon$   
→ useless to approximate  $W_p$

*Proof:*

Variation semi-norm:  $\|h\|_V \triangleq \sup(f) - \inf(f)$

Birkhoff's contraction theorem:

$f \mapsto f^{c,\varepsilon}$  is contractant for  $\|\cdot\|_V$   
 $\implies \|f_k - f^*\|_V = O(1 - e^{-\frac{\|d\|_\infty^p}{\varepsilon}})^k$

# Hilbert vs Mirror Analysis

Theorem:

$$|W^{(k)}(\alpha, \beta) - W_p^\varepsilon(\alpha, \beta)^p| \leq \begin{cases} C(1 - e^{-\frac{\|d\|_\infty^p}{\varepsilon}})^k \\ \frac{\|d\|_\infty^{2p}}{\varepsilon k} \end{cases}$$

# Hilbert vs Mirror Analysis

Theorem:

$$|\mathbf{W}^{(k)}(\alpha, \beta) - \mathbf{W}_p^\varepsilon(\alpha, \beta)^p| \leq \begin{cases} C(1 - e^{-\frac{\|d\|_\infty^p}{\varepsilon}})^k \\ \frac{\|d\|_\infty^{2p}}{\varepsilon k} \end{cases}$$

[Altschuler et al 2017]

$$\begin{aligned} \Delta_k &\stackrel{\text{def.}}{=} \mathbf{W}_{\varepsilon, p}^p(\alpha, \beta) - \mathbf{W}^{(k)}(\alpha, \beta) \\ \Delta_k - \Delta_{k+1} &= \varepsilon \text{KL}(\alpha | (\pi_k)_1) + \varepsilon \text{KL}(\beta | (\pi_k)_2) \\ &\geq \varepsilon \|\alpha - (\pi_k)_1\|_1^2 + \varepsilon \|\beta - (\pi_k)_2\|_1^2 \quad (\text{ Pinsker }) \\ &\geq \frac{\varepsilon}{\|d\|_\infty^{2p}} \Delta_k^2 \quad (\|f_k\|_V, \|g_k\|_V \leq \|d\|_\infty^p) \end{aligned}$$

# Hilbert vs Mirror Analysis

*Theorem:*  $|W^{(k)}(\alpha, \beta) - W_p^\varepsilon(\alpha, \beta)^p| \leq \begin{cases} C(1 - e^{-\frac{\|d\|_\infty^p}{\varepsilon}})^k \\ \frac{\|d\|_\infty^{2p}}{\varepsilon k} \end{cases}$

[Altschuler et al 2017]

$$\begin{aligned} \Delta_k &\stackrel{\text{def.}}{=} W_{\varepsilon, p}^p(\alpha, \beta) - W^{(k)}(\alpha, \beta) \\ \Delta_k - \Delta_{k+1} &= \varepsilon \text{KL}(\alpha | (\pi_k)_1) + \varepsilon \text{KL}(\beta | (\pi_k)_2) \\ &\geq \varepsilon \|\alpha - (\pi_k)_1\|_1^2 + \varepsilon \|\beta - (\pi_k)_2\|_1^2 \quad (\text{ Pinsker }) \\ &\geq \frac{\varepsilon}{\|d\|_\infty^{2p}} \Delta_k^2 \quad (\|f_k\|_V, \|g_k\|_V \leq \|d\|_\infty^p) \end{aligned}$$

for  $\alpha = \frac{1}{n} \sum_i \delta_{x_i}$

*Proposition:*

$$|W_p^{\varepsilon, p} - W_p^p| \leq \varepsilon \log(n)$$

$$\varepsilon = \frac{\delta}{\log(n)}$$

$|W^{(k)} - W_p^p| \leq \delta$   
in  $n^2 \log(n) \|d\|_\infty^{2p} / \varepsilon^2$   
operations

# Overview

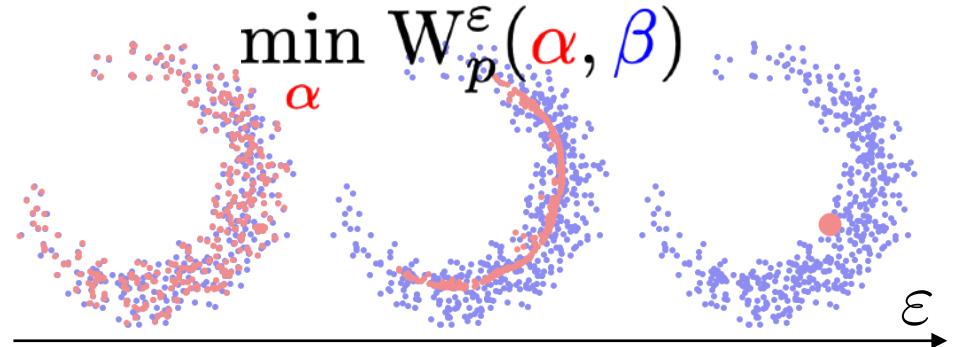
---

- Entropic Regularization and Sinkhorn
- Convergence Analysis
- **Sinkhorn Divergences**
- Generative Model Fitting

# Kernel norms and MMDs

$$W_p^\varepsilon(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d(x, y)^p d\pi(x, y) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta)$$

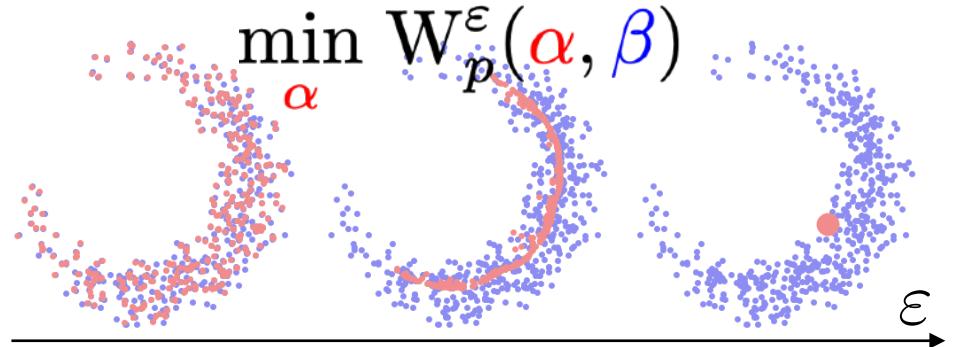
Problem:  $W_p^\varepsilon(\alpha, \alpha) \neq 0$



# Kernel norms and MMDs

$$W_p^\varepsilon(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d(x, y)^p d\pi(x, y) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta)$$

Problem:  $W_p^\varepsilon(\alpha, \alpha) \neq 0$



Prop.:  $\pi^{(\varepsilon)} \xrightarrow{\varepsilon \rightarrow +\infty} \alpha \otimes \beta$

$$W_p^\varepsilon(\alpha, \beta)^p \xrightarrow{\varepsilon \rightarrow +\infty} -\langle \alpha, \beta \rangle_k$$

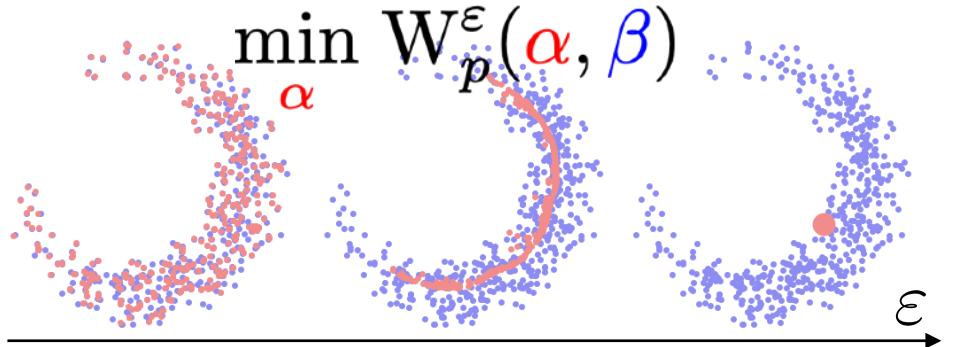
for  $k(x, y) = -d(x, y)^p$  and

$$\langle \alpha, \beta \rangle_k \stackrel{\text{def.}}{=} \int k(x, y) d\alpha(x) d\beta(y)$$

# Kernel norms and MMDs

$$W_p^\varepsilon(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d(x, y)^p d\pi(x, y) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta)$$

Problem:  $W_p^\varepsilon(\alpha, \alpha) \neq 0$



Prop.:  $\pi^{(\varepsilon)} \xrightarrow{\varepsilon \rightarrow +\infty} \alpha \otimes \beta$

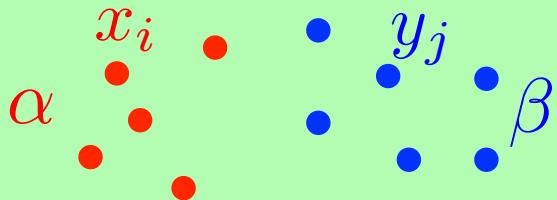
$$W_p^\varepsilon(\alpha, \beta)^p \xrightarrow{\varepsilon \rightarrow +\infty} -\langle \alpha, \beta \rangle_k$$

for  $k(x, y) = -d(x, y)^p$  and

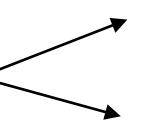
$$\langle \alpha, \beta \rangle_k \stackrel{\text{def.}}{=} \int k(x, y) d\alpha(x) d\beta(y)$$

Kernel norms (MMD):  $\|\alpha - \beta\|_k^2 \stackrel{\text{def.}}{=} \langle \alpha - \beta, \alpha - \beta \rangle_k$

$$\|\alpha - \beta\|_k^2 = \frac{1}{n^2} \sum_{i, i'} k(x_i, x_{i'}) + \frac{1}{m^2} \sum_{j, j'} k(y_j, y_{j'}) - \frac{2}{nm} \sum_{i, j} k(x_i, y_j)$$



Arthur  
Gretton

$k$  must be:  conditionally positive  
universal

# Sinkhorn Divergences

Sinkhorn Divergences:

$$\overline{W}_p^\varepsilon(\alpha, \beta)^p \stackrel{\text{def.}}{=} W_p^\varepsilon(\alpha, \beta)^p - \frac{1}{2} W_p^\varepsilon(\alpha, \alpha)^p - \frac{1}{2} W_p^\varepsilon(\beta, \beta)^p$$

[Ramdas, García Trillos, Cuturi, 2017]

# Sinkhorn Divergences

Sinkhorn Divergences:

$$\overline{W}_p^\varepsilon(\alpha, \beta)^p \stackrel{\text{def.}}{=} W_p^\varepsilon(\alpha, \beta)^p - \frac{1}{2} W_p^\varepsilon(\alpha, \alpha)^p - \frac{1}{2} W_p^\varepsilon(\beta, \beta)^p$$

[Ramdas, García Trillos, Cuturi, 2017]

*Theorem:*  $W_p(\alpha, \beta)^p \xleftarrow[\text{Léonard 2012}]{\varepsilon \rightarrow 0} \overline{W}_p^\varepsilon(\alpha, \beta)^p \xrightarrow[\text{[Ramdas, García Trillos, Cuturi, 2017]}]{\varepsilon \rightarrow +\infty} \frac{1}{2} \|\alpha - \beta\|_{-d^p}^2$

# Sinkhorn Divergences

Sinkhorn Divergences:

$$\overline{W}_p^\varepsilon(\alpha, \beta)^p \stackrel{\text{def.}}{=} W_p^\varepsilon(\alpha, \beta)^p - \frac{1}{2} W_p^\varepsilon(\alpha, \alpha)^p - \frac{1}{2} W_p^\varepsilon(\beta, \beta)^p$$

[Ramdas, García Trillo, Cuturi, 2017]

*Theorem:*  $W_p(\alpha, \beta)^p \xleftarrow[\text{[Léonard 2012]}]{\varepsilon \rightarrow 0} \overline{W}_p^\varepsilon(\alpha, \beta)^p \xrightarrow[\text{[Ramdas, García Trillo, Cuturi, 2017]}]{\varepsilon \rightarrow +\infty} \frac{1}{2} \|\alpha - \beta\|_{-d^p}^2$

*Key problem:* when is  $k(x, y) = -d(x, y)^p$  a universal conditionaly positive kernel?

