

Robust shape matching with Optimal Transport

Jean Feydy

BIRS, Banff seminar 18w5151 – 13th December, 2018

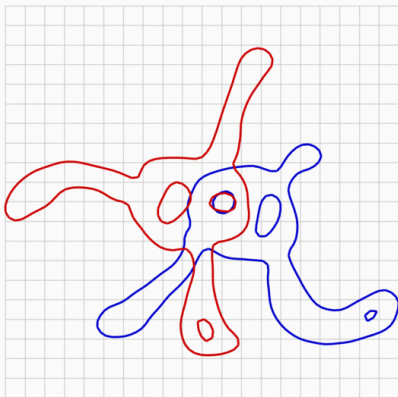
Écoles Normales Supérieures de Paris et Paris-Saclay

Collaboration with B. Charlier, J. Glaunès (KeOps library);

S.-i. Amari, G. Peyré, T. Séjourné, A. Trounev, F.-X. Vialard (OT theory)

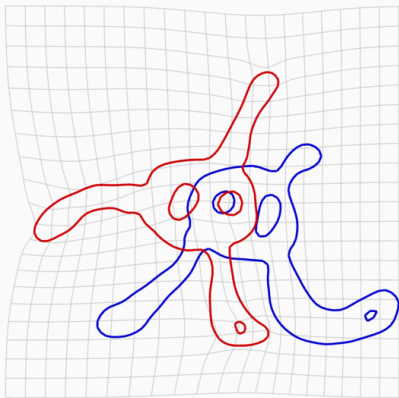
What is shape matching?

Source **A**, target **B**,



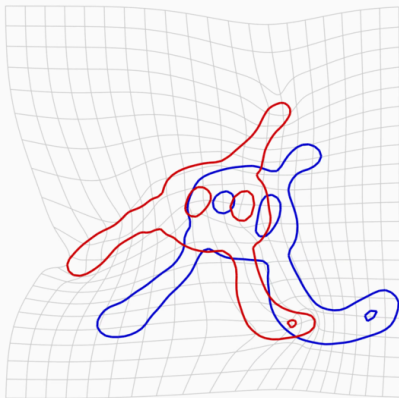
What is shape matching?

Source **A**, target **B**, mapping φ



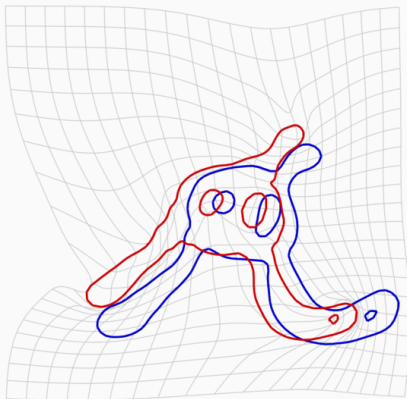
What is shape matching?

Source **A**, target **B**, mapping φ



What is shape matching?

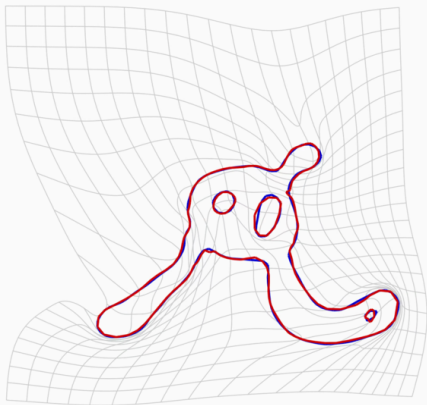
Source A , target B , mapping φ



What is shape matching?

Source A , target B , mapping φ

$$A \xrightarrow[\text{Model}]{\varphi} \varphi(A) = A' \xleftrightarrow[\text{Loss}]{\leftarrow} B$$



A good Loss function is a guarantee of robustness

Iterative Matching Algorithm

- 1: $A' \leftarrow A$
 - 2: **repeat**
 - 3: $L, v \leftarrow \text{Loss}(A', B), -\partial_{A'} \text{Loss}(A', B)$
 - 4: $A' \leftarrow A' + \text{Model}(v)$
 - 5: **until** $L < \text{tol}$
 Output: deformed shape $A' = \varphi(A)$.
-

A good Loss function is a guarantee of robustness

Iterative Matching Algorithm

- 1: $A' \leftarrow A$
 - 2: **repeat**
 - 3: $L, \nu \leftarrow \text{Loss}(A', B), -\partial_{A'} \text{Loss}(A', B)$
 - 4: $A' \leftarrow A' + \text{Model}(\nu)$
 - 5: **until** $L < \text{tol}$
- Output:** deformed shape $A' = \varphi(A)$.
-

“Model” encodes the **prior knowledge** on admissible deformations:

- *smoothing* convolution
- LDDMM/SVF *backprop* + regularization + *shooting*
- *trained* neural network

A good Loss function is a guarantee of robustness

Iterative Matching Algorithm

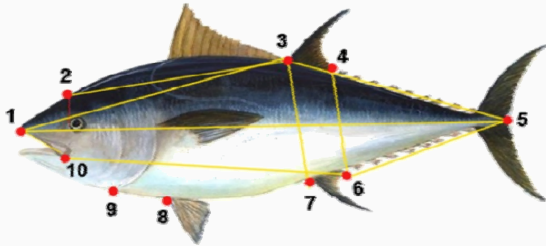
- 1: $A' \leftarrow A$
 - 2: **repeat**
 - 3: $L, \nu \leftarrow \text{Loss}(A', B), -\partial_{A'} \text{Loss}(A', B)$
 - 4: $A' \leftarrow A' + \text{Model}(\nu)$
 - 5: **until** $L < \text{tol}$
- Output:** deformed shape $A' = \varphi(A)$.
-

