

*Mind the entropy bias in regularized OT :*

# Debiased Sinkhorn barycenters

Hicham Janati — Marco Cuturi — Alexandre Gramfort

Inria Saclay / ENSAE

Google Brain / ENSAE

Inria Saclay



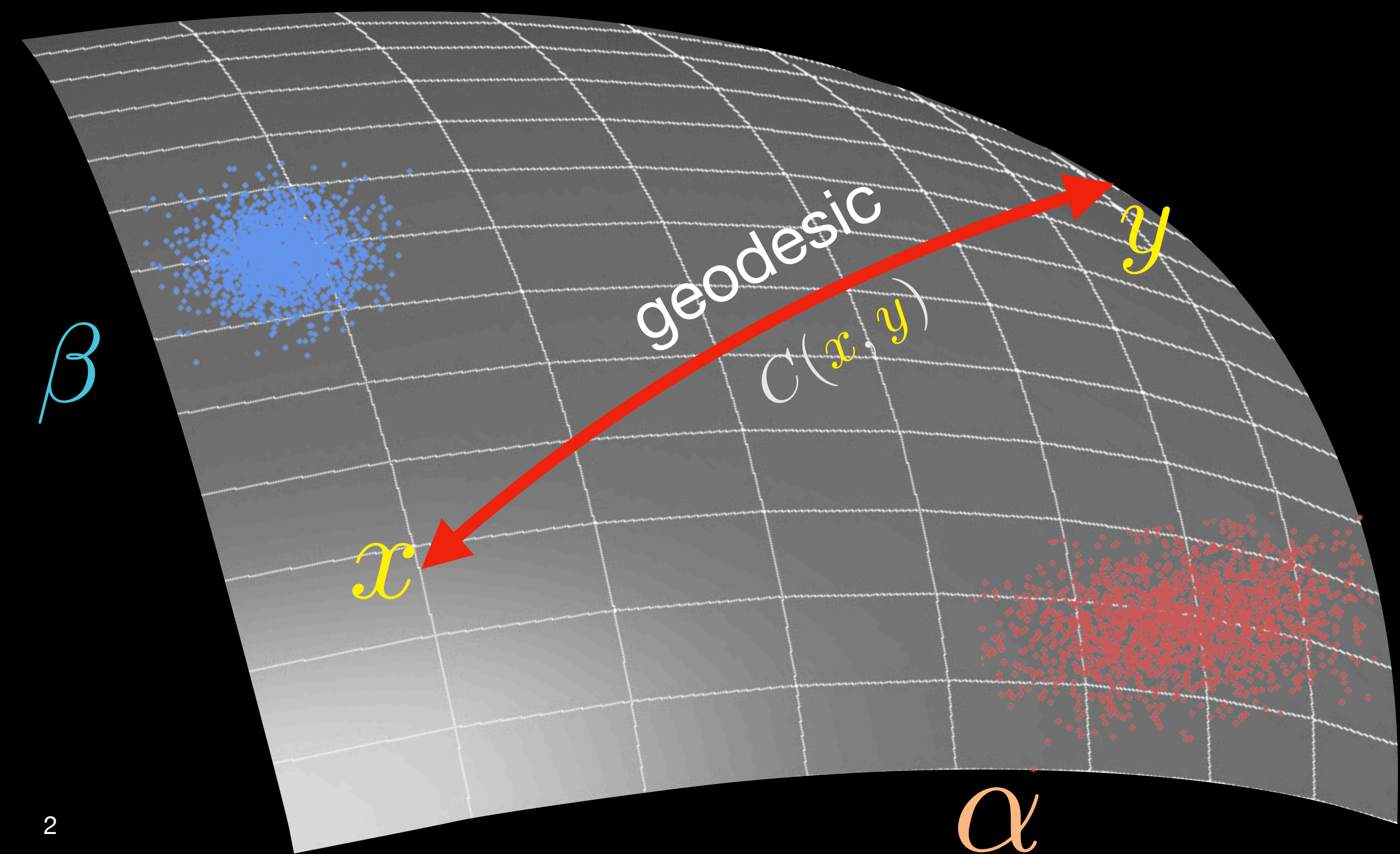
European  
Research  
Council

a metric space  $(\mathcal{X}, C)$

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

“unit transport cost”

