# Robust shape matching
# with Optimal Transport

Jean Feydy
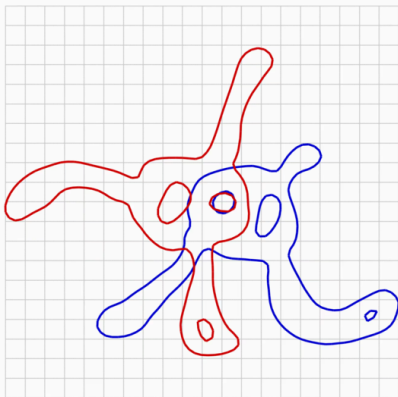
GTTI, ENS Cachan — 13th February, 2019

Écoles Normales Supérieures de Paris et Paris-Saclay
Collaboration with B. Charlier, J. Glaunès (KeOps library);
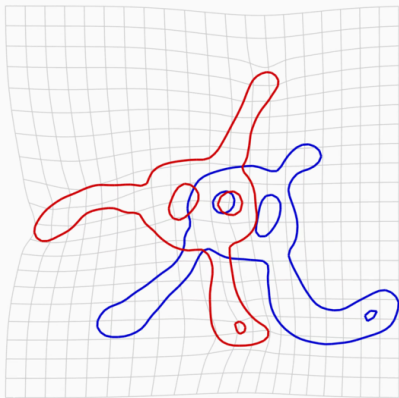S.-i. Amari, G. Peyré, T. Séjourné, A. Trouvé, F.-X. Vialard (OT theory)
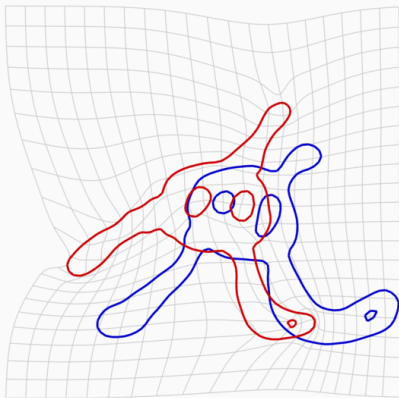
Source *A*, target *B*,

Source $A$, target $B$, mapping $\varphi$

Source *A*, target *B*, mapping $\varphi$

Source *A*, target *B*, mapping $\varphi$

Source $A$, target $B$, mapping $\varphi$

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

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

"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
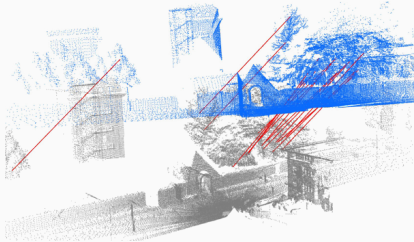
<div align="center">

**Iterative** Matching Algorithm

</div>

1: $A' \leftarrow A$

2: **repeat**

3:   $L \,, \, v \; \leftarrow \; \mathrm{Loss}(A', B) \,, \; -\partial_{A'}\mathrm{Loss}(A', B)$

4:   $A' \; \leftarrow \; A' \, + \, \mathrm{Model}(v)$

5: **until**   $L < \mathrm{tol}$

  **Output :** deformed shape   $A' = \varphi(A)$.

"Model" encodes the **prior knowledge** on admissible deformations:

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

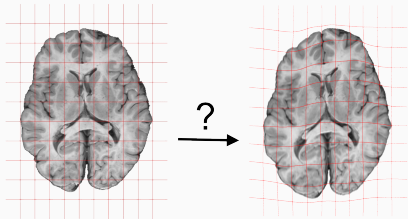   $\Rightarrow$   The *raw* Loss gradient $v$ is what **drives** the registration

From the documentation of the
Point Cloud Library.

- $\varphi$ is **rigid** or affine
- Occlusions
- Outliers

From Marc Niethammer's
Quicksilver slides.

- $\varphi$ is a spline or a **diffeomorphism**
- Ill-posed problem
- Some occlusions

*Wasserstein Auto-Encoders,*
Tolstikhin et al., 2018.