“Model” encodes the **prior knowledge** on admissible deformations:

- *smoothing* convolution
- LDDMM/SVF *backprop* + regularization + *shooting*
- *trained* neural network

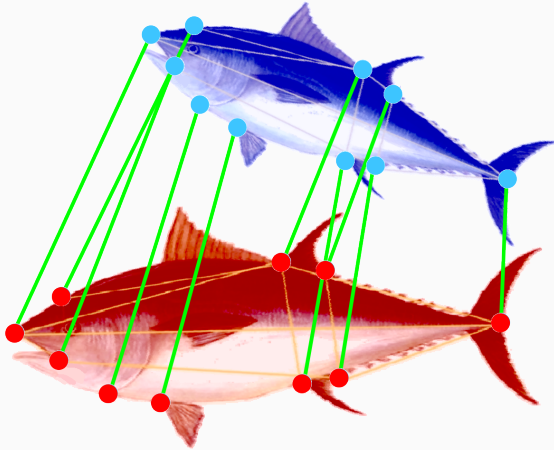
\Rightarrow The *raw* Loss gradient ν is what **drives** the registration

On labeled shapes, use a spring energy



Anatomical landmarks from *A morphometric approach for the analysis of body shape in bluefin tuna*, Addis et al., 2009.

On labeled shapes, use a spring energy



Anatomical landmarks from *A morphometric approach for the analysis of body shape in bluefin tuna*, Addis et al., 2009.

Encoding unlabeled shapes as measures

Let's enforce sampling invariance:

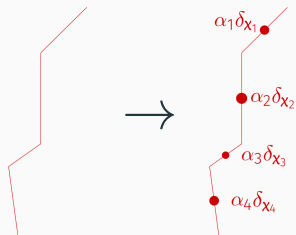
$$A \longrightarrow \alpha = \sum_{i=1}^N \alpha_i \delta_{x_i}, \quad B \longrightarrow \beta = \sum_{j=1}^M \beta_j \delta_{y_j}.$$

Encoding unlabeled shapes as measures

Let's enforce sampling invariance:

$$A \longrightarrow \alpha = \sum_{i=1}^N \alpha_i \delta_{x_i},$$

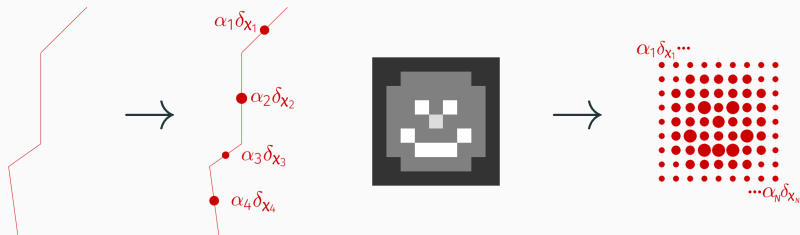
$$B \longrightarrow \beta = \sum_{j=1}^M \beta_j \delta_{y_j}.$$



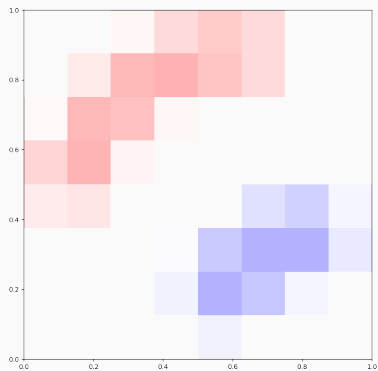
Encoding unlabeled shapes as measures

Let's enforce sampling invariance:

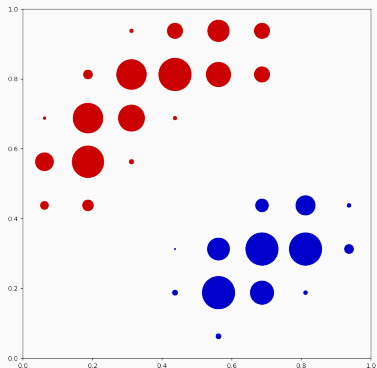
$$A \longrightarrow \alpha = \sum_{i=1}^N \alpha_i \delta_{x_i}, \quad B \longrightarrow \beta = \sum_{j=1}^M \beta_j \delta_{y_j}.$$



A baseline setting: density registration

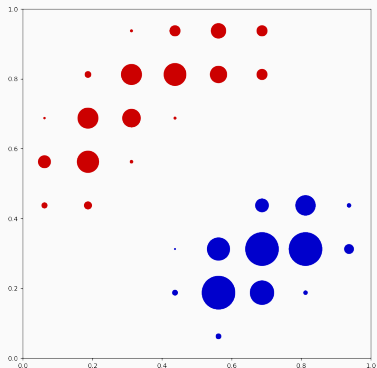


A baseline setting: density registration



$$\alpha = \sum_{i=1}^N \alpha_i \delta_{x_i}, \quad \beta = \sum_{j=1}^M \beta_j \delta_{y_j}.$$

