

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

# Learning with Entropy-Regularized Optimal Transport

Aude Genevay

MIT CSAIL

Sept. 23rd 2020

*Joint work with Gabriel Peyré, Marco Cuturi, Francis Bach, Lénaïc Chizat*

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

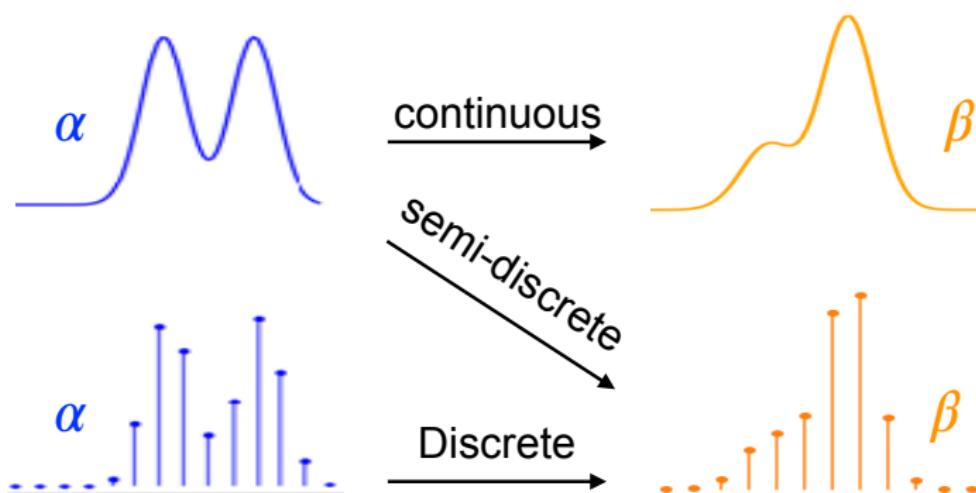
Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Comparing Probability Measures



Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Discrete Setting (Quantization)

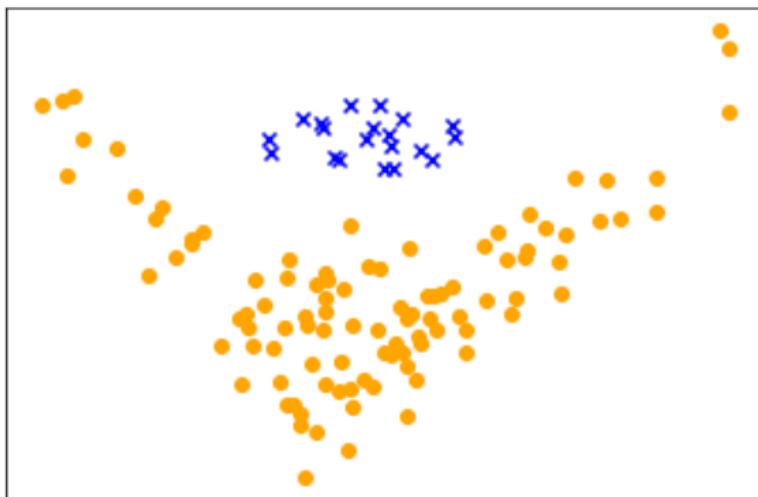


Figure 1 –  $\min_{(\mathbf{x}_1, \dots, \mathbf{x}_k)} \mathcal{D}\left(\frac{1}{k} \sum_{i=1}^k \delta \mathbf{x}_i, \frac{1}{n} \sum_{j=1}^n \delta \mathbf{y}_j\right)$

Distances

oooooooooooooooo

Entropic Regularization

o  
oooo  
o  
oooooo

Sinkhorn Divergences

o  
ooooooo  
ooooooo

Conclusion

oo

## Discrete Setting (Quantization)

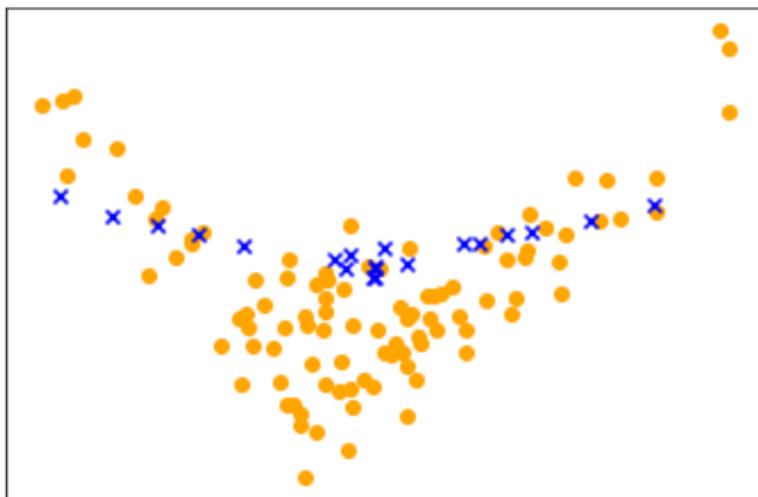


Figure 1 –  $\min_{(x_1, \dots, x_k)} \mathcal{D}\left(\frac{1}{k} \sum_{i=1}^k \delta x_i, \frac{1}{n} \sum_{j=1}^n \delta y_j\right)$

Distances

oooooooooooooooo

Entropic Regularization

o  
oooo  
o  
oooooo

Sinkhorn Divergences

o  
ooooooo  
ooooooo

Conclusion

oo

## Discrete Setting (Quantization)

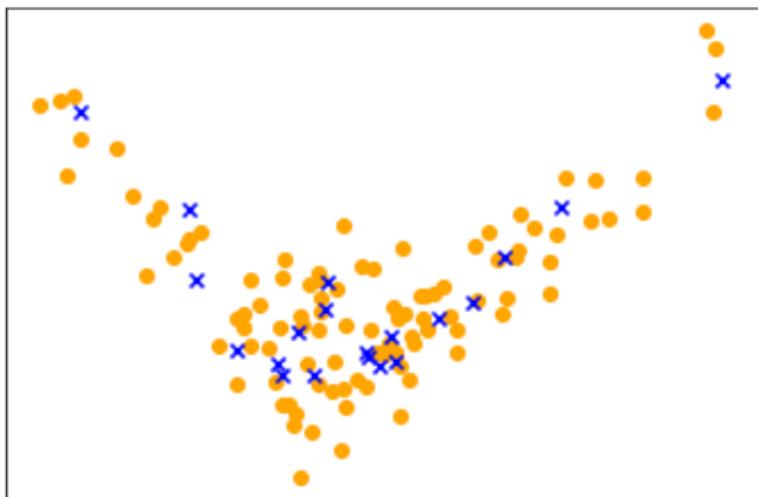


Figure 1 –  $\min_{(x_1, \dots, x_k)} \mathcal{D}\left(\frac{1}{k} \sum_{i=1}^k \delta x_i, \frac{1}{n} \sum_{j=1}^n \delta y_j\right)$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Discrete Setting (Quantization)

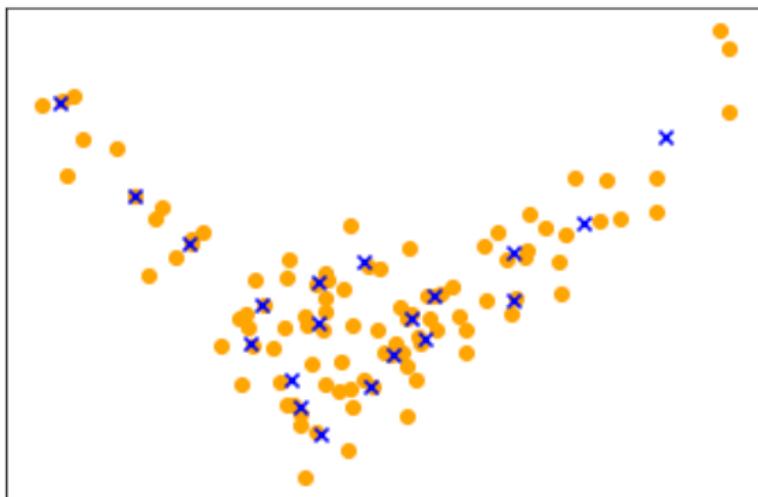


Figure 1 –  $\min_{(x_1, \dots, x_k)} \mathcal{D}\left(\frac{1}{k} \sum_{i=1}^k \delta x_i, \frac{1}{n} \sum_{j=1}^n \delta y_j\right)$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Semi-discrete Setting (Density Fitting)