$$\alpha, \beta \in \mathcal{P}(\mathcal{X})$$



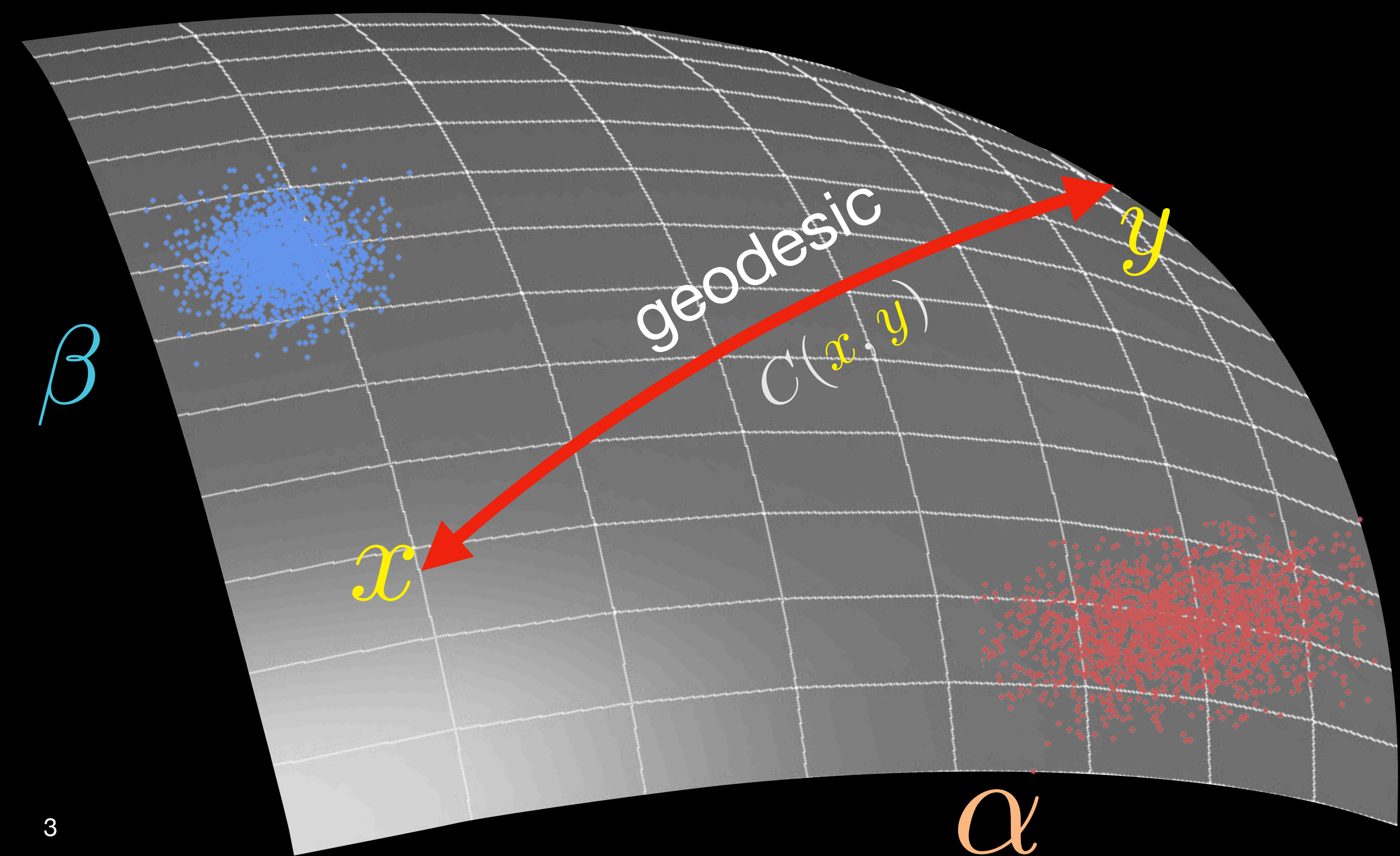
a metric space  $(\mathcal{X}, C)$

$$\text{OT}(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\substack{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) \\ \pi_1 = \alpha, \pi_2 = \beta}} \int_{\mathcal{X} \times \mathcal{X}} C(x, y) d\pi(x, y)$$

“Wasserstein distance”

Linear program  
 $O(n^3)$

$$\alpha, \beta \in \mathcal{P}(\mathcal{X})$$



## Entropy regularized OT

$$\text{OT}_{\varepsilon}^{m_1, m_2}(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\substack{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) \\ \pi_1 = \alpha, \pi_2 = \beta}} \int_{\mathcal{X} \times \mathcal{X}} C(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | m_1 \otimes m_2)$$

regularizer: relative entropy

 $m_1, m_2 \in \mathcal{P}(\mathcal{X})$  reference measures in  $\mathcal{P}(\mathcal{X})$ 

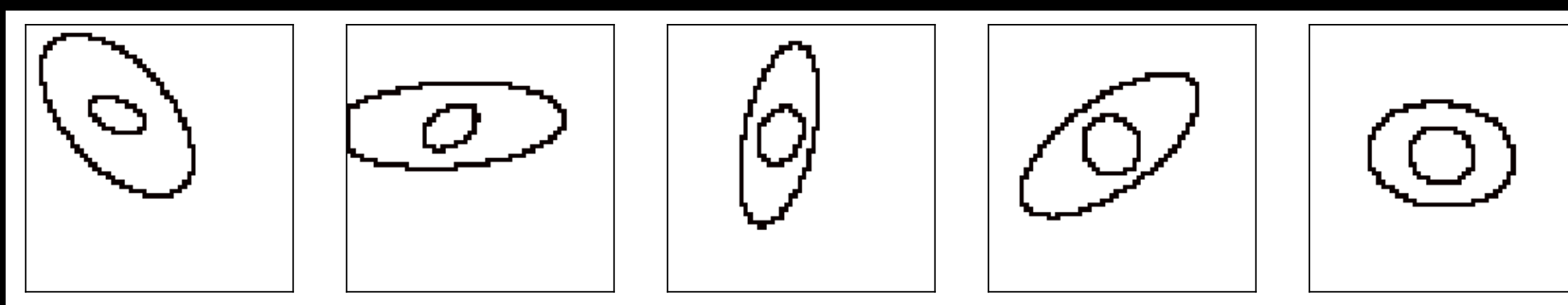
- + Computational cost: cubic  $\rightarrow$  quadratic
- + Breaks the curse of dimension