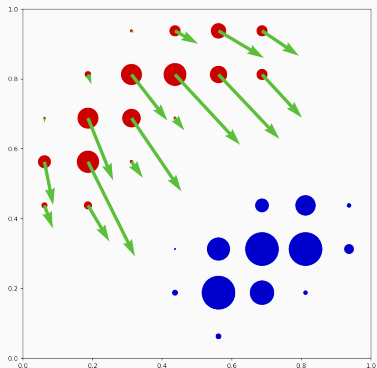
A baseline setting: density registration



$$\alpha = \sum_{i=1}^N \alpha_i \delta_{x_i}, \quad \beta = \sum_{j=1}^M \beta_j \delta_{y_j}.$$

$$\sum_{i=1}^N \alpha_i = 1 = \sum_{j=1}^M \beta_j$$

A baseline setting: density registration

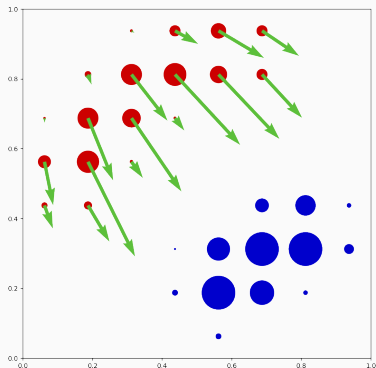


$$\alpha = \sum_{i=1}^N \alpha_i \delta_{x_i}, \quad \beta = \sum_{j=1}^M \beta_j \delta_{y_j}.$$

$$\sum_{i=1}^N \alpha_i = 1 = \sum_{j=1}^M \beta_j$$

Display $v = -\nabla_{x_i} d(\alpha, \beta)$.

A baseline setting: density registration



$$\alpha = \sum_{i=1}^N \alpha_i \delta_{x_i}, \quad \beta = \sum_{j=1}^M \beta_j \delta_{y_j}.$$

$$\sum_{i=1}^N \alpha_i = 1 = \sum_{j=1}^M \beta_j$$

Display $v = -\nabla_{x_i} d(\alpha, \beta)$.

Seamless extensions to:

- $\sum_i \alpha_i \neq \sum_j \beta_j$, outliers [Chizat et al., 2018],
- curves and surfaces [Kaltenmark et al., 2017],
- variable weights α_i .

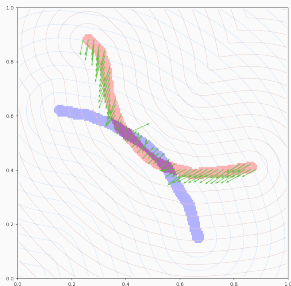
1. Computing fidelities between **measures**

1. Computing fidelities between **measures**
2. What's **new**, in 2018?

1. Computing fidelities between **measures**
2. What's **new**, in 2018?
3. Efficient GPU routines: **KeOps**

A simple formula: Hausdorff distance (aka. ICP, \simeq GMM-MLE)

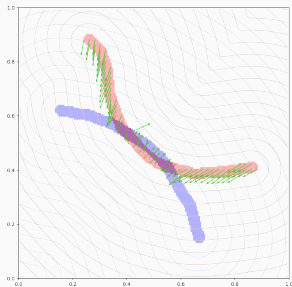
Define the fields $a(x) = d(x, \text{supp}(\alpha)) = \min_i \|x_i - x\|$,
 $b(x) = d(x, \text{supp}(\beta)) = \min_j \|x - y_j\|$,



A simple formula: Hausdorff distance (aka. ICP, \simeq GMM-MLE)

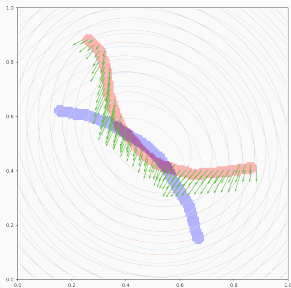
Define the fields $a(x) = d(x, \text{supp}(\alpha)) = \min_i \|x_i - x\|$,
 $b(x) = d(x, \text{supp}(\beta)) = \min_j \|x - y_j\|$,

$$\begin{aligned} \text{Loss}(\alpha, \beta) &= \frac{1}{2} \langle \alpha, b \rangle + \frac{1}{2} \langle \beta, a \rangle \\ &= \frac{1}{2} \sum_i \alpha_i \cdot \min_j \|x_i - y_j\| + \frac{1}{2} \sum_j \beta_j \cdot \min_i \|x_i - y_j\| \end{aligned}$$



A simple formula: Kernel norms (aka. MMD)

Define the fields

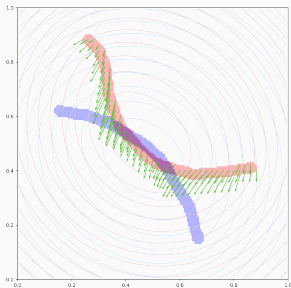
$$a(x) = \sum_i \alpha_i \|x - x_i\| = (\|\cdot\| \star \alpha)(x),$$
$$b(x) = \sum_j \beta_j \|x - y_j\| = (\|\cdot\| \star \beta)(x),$$


A simple formula: Kernel norms (aka. MMD)