*Proposition:*  $\|\cdot\|_{-\|\cdot\|^p}$  is a norm for  $0 < p < 2$ .

For  $p = 1$ :  $\dot{H}^{-\frac{d+1}{2}}(\mathbb{R}^d)$  Sobolev norm;

For  $p = 2$ :  $\|\xi\|_{-\|\cdot\|^2}^2 = |\int x d\xi(x)|^2$

# Sinkhorn Divergences Positivity

$$\overline{W}_p^\varepsilon(\alpha, \beta)^p \stackrel{\text{def.}}{=} W_p^\varepsilon(\alpha, \beta)^p - \frac{1}{2} W_p^\varepsilon(\alpha, \alpha)^p - \frac{1}{2} W_p^\varepsilon(\beta, \beta)^p$$

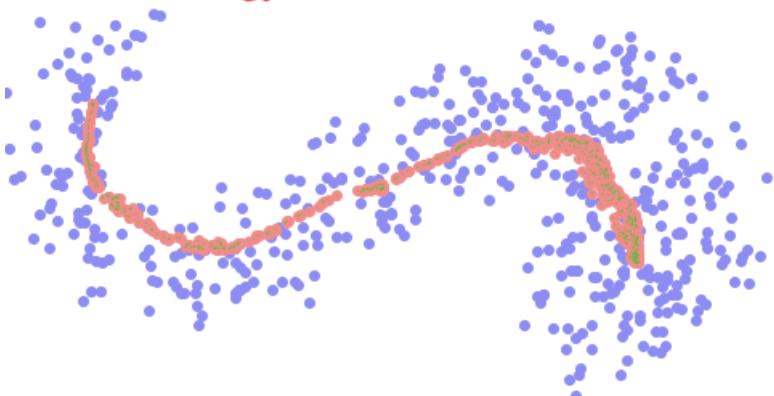
↓ concave      ↓ concave

*Theorem:* [Feydy, Séjourné, P, Vialard, Trouvé, Amari 2018]

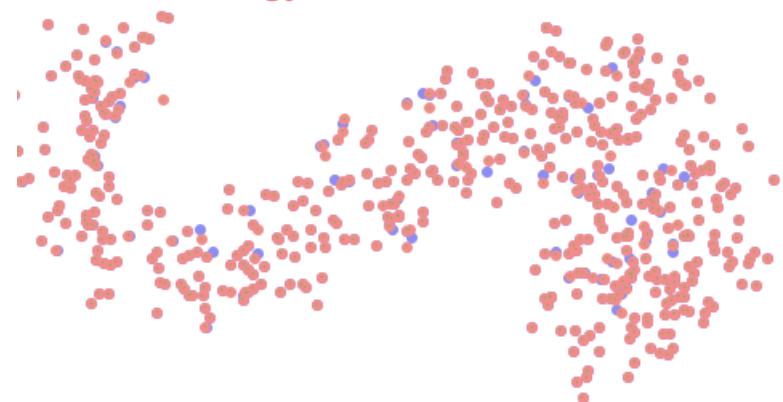
If  $e^{-\frac{d^p}{\varepsilon}}$  is positive:

$$\begin{aligned}\overline{W}_p^\varepsilon &\geq 0 \text{ and } \overline{W}_p^\varepsilon(\cdot, \beta)^p \text{ is convex.} \\ \overline{W}_p^\varepsilon(\alpha_n, \beta) \rightarrow 0 &\iff \alpha_n \xrightarrow{\text{weak*}} \beta\end{aligned}$$

$$\min_{\alpha} W_p^\varepsilon(\alpha, \beta)$$



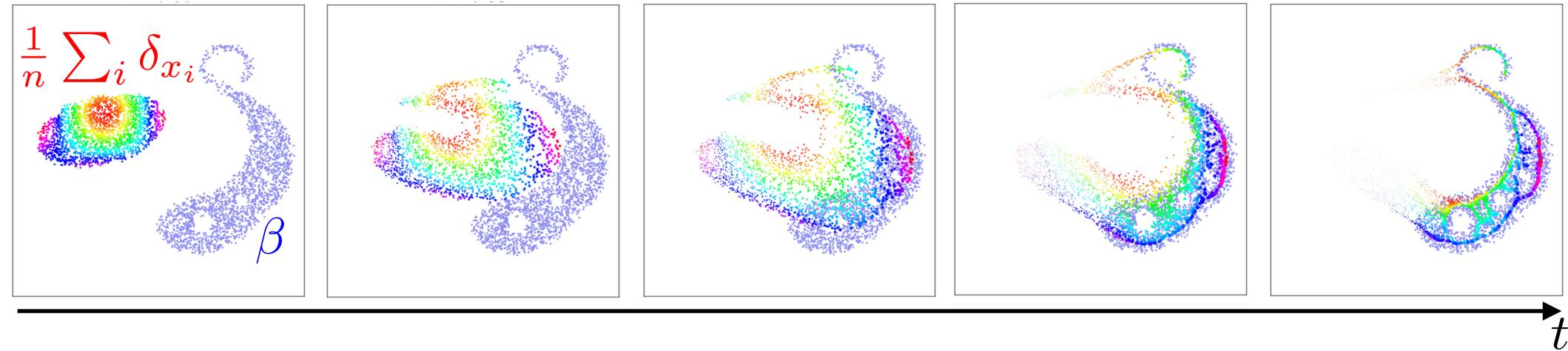
$$\min_{\alpha} \overline{W}_p^\varepsilon(\alpha, \beta)$$



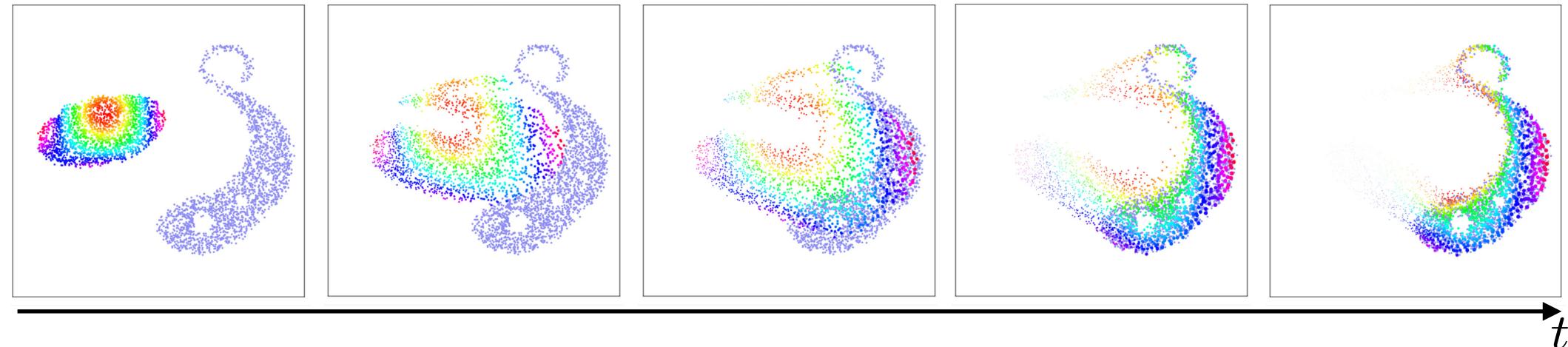
# Wasserstein Gradient Flows

$$\min_{x=(x_i)_i} \mathcal{E}(x) \triangleq W_2^\varepsilon\left(\frac{1}{n} \sum_i \delta_{x_i}, \beta\right)$$

Wasserstein flow:  $\frac{dx(t)}{dt} = -\nabla \mathcal{E}(x(t))$



$$\min_{x=(x_i)_i} f(x) \triangleq \bar{W}_2^\varepsilon\left(\frac{1}{n} \sum_i \delta_{x_i}, \beta\right)$$