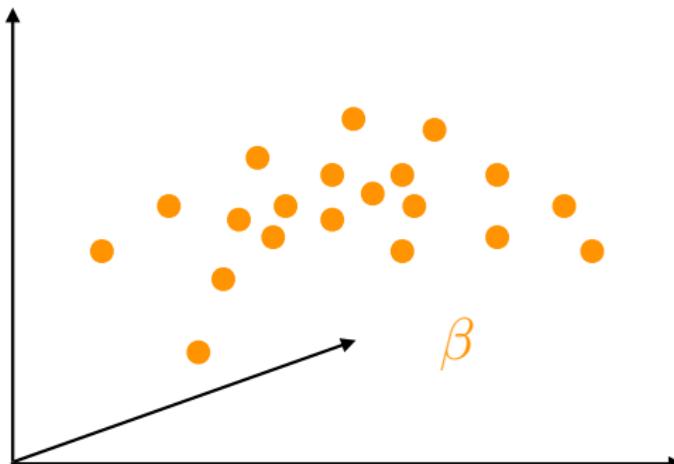


Figure 2 –  $\min_{\theta} \mathcal{D}(\alpha_{\theta}, \beta)$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Semi-discrete Setting (Density Fitting)

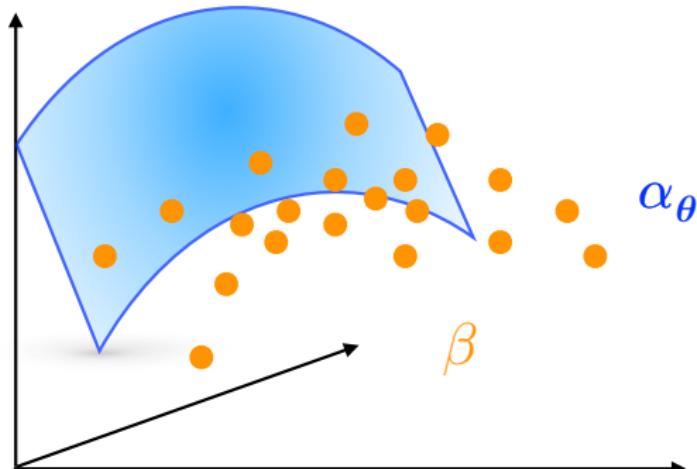


Figure 2 –  $\min_{\theta} \mathcal{D}(\alpha_{\theta}, \beta)$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Semi-discrete Setting (Density Fitting)

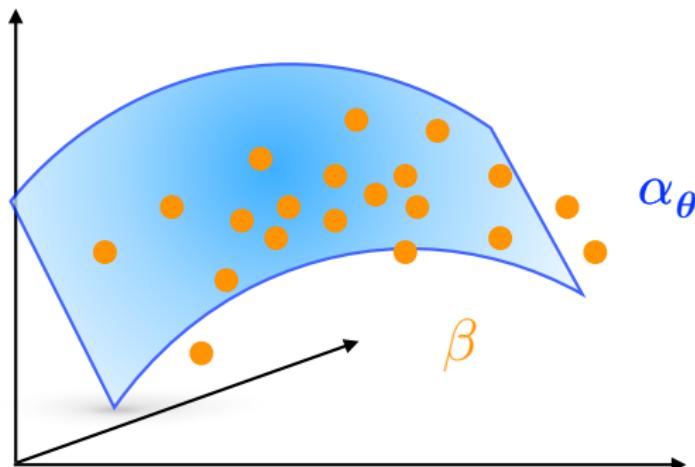


Figure 2 –  $\min_{\theta} \mathcal{D}(\alpha_{\theta}, \beta)$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Semi-discrete Setting (Density Fitting)

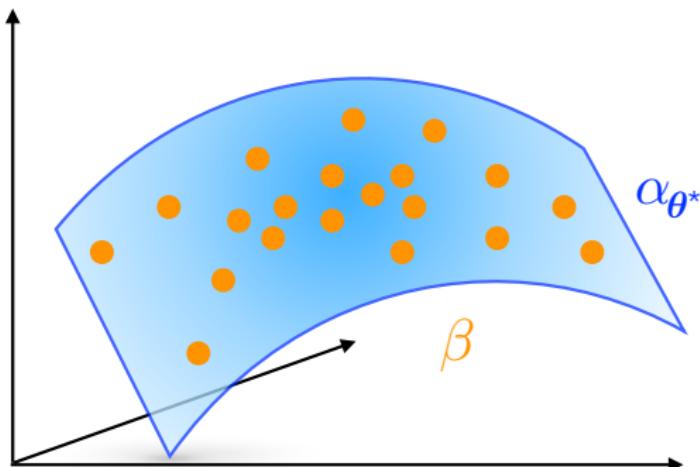


Figure 2 –  $\min_{\theta} \mathcal{D}(\alpha_{\theta}, \beta)$

**Distances**

●ooooooooooooooo

**Entropic Regularization**

○  
○○○○  
○  
○○○○○○

**Sinkhorn Divergences**

○  
○○○○○○  
○○○○○○○○

**Conclusion**

○○

- 1 Notions of Distance between Measures
- 2 Entropic Regularization of Optimal Transport
- 3 Sinkhorn Divergences : Interpolation between OT and MMD
- 4 Conclusion

## Distances

## Entropic Regularization

A 4x3 grid of 12 small circles, arranged in four rows and three columns.

## Sinkhorn Divergences

10

## Conclusion

$\varphi$ -divergences (Czisar '63)

## Definition ( $\varphi$ -divergence)

Let  $\varphi$  convex l.s.c. function such that  $\varphi(1) = 0$ , the  $\varphi$ -divergence  $D_\varphi$  between two measures  $\alpha$  and  $\beta$  is defined by :

$$D_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha(x)}{d\beta(x)}\right) d\beta(x).$$

## Example (Kullback Leibler Divergence)

$$D_{KL}(\alpha \mid\beta) = \int_{\mathcal{X}} \log \left( \frac{d\alpha}{d\beta}(x) \right) d\alpha(x) \quad \leftrightarrow \quad \varphi(x) = x \log(x)$$

## Distances

○○●○○○○○○○○○○○○○○○○

Entropic Regularization

1

## Sinkhorn Divergences

10

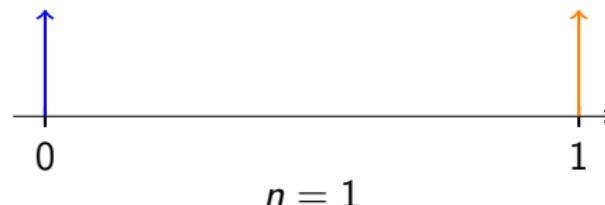
## Conclusion

6

## Weak Convergence of measures

## Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



Distances

○○●○○○○○○○○○○○○○○

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

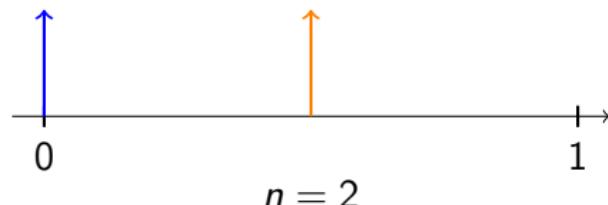
Conclusion

○○

## Weak Convergence of measures

### Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



## Distances

○●○○○○○○○○○○○○○○○○

Entropic Regularization

1

Sinkhorn Divergences

A 3x5 grid of 15 small circles, arranged in three rows and five columns.