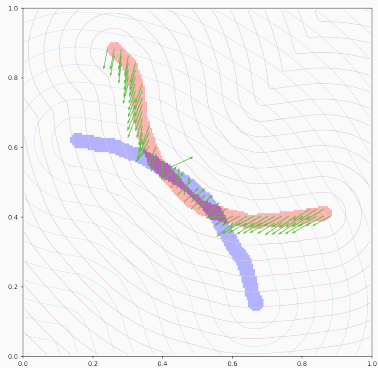
Define the fields $a(x) = \sum_i \alpha_i \|x - x_i\| = (\|\cdot\| \star \alpha)(x)$,

$$b(x) = \sum_j \beta_j \|x - y_j\| = (\|\cdot\| \star \beta)(x),$$

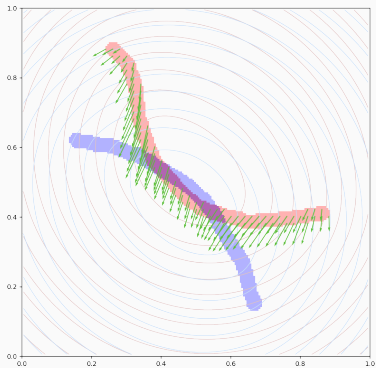
$$\begin{aligned} \text{Loss}(\alpha, \beta) &= \frac{1}{2} \langle \alpha - \beta, b - a \rangle = \sum_i \sum_j \alpha_i \beta_j \|x_i - y_j\| \\ &\quad - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \|x_i - x_j\| - \frac{1}{2} \sum_i \sum_j \beta_i \beta_j \|y_i - y_j\| \end{aligned}$$



The Hausdorff distance is local, the Energy Distance is global

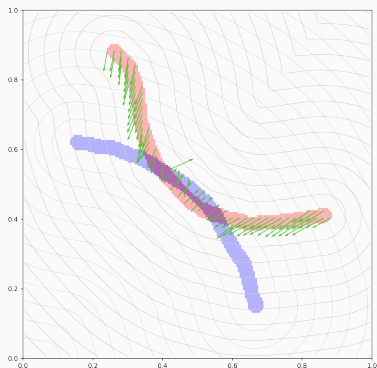


Hausdorff, min

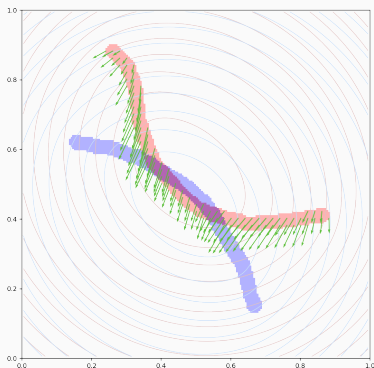


Kernel, Σ

The Hausdorff distance is local, the Energy Distance is global



Hausdorff, min

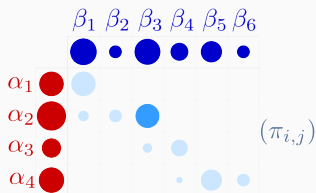
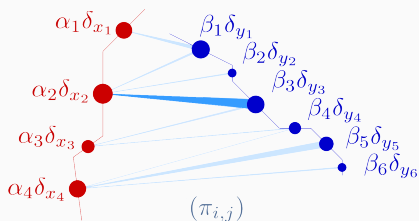


Kernel, Σ

\implies Can we get the best of both worlds?

Computational Optimal Transport

The Optimal Transport problem



Minimize over N -by- M matrices
(transport plans) π :

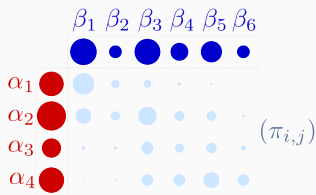
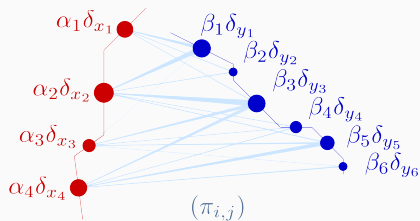
$$\text{OT}(\alpha, \beta) = \min_{\pi} \underbrace{\sum_{i,j} \pi_{i,j} \cdot |x_i - y_j|^2}_{\text{transport cost}}$$

subject to $\pi_{i,j} \geq 0$,

$$\sum_j \pi_{i,j} = \alpha_i, \quad \sum_i \pi_{i,j} = \beta_j.$$

\Rightarrow Hungarian method in $O(N^3)$.

Entropic regularization = add temperature, blur the transport plan



For $\varepsilon > 0$:

$$\text{OT}_\varepsilon(\alpha, \beta) = \min_{\pi} \underbrace{\sum_{i,j} \pi_{i,j} \cdot |\mathbf{x}_i - \mathbf{y}_j|^2}_{\text{transport cost}} + \varepsilon \underbrace{\sum_{i,j} \pi_{i,j} \cdot \log \frac{\pi_{i,j}}{\alpha_i \beta_j}}_{\text{entropic barrier}}$$

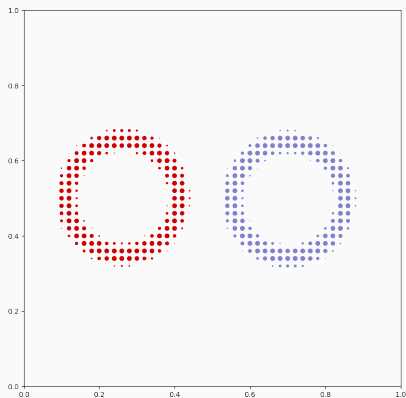
subject to

$$\sum_j \pi_{i,j} = \alpha_i, \quad \sum_i \pi_{i,j} = \beta_j.$$

⇒ Sinkhorn algorithm (GPU).