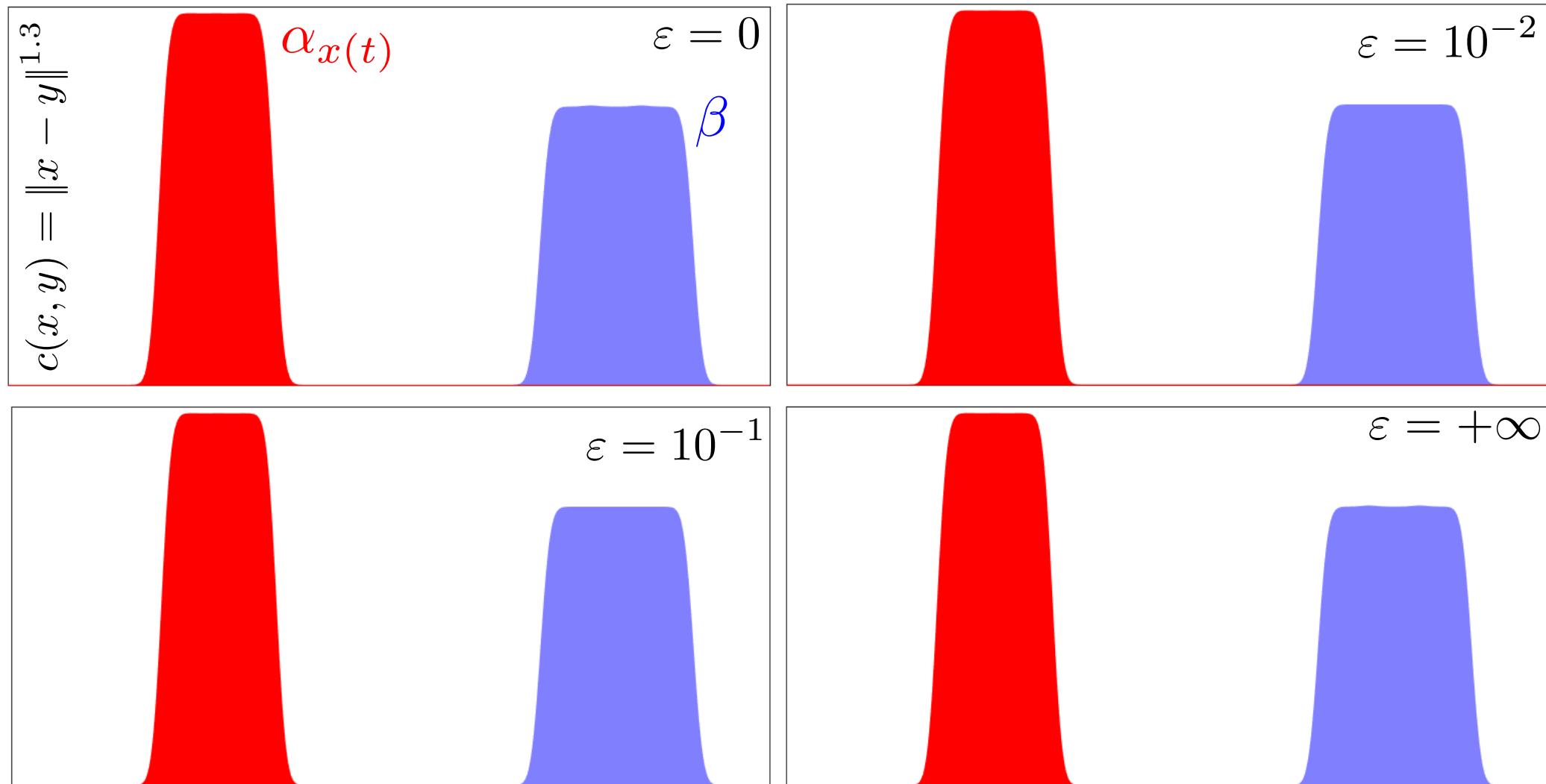


# Wasserstein Gradient Flows

$$\frac{dx(t)}{dt} = -\nabla \mathcal{E}(x(t))$$

$$\mathcal{E}(x) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon, 1.3}^{1.3}(\alpha_x, \beta)$$

$$\alpha_x \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$

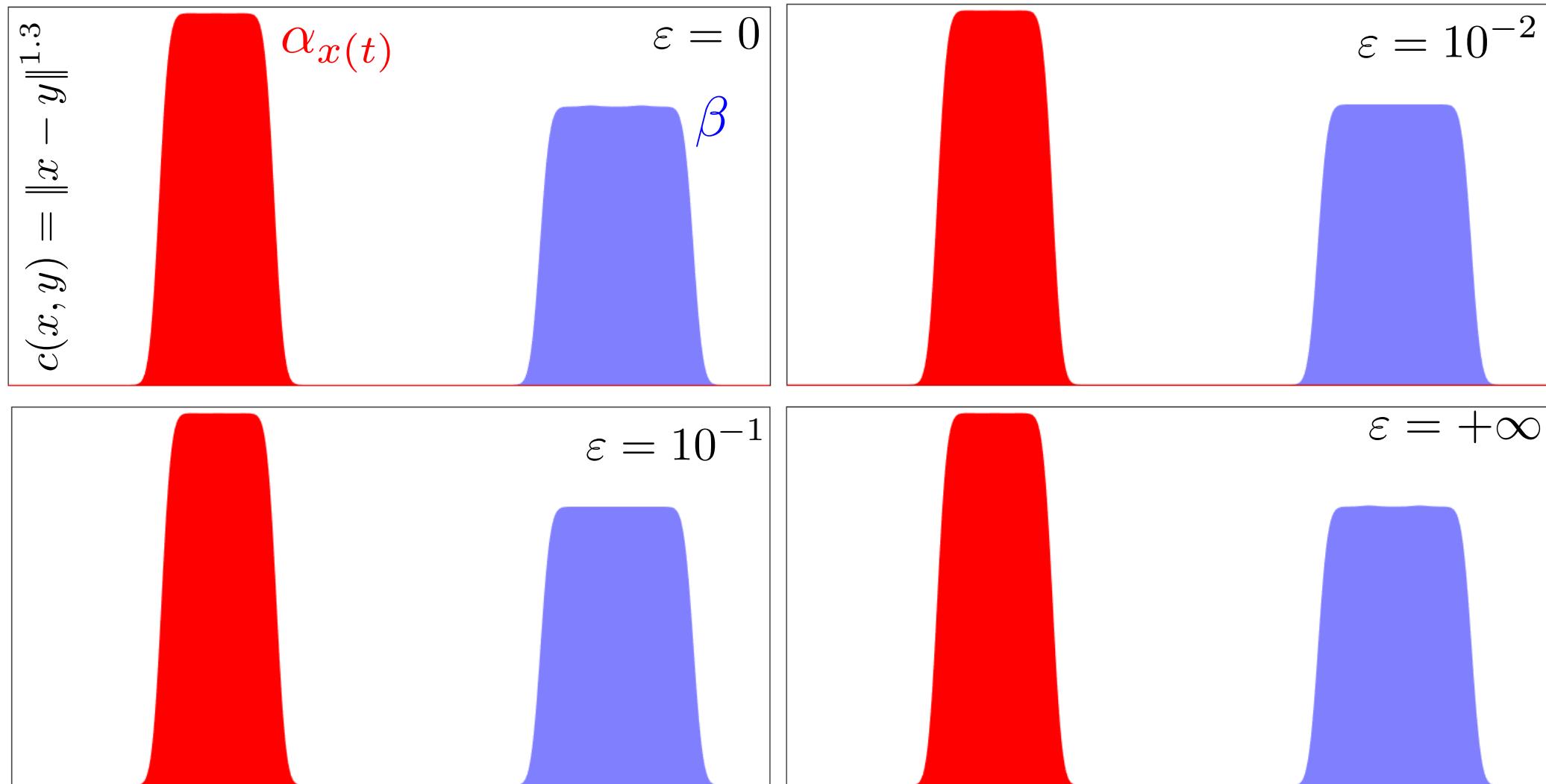


# Wasserstein Gradient Flows

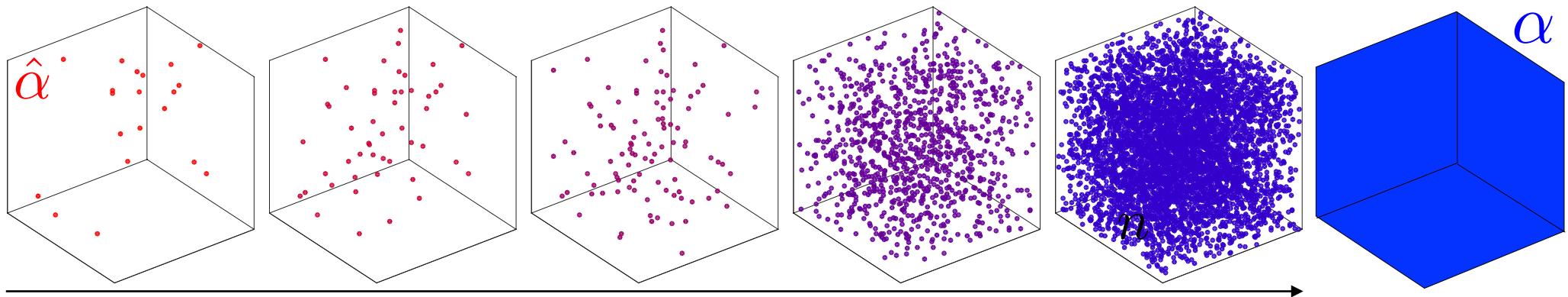
$$\frac{dx(t)}{dt} = -\nabla \mathcal{E}(x(t))$$

$$\mathcal{E}(x) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon, 1.3}^{1.3}(\alpha_x, \beta)$$

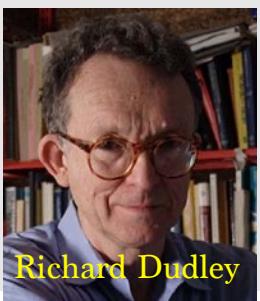
$$\alpha_x \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$



# Sample Complexity



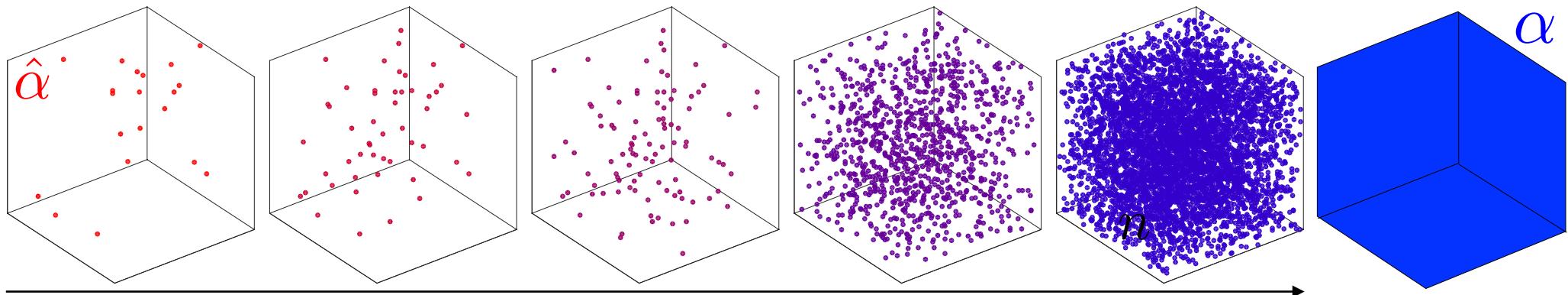
Theorem:



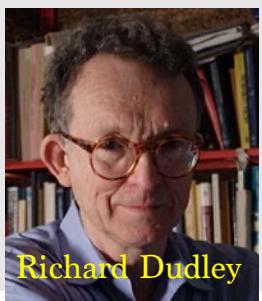
Richard Dudley

$$\mathbb{E}(|W_p(\hat{\alpha}, \hat{\beta}) - W_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}})$$

# Sample Complexity



Theorem:

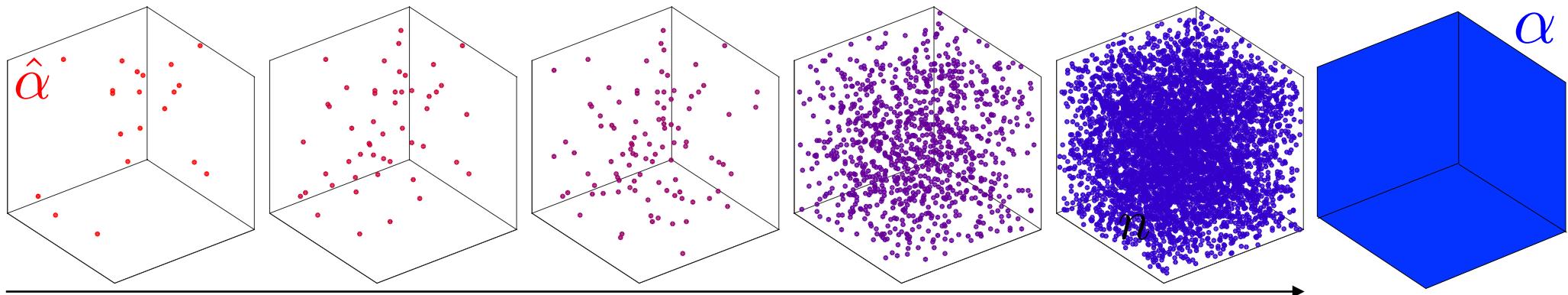


Richard Dudley

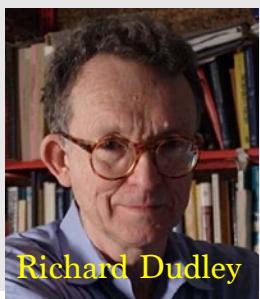
$$\mathbb{E}(|W_p(\hat{\alpha}, \hat{\beta}) - W_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}}) \rightarrow \text{if } \alpha \neq \beta$$
$$\mathbb{E}(|W_2(\hat{\alpha}, \hat{\beta}) - W_2(\alpha, \beta)|) = O(n^{-\frac{2}{d}})$$

[Chizat, Roussillon, Léger, Vialard, P. 2020]

# Sample Complexity



Theorem:



Richard Dudley

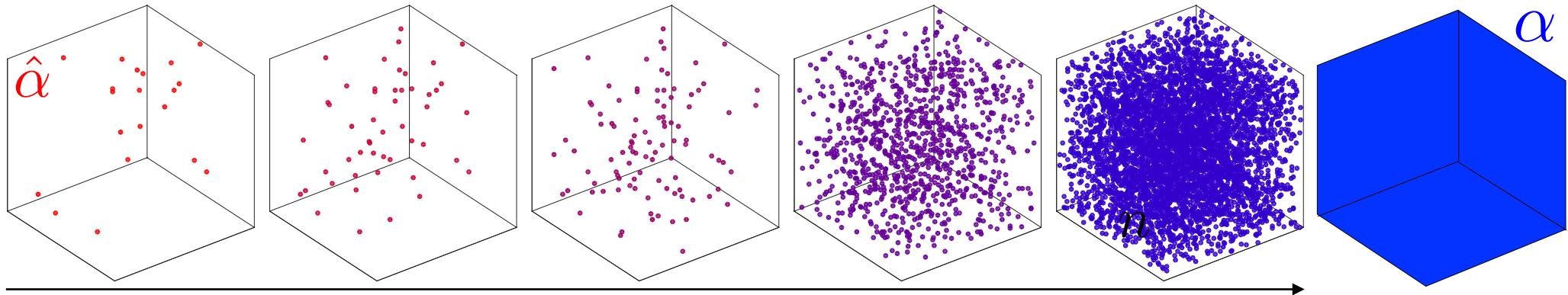
$$\mathbb{E}(|W_p(\hat{\alpha}, \hat{\beta}) - W_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}}) \rightarrow \text{if } \alpha \neq \beta$$

$$\mathbb{E}(|W_2(\hat{\alpha}, \hat{\beta}) - W_2(\alpha, \beta)|) = O(n^{-\frac{2}{d}})$$

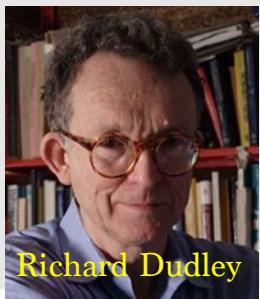
[Chizat, Roussillon, Léger, Vialard, P. 2020]

$$\mathbb{E}(|\|\hat{\alpha} - \hat{\beta}\|_k - \|\alpha - \beta\|_k|) = O(n^{-\frac{1}{2}})$$

# Sample Complexity



Theorem:



Richard Dudley

$$\mathbb{E}(|W_p(\hat{\alpha}, \hat{\beta}) - W_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}}) \rightarrow \text{if } \alpha \neq \beta$$

[Chizat, Roussillon, Léger, Vialard, P. 2020]

$$\mathbb{E}(|\|\hat{\alpha} - \hat{\beta}\|_k - \|\alpha - \beta\|_k|) = O(n^{-\frac{1}{2}})$$

$$\mathbb{E}(|\overline{W}_p^\varepsilon(\hat{\alpha}, \hat{\beta}) - \overline{W}_p^\varepsilon(\alpha, \beta)|) = O(\varepsilon^{-\frac{d}{2}} n^{-\frac{1}{2}})$$

[Genevay, Bach, P, Cuturi, 2019]

# Sinkhorn Divergence Estimator

*Theorem:* [S. Pal, 2019]

$$W_2^\varepsilon(\alpha, \beta)^2 - W_2(\alpha, \beta)^2 = d\varepsilon \log \sqrt{2\pi\varepsilon} + \varepsilon(H(\alpha) + H(\beta)) + o(\varepsilon)$$

# Sinkhorn Divergence Estimator

*Theorem:* [S. Pal, 2019]

$$W_2^\varepsilon(\alpha, \beta)^2 - W_2(\alpha, \beta)^2 = d\varepsilon \log \sqrt{2\pi\varepsilon} + \varepsilon(H(\alpha) + H(\beta)) + o(\varepsilon)$$

Displacement interpolation:  $\rho_t : \alpha \rightarrow \beta$

Fisher information:  $I(\alpha, \beta) \stackrel{\text{def.}}{=} \int_0^1 \int_{\mathbb{R}^d} \|\nabla \log(\rho_t(x))\|^2 \rho_t(x) dx dt$

*Theorem:* [G. Conforti and L. Tamanini 2019]

$$\overline{W}_2^\varepsilon(\alpha, \beta)^2 - W_2(\alpha, \beta)^2 = \varepsilon^2 \left( I(\alpha, \beta) - \frac{I(\alpha, \alpha)}{2} - \frac{I(\beta, \beta)}{2} \right) + o(\varepsilon^2)$$

# Sinkhorn Divergence Estimator

Theorem: [S. Pal, 2019]

$$W_2^\varepsilon(\alpha, \beta)^2 - W_2(\alpha, \beta)^2 = d\varepsilon \log \sqrt{2\pi\varepsilon} + \varepsilon(H(\alpha) + H(\beta)) + o(\varepsilon)$$

Displacement interpolation:  $\rho_t : \alpha \rightarrow \beta$

Fisher information:  $I(\alpha, \beta) \stackrel{\text{def.}}{=} \int_0^1 \int_{\mathbb{R}^d} \|\nabla \log(\rho_t(x))\|^2 \rho_t(x) dx dt$

Theorem: [G. Conforti and L. Tamanini 2019]

$$\overline{W}_2^\varepsilon(\alpha, \beta)^2 - W_2(\alpha, \beta)^2 = \varepsilon^2 \left( I(\alpha, \beta) - \frac{I(\alpha, \alpha)}{2} - \frac{I(\beta, \beta)}{2} \right) + o(\varepsilon^2)$$

$$\mathbb{E}(|\overline{W}_2^\varepsilon(\hat{\alpha}, \hat{\beta})^2 - \overline{W}_2^\varepsilon(\alpha, \beta)^2|) = O(\varepsilon^{-\frac{d}{2}} n^{-\frac{1}{2}})$$

$$\varepsilon = n^{-\frac{1}{d}}$$
$$\alpha \neq \beta$$

$$\mathbb{E}(|\overline{W}_2^\varepsilon(\hat{\alpha}, \hat{\beta}) - W_2(\alpha, \beta)|) = O(n^{-\frac{2}{d}})$$

# Leveraging smoothness in high dimension?

$(\frac{d\alpha}{dx}, \frac{d\beta}{dy})$  with  $d$  derivatives

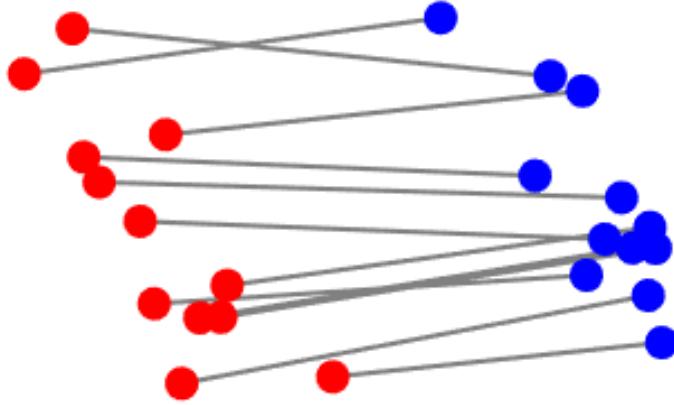
$$\begin{array}{l} \mathbb{E}|\hat{W} - W| \leq C(d)n^{-1/2} \\ \text{Complexity } O(n^2) \end{array}$$

# Leveraging smoothness in high dimension?

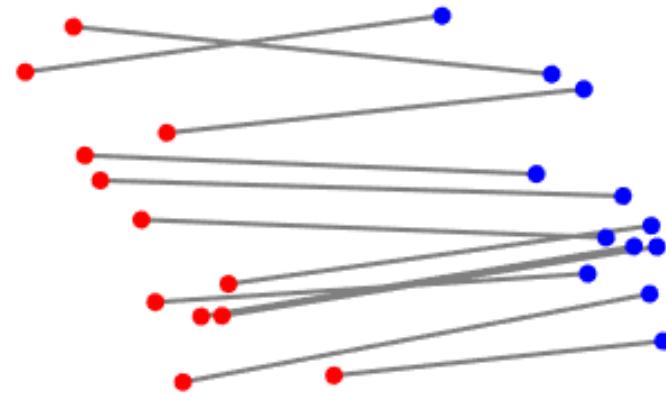
$(\frac{d\alpha}{dx}, \frac{d\beta}{dy})$  with  $d$  derivatives

$$\begin{aligned} & \mathbb{E}|\hat{W} - W| \leq C(d)n^{-1/2} \\ & \text{Complexity } O(n^2) \end{aligned}$$

$\overline{W}_2^\varepsilon(\hat{\alpha}, \hat{\beta})$  is a bad proxy for  $W_2(G_\varepsilon \star \hat{\alpha}, G_\varepsilon \star \hat{\beta}) \dots$



$$\overline{W}_2^\varepsilon(\hat{\alpha}, \hat{\beta})$$



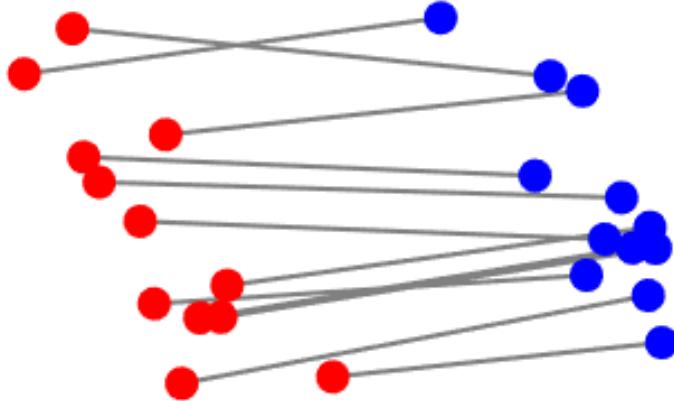
$$W_2(G_\varepsilon \star \hat{\alpha}, G_\varepsilon \star \hat{\beta})$$