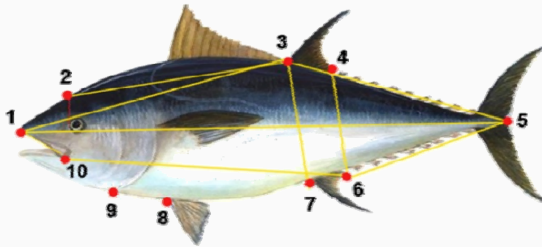- $\varphi$ is a **neural network**
- Very weak regularization
- High-dimensional space

*Wasserstein Auto-Encoders,*
Tolstikhin et al., 2018.

- $\varphi$ is a **neural network**
- Very weak regularization
- High-dimensional space
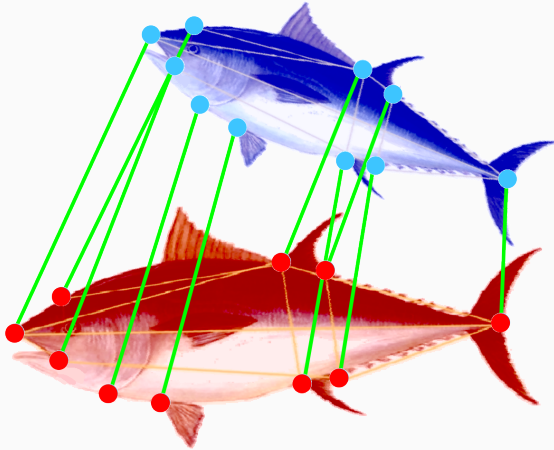
Which **Loss** function
should we use?

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

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

Let's enforce sampling invariance:

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

Let's enforce sampling invariance:

$$A \; \longrightarrow \; \alpha \; = \; \sum_{i=1}^{N} \alpha_i \delta_{x_i} \, , \qquad B \; \longrightarrow \; \beta \; = \; \sum_{j=1}^{M} \beta_j \delta_{y_j} \, .$$
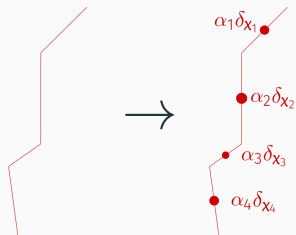
Let's enforce sampling invariance:

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

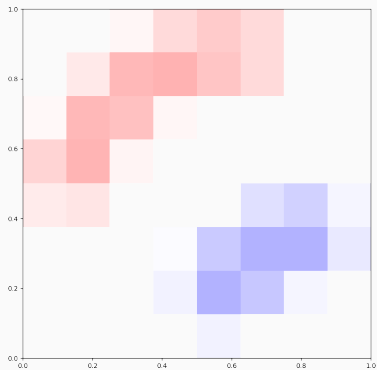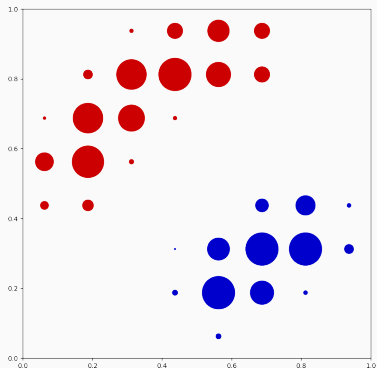$$\alpha = \sum_{i=1}^{N} \alpha_i \delta_{x_i}, \quad \beta = \sum_{j=1}^{M} \beta_j \delta_{y_j}.$$

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

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

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

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

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

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

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

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

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 $\alpha_i$.

8

Computing fidelities between **measures**:

1. **Computer graphics**: weighted Hausdorff distance
2. **Statistics**: kernel distances
3. **Optimal Transport**: Wasserstein distance
   $\simeq$ Robust Point Matching

Computing fidelities between **measures**:
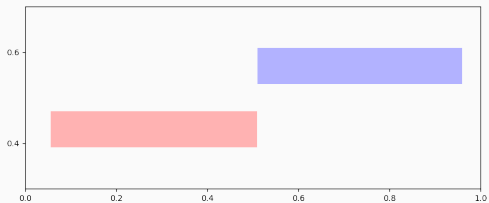
1. **Computer graphics**: weighted Hausdorff distance
2. **Statistics**: kernel distances
3. **Optimal Transport**: Wasserstein distance
$\simeq$ Robust Point Matching
4. What's **new**, in 2019?