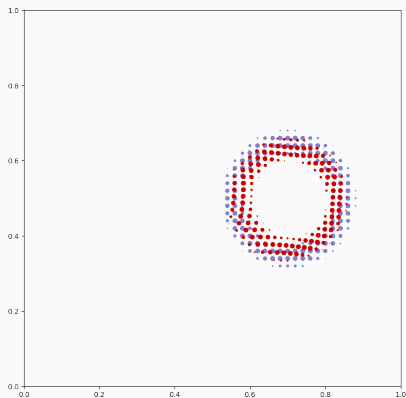
Problem : if $\varepsilon > 0$, OT_ε is not a valid divergence

Registrating circles, $C(x,y) = \|x - y\|^2$, $\sqrt{\varepsilon} = 0.1$:



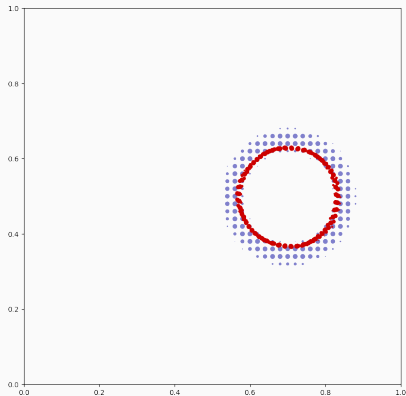
Problem : if $\varepsilon > 0$, OT_ε is not a valid divergence

Registrating circles, $C(x,y) = \|x - y\|^2$, $\sqrt{\varepsilon} = 0.1$:



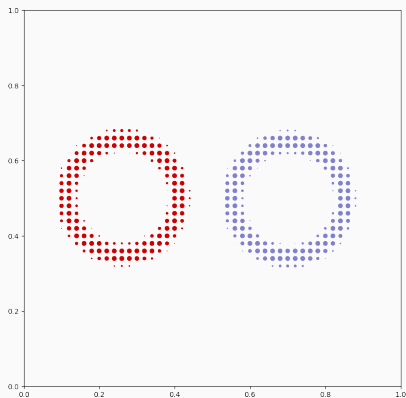
Problem : if $\varepsilon > 0$, OT_ε is not a valid divergence

Registrating circles, $C(x,y) = \|x - y\|^2$, $\sqrt{\varepsilon} = 0.1$:



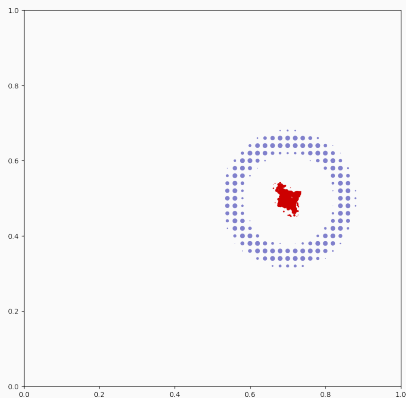
Problem : if $\varepsilon > 0$, OT_ε is not a valid divergence

Registrating circles, $C(x,y) = \|x - y\|^2$, $\sqrt{\varepsilon} = 0.2$:



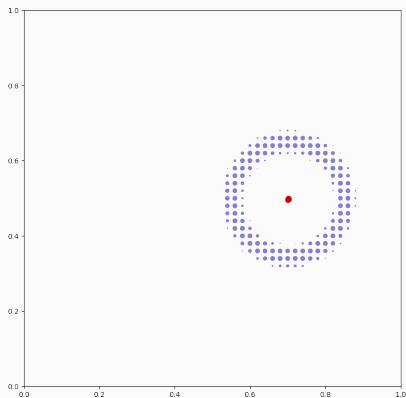
Problem : if $\varepsilon > 0$, OT_ε is not a valid divergence

Registrating circles, $C(x,y) = \|x - y\|^2$, $\sqrt{\varepsilon} = 0.2$:



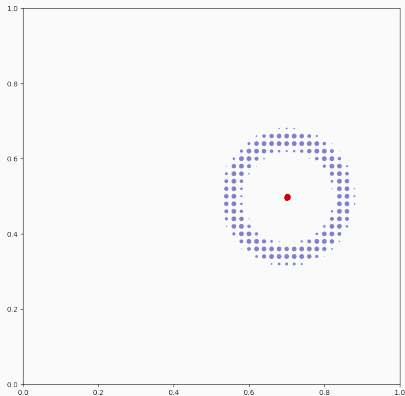
Problem : if $\varepsilon > 0$, OT_ε is not a valid divergence

Registrating circles, $C(x,y) = \|x - y\|^2$, $\sqrt{\varepsilon} = 0.2$:



Problem : if $\varepsilon > 0$, OT_ε is not a valid divergence

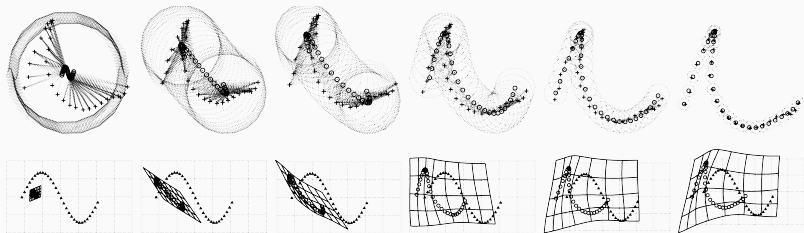
Registrating circles, $C(x,y) = \|x - y\|^2$, $\sqrt{\varepsilon} = 0.2$:



Bad news: for $0 < \varepsilon \leq +\infty$, we converge towards α such that

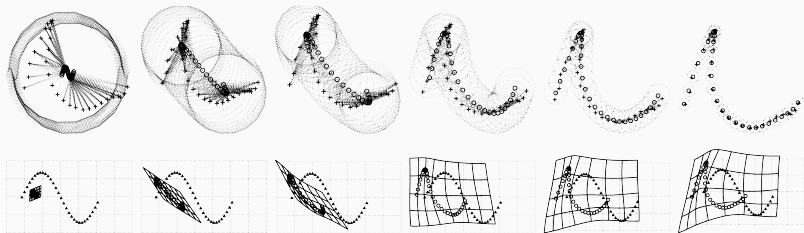
$$OT_\varepsilon(\alpha, \beta) < OT_\varepsilon(\beta, \beta).$$

Standard solution: use an annealing scheme



TPS-RPM algorithm, Chui and Rangarajan, CVPR 2000

Standard solution: use an annealing scheme



TPS-RPM algorithm, Chui and Rangarajan, CVPR 2000

⇒ **Expensive** and cumbersome workaround,
with parameters to tune.

A new idea in 2017 : un-biased Sinkhorn divergences

$$\text{OT}_\varepsilon(\alpha, \beta) = \min_{\pi} \langle \pi, \mathbf{C} \rangle + \varepsilon \text{KL}(\pi, \alpha \otimes \beta) \longrightarrow \text{Fuzzy assignment}$$

s.t. $\pi \mathbf{1} = \alpha, \quad \pi^T \mathbf{1} = \beta$

A new idea in 2017 : un-biased Sinkhorn divergences

$$\text{OT}_\varepsilon(\alpha, \beta) = \min_{\pi} \langle \pi, \mathbf{C} \rangle + \varepsilon \text{KL}(\pi, \alpha \otimes \beta) \longrightarrow \text{Fuzzy assignment}$$

s.t. $\pi \mathbf{1} = \alpha, \quad \pi^T \mathbf{1} = \beta$

$$\text{OT}_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \rightarrow +\infty} \langle \alpha \otimes \beta, \mathbf{C} \rangle = \langle \alpha, \mathbf{C} \star \beta \rangle$$

A new idea in 2017 : un-biased Sinkhorn divergences

$$\text{OT}_\varepsilon(\alpha, \beta) = \min_{\pi} \langle \pi, \mathbf{C} \rangle + \varepsilon \text{KL}(\pi, \alpha \otimes \beta) \longrightarrow \text{Fuzzy assignment}$$

s.t. $\pi \mathbf{1} = \alpha, \quad \pi^T \mathbf{1} = \beta$

$$\text{OT}_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \rightarrow +\infty} \langle \alpha \otimes \beta, \mathbf{C} \rangle = \langle \alpha, \mathbf{C} \star \beta \rangle$$

Define the **Sinkhorn divergence** [Raudas et al., 2017]:

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

A new idea in 2017 : un-biased Sinkhorn divergences

$$\text{OT}_\varepsilon(\alpha, \beta) = \min_{\pi} \langle \pi, C \rangle + \varepsilon \text{KL}(\pi, \alpha \otimes \beta) \longrightarrow \text{Fuzzy assignment}$$

s.t. $\pi \mathbf{1} = \alpha, \quad \pi^T \mathbf{1} = \beta$

$$\text{OT}_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \rightarrow +\infty} \langle \alpha \otimes \beta, C \rangle = \langle \alpha, C \star \beta \rangle$$

Define the **Sinkhorn divergence** [Raudas et al., 2017]:

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

$$\text{Wasserstein}_{+C}(\alpha, \beta) \xleftarrow{\varepsilon \rightarrow 0} S_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \rightarrow +\infty} \text{Kernel}_{-C}(\alpha, \beta)$$

A new idea in 2017 : un-biased Sinkhorn divergences

$$\text{OT}_\varepsilon(\alpha, \beta) = \min_{\pi} \langle \pi, \mathbf{C} \rangle + \varepsilon \text{KL}(\pi, \alpha \otimes \beta) \longrightarrow \text{Fuzzy assignment}$$

s.t. $\pi \mathbf{1} = \alpha, \quad \pi^T \mathbf{1} = \beta$

$$\text{OT}_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \rightarrow +\infty} \langle \alpha \otimes \beta, \mathbf{C} \rangle = \langle \alpha, \mathbf{C} \star \beta \rangle$$

Define the **Sinkhorn divergence** [Ramdas et al., 2017]:

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

$$\text{Wasserstein}_{+\mathbf{C}}(\alpha, \beta) \xleftarrow{\varepsilon \rightarrow 0} S_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \rightarrow +\infty} \text{Kernel}_{-\mathbf{C}}(\alpha, \beta)$$

In practice, S_ε is “good enough” for ML applications

[Genevay et al., 2018, Salimans et al., 2018, Sanjabi et al., 2018].