- Entropy bias
- Not a metric

$$\arg \min_{\alpha} \text{OT}_{\varepsilon}(\alpha, \beta) \neq \beta$$

## Entropy OT barycenters and bias

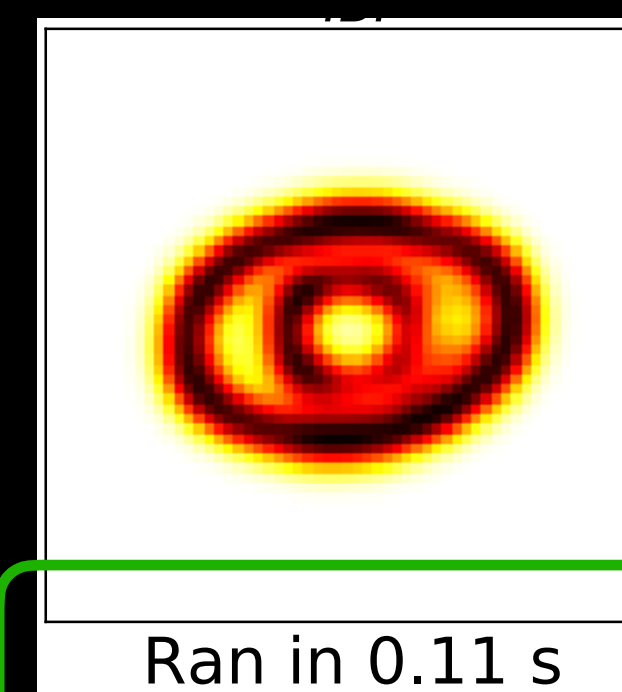
$$\mathcal{X} \text{ finite} \quad m_1, m_2 = \mathcal{U}(\mathcal{X})$$

 $\alpha_1, \alpha_2 \dots$ 

 $\alpha$ 

$\varepsilon = 0.5$

$\varepsilon = 0$

$$\arg \min_{\alpha} \sum_{k=1}^K w_k \text{OT}_{\varepsilon}(\alpha_k, \alpha)$$



Using Sinkhorn's algorithm

(Benamou et al, 15')

## Entropy regularized OT

$$\text{OT}_{\varepsilon}^{m_1, m_2}(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\substack{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) \\ \pi_1 = \alpha, \pi_2 = \beta}} \int_{\mathcal{X} \times \mathcal{X}} C(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | m_1 \otimes m_2)$$

regularizer: relative entropy

 $m_1, m_2 \in$  reference measures in  $\mathcal{P}(\mathcal{X})$ 

- + Computational cost: cubic  $\rightarrow$  quadratic
- + Breaks the curse of dimension

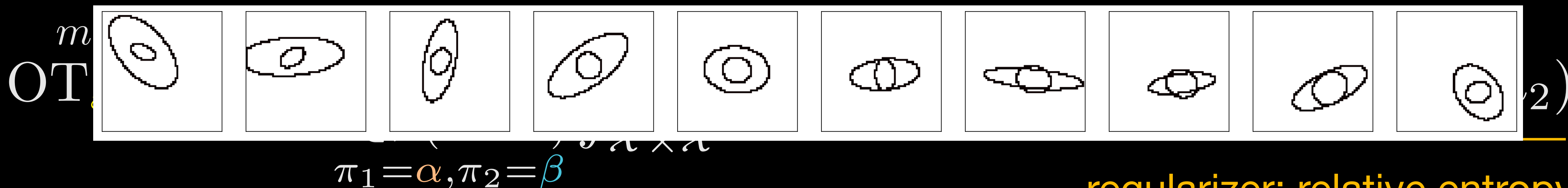
- Entropy bias
- Not a metric

(Feydy et al, 19):  $\mathcal{X}$  is compact  $m_1 = \alpha$   $m_2 = \beta$ 

$$S_{\varepsilon}(\alpha, \beta) \stackrel{\text{def}}{=} \text{OT}_{\varepsilon}(\alpha, \beta) - \frac{1}{2}(\text{OT}_{\varepsilon}(\alpha, \alpha) + \text{OT}_{\varepsilon}(\beta, \beta))$$