# Leveraging smoothness in high dimension?

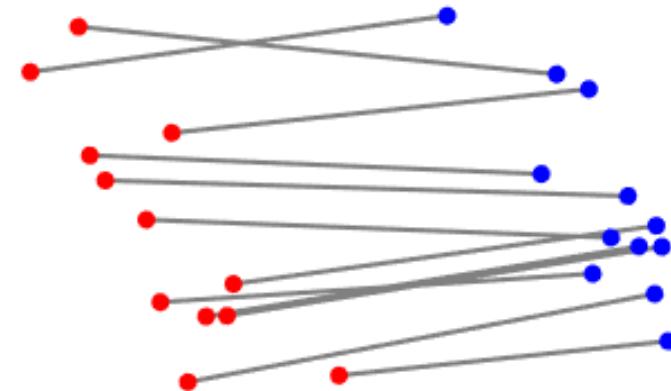
$(\frac{d\alpha}{dx}, \frac{d\beta}{dy})$  with  $d$  derivatives

$$\begin{aligned} & \mathbb{E}|\hat{W} - W| \leq C(d)n^{-1/2} \\ & \text{Complexity } O(n^2) \end{aligned}$$

$\overline{W}_2^\varepsilon(\hat{\alpha}, \hat{\beta})$  is a bad proxy for  $W_2(G_\varepsilon \star \hat{\alpha}, G_\varepsilon \star \hat{\beta}) \dots$



$$\overline{W}_2^\varepsilon(\hat{\alpha}, \hat{\beta})$$



$$W_2(G_\varepsilon \star \hat{\alpha}, G_\varepsilon \star \hat{\beta})$$

Recent breakthrough: [Vacher, Muzellec, Rudi, Bach, Vialard, 2021]  
→ SDP programming.

# Overview

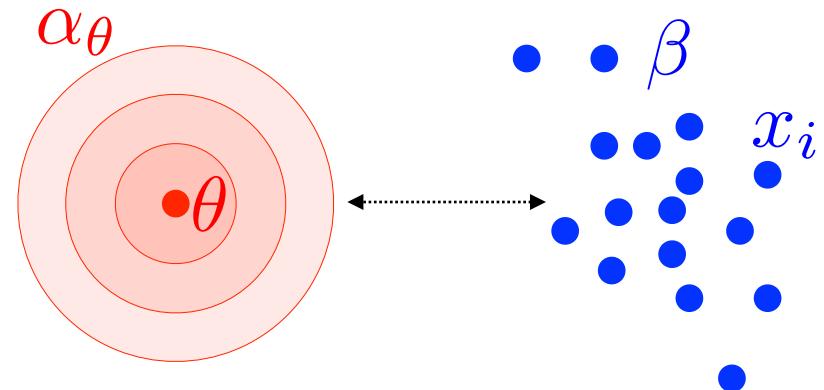
---

- Entropic Regularization and Sinkhorn
- Convergence Analysis
- Sinkhorn Divergences
- **Generative Model Fitting**

# Density Fitting and Generative Models

*Observations:*  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

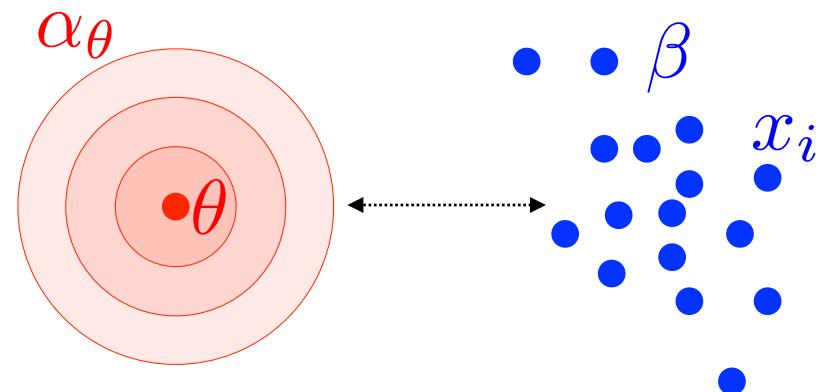
*Parametric model:*  $\theta \mapsto \alpha_\theta$



# Density Fitting and Generative Models

Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \alpha_\theta$



Density fitting:  $d\alpha_\theta(x) = \rho_\theta(x)dx$

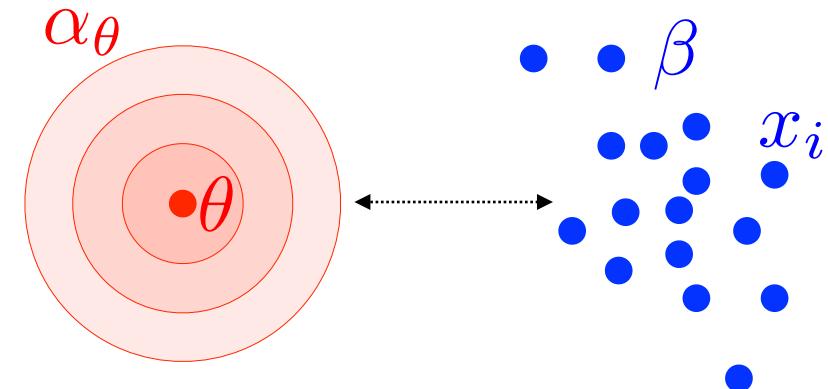
$$\min_{\theta} - \sum_i \log(\rho_\theta(x_i)) \xrightarrow{n \rightarrow +\infty} \text{KL}(\beta | \alpha_\theta)$$

Maximum likelihood (MLE)

# Density Fitting and Generative Models

Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \alpha_\theta$



Density fitting:  $d\alpha_\theta(x) = \rho_\theta(x)dx$

$$\min_{\theta} - \sum_i \log(\rho_\theta(x_i)) \xrightarrow{n \rightarrow +\infty} \text{KL}(\beta | \alpha_\theta)$$

Maximum likelihood (MLE)

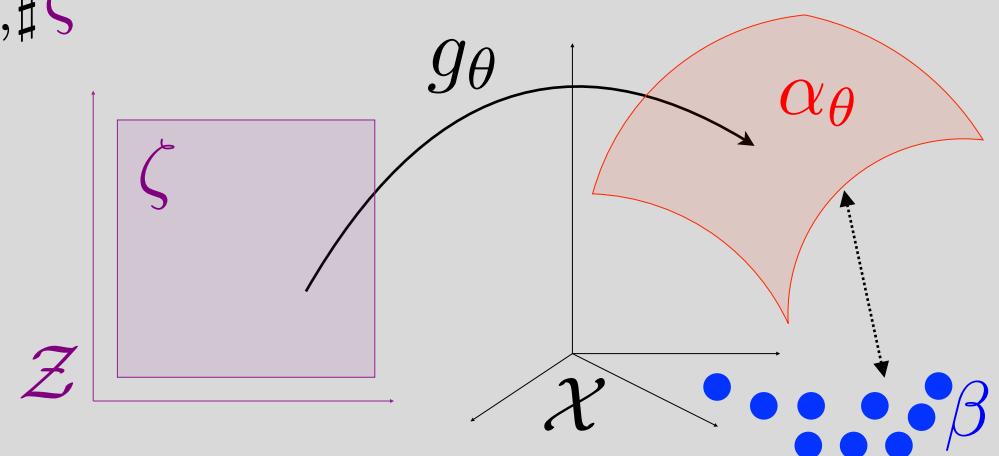
Generative model fit:  $\alpha_\theta = g_{\theta, \sharp} \zeta$

$$\text{KL}(\beta | \alpha_\theta) = +\infty$$

→ MLE undefined.

→ Need a weaker metric.

$$\min_{\theta} \overline{W}_{\varepsilon, p}^p(\alpha_\theta, \beta)$$



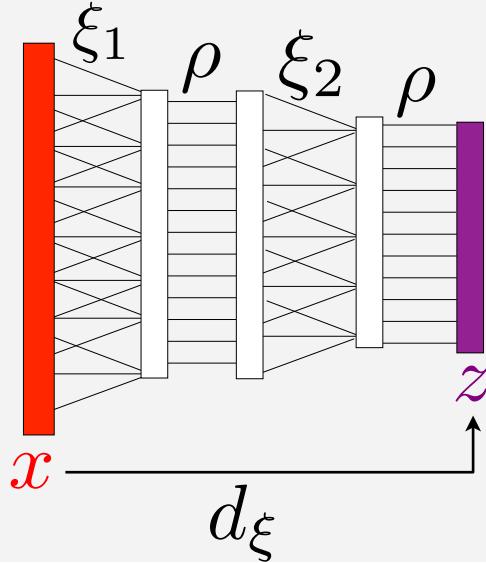
# Deep Discriminative vs Generative Models

Deep networks:

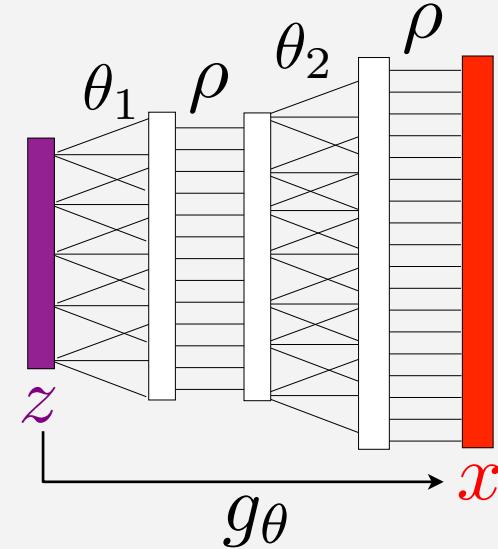
$$d_\xi(\textcolor{red}{x}) = \rho(\xi_K(\dots \rho(\xi_2(\rho(\xi_1(\textcolor{red}{x}) \dots)$$

$$g_\theta(\textcolor{violet}{z}) = \rho(\theta_K(\dots \rho(\theta_2(\rho(\theta_1(\textcolor{violet}{z}) \dots)$$