# The weighted Hausdorff distance: Iterative Closest Point algorithm

*p*-Hausdorff distance:

$$\text{Loss}(\alpha, \beta) \;=\; \tfrac{1}{2} \sum_i \alpha_i \cdot \min_j \|x_i - y_j\|^p$$

$p$-Hausdorff distance:

$$\text{Loss}(\alpha, \beta) \;=\; \tfrac{1}{2} \sum_i \alpha_i \cdot \min_j \|x_i - y_j\|^p \;\;+\;\; \tfrac{1}{2} \sum_j \beta_j \cdot \min_i \|x_i - y_j\|^p$$

$p$-Hausdorff distance:

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

with $a(x) = \text{d}(x, \text{supp}(\alpha))^p$
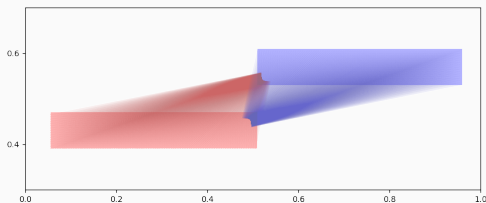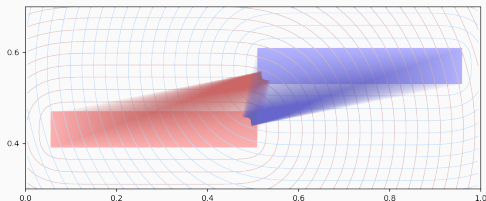$\quad\quad b(x) = \text{d}(x, \text{supp}(\beta))^p$

10

# The weighted Hausdorff distance



$p$-Hausdorff distance:

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

with $a(x) = d(x, \text{supp}(\alpha))^p$
$\quad\quad\; b(x) = d(x, \text{supp}(\beta))^p$

$p$-Hausdorff distance:

$$\text{Loss}(\alpha, \beta) \;=\; \tfrac{1}{2}\textstyle\sum_i \alpha_i \cdot \min_j \|x_i - y_j\|^p \;+\; \tfrac{1}{2}\textstyle\sum_j \beta_j \cdot \min_i \|x_i - y_j\|^p$$

$$= \tfrac{1}{2}\langle\, \alpha \,,\, b \,\rangle \;+\; \tfrac{1}{2}\langle\, \beta \,,\, a \,\rangle$$

$$= \tfrac{1}{2}\langle\, \alpha \,,\, b - a \,\rangle \;+\; \tfrac{1}{2}\langle\, \beta \,,\, a - b \,\rangle$$

with $a(x) = \mathrm{d}(\,x\,,\,\mathrm{supp}(\alpha))^p$
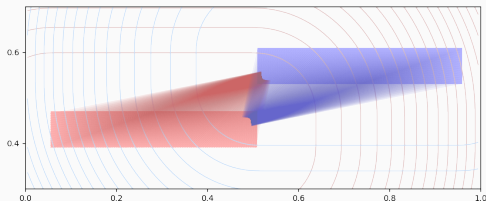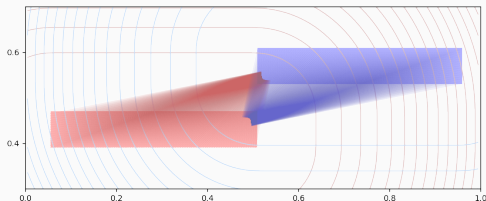$\phantom{with\ } b(x) = \mathrm{d}(\,x\,,\,\mathrm{supp}(\beta))^p$

$p$-Hausdorff distance:

$$\text{Loss}(\alpha, \beta) \;=\; \tfrac{1}{2}\sum_i \alpha_i \cdot \min_j \|x_i - y_j\|^p \quad + \quad \tfrac{1}{2}\sum_j \beta_j \cdot \min_i \|x_i - y_j\|^p$$

$$= \quad \tfrac{1}{2}\langle\, \alpha\,,\, b\,\rangle \qquad\qquad + \qquad\qquad \tfrac{1}{2}\langle\, \beta\,,\, a\,\rangle$$

$$= \quad \tfrac{1}{2}\langle\, \alpha\,,\, b - a\,\rangle \qquad\quad + \qquad\quad \tfrac{1}{2}\langle\, \beta\,,\, a - b\,\rangle$$

$$= \quad \tfrac{1}{2}\langle\, \alpha - \beta\,,\, b - a\,\rangle$$

with $a(x) = \mathrm{d}(x, \mathrm{supp}(\alpha))^p$
$\phantom{with}\; b(x) = \mathrm{d}(x, \mathrm{supp}(\beta))^p$