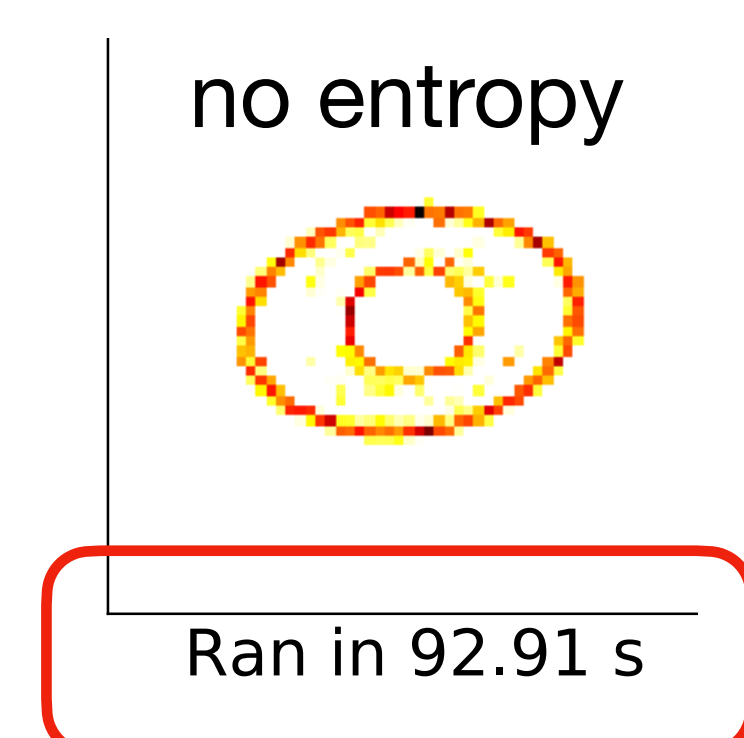
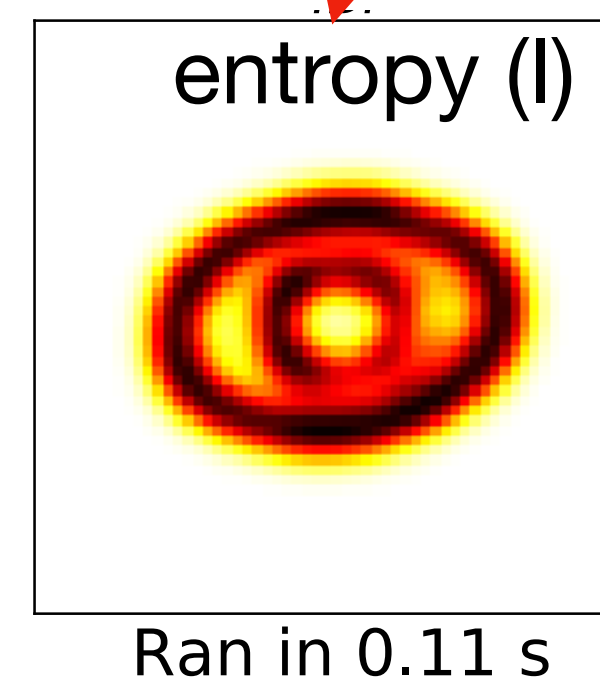
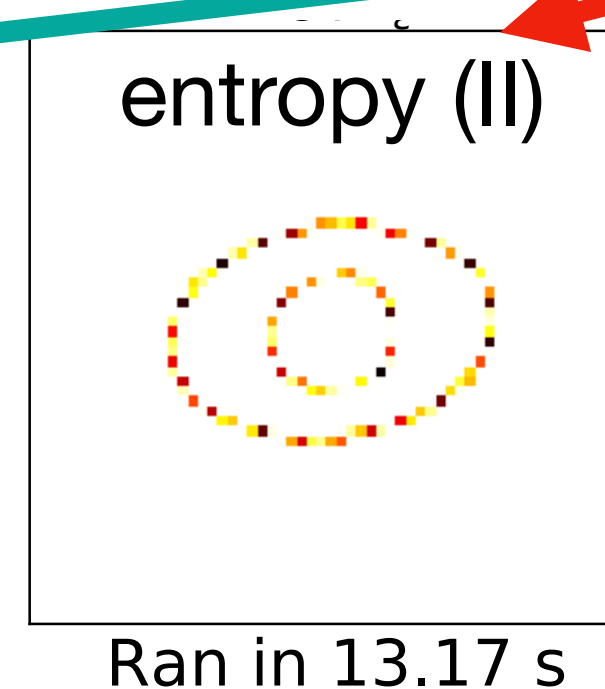
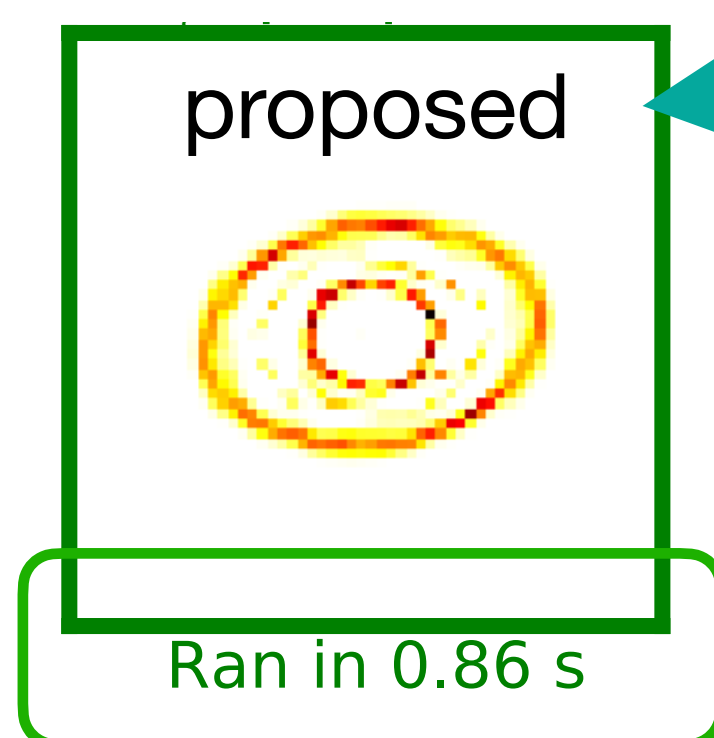
- ~~- Entropy bias~~
- +  $S_{\varepsilon}(\alpha, \beta) \geq 0$ ;
- +  $S_{\varepsilon}(\alpha, \beta) = 0 \Rightarrow \alpha = \beta$