Discriminative



Generative



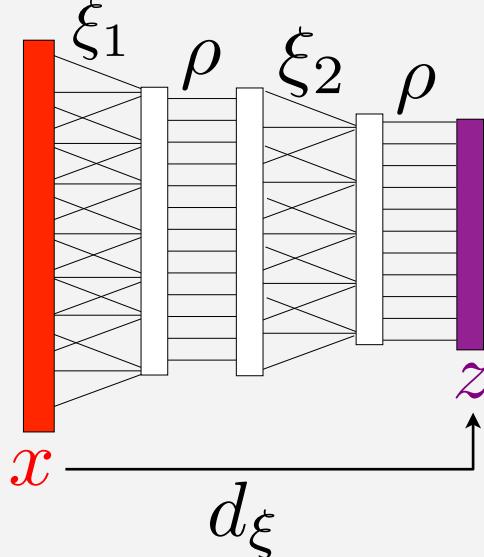
# Deep Discriminative vs Generative Models

Deep networks:

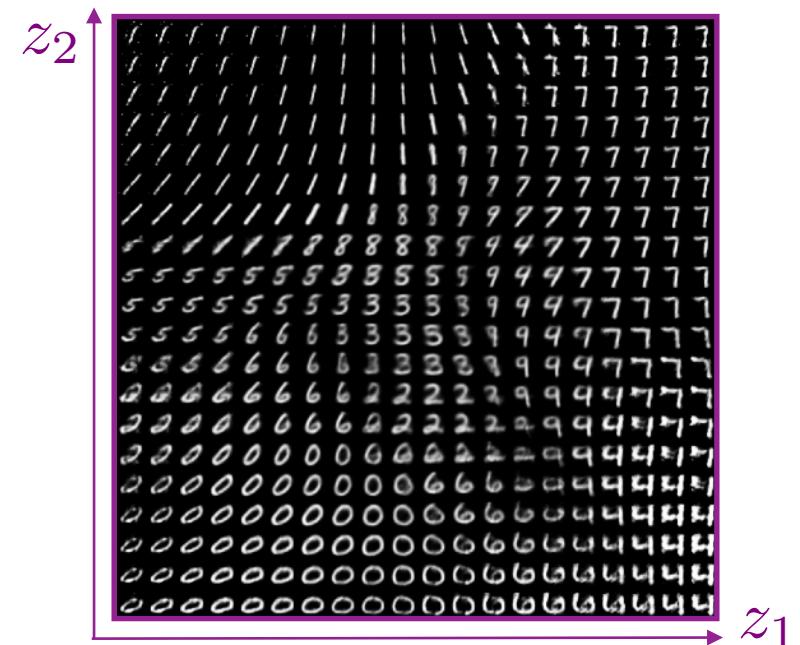
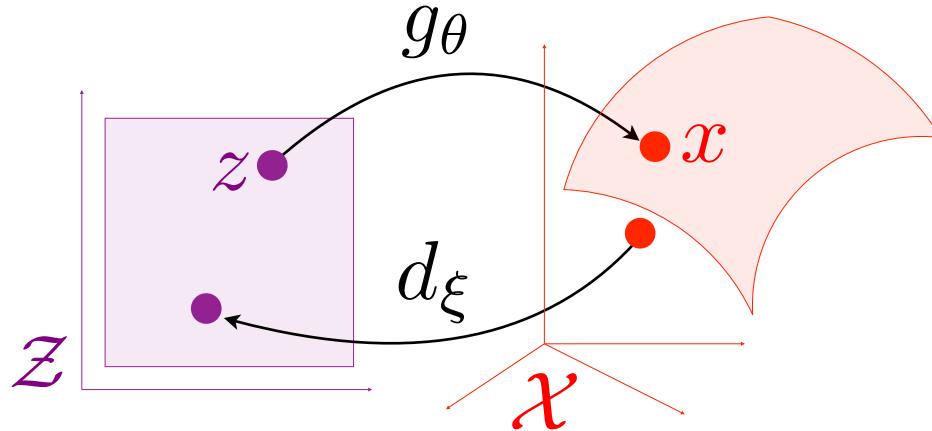
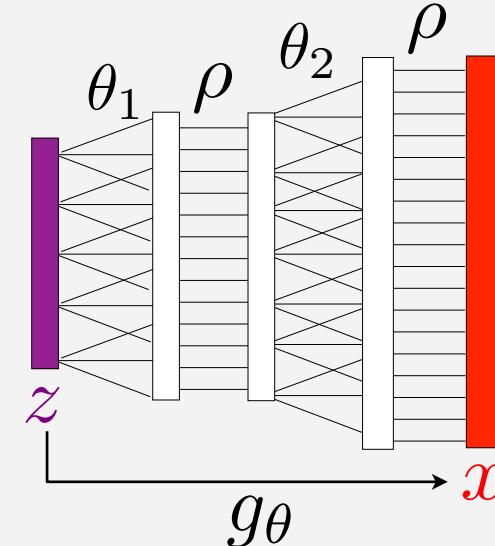
$$d_\xi(\mathbf{x}) = \rho(\xi_K(\dots \rho(\xi_2(\rho(\xi_1(\mathbf{x}) \dots)$$

$$g_\theta(\mathbf{z}) = \rho(\theta_K(\dots \rho(\theta_2(\rho(\theta_1(\mathbf{z}) \dots)$$

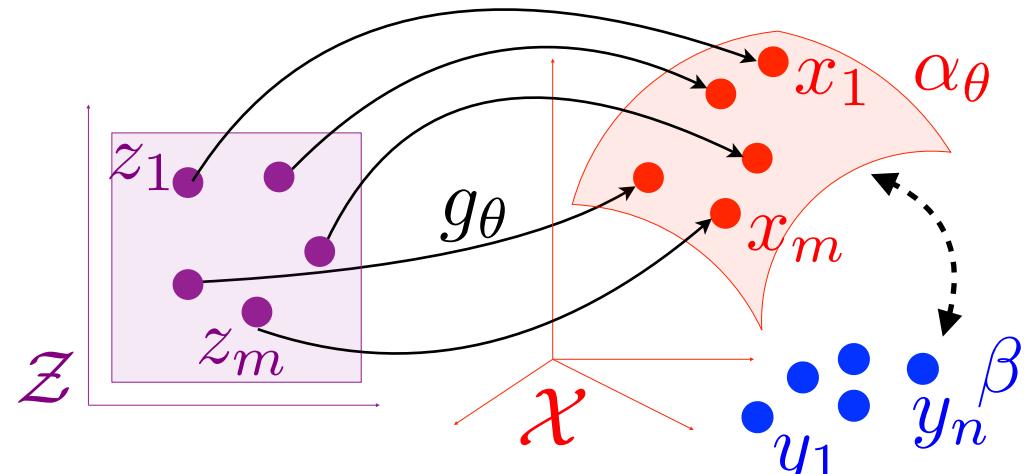
Discriminative



Generative



# Training Architecture



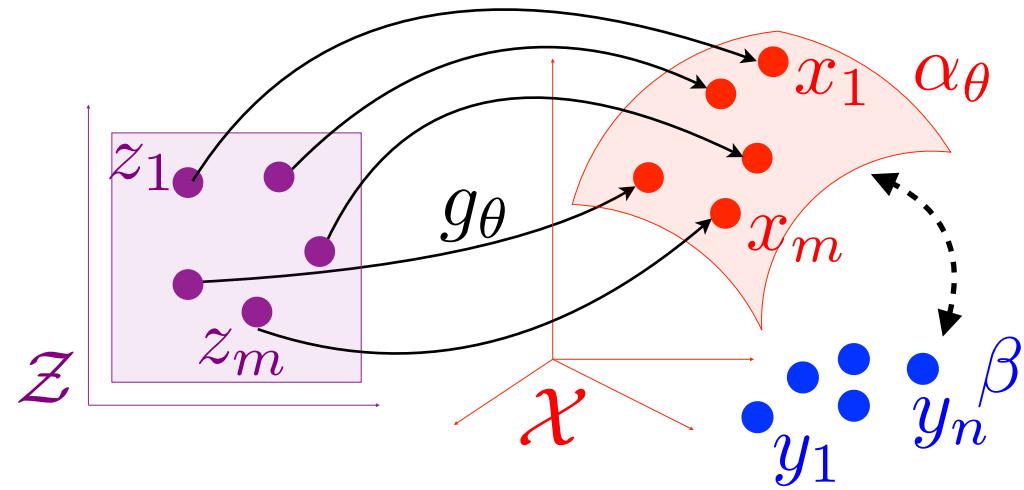
$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{\mathbf{W}}_{\varepsilon,p}^p(\alpha_\theta, \beta)$$

Stochastic gradient descent

$$\theta \leftarrow \theta - \tau \nabla \hat{\mathcal{E}}(\theta)$$

$$\hat{\mathcal{E}}(\theta) \stackrel{\text{def.}}{=} \overline{\mathbf{W}}_{\varepsilon,p}^p\left(\frac{1}{m} \sum_i \delta_{g_\theta(z_i)}, \beta\right)$$

# Training Architecture

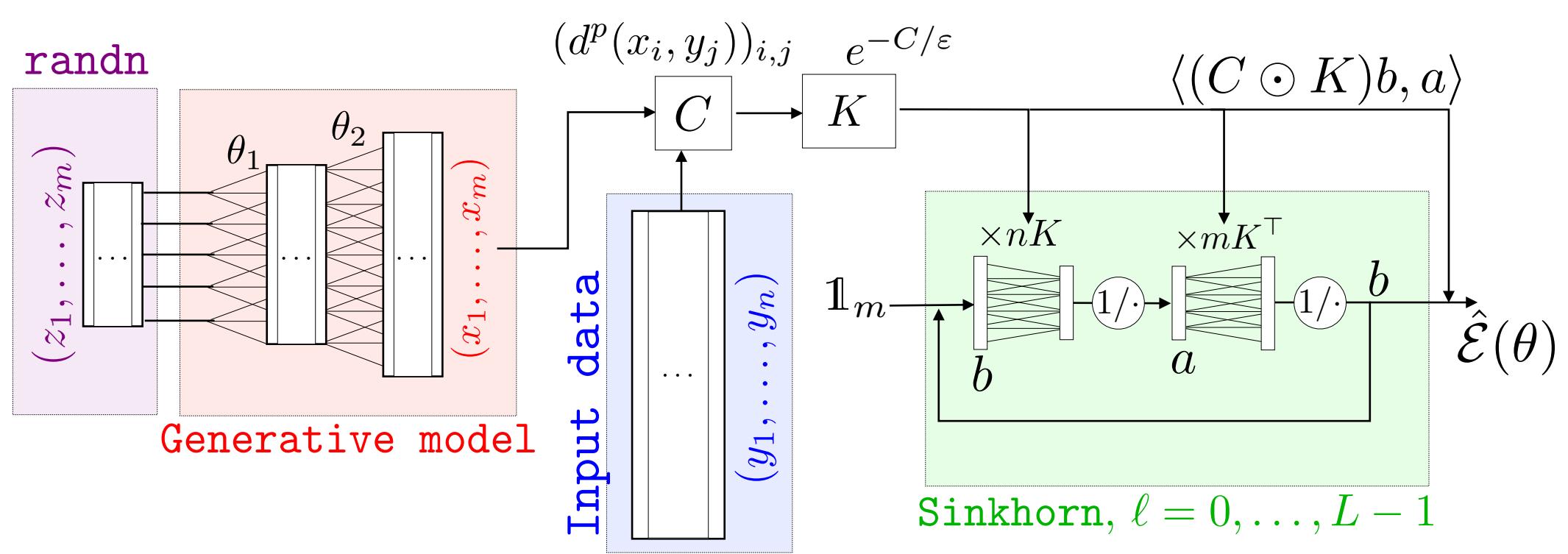


$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{\mathbf{W}}_{\varepsilon,p}^p(\alpha_\theta, \beta)$$