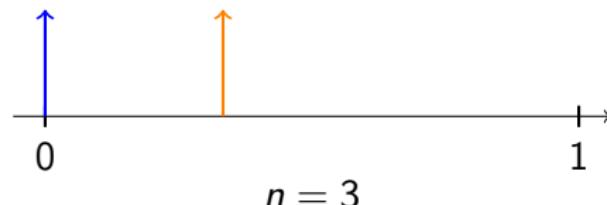
## Conclusion

○○

## Weak Convergence of measures

## Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



Distances

○○●○○○○○○○○○○○○○○○○

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

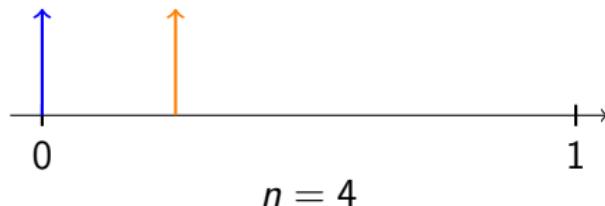
Conclusion

○○

## Weak Convergence of measures

### Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n}$  :  $D_{KL}(\alpha_n | \alpha) = +\infty$ .



## Distances

○●○○○○○○○○○○○○○○○○

Entropic Regularization

1

## Sinkhorn Divergences

10

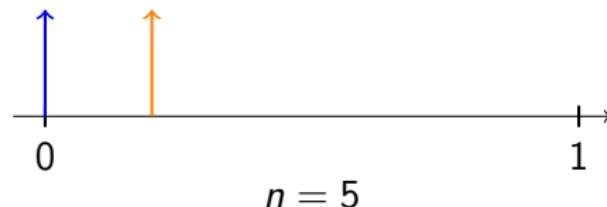
## Conclusion

○○

## Weak Convergence of measures

## Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



## Distances

○●○○○○○○○○○○○○○○○○

Entropic Regularization

1

## Sinkhorn Divergences

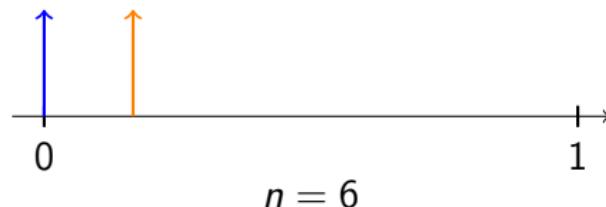
## Conclusion

○○

## Weak Convergence of measures

## Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



## Distances

○●○○○○○○○○○○○○○○○○

Entropic Regularization

1

## Sinkhorn Divergences

10

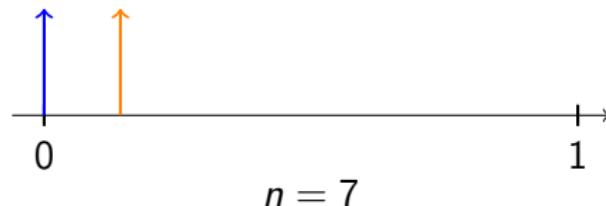
## Conclusion

○○

## Weak Convergence of measures

## Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



Distances

○○●○○○○○○○○○○○○○○○○

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

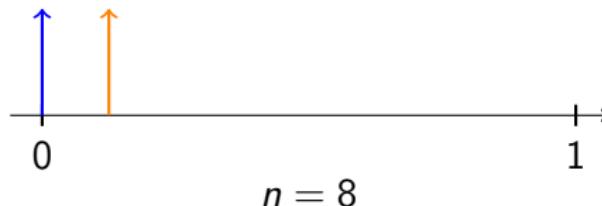
Conclusion

○○

## Weak Convergence of measures

### Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n}$  :  $D_{KL}(\alpha_n | \alpha) = +\infty$ .



## Distances

○●○○○○○○○○○○○○○○○○

Entropic Regularization

1

## Sinkhorn Divergences

10

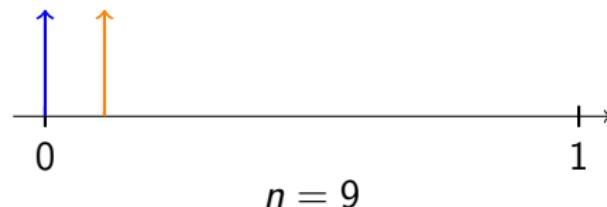
### Conclusion

○○

## Weak Convergence of measures

## Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



## Distances

○●○○○○○○○○○○○○○○○○

Entropic Regularization

A 2x4 grid of circles, representing a 2x4 matrix.

Sinkhorn Divergences

A 3x5 grid of 15 small circles, arranged in three rows and five columns.

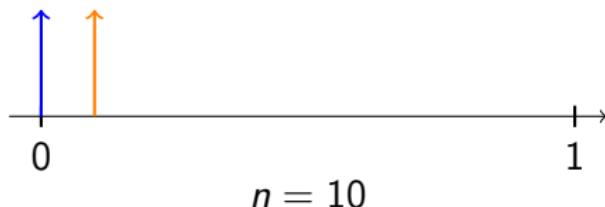
## Conclusion

○○

## Weak Convergence of measures

## Example

On  $\mathbb{R}$ ,  $\alpha = \delta_0$  and  $\alpha_n = \delta_{1/n} : D_{KL}(\alpha_n | \alpha) = +\infty$ .



## Definition (Weak Convergence)

$\alpha_n$  weakly converges to  $\alpha$ , ( denoted  $\alpha_n \rightharpoonup \alpha$ )

$$\Leftrightarrow \int f(x) d\alpha_{\textcolor{brown}{n}}(x) \rightarrow \int f(x) d\alpha(x) \quad \forall f \in \mathcal{C}_b(\mathcal{X}).$$

Let  $\mathcal{D}$  distance between measures ,  $\mathcal{D}$  metrises weak convergence IFF  $(\mathcal{D}(\alpha_n, \alpha) \rightarrow 0 \Leftrightarrow \alpha_n \rightharpoonup \alpha)$ .

## Distances

A horizontal row of 15 small circles. The 6th circle from the left is filled with a dark blue color, while all other circles are unfilled.

## Entropic Regularization

10

## Sinkhorn Divergences

A 3x5 grid of 15 small circles, arranged in three rows and five columns.

## Conclusion

00

# Integral Probability Metrics

Definition (Integral Probability Metric, Muller '97)

$$d_{\mathcal{F}}(\alpha, \beta) \stackrel{\text{def.}}{=} \left( \sup_{f \in \mathcal{F}} |\mathbb{E}_{\alpha}(f(X)) - \mathbb{E}_{\beta}(f(Y))| \right)$$

## Examples :

- Total Variation :  $\mathcal{F} = \{f \mid \|f\|_\infty \leq 1\}$  (upper-bounded by 1)
  - 1-Wasserstein :  $\mathcal{F} = \{f \mid \|f\|_{Lip} \leq 1\}$  (1-Lipschitz)
  - MMD :  $\mathcal{F} = \{f \mid \|f\|_{\mathcal{H}} \leq 1\}$  (unit ball of RKHS  $\mathcal{H}$ )

Distances

○○○●○○○○○○○○○○

Entropic Regularization

- 
- 
- 
- 

Sinkhorn Divergences

- 
- 
- 

Conclusion

○○

## Maximum Mean Discrepancies (Gretton '06)

### Definition (RKHS)

Let  $\mathcal{H}$  Hilbert space with kernel  $k$ , then  $\mathcal{H}$  is a Reproducing Kernel Hilbert Space (RKHS) IFF :

- ①  $\forall x \in \mathcal{X}, \quad k(x, \cdot) \in \mathcal{H},$
- ②  $\forall f \in \mathcal{H}, \quad f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{H}}.$

$$\begin{aligned} MMD_k^2(\alpha, \beta) &\stackrel{\text{def.}}{=} \left( \sup_{\{f \mid \|f\|_{\mathcal{H}} \leq 1\}} |\mathbb{E}_{\alpha}(f(X)) - \mathbb{E}_{\beta}(f(Y))| \right)^2 \\ &= \mathbb{E}_{\alpha \otimes \alpha}[k(X, X')] + \mathbb{E}_{\beta \otimes \beta}[k(Y, Y')] \\ &\quad - 2\mathbb{E}_{\alpha \otimes \beta}[k(X, Y)]. \end{aligned}$$