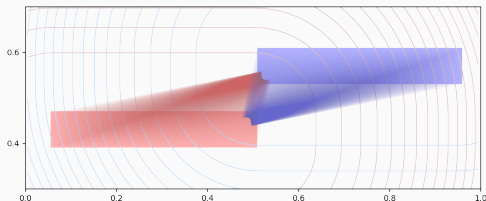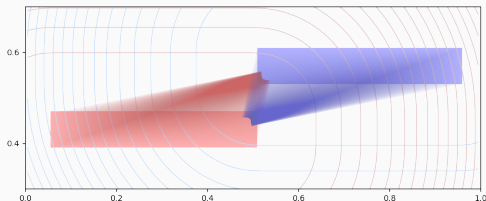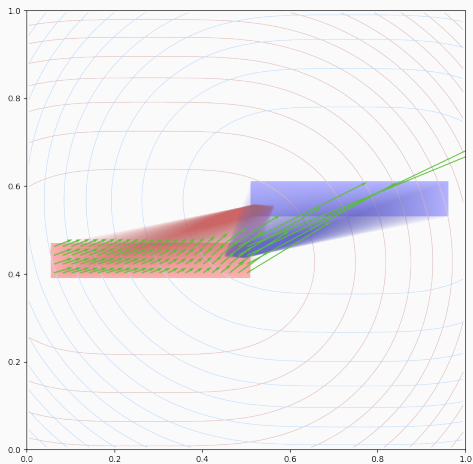$p$-Hausdorff distance:

$$\text{Loss}(\alpha, \beta) \;=\; \tfrac{1}{2}\sum_i \alpha_i \cdot \min_j \|x_i - y_j\|^p \;+\; \tfrac{1}{2}\sum_j \beta_j \cdot \min_i \|x_i - y_j\|^p$$

$$=\; \tfrac{1}{2}\langle\, \alpha \,,\, b \,\rangle \;+\; \tfrac{1}{2}\langle\, \beta \,,\, a \,\rangle$$

$$=\; \tfrac{1}{2}\langle\, \alpha \,,\, b - a \,\rangle \;+\; \tfrac{1}{2}\langle\, \beta \,,\, a - b \,\rangle$$

$$=\; \tfrac{1}{2}\langle\, \alpha - \beta \,,\, b - a \,\rangle$$

$$\left. \begin{array}{llll} \text{with } a(x) &=& \mathrm{d}(\,x\,,\,\mathrm{supp}(\alpha))^p &\simeq\; -\log(k \star \alpha) \\ b(x) &=& \mathrm{d}(\,x\,,\,\mathrm{supp}(\beta))^p &\simeq\; -\log(k \star \beta) \end{array} \right\} \text{ GMM log-likelihoods}$$

10

$$\text{Loss}(\alpha, \beta) \quad = \quad \tfrac{1}{2}\langle\, \alpha,\ b - a\,\rangle \ + \ \tfrac{1}{2}\langle\, \beta,\ a - b\,\rangle$$



11

$$\text{Loss}(\alpha, \beta) \quad = \quad \tfrac{1}{2} \langle\, \alpha,\ b - a \,\rangle \ + \ \tfrac{1}{2} \langle\, \beta,\ a - b \,\rangle$$

$$\text{Loss}(\alpha, \beta) \quad = \quad \tfrac{1}{2}\langle\, \alpha,\ b - a\,\rangle \ + \ \tfrac{1}{2}\langle\, \beta,\ a - b\,\rangle$$

$$\text{Loss}(\alpha, \beta) \quad = \quad \tfrac{1}{2}\langle\, \alpha,\; b - a \,\rangle \; + \; \tfrac{1}{2}\langle\, \beta,\; a - b \,\rangle$$

$t = .00$

$t = .25$

$t = .50$

$t = 1.00$

$t = 5.00$

$t = 10.00$

An idea from statistics:
Kernel distances

Raw signal ($\alpha - \beta$).

Blurred signal $g \star (\alpha - \beta)$.

Choose a symmetric blurring function $g$, a **kernel** $k = g \star g$:

$$d_k(\alpha, \beta) = \| g \star \alpha - g \star \beta \|_{L^2}^2$$

Blurred signal $g \star (\alpha - \beta)$.