## Entropy regularized OT



## Contributions:

- 1)  $\mathcal{S}_\varepsilon$  and  $\text{OT}_\varepsilon$  are convex and differentiable on non compact spaces ( $\mathcal{X} = \mathbb{R}^d$ )
- 2) Depending on  $m_1, m_2$  the entropy bias can be blurring or shrinking
- 3) a Fast modified Sinkhorn algorithm to compute debiased barycenters



## Entropy regularized OT

$$\text{OT}_{\varepsilon}^{m_1, m_2}(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\substack{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) \\ \pi_1 = \alpha, \pi_2 = \beta}} \int_{\mathcal{X} \times \mathcal{X}} C(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | m_1 \otimes m_2)$$

regularizer: relative entropy

$m_1, m_2 \in$  reference measures in  $\mathcal{P}(\mathcal{X})$

$$\text{KL}(\pi | r) \stackrel{\text{def}}{=} \int_{\mathcal{X} \times \mathcal{X}} \log \left( \frac{d\pi}{dr} \right) d\pi \quad \text{if } \pi \ll r \quad \text{else } +\infty$$



(0) Discrete case  $\mathcal{X} = \{x_1, \dots, x_n\}$ 

$$\alpha, \beta \in \Delta_n \stackrel{\text{def}}{=} \left\{ x \in \mathbb{R}_+^d \mid \sum_{i=1}^n x_i = 1 \right\}$$

$$\text{OT}_\varepsilon^{\mathcal{U}}(\alpha, \beta) = \min_{\substack{\pi \in \mathbb{R}_+^{n \times n} \\ \pi \mathbf{1} = \alpha, \pi^\top \mathbf{1} = \beta}} \langle (C(x_i, x_j)_{ij}), \pi \rangle + \varepsilon \langle \pi, \log(\pi) - \mathbf{1} \rangle$$

(Cuturi, Neurips 13')

(I) Lebesgue continuous case

$$\text{OT}_\varepsilon^{\mathcal{L}}(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\substack{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) \\ \pi_1 = \alpha, \pi_2 = \beta}} \int_{\mathcal{X} \times \mathcal{X}} C(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \mathcal{L} \otimes \mathcal{L})$$

(II) General case

$$\text{OT}_\varepsilon^{\otimes}(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\substack{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) \\ \pi_1 = \alpha, \pi_2 = \beta}} \int_{\mathcal{X} \times \mathcal{X}} C(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta)$$

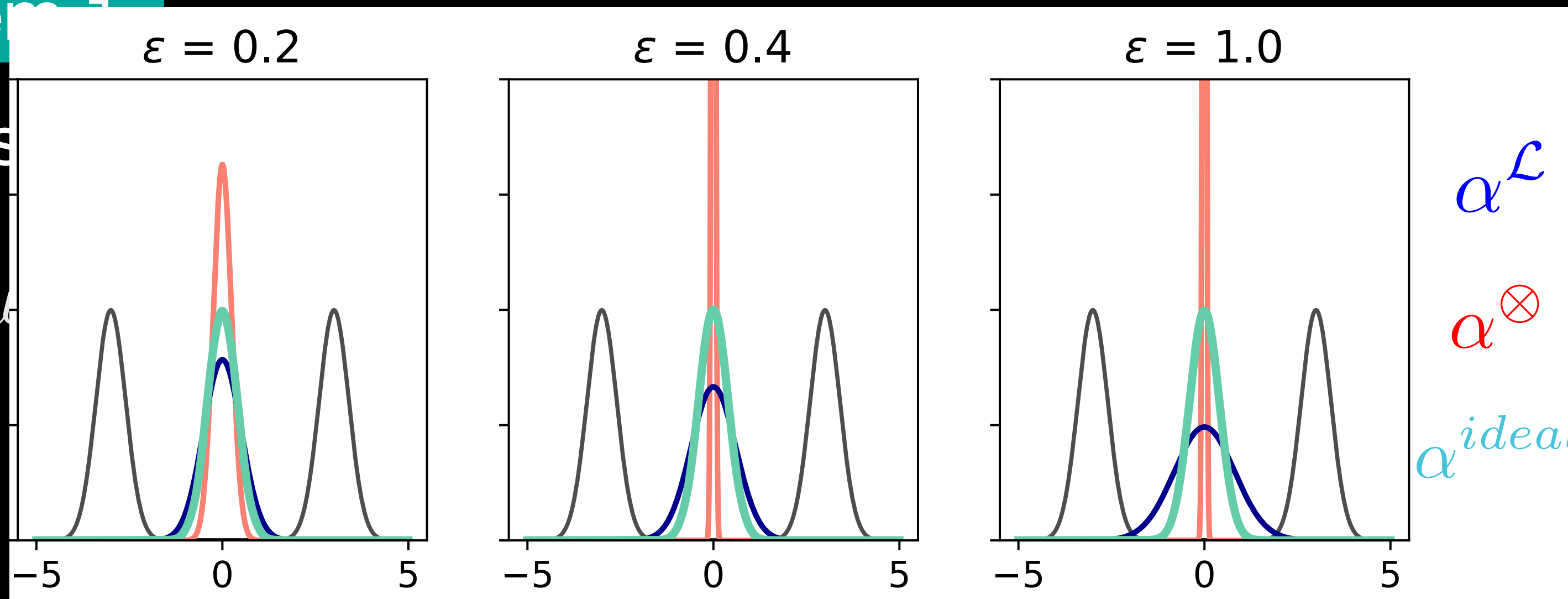
(Genevay 18', Feydy 19')