Distances

oooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

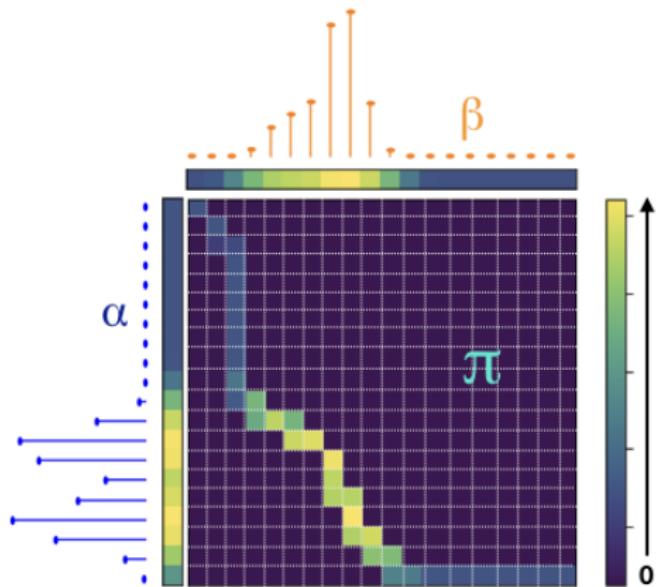
○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Optimal Transport (Monge 1781, Kantorovitch '42)

- $c(x, y)$  : cost of moving a unit of mass from  $x$  to  $y$
- $\pi(x, y)$  (coupling) : how much mass moves from  $x$  to  $y$



Distances

oooooooo●oooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○○○

Conclusion

○○

## The Wasserstein Distance

Minimal cost of moving all the mass from  $\alpha$  to  $\beta$ ?

Let  $\alpha \in \mathcal{M}_+^1(\mathcal{X})$  and  $\beta \in \mathcal{M}_+^1(\mathcal{Y})$ ,

$$W_c(\alpha, \beta) = \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) \quad (\mathcal{P})$$

For  $c(x, y) = \|x - y\|_2^p$ ,  $W_c(\alpha, \beta)^{1/p}$  is the **p-Wasserstein distance**.

Distances

oooooooo●oooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Optimal Transport vs. MMD

MMD

OT

sample complexity

$$(\frac{1}{\sqrt{n}})$$

$O(n^{-1/d})$   
(curse of dimension)

computation

$$O(n^2)$$

$$O(n^3 \log(n))$$

Distances

oooooooooooo

Entropic Regularization

o  
oooo  
o  
oooooo

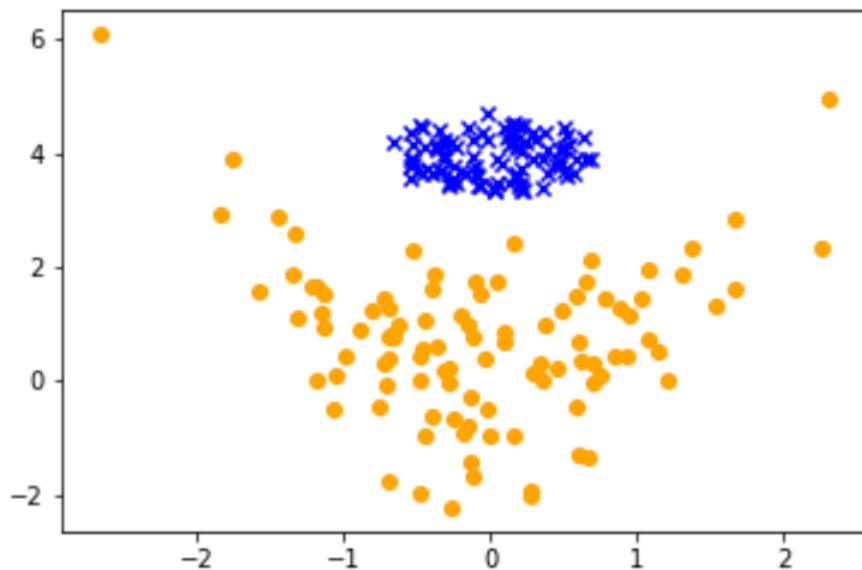
Sinkhorn Divergences

o  
ooooooo  
oooooooooooo

Conclusion

oo

## Simple example



$$\min_{(x_1, \dots, x_n)} \mathcal{D}\left(\frac{1}{n} \sum_{i=1}^n \delta x_i, \frac{1}{n} \sum_{j=1}^n \delta y_j\right)$$

**Distances**

oooooooooooo●oooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Gradient descent for $MMD$ (gaussian kernel)

**Distances**

oooooooooooo●ooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

Gradient descent for *MMD* (energy  
distance)

**Distances**

oooooooooooo●oooo

**Entropic Regularization**

○  
oooo  
○  
oooooo

**Sinkhorn Divergences**

○  
ooooooo  
ooooooo

**Conclusion**

oo

## Gradient descent for $OT$

Distances

oooooooooooo●ooo

Entropic Regularization

○  
○○○○  
○  
○○○○○

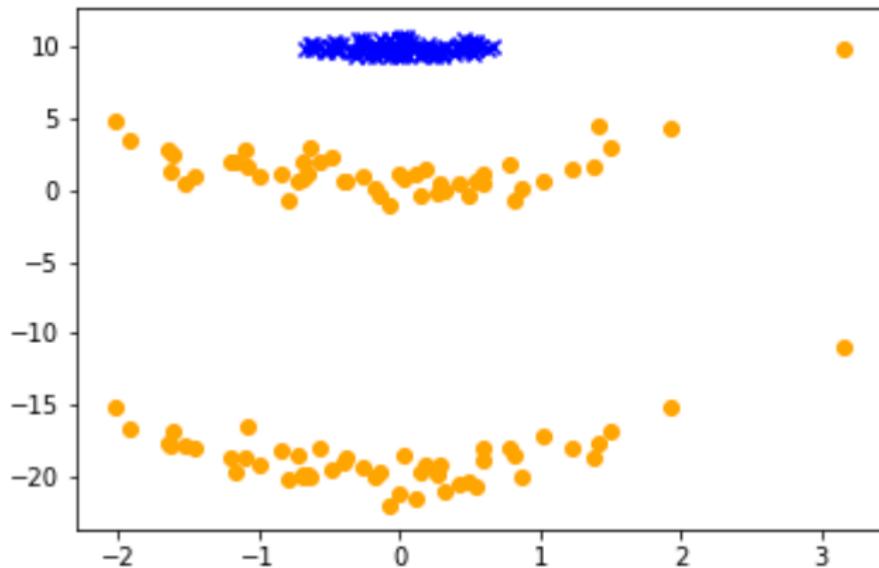
Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Another example



$$\min_{(x_1, \dots, x_n)} \mathcal{D}\left(\frac{1}{n} \sum_{i=1}^n \delta x_i, \frac{1}{n} \sum_{j=1}^n \delta y_j\right)$$

**Distances**

oooooooooooo●oo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Gradient descent for *MMD*

**Distances**

oooooooooooooo●○

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## Gradient descent for $OT$

Distances

oooooooooooooo●

Entropic Regularization

○  
○○○  
○  
○○○○○

Sinkhorn Divergences

○  
○○○○○  
○○○○○○○

Conclusion

○○

## Optimal Transport vs. MMD

MMD

OT

sample complexity

$$(\frac{1}{\sqrt{n}})$$

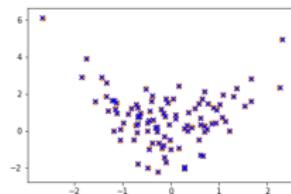
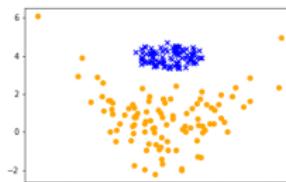
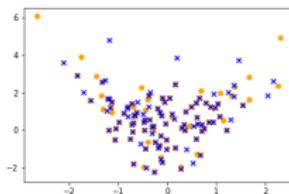
$O(n^{-1/d})$   
(curse of dimension)

computation

$$O(n^2)$$

$$O(n^3 \log(n))$$

better gradients!