Choose a symmetric blurring function $g$, a **kernel** $k = g \star g$:

$$d_k(\alpha, \beta) = \| g \star \alpha - g \star \beta \|_{L^2}^2$$
$$= \langle\, \alpha - \beta \mid k \star (\alpha - \beta)\,\rangle$$

Blurred signal $g \star (\alpha - \beta)$.

Choose a symmetric blurring function $g$, a **kernel** $k = g \star g$:

$$
\begin{aligned}
\mathrm{d}_k(\alpha, \beta) &= \| \, g \star \alpha - g \star \beta \, \|_{L^2}^2 \\
&= \langle \, \alpha - \beta \mid k \star (\alpha - \beta) \, \rangle \\
&= -2 \sum_{i,j} k(x_i, y_j) \, \alpha_i \, \beta_j + \cdots
\end{aligned}
$$

Blurred signal $g \star (\alpha - \beta)$.

Choose a symmetric blurring function $g$, a **kernel** $k = g \star g$:

$$
\begin{aligned}
d_k(\alpha, \beta) &= \| g \star \alpha - g \star \beta \|_{L^2}^2 \\
&= \langle \alpha - \beta \mid k \star (\alpha - \beta) \rangle \\
&= -2 \sum_{i,j} k(x_i, y_j)\, \alpha_i\, \beta_j + \cdots \\
&= \langle \alpha - \beta \mid b^k - a^k \rangle
\end{aligned}
$$

with $a^k = -k \star \alpha,\ b^k = -k \star \beta$.

**Kernel** distances, aka. **blurred** SSDs:

$$\text{choose} \qquad a(x) \;=\; -(k \star \alpha)(x) \;=\; -\sum_i \alpha_i \, k(x, x_i)$$

$$\text{and use} \qquad \tfrac{1}{2}\langle\, \alpha - \beta \,,\, b - a \,\rangle \;=\; \tfrac{1}{2}\langle\, \alpha - \beta \,,\, k \star (\alpha - \beta) \,\rangle.$$

**Kernel** distances, aka. **blurred** SSDs:

$$\text{choose} \qquad a(x) \;=\; -(k \star \alpha)(x) \;=\; -\sum_i \alpha_i \, k(x, x_i)$$

$$\text{and use} \qquad \tfrac{1}{2}\langle\, \alpha - \beta \,,\; b - a \,\rangle \;=\; \tfrac{1}{2}\langle\, \alpha - \beta \,,\; k \star (\alpha - \beta) \,\rangle.$$

The **Energy Distance:** an underrated kernel, $k(x, y) \;=\; -\|x - y\|$.

$$a(x) \;=\; \sum_i \alpha_i \|x - x_i\| \qquad \text{instead of} \quad a(x) \;=\; \min_i \|x - x_i\|$$

$$b(x) \;=\; \sum_j \beta_j \|x - y_j\| \qquad \text{instead of} \quad b(x) \;=\; \min_j \|x - y_j\|.$$

**Kernel** distances, aka. **blurred** SSDs:

$$\text{choose} \qquad a(x) \;=\; -(k \star \alpha)(x) \;=\; -\sum_i \alpha_i \, k(x, x_i)$$

$$\text{and use} \qquad \tfrac{1}{2}\langle\, \alpha - \beta \,,\, b - a \,\rangle \;=\; \tfrac{1}{2}\langle\, \alpha - \beta \,,\, k \star (\alpha - \beta) \,\rangle.$$

The **Energy Distance:** an underrated kernel, $k(x, y) \;=\; -\|x - y\|$.

$$a(x) \;=\; \sum_i \alpha_i \|x - x_i\| \qquad \text{instead of} \quad a(x) \;=\; \min_i \|x - x_i\|$$

$$b(x) \;=\; \sum_j \beta_j \|x - y_j\| \qquad \text{instead of} \quad b(x) \;=\; \min_j \|x - y_j\|.$$

$$\text{Loss}(\alpha, \beta) \;=\; \sum_i \sum_j \alpha_i \beta_j \,\|x_i - y_j\|$$

$$- \; \tfrac{1}{2}\sum_i \sum_j \alpha_i \alpha_j \,\|x_i - x_j\| \;-\; \tfrac{1}{2}\sum_i \sum_j \beta_i \beta_j \,\|y_i - y_j\|$$