Theorem (F., Séjourné, Vialard, Amari, Trouvé, Peyré; 2018)

For all probability measures α, β and regularization $\varepsilon > 0$:

Theorem (F., Séjourné, Vialard, Amari, Trouvé, Peyré; 2018)

For all probability measures α, β and regularization $\varepsilon > 0$:

$$0 \leq S_\varepsilon(\alpha, \beta) \quad \text{with equality iff. } \alpha = \beta$$

Theorem (F., Séjourné, Vialard, Amari, Trouvé, Peyré; 2018)

For all probability measures α, β and regularization $\varepsilon > 0$:

$$0 \leq S_\varepsilon(\alpha, \beta) \quad \text{with equality iff. } \alpha = \beta$$

$\alpha \mapsto S_\varepsilon(\alpha, \beta)$ is convex, differentiable and metrizes $\alpha \rightarrow \beta$

In our paper(s): theoretical guarantees

Theorem (F., Séjourné, Vialard, Amari, Trouvé, Peyré; 2018)

For all probability measures α, β and regularization $\varepsilon > 0$:

$$0 \leq S_\varepsilon(\alpha, \beta) \quad \text{with equality iff. } \alpha = \beta$$

$\alpha \mapsto S_\varepsilon(\alpha, \beta)$ is convex, differentiable and metrizes $\alpha \rightarrow \beta$

These results can be generalized to arbitrary **feature** spaces

– e.g. (position, orientation, curvature).

In our paper(s): theoretical guarantees

Theorem (F., Séjourné, Vialard, Amari, Trouvé, Peyré; 2018)

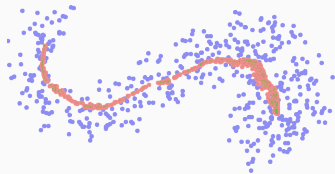
For all probability measures α, β and regularization $\varepsilon > 0$:

$$0 \leq S_\varepsilon(\alpha, \beta) \quad \text{with equality iff. } \alpha = \beta$$

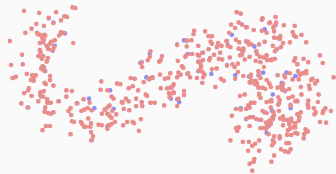
$\alpha \mapsto S_\varepsilon(\alpha, \beta)$ is convex, differentiable and metrizes $\alpha \rightarrow \beta$

These results can be generalized to arbitrary **feature** spaces

– e.g. (position, orientation, curvature).



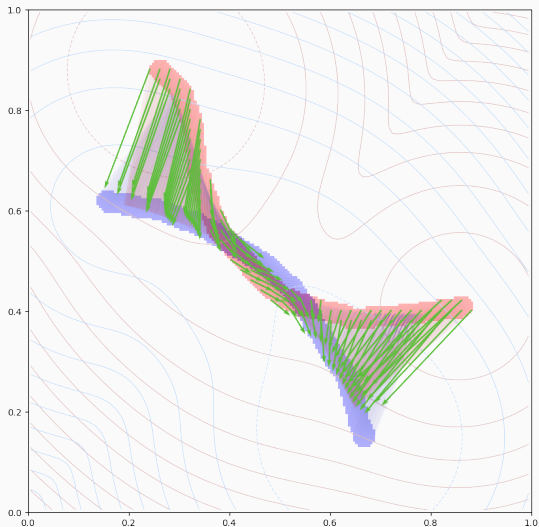
Loss = OT_ε



Loss = S_ε

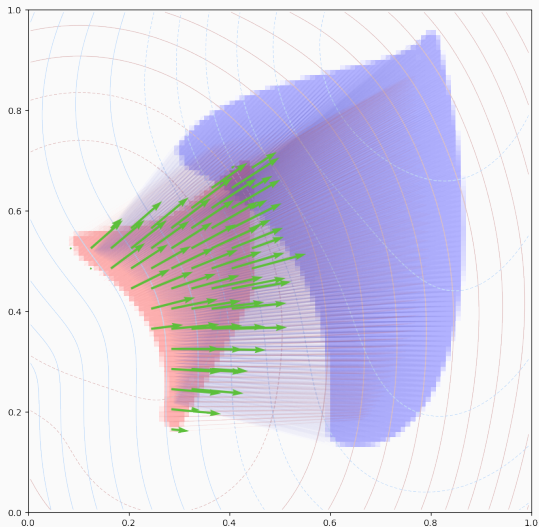
In practice

The ε -Sinkhorn divergence; with $\|x - y\|^2$ and $\sqrt{\varepsilon} = .1$



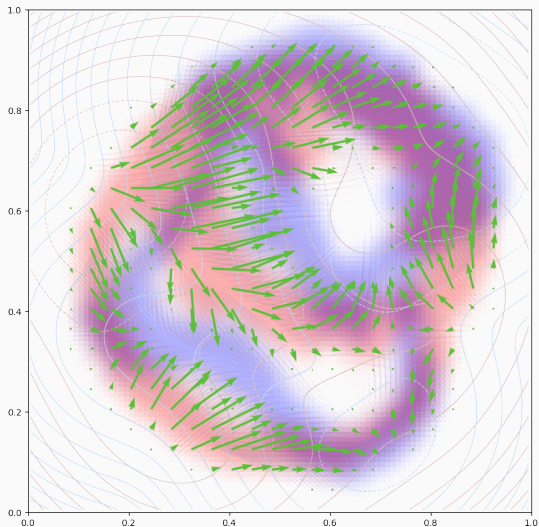
A high-quality gradient.

The ε -Sinkhorn divergence; with $\|x - y\|^2$ and $\sqrt{\varepsilon} = .1$



A high-quality gradient.