$\min_{(x_1, \dots, x_k)} \mathcal{D}(\frac{1}{k} \sum_{i=1}^k \delta x_i, \frac{1}{n} \sum_{j=1}^n \delta y_j)$  after 200 steps of grad. descent.

Distances  
ooooooooooooooo

Entropic Regularization



Sinkhorn Divergences  
o ooooooo oooooooo

Conclusion  
oo

## 1 Notions of Distance between Measures

## 2 Entropic Regularization of Optimal Transport

The basics

A magic regularizing tool

Sample Complexity

## 3 Sinkhorn Divergences : Interpolation between OT and MMD

## 4 Conclusion

Distances

oooooooooooooooo

Entropic Regularization



Sinkhorn Divergences

A legend for Sinkhorn Divergences. It shows two open circles, followed by two horizontal lines of open circles.

Conclusion

oo

The basics

## Entropic Regularization (Cuturi '13)

Let  $\alpha \in \mathcal{M}_+^1(\mathcal{X})$  and  $\beta \in \mathcal{M}_+^1(\mathcal{Y})$ ,

$$W_c(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) \quad (\mathcal{P})$$

Distances

oooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Sinkhorn Divergences

- 
- 
- 

Conclusion

○○

The basics

## Entropic Regularization (Cuturi '13)

Let  $\alpha \in \mathcal{M}_+^1(\mathcal{X})$  and  $\beta \in \mathcal{M}_+^1(\mathcal{Y})$ ,

$$W_{c,\varepsilon}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) + \varepsilon H(\pi | \alpha \otimes \beta), \quad (\mathcal{P}_\varepsilon)$$

where

$$H(\pi | \alpha \otimes \beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X} \times \mathcal{Y}} \log \left( \frac{d\pi(x, y)}{d\alpha(x)d\beta(y)} \right) d\pi(x, y).$$

relative entropy of the transport plan  $\pi$  with respect to the product measure  $\alpha \otimes \beta$ .

Distances

oooooooooooooooo

Entropic Regularization

○  
●●○  
○  
○○○○○

Sinkhorn Divergences

○  
○○○○○  
○○○○○○○

Conclusion

○○

The basics

## Entropic Regularization

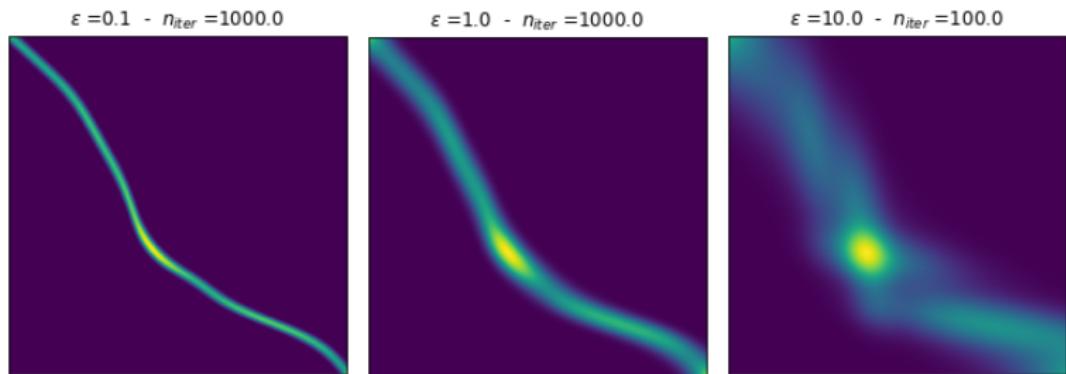


Figure 3 – Influence of the regularization parameter  $\varepsilon$  on the transport plan  $\pi$ .

Distances

oooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Sinkhorn Divergences

- 
- 
- 

Conclusion

○○

The basics

## Dual Formulation

Contrary to standard OT, no constraint on the dual problem :

$$W_c(\alpha, \beta) = \max_{\substack{u \in \mathcal{C}(\mathcal{X}) \\ v \in \mathcal{C}(\mathcal{Y})}} \int_{\mathcal{X}} u(x) d\alpha(x) + \int_{\mathcal{Y}} v(y) d\beta(y) \quad (\mathcal{D})$$

such that  $\{u(x) + v(y) \leq c(x, y) \forall (x, y) \in \mathcal{X} \times \mathcal{Y}\}$

Distances

oooooooooooooooo

Entropic Regularization

Sinkhorn Divergences

Conclusion

oo

The basics

## Dual Formulation

Contrary to standard OT, no constraint on the dual problem :

$$\begin{aligned}
 W_{c,\varepsilon}(\alpha, \beta) = & \max_{\substack{\mathbf{u} \in \mathcal{C}(\mathcal{X}) \\ \mathbf{v} \in \mathcal{C}(\mathcal{Y})}} \int_{\mathcal{X}} \mathbf{u}(x) d\alpha(x) + \int_{\mathcal{Y}} \mathbf{v}(y) d\beta(y) \\
 & - \varepsilon \int_{\mathcal{X} \times \mathcal{Y}} e^{\frac{\mathbf{u}(x) + \mathbf{v}(y) - c(x,y)}{\varepsilon}} d\alpha(x) d\beta(y) + \varepsilon.
 \end{aligned}$$

Distances

○○○○○○○○○○○○○○○○

Entropic Regularization

○  
○○○●  
○  
○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

The basics

## Sinkhorn's Algorithm

Iterative algorithm : alternate between optimizing over  $u$  with fixed  $v$  and optimizing over  $v$  with fixed  $u$ .

Distances

oooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Sinkhorn Divergences

- 
- 
- 

Conclusion

○○

The basics

## Sinkhorn's Algorithm

Iterative algorithm : alternate between optimizing over  $\mathbf{u}$  with fixed  $\mathbf{v}$  and optimizing over  $\mathbf{v}$  with fixed  $\mathbf{u}$ .

### Sinkhorn's Algorithm

Let  $K_{ij} = e^{-\frac{c(x_i, y_j)}{\varepsilon}}$ ,  $\mathbf{a} = e^{\frac{\mathbf{u}}{\varepsilon}}$ ,  $\mathbf{b} = e^{\frac{\mathbf{v}}{\varepsilon}}$ .

$$\mathbf{a}^{(\ell+1)} = \frac{1}{K(\mathbf{b}^{(\ell)} \odot \boldsymbol{\beta})} \quad ; \quad \mathbf{b}^{(\ell+1)} = \frac{1}{K^T(\mathbf{a}^{(\ell+1)} \odot \boldsymbol{\alpha})}$$

Complexity of each iteration :  $O(n^2)$ ,

Linear convergence, constant degrades when  $\varepsilon \rightarrow 0$ .

Distances

ooooooooooooooo

Entropic Regularization

○  
○○○○  
●  
○○○○○

Sinkhorn Divergences

○  
○○○○○  
○○○○○○○

Conclusion

○○

A magic regularizing tool

## Differentiable approximation of OT

**Bonus : Sinkhorn procedure is fully differentiable with auto-diff tools (e.g TensorFlow)  $\Rightarrow$  yields a differentiable approximation of OT !**

Some applications :

- Differentiable sorting (Cuturi et al '19)
- Differentiable (or 'soft') assignments
- Differentiable clustering (G. et al '19)
- Learning with a regularized Wasserstein loss  
( $\rightarrow$  more on that later...)