# Quantifying the entropy bias

## Theorem 1

Let  $\mathcal{G}$  the set of...

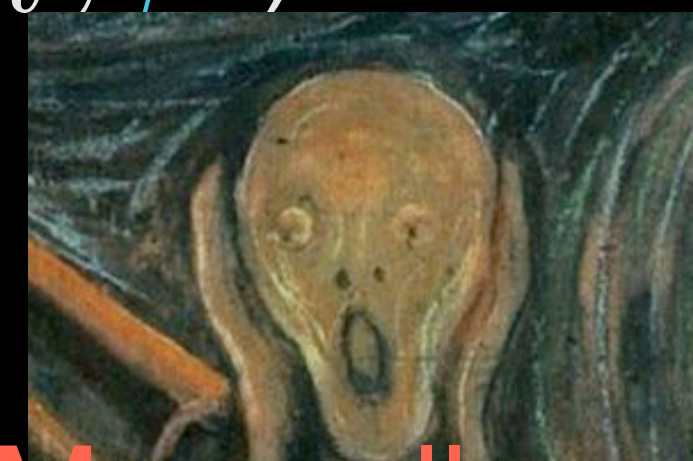
$$\bar{\mu} = \sum_{k=1}^K w_k \mu_k$$



$\alpha^{\mathcal{L}} \sim \mathcal{N}(\mu_k, \sigma^2)$   
 $\alpha^{\otimes}$   
 $\alpha^{ideal}$   
 smoothing

$$\alpha^{\mathcal{L}} \stackrel{\text{def}}{=} \arg \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k \text{OT}_{\epsilon}^{\mathcal{L}}(\alpha_k, \beta) \sim \mathcal{N}\left(\bar{\mu}, \sigma^2 + \frac{\epsilon^2}{2}\right)$$

(II) General case



May collapse to a dirac

$$\alpha^{\otimes} \stackrel{\text{def}}{=} \arg \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k \text{OT}_{\epsilon}^{\otimes}(\alpha_k, \beta) \sim \mathcal{N}\left(\bar{\mu}, \left(\sigma^2 - \frac{\epsilon^2}{2}\right)_+\right)$$

Entropy shrinking

# Quantifying the **entropy** bias

## Theorem 1

Let  $\mathcal{G}$  the set of sub-Gaussian measures in  $\mathbb{R}$  and  $\alpha_k \sim \mathcal{N}(\mu_k, \sigma^2)$

“Entropic OT is maximum likelihood deconvolution”

$$Y = X + \sigma^2 Z, \quad Z \sim \mathcal{N}(0, \text{Id})$$

$$P_X = \arg \min_P \text{OT}_{\sigma^2}^{\otimes} \left( P, \sum_{i=1}^n \frac{1}{n} \delta_{y_i} \right) \quad (\text{Rigollet \& Weed, 18'})$$

$$\alpha^{\mathcal{L}} \stackrel{\text{def}}{=} \arg \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k \text{OT}_{\epsilon}^{\mathcal{L}}(\alpha_k, \beta) \sim \mathcal{N} \left( \bar{\mu}, \sigma^2 + \frac{\epsilon^2}{2} \right)$$

(II) General case

$$\alpha^{\otimes} \stackrel{\text{def}}{=} \arg \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k \text{OT}_{\epsilon}^{\otimes}(\alpha_k, \beta) \sim \mathcal{N} \left( \bar{\mu}, \left( \sigma^2 - \frac{\epsilon^2}{2} \right)_+ \right)$$

Entropy shrinking

Yet, useful

# Fixing the entropy bias

(Feydy et al, 19):  $\mathcal{X}$  is compact

$$S_\varepsilon(\alpha, \beta) \stackrel{\text{def}}{=} \text{OT}_\varepsilon^\otimes(\alpha, \beta) - \frac{1}{2}(\text{OT}_\varepsilon^\otimes(\alpha, \alpha) + \text{OT}_\varepsilon^\otimes(\beta, \beta))$$

~~- Entropy bias~~

+  $S_\varepsilon(\alpha, \beta) \geq 0$ ;

+  $S_\varepsilon(\alpha, \beta) = 0 \Rightarrow \alpha = \beta$

## Theorem

Why?