Stochastic gradient descent

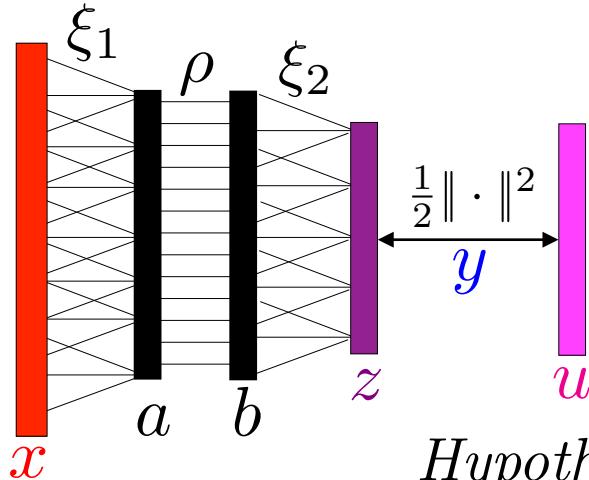
$$\theta \leftarrow \theta - \tau \nabla \hat{\mathcal{E}}(\theta)$$

$$\hat{\mathcal{E}}(\theta) \stackrel{\text{def.}}{=} \overline{\mathbf{W}}_{\varepsilon,p}^p\left(\frac{1}{m} \sum_i \delta_{g_\theta(z_i)}, \beta\right)$$



# Automatic Differentiation

**Setup:**  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  computable in  $K$  operations.



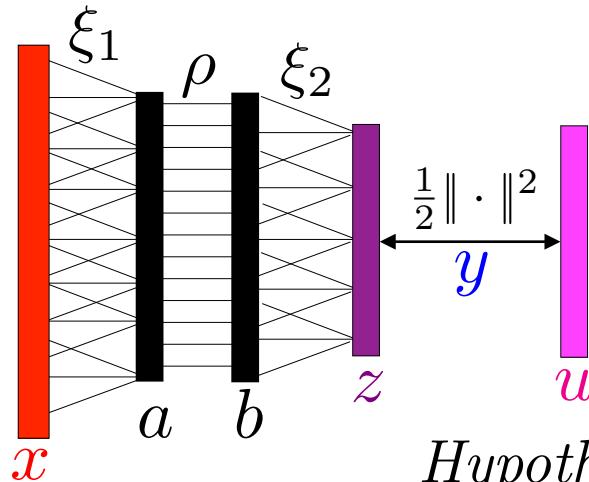
```
function y = E(x)
  a = xi1*x
  b = rho(a)
  z = xi2*b
  y = 1/2*norm(z-u)^2
```

*Hypothesis:* elementary operations ( $a \times b, \log(a), \sqrt{a} \dots$ )  
and their derivatives cost  $O(1)$ .

**Question:** What is the complexity of computing  $\nabla \mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ?

# Automatic Differentiation

**Setup:**  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  computable in  $K$  operations.



```
function y = E(x)
    a = xi1*x
    b = rho(a)
    z = xi2*b
    y = 1/2*norm(z-u)^2
```

*Hypothesis:* elementary operations ( $a \times b, \log(a), \sqrt{a} \dots$ )  
and their derivatives cost  $O(1)$ .

**Question:** What is the complexity of computing  $\nabla \mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ?

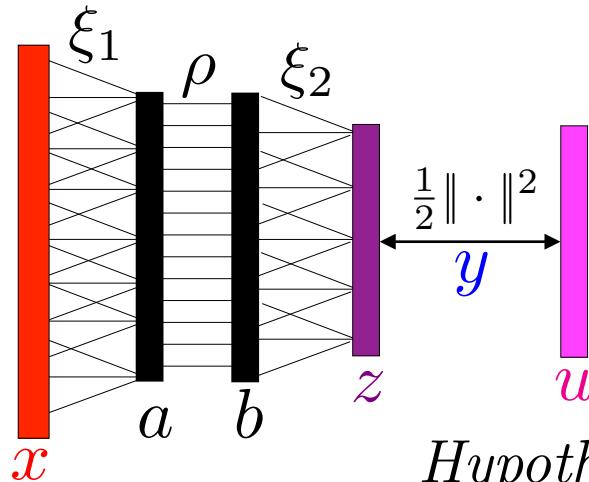
Finite differences:

$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon} (\mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta), \dots, \mathcal{E}(\theta + \varepsilon \delta_n) - \mathcal{E}(\theta))$$

$K(n+1)$  operations, intractable for large  $n$ .

# Automatic Differentiation

**Setup:**  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  computable in  $K$  operations.



```
function y = E(x)
    a = xi1*x
    b = rho(a)
    z = xi2*b
    y = 1/2*norm(z-u)^2
```

```
function dx = nablaE(x)
    dz = z-u
    db = xi2'*dz
    da = diag(dphi(a)) * db
    dx = xi1'*da
```

*Hypothesis:* elementary operations ( $a \times b, \log(a), \sqrt{a} \dots$ )  
and their derivatives cost  $O(1)$ .

**Question:** What is the complexity of computing  $\nabla \mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ?

Finite differences:

$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon} (\mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta), \dots, \mathcal{E}(\theta + \varepsilon \delta_n) - \mathcal{E}(\theta))$$

$K(n+1)$  operations, intractable for large  $n$ .

*Theorem:* there is an algorithm to compute  $\nabla \mathcal{E}$   
in  $O(K)$  operations. [Seppo Linnainmaa, 1970]

This algorithm is reverse mode automatic differentiation



Seppo  
Linnainmaa

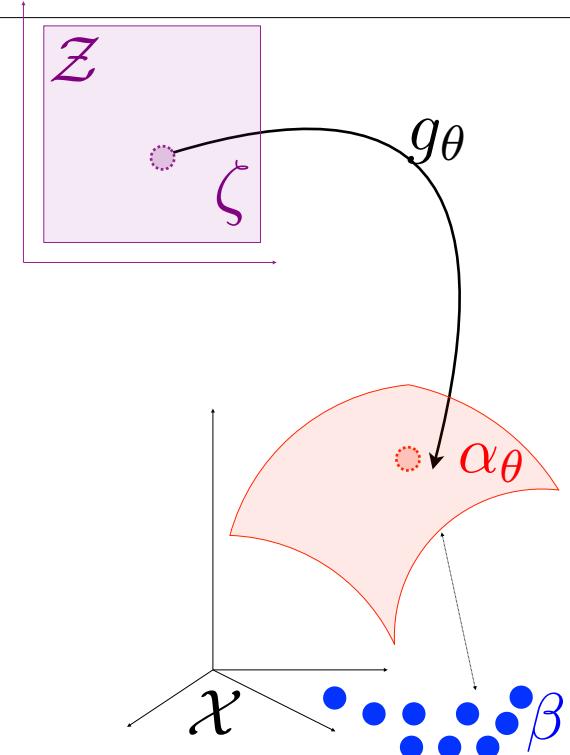
# Examples of Images Generation

Inputs  $\beta$

3	4	2	1	9	5	6	2	1
8	9	1	2	5	0	0	6	6
6	7	0	1	6	3	6	3	7
3	7	7	9	4	6	6	1	8
2	9	3	4	3	9	8	7	2
1	5	9	8	3	6	5	7	2
9	3	1	9	1	5	8	0	8
5	6	2	6	8	5	8	8	9
3	7	7	0	9	4	8	5	4

Generated  $\alpha_\theta$

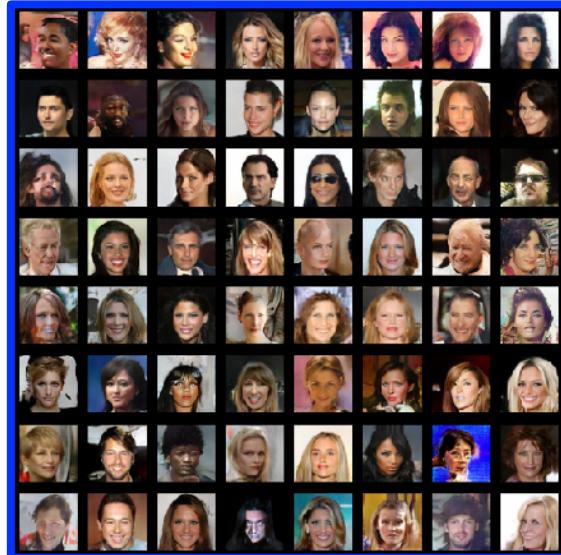
9	4	7	3	3	7	6	8
5	5	1	0	8	1	2	0
5	4	0	8	0	0	7	9
8	8	6	0	7	2	4	7
3	9	0	6	1	9	1	8
4	2	6	7	9	3	6	7
8	7	0	8	4	8	5	7
2	6	0	5	3	4	0	3



# Examples of Images Generation

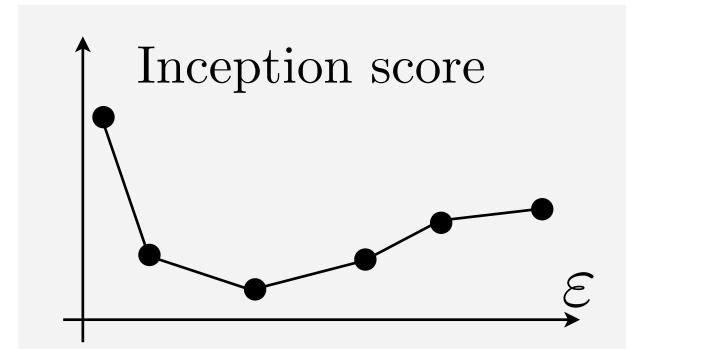
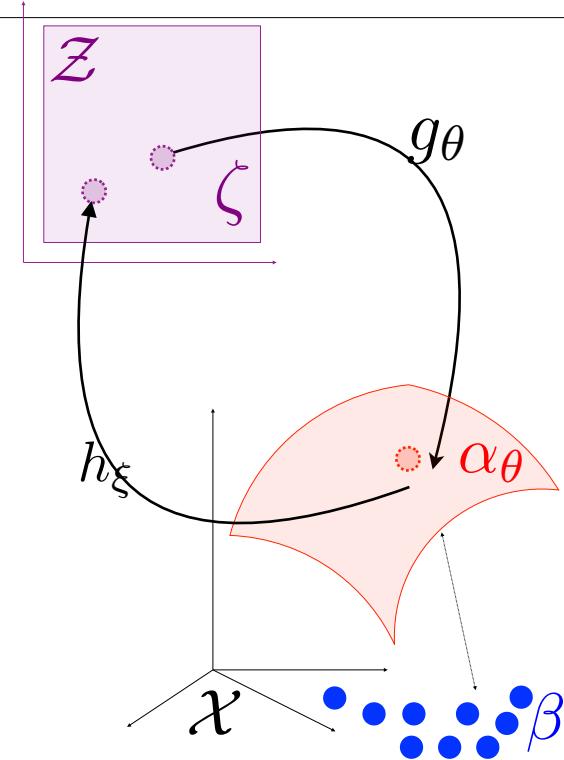
Inputs  $\beta$