Distances

oooooooooooooooo

Entropic Regularization

Sample Complexity

Sinkhorn Divergences

Conclusion

oo

## The 'sample complexity'

### Informal Definition

*Given a distance between measures , its **sample complexity** corresponds to the error made when approximating this distance with samples of the measures.*

Known cases :

- OT :  $\mathbb{E}|W(\alpha, \beta) - W(\hat{\alpha}_n, \hat{\beta}_n)| = O(n^{-1/d})$   
 $\Rightarrow$  curse of dimension (Dudley '84, Weed and Bach '18)
- MMD :  $\mathbb{E}|MMD(\alpha, \beta) - MMD(\hat{\alpha}_n, \hat{\beta}_n)| = O(\frac{1}{\sqrt{n}})$   
 $\Rightarrow$  independent of dimension (Gretton '06)

*What about  $\mathbb{E}|W_\varepsilon(\alpha, \beta) - W_\varepsilon(\hat{\alpha}_n, \hat{\beta}_n)|$  ?*

Distances

oooooooooooooooo

Entropic Regularization

Sample Complexity

Sinkhorn Divergences

Conclusion

oo

## 'Sample Complexity' of $W_\varepsilon$ .

Theorem (G., Chizat, Bach, Cuturi, Peyré '19) (Mena+Weed '19)

Let  $\mathcal{X}, \mathcal{Y} \subset \mathbb{R}^d$  bounded , and  $c \in \mathcal{C}^\infty$ . Then

$$\mathbb{E}|W_\varepsilon(\alpha, \beta) - W_\varepsilon(\hat{\alpha}_n, \hat{\beta}_n)| = \begin{cases} O\left(\frac{1}{\varepsilon^{\lfloor d/2 \rfloor} \sqrt{n}}\right) & \text{when } \varepsilon \rightarrow 0, \\ O\left(\frac{1}{\sqrt{n}}\right) & \text{when } \varepsilon \rightarrow +\infty. \end{cases}$$

where constants depend on  $|\mathcal{X}|$ ,  $|\mathcal{Y}|$ ,  $d$ , and  $\|c^{(k)}\|_\infty$  for  $k = 0 \dots \lfloor d/2 \rfloor + 1$ .

→ A large enough regularization breaks the curse of dimension.

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○●○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

Sample Complexity

Gradient descent for  $W_\varepsilon$ ,  $\varepsilon = 1$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○●○○

Sample Complexity

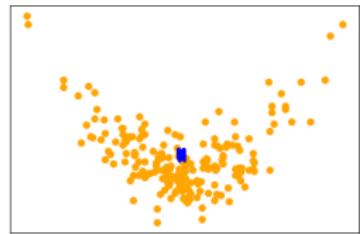
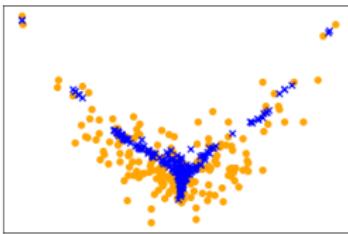
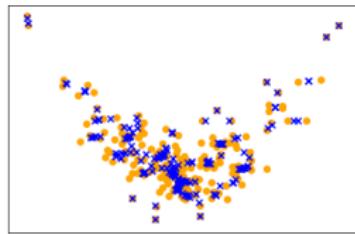
Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

## The effect of entropy



$$\varepsilon = 0.1$$

$$\varepsilon = 1$$

$$\varepsilon = 10$$

Distances

ooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○●○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○○

Sample Complexity

## The effect of entropy

Entropic Transport is Maximum Likelihood under Gaussian noise (Rigollet+Weed '18)

Consider a noisy sample  $(x_1, \dots, x_n) \sim \alpha_\theta + \mathcal{N}(0, \varepsilon)$ ,

$$\hat{\theta}^{MLE} = \operatorname{argmin}_\theta W_\varepsilon(\alpha_\theta, \frac{1}{n} \sum_{i=1}^n \delta_{x_i})$$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○●

Sample Complexity

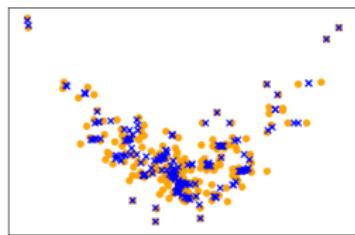
Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

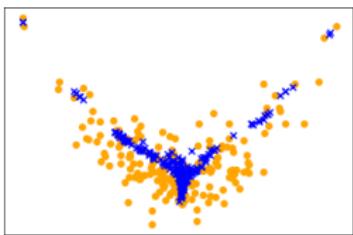
Conclusion

○○

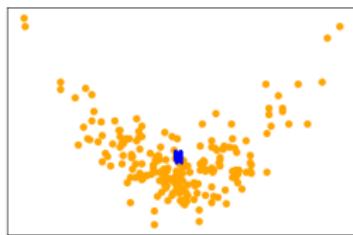
## The effect of entropy



$$\varepsilon = 0.1$$



$$\varepsilon = 1$$



$$\varepsilon = 10$$

Distances

ooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

●  
○○○○○○  
○○○○○○○○

Conclusion

○○

- ① Notions of Distance between Measures
- ② Entropic Regularization of Optimal Transport
- ③ Sinkhorn Divergences : Interpolation between OT and MMD
  - Definition and properties
  - Learning with Sinkhorn Divergences
- ④ Conclusion

Distances

oooooooooooooooo

Entropic Regularization

Sinkhorn Divergences

Conclusion

oo

Definition and properties

## Sinkhorn Divergences

**Issue of regularized Wass. Distance :**  $W_{c,\varepsilon}(\alpha, \alpha) \neq 0$

**Proposed Solution :** introduce corrective terms to ‘debias’  
regularized Wasserstein distance

### Definition (Sinkhorn Divergences)

Let  $\alpha \in \mathcal{M}_+^1(\mathcal{X})$  and  $\beta \in \mathcal{M}_+^1(\mathcal{Y})$ ,

$$SD_{c,\varepsilon}(\alpha, \beta) \stackrel{\text{def.}}{=} W_{c,\varepsilon}(\alpha, \beta) - \frac{1}{2} W_{c,\varepsilon}(\alpha, \alpha) - \frac{1}{2} W_{c,\varepsilon}(\beta, \beta),$$

Distances  
oooooooooooo

Entropic Regularization  
○  
○○○○  
○  
○○○○○

Sinkhorn Divergences  
○  
○●○○○○  
○○○○○○○○

Conclusion  
○○

Definition and properties

## Interpolation Property

Theorem (G., Peyré, Cuturi '18), (Ramdas and al. '17)

Sinkhorn Divergences have the following asymptotic behavior :

$$\text{when } \varepsilon \rightarrow 0, \quad SD_{c,\varepsilon}(\alpha, \beta) \rightarrow W_c(\alpha, \beta), \quad (1)$$

$$\text{when } \varepsilon \rightarrow +\infty, \quad SD_{c,\varepsilon}(\alpha, \beta) \rightarrow \frac{1}{2} MMD_{-c}^2(\alpha, \beta). \quad (2)$$

Remark : To get an MMD,  $-c$  must be positive definite. For  $c = \|\cdot\|_2^p$  with  $0 < p < 2$ , the MMD is called **Energy Distance**.

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○●○○○○  
○○○○○○○○

Conclusion

○○

Definition and properties

## Gradient descent for Energy Distance, $p = 1$

Distances

oooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○●○○○  
○○○○○○○○

Conclusion

○○

Definition and properties

## Gradient descent for Energy Distance, $p = 2$

Distances

○○○○○○○○○○○○○○○○

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○●○○  
○○○○○○○○

Conclusion

○○

Definition and properties

Gradient descent for  $SD_\varepsilon$ ,  $\varepsilon = 1$

Distances

○○○○○○○○○○○○○○○○

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○●○  
○○○○○○○○

Conclusion

○○

Definition and properties

Gradient descent for  $SD_\varepsilon$ ,  $\varepsilon = 1$

## Distances

○○○○○○○○○○○○○○

Entropic Regularization

1

## Sinkhorn Divergences

1

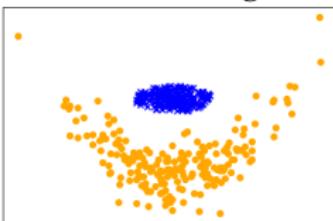
## Conclusion

○○

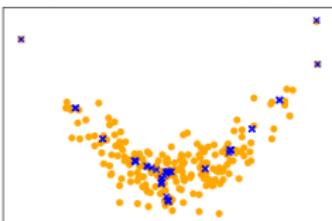
## Definition and properties

## Summary

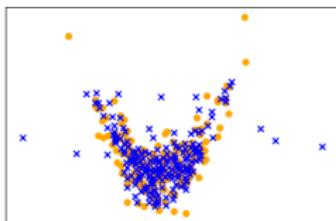
### *Initial Setting*



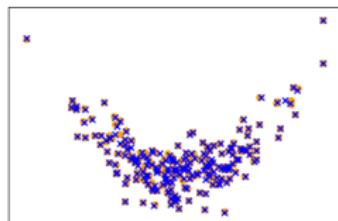
$$W_{c,\varepsilon} - \varepsilon = 1, c = \|\cdot\|_2^{1.5}$$