Let  $\mathcal{G}$  the set of sub-Gaussian measures in  $\mathbb{R}$  and  $\alpha_k \sim \mathcal{N}(\mu_k, \sigma^2)$

$\bar{\mu} = \sum_{k=1}^K w_k \mu_k$   $C(x, y) = (x - y)^2$  then:

$$\alpha^S \stackrel{\text{def}}{=} \arg \min_{\beta \in \mathcal{G}} \sum_{k=1}^K w_k S_\varepsilon(\alpha_k, \beta) \sim \mathcal{N}(\bar{\mu}, \sigma^2)$$

To prove the convexity and differentiability of  $S_\varepsilon$  and  $\text{OT}_\varepsilon$  and characterize the barycentric optimality

# Barycentric Algorithms: two different views

## Lagrangian (Free supports)



*Distributions represented  
by point clouds (positions, weights)*

- + scalable for sparse distributions in high dimensions
- Optimization over weights and **positions**: requires several Sinkhorn loops

## Eulerian (Fixed supports)



*Distributions represented  
by histograms on a fixed grid*

- + Parallelizable on regular grids (e.g. images)
- Not scalable memory footprint

*Done by (Luise, 19')*

## Sinkhorn algorithm and IBP

$$\alpha^{\mathcal{U}} \stackrel{\text{def}}{=} \arg \min_{\beta \in \Delta_n} \sum_{k=1}^K w_k \text{OT}_{\varepsilon}^{\mathcal{U}}(\alpha_k, \beta) \quad \text{can be written as a KL projection:}$$

$$\min_{\substack{\pi_1, \dots, \pi_K \in \mathbb{R}_+^{n \times n} \\ \pi_k \mathbf{1} = \alpha_k \\ \pi_1^{\top} \mathbf{1} = \dots = \pi_K^{\top} \mathbf{1}}} \sum_{k=1}^K w_k \text{KL}(\pi_k | e^{-\frac{(C_{ij})}{\varepsilon}})$$

Solved using  
 “Iterative Bregman Projections” (IBP)  
 a.k.a. generalized Sinkhorn

Not possible with the product measure as reference :(

$$\alpha^{\otimes} \stackrel{\text{def}}{=} \arg \min_{\beta \in \Delta_n} \sum_{k=1}^K w_k \text{OT}_{\varepsilon}^{\otimes}(\alpha_k, \beta)$$

What about

$S_{\varepsilon}$  ?

# Debiased Sinkhorn algorithm and IBP

$\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X})$ ,  $\pi_1 = \alpha$  and  $\pi_2 = \beta$  then:

$$\text{KL}(\pi | m_1 \otimes m_2) = \text{KL}(\pi | \alpha \otimes \beta) + \text{KL}(\alpha | m_1) + \text{KL}(\beta | m_2)$$

(Di marino, 19')

Then  $S_\varepsilon$  is independent of the reference measure:

If discrete:  $S_\varepsilon^{\mathcal{U}}(\alpha, \beta) = S_\varepsilon^{\otimes}(\alpha, \beta)$

If Lebesgue-continuous:  $S_\varepsilon^{\mathcal{L}}(\alpha, \beta) = S_\varepsilon^{\otimes}(\alpha, \beta)$

# Debiased Sinkhorn algorithm and IBP

Using  $S_{\varepsilon}^{\mathcal{U}}(\alpha, \beta) = S_{\varepsilon}^{\otimes}(\alpha, \beta)$

$$\alpha^S \stackrel{\text{def}}{=} \arg \min_{\beta \in \Delta_n} \sum_{k=1}^K w_k S_{\varepsilon}^{\mathcal{U}}(\alpha_k, \beta) \quad \text{equivalent to:}$$