3	4	2	1	9	5	6	2	1
8	9	1	2	5	0	0	6	6
6	7	0	1	6	3	6	3	7
3	7	7	9	4	6	6	1	8
2	9	3	4	3	9	8	7	2
1	5	9	8	3	6	5	7	2
9	3	1	9	1	5	8	0	8
5	6	2	6	8	5	8	8	9
3	7	7	0	9	4	8	5	4



Generated  $\alpha_\theta$

9	4	7	3	3	9	6	8
5	5	1	0	8	1	2	0
5	4	0	8	0	0	5	9
8	2	6	0	7	2	4	7
3	9	0	6	1	9	1	8
4	2	6	3	9	3	6	2
9	7	0	8	4	8	5	7
2	6	0	5	3	4	0	3



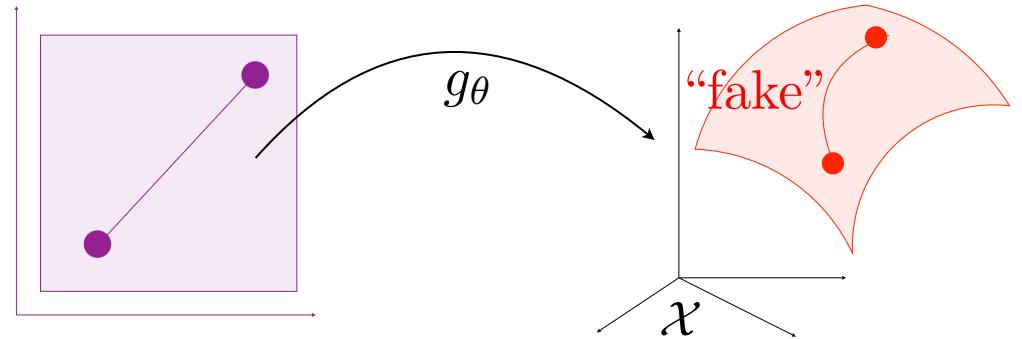
- Need to learn the metric  $d(x, y) = \|h_\xi(x) - h_\xi(y)\|$  (GANs)
- Influence of  $\epsilon$ ?
- Performance evaluation of generative models is an open problem.

Ian Goodfellow



*Progressive Growing of GANs for Improved Quality, Stability, and Variation*

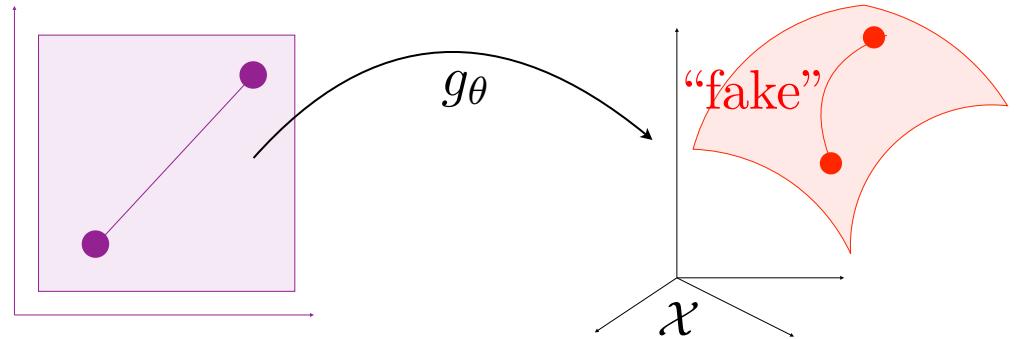
Tero Karras, Timo Aila, Samuli Laine,  
Jaakko Lehtinen, ICLR 2018



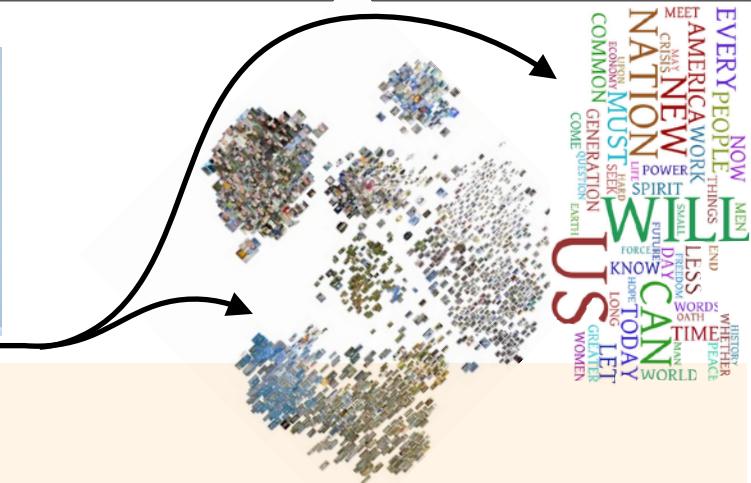


*Progressive Growing of GANs for Improved  
Quality, Stability, and Variation*

Tero Karras, Timo Aila, Samuli Laine,  
Jaakko Lehtinen, ICLR 2018



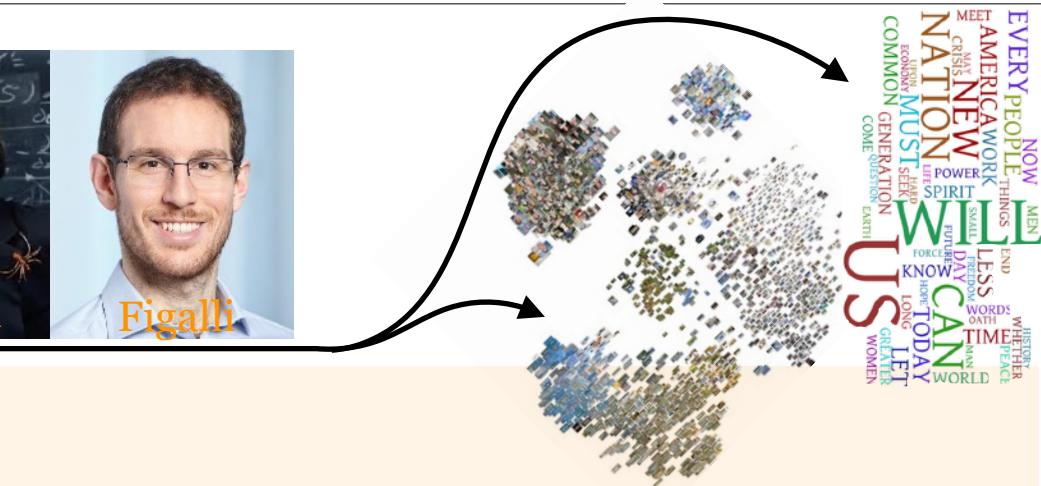
# Open Problems



Toward high-dimensional OT:

- Scalable geometrical loss functions in high dimension?
- Performance quality measures for unsupervised learning?

# Open Problems

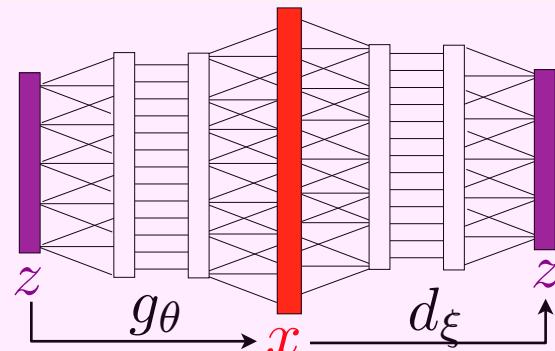


Toward high-dimensional OT:

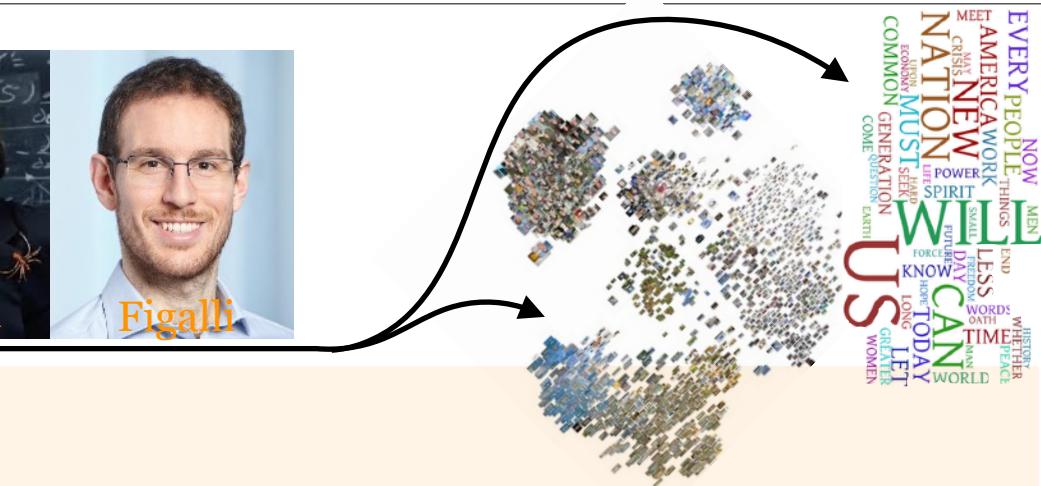
- Scalable geometrical loss functions in high dimension?
- Performance quality measures for unsupervised learning?

Metric learning for OT:

- Adversarial training to leverage multi scale priors?



# Open Problems

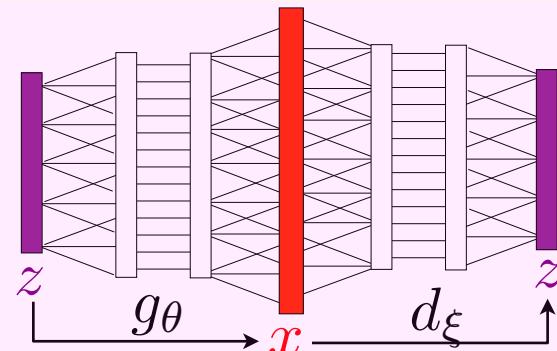


Toward high-dimensional OT:

- Scalable geometrical loss functions in high dimension?
- Performance quality measures for unsupervised learning?

Metric learning for OT:

- Adversarial training to leverage multi scale priors?



Beyond comparing measures:

- Learning for surfaces, graphs, metric spaces?
- Using Gromov-Wasserstein geometry?