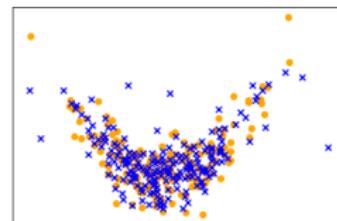
$$ED_p - p = 1.5$$



$$SD_{c,\varepsilon} - \varepsilon = 1, c = \|\cdot\|_2^{1.5}$$



$$SD_{c,\varepsilon} - \varepsilon = 10^2, c = \|\cdot\|_2^{1.5}$$



**Figure 4 – Goal :** Recover the positions of the Diracs with gradient descent. Orange circles : target distribution  $\beta$ , blue crosses : parametric model after convergence  $\alpha_{\theta^*}$ . Upper right : initial setting  $\alpha_{\theta_0}$ .

Distances

○○○○○○○○○○○○○○○○

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

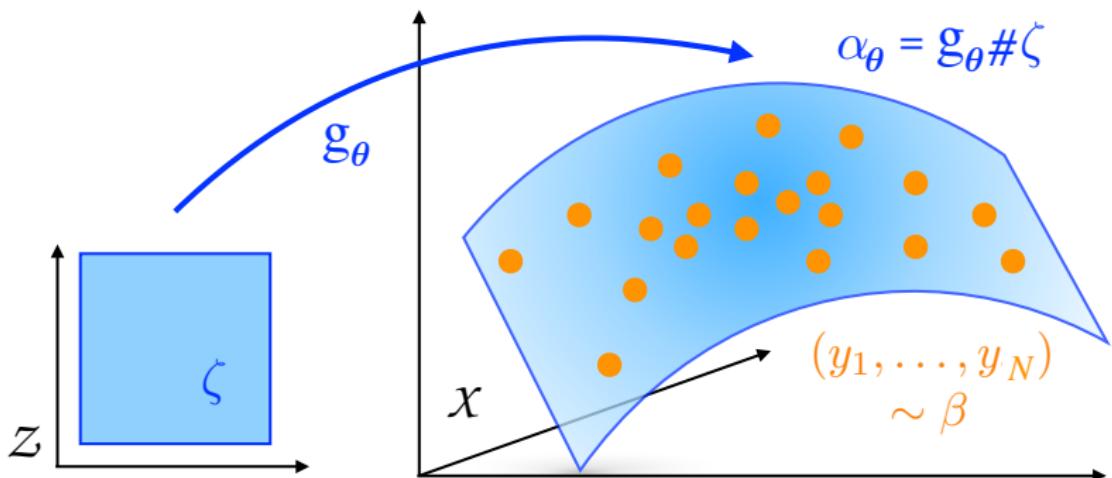
○  
○○○○○○  
●○○○○○○○○

Conclusion

○○

Learning

## Generative Models



Distances

ooooooooooooooo

Entropic Regularization

Learning

Sinkhorn Divergences

Conclusion

oo

## Problem Formulation

- $\beta$  the **unknown** measure of the data : finite number of samples  $(y_1, \dots, y_N) \sim \beta$
- $\alpha_\theta$  the parametric model of the form  $\alpha_\theta \stackrel{\text{def.}}{=} g_\theta \# \zeta$  : to sample  $x \sim \alpha_\theta$ , draw  $z \sim \zeta$  and take  $x = g_\theta(z)$ .

We are looking for the optimal parameter  $\theta^*$  defined by

$$\theta^* \in \operatorname{argmin}_\theta SD_{c,\varepsilon}(\alpha_\theta, \beta)$$

*NB :  $\alpha_\theta$  and  $\beta$  are only known via their samples.*

Distances

ooooooooooooooo

Entropic Regularization

Learning

Sinkhorn Divergences

Conclusion

oo

## The Optimization Procedure

We want to solve by gradient descent

$$\min_{\theta} SD_{c,\varepsilon}(\alpha_\theta, \beta)$$

At each descent step  $k$  :

- we approximate  $SD_{c,\varepsilon}(\alpha_{\theta(k)}, \beta)$  via
  - minibatches,
  - fixed number of Sinkhorn iterations
- we compute the gradient  $\nabla_{\theta} SD_{c,\varepsilon}^{(L)}(\hat{\alpha}_{\theta(k)}, \hat{\beta})$  by backpropagation (with automatic differentiation library)
- we do an update  $\theta^{(k+1)} = \theta^{(k)} - C_k \nabla_{\theta} SD_{c,\varepsilon}^{(L)}(\hat{\alpha}_{\theta(k)}, \hat{\beta})$

Distances

oooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Sinkhorn Divergences

- 
- 
- 

Conclusion

○○

Learning

## Computing the Gradient in Practice

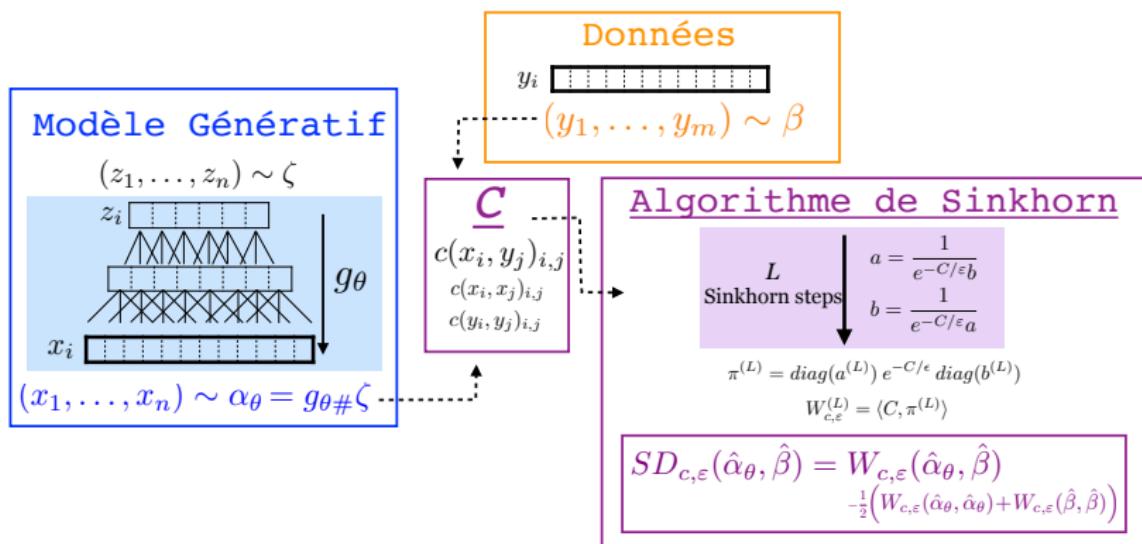


Figure 5 – Scheme of the approximation of the Sinkhorn Divergence from samples (here,  $g_\theta : z \mapsto x$  is represented as a 2-layer NN).