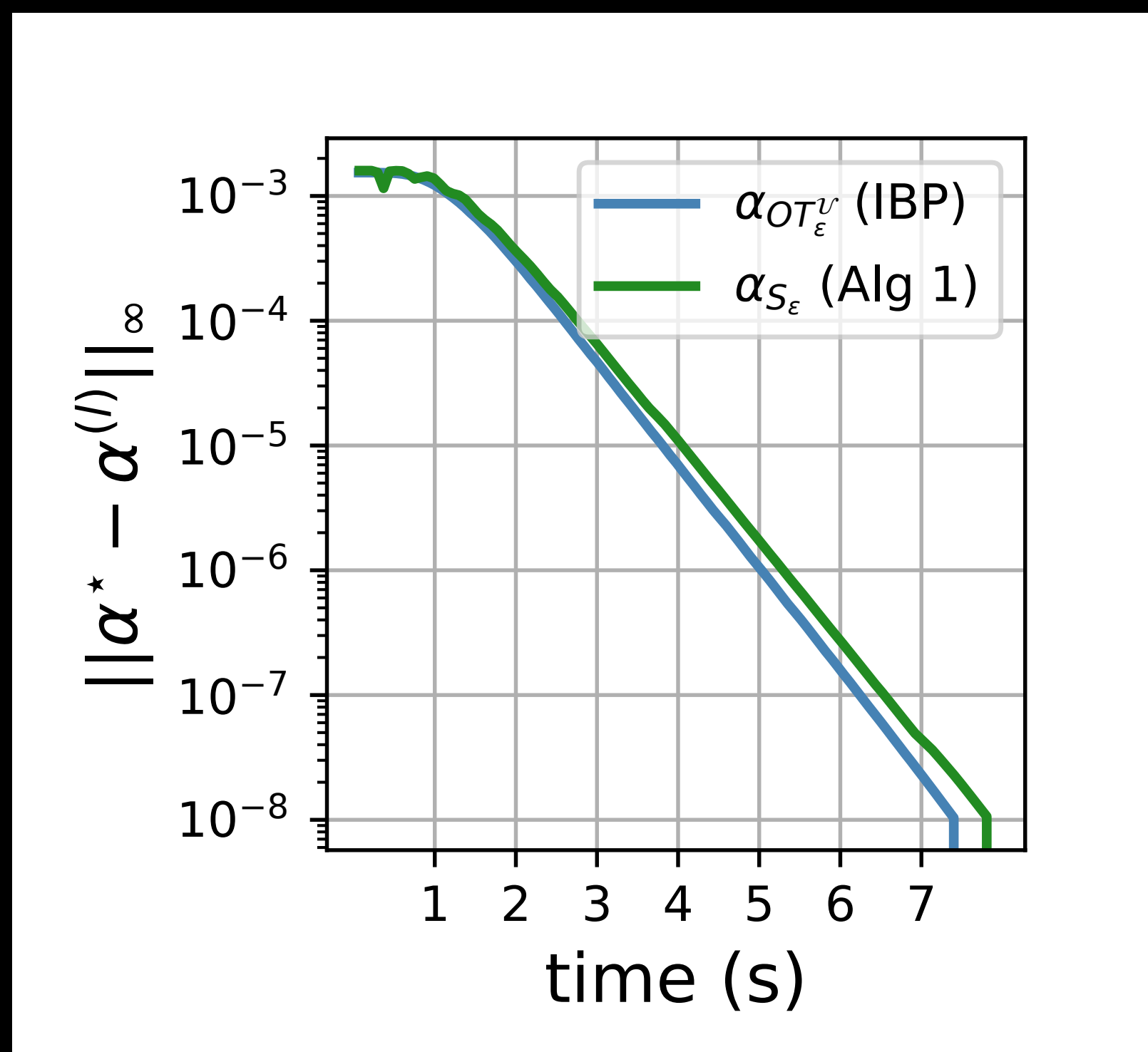
$$\min_{\substack{\pi_1, \dots, \pi_K \in \mathbb{R}_+^{n \times n} \\ \pi_k \mathbf{1} = \alpha_k \\ \pi_1^{\top} \mathbf{1} = \dots = \pi_K^{\top} \mathbf{1} \\ d \in \mathbb{R}_+^n}} \sum_{k=1}^K w_k \text{KL}(\pi_k | e^{-\frac{(C_{ij})}{\varepsilon}} d_j) + \frac{\varepsilon}{2} \langle d - \mathbf{1}, e^{-\frac{(C_{ij})}{\varepsilon}} (d - \mathbf{1}) \rangle$$

Alternating block minimization:

- 1) IBP loop
- 2) Symmetric Sinkhorn



# Debiased Sinkhorn algorithm and IBP



## Algorithm 1 Debiased Sinkhorn Barycenter

**Input:**  $\alpha_1, \dots, \alpha_K, \mathbf{K} = e^{-\frac{c}{\epsilon}}$

**Output:**  $\alpha_{S_\epsilon}$

Initialize all scalings  $(b_k), d$  to  $\mathbb{1}$ ,

**repeat**

**for**  $k = 1$  **to**  $K$  **do**

$$a_k \leftarrow \left( \frac{\alpha_k}{\mathbf{K} b_k} \right)$$

**end for**

$$\alpha \leftarrow d \odot \prod_{k=1}^K (\mathbf{K}^\top a_k)^{w_k}$$

**for**  $k = 1$  **to**  $K$  **do**

$$b_k \leftarrow \left( \frac{\alpha}{\mathbf{K}^\top a_k} \right)$$

**end for**

$$d \leftarrow \sqrt{d \odot \left( \frac{\alpha}{\mathbf{K} d} \right)}$$

**until convergence**

# Take home message

- The entropic barycenter of univariate Gaussians is Gaussian
- The entropic bias can be smoothing or shrinking
- The debiased barycenter can be computed on a GPU-friendly modified Sinkhorn algorithm