The ε -Sinkhorn divergence; with $\|x - y\|^2$ and $\sqrt{\varepsilon} = .1$



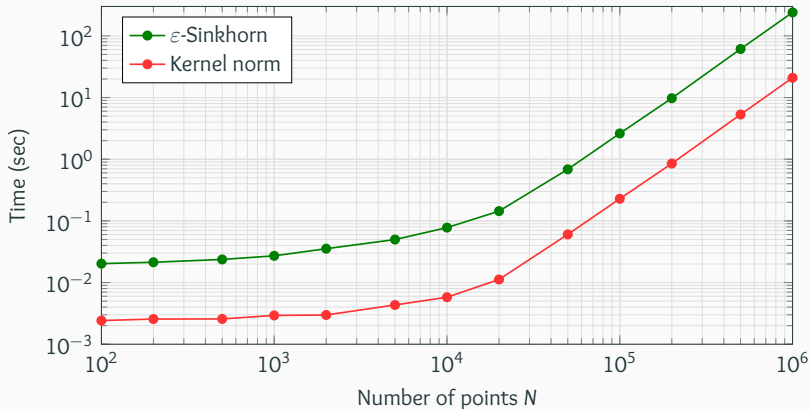
A high-quality gradient?

(Data from the Spectral Log-Demons paper.)

Kernel OPERATIONS, with autodiff, without memory overflows

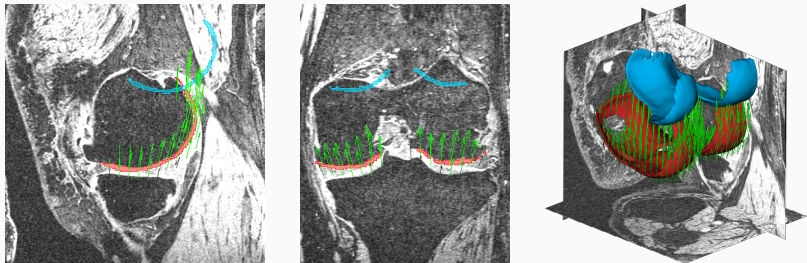
⇒ pip install pykeops ⇐
(Thanks Benjamin and Joan!)

Fidelity + gradient with N vertices on a **high-end GPU (Tesla P100)**



We provide a reference PyTorch implementation

github.com/jeanfeydy/global-divergences.



Gradient of the Energy Distance, computed in 0.5s on my laptop.

Data from the OsteoArthritis Initiative:

52,319 and 34,966 voxels out of a 192-192-160 volume.

Global, **geometry-aware** loss functions are easy to compute.

Global, **geometry-aware** loss functions are easy to compute.

- Try using $k(x,y) = -\|x - y\|$!

Global, **geometry-aware** loss functions are easy to compute.

- Try using $k(x,y) = -\|x - y\|$!
- Remove the **entropic bias** from the SoftAssign algorithm!

Global, **geometry-aware** loss functions are easy to compute.

- Try using $k(x,y) = -\|x - y\|$!
- Remove the **entropic bias** from the SoftAssign algorithm!
- Sinkhorn = Hausdorff + mass **spreading** constraint
 - \simeq best you can do without topology or landmarks
 - \simeq 20-50 convolutions through the data
 - Is it worth it?

Open questions:

- Rigorous link with the **auction algorithm**?

Open questions:

- Rigorous link with the **auction algorithm**?
- Link between S_ε and Sobolev **distances**?

Open questions:

- Rigorous link with the **auction algorithm**?
- Link between S_ε and Sobolev **distances**?
- What about **multiscale** schemes?

Open questions:

- Rigorous link with the **auction algorithm**?
- Link between S_ε and Sobolev **distances**?
- What about **multiscale** schemes?
- Interest in the **CVPR/SIGGRAPH** communities?

Thank you for your attention.

Any questions ?

Our papers:



- *Global divergences between measures: from Hausdorff distance to Optimal Transport*, F., Trouvé, 2018



Our papers:




- *Global divergences between measures: from Hausdorff distance to Optimal Transport*, F., Trouvé, 2018
- *Sinkhorn entropies and divergences*, F., Séjourné, Vialard, Amari, Trouvé, Peyré, 2018

Our papers:

- *Global divergences between measures: from Hausdorff distance to Optimal Transport*, F., Trouvé, 2018
- *Sinkhorn entropies and divergences*, F., Séjourné, Vialard, Amari, Trouvé, Peyré, 2018
- *Optimal Transport for diffeomorphic registration*, F., Charlier, Vialard, Peyré, 2017

-  Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. (2018).
Unbalanced optimal transport: Dynamic and kantorovich formulations.
Journal of Functional Analysis, 274(11):3090–3123.
-  Cuturi, M. (2013).
Sinkhorn distances: Lightspeed computation of optimal transport.
In *Advances in neural information processing systems*, pages 2292–2300.

-  Genevay, A., Peyre, G., and Cuturi, M. (2018).
Learning generative models with sinkhorn divergences.
In Storkey, A. and Perez-Cruz, F., editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1608–1617. PMLR.
-  Kaltenmark, I., Charlier, B., and Charon, N. (2017).
A general framework for curve and surface comparison and registration with oriented varifolds.
In *Computer Vision and Pattern Recognition (CVPR)*.

-  Peyré, G. and Cuturi, M. (2018).
Computational optimal transport.
arXiv preprint arXiv:1803.00567.
-  Ramdas, A., Trillos, N. G., and Cuturi, M. (2017).
On wasserstein two-sample testing and related families of nonparametric tests.
Entropy, 19(2).
-  Salimans, T., Zhang, H., Radford, A., and Metaxas, D. (2018).
Improving GANs using optimal transport.
arXiv preprint arXiv:1803.05573.



Sanjabi, M., Ba, J., Razaviyayn, M., and Lee, J. D. (2018).
**On the convergence and robustness of training GANs with
regularized optimal transport.**
arXiv preprint arXiv:1802.08249.