Distances

oooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Sinkhorn Divergences

- 
- 
- 

Conclusion

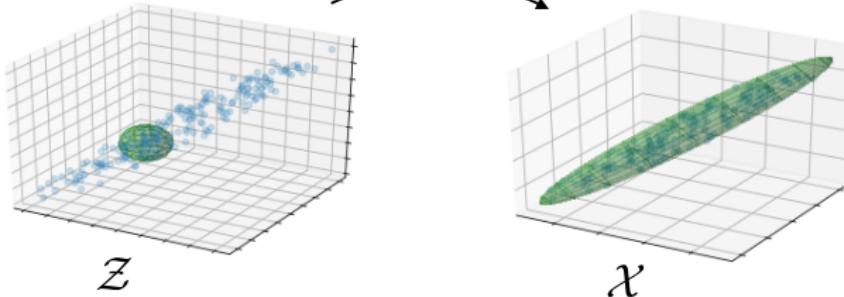
○○

Learning

## Toy Example : 3D ellipse

- Latent distribution : Uniform over unit ball in  $\mathbb{R}^3$
- Pushforward function :  $g_\theta(z) = Az + \omega$
- Parameters :  $\theta \stackrel{\text{def.}}{=} (A, \omega)$  where  $A$  covariance of the ellipse,  $\omega$  its center

$$x = Az + \omega$$



Distances

oooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Learning

Sinkhorn Divergences

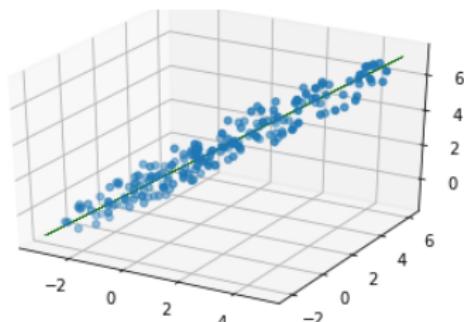
- 
- 
- 

Conclusion

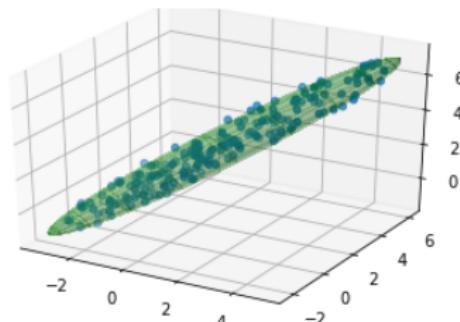
○○

## Empirical Results

$$W_{c,\varepsilon} - \varepsilon = 1, c = \|\cdot\|_2^2$$



$$SD_{c,\varepsilon} - \varepsilon = 1, c = \|\cdot\|_2^2$$



**Figure 6 –** Influence of the ‘debiasing’ of the Sinkhorn Divergence ( $SD_{\varepsilon}$ ) compared to regularized OT ( $W_{\varepsilon}$ ).

Distances

oooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Learning

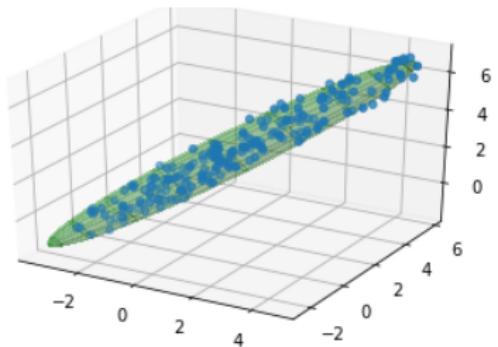
Sinkhorn Divergences

- 
- 
- 

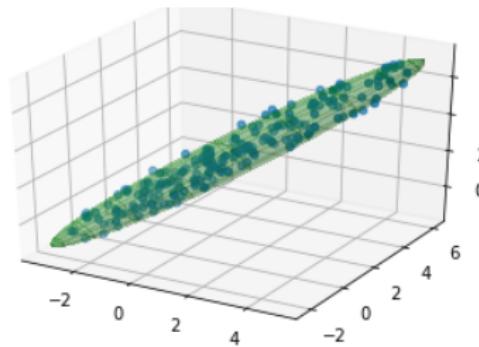
Conclusion

○○

$$ED_p - p = 1.5$$



$$SD_{c,\varepsilon} - \varepsilon = 1, c = \|\cdot\|_2^2$$



$ED_p$		
1.5,-		
3.12	1.74	2.08
2.25	2.83	2.09
2.30	1.74	3.07
( 0.63 , 1.75 , 2.75 )		

ground truth		
3	2	2
2	3	2
2	2	3
(1,2,3)		

$SD_{c,\varepsilon}$		
2, 1		
2.90	1.96	2.13
2.02	3.03	2.10
2.06	1.95	3.03
(0.94 , 1.96 , 2.90)		

Figure 7 – Comparison of the Sinkhorn Divergence ( $SD_{c,\varepsilon}$ ) and Energy Distance ( $ED_p$ ) on the ellipse fitting task.

Distances  
ooooooooooooooo

Entropic Regularization  
○  
○○○○  
○  
○○○○○

Sinkhorn Divergences  
○  
○○○○○○  
○○○○○○●○

Conclusion  
○○

Learning

## Learning the cost function

In high dimension (e.g. images), the Euclidean distance is not relevant → choosing the cost  $c$  is a complex problem.

**Idea** : the cost should yield high values for the Sinkhorn Divergence when  $\alpha_\theta \neq \beta$  to differentiate between synthetic samples (from  $\alpha_\theta$ ) and 'real' data (from  $\beta$ ). (Li and al '18)

We learn a parametric cost of the form :

$$c_\varphi(x, y) \stackrel{\text{def.}}{=} \|f_\varphi(x) - f_\varphi(y)\|^p \quad \text{where} \quad f_\varphi : \mathcal{X} \rightarrow \mathbb{R}^{d'},$$

The optimization problem becomes a min-max on  $(\theta, \varphi)$

$$\min_{\theta} \max_{\varphi} SD_{c_\varphi, \varepsilon}(\alpha_\theta, \beta)$$

→ GAN-type problem, cost  $c$  acts as a discriminator.

Distances

oooooooooooooooooooo

Entropic Regularization

- 
- 
- 
- 

Learning

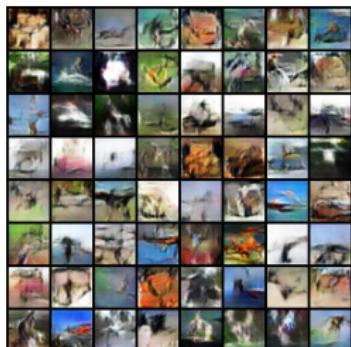
Sinkhorn Divergences

- 
- 
- 

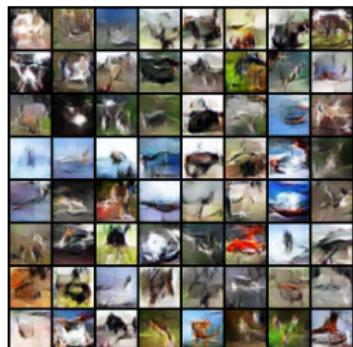
Conclusion

○○

## Empirical Results - CIFAR10



(a) MMD

(b)  $\varepsilon = 100$ (c)  $\varepsilon = 1$ 

MMD (Gaussian)

 $\varepsilon = 100$  $\varepsilon = 10$  $\varepsilon = 1$  $4.56 \pm 0.07$  $4.81 \pm 0.05$  $4.79 \pm 0.13$  $4.43 \pm 0.07$ 

**Table 1 – Inception Scores on CIFAR10 (same setting as MMD-GAN paper (Li et al. '18)).**

Distances

oooooooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

●○

- 1 Notions of Distance between Measures
- 2 Entropic Regularization of Optimal Transport
- 3 Sinkhorn Divergences : Interpolation between OT and MMD
- 4 Conclusion

Distances

ooooooooooooooo

Entropic Regularization

○  
○○○○  
○  
○○○○○

Sinkhorn Divergences

○  
○○○○○○  
○○○○○○○○

Conclusion

○●

## Take Home Message

Sinkhorn Divergences are a great notion of distance between measures !

- 'debias' regularized Wasserstein Distance
- interpolate between OT (small  $\varepsilon$ ) and MMD (large  $\varepsilon$ ) and get the best of both worlds :
  - inherit geometric properties from OT
  - break curse of dimension for  $\varepsilon$  large enough
- fast algorithms for implementation in ML tasks