**Kernel** distances, aka. **blurred** SSDs:

$$\text{choose} \quad a(x) \; = \; -(k \star \alpha)(x) \; = \; -\sum_i \alpha_i \, k(x, x_i)$$

$$\text{and use} \quad \tfrac{1}{2}\langle \, \alpha - \beta \, , \, b - a \, \rangle \; = \; \tfrac{1}{2}\langle \, \alpha - \beta \, , \, k \star (\alpha - \beta) \, \rangle.$$

The **Energy Distance:** an underrated kernel, $k(x, y) \; = \; -\|x - y\|$.

$$a(x) \; = \; \sum_i \alpha_i \|x - x_i\| \quad \text{instead of} \quad a(x) \; = \; \min_i \|x - x_i\|$$

$$b(x) \; = \; \sum_j \beta_j \|x - y_j\| \quad \text{instead of} \quad b(x) \; = \; \min_j \|x - y_j\|.$$

$$\text{Loss}(\alpha, \beta) \; = \; \sum_i \sum_j \alpha_i \beta_j \|x_i - y_j\| \; \simeq \; \textbf{Electrostatic Energy}$$

$$- \; \tfrac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \|x_i - x_j\| \; - \; \tfrac{1}{2} \sum_i \sum_j \beta_i \beta_j \|y_i - y_j\|$$

$t = .00$

$$t = .25$$

$t = .50$

$t = 1.00$

$t = 5.00$

$t = 10.00$

Hausdorff, min

Kernel, $\sum$

Hausdorff, min                              Kernel, $\sum$

$\implies$ Can we get the best of both worlds?

# An idea from Optimal Transport theory: The SoftAssign algorithm

Minimize over $N$-by-$M$ matrices (transport plans) $\pi$ :

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

subject to $\quad \pi_{i,j} \geqslant 0,$

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

17

With $C(x_i, y_j) = \|x_i - y_j\|^p$,

$$OT(\alpha, \beta) = \min_\pi \langle \pi, C \rangle \qquad\qquad \longrightarrow \text{Assignment}$$
$$\text{s.t. } \pi \geqslant 0, \qquad \pi \mathbf{1} = \alpha, \qquad \pi^\mathsf{T} \mathbf{1} = \beta$$

With $C(x_i, y_j) = \|x_i - y_j\|^p$,

$$OT(\alpha, \beta) = \min_{\pi} \langle \pi, C \rangle \qquad\qquad \longrightarrow \text{Assignment}$$

$$\text{s.t. } \pi \geqslant 0, \qquad \pi\, \mathbf{1} = \alpha, \qquad \pi^{\mathsf{T}}\mathbf{1} = \beta$$

$$= \max_{f,g} \langle \alpha, f \rangle + \langle \beta, g \rangle \qquad\qquad \longrightarrow \text{FedEx}$$

$$\text{s.t.} \qquad f(x_i) + g(y_j) \leqslant C(x_i, y_j),$$

With $C(x_i, y_j) = \|x_i - y_j\|^p$,

$$OT(\alpha, \beta) = \min_{\pi} \langle \pi, C \rangle \qquad\qquad \longrightarrow \text{Assignment}$$

$$\text{s.t. } \pi \geqslant 0, \qquad \pi\, \mathbf{1} = \alpha, \qquad \pi^\mathsf{T} \mathbf{1} = \beta$$

$$= \max_{f,g} \langle \alpha, f \rangle + \langle \beta, g \rangle \qquad\qquad \longrightarrow \text{FedEx}$$

$$\text{s.t. } \qquad f(x_i) + g(y_j) \leqslant C(x_i, y_j),$$

$\implies$ **Combinatorial** problem on the simplex

With $C(x_i, y_j) = \|x_i - y_j\|^p$,

$$OT(\alpha, \beta) = \min_{\pi} \langle \pi, C \rangle \qquad\qquad \longrightarrow \text{Assignment}$$
$$\text{s.t. } \pi \geqslant 0, \qquad \pi\mathbf{1} = \alpha, \qquad \pi^{\mathsf{T}}\mathbf{1} = \beta$$

$$= \max_{f,g} \langle \alpha, f \rangle + \langle \beta, g \rangle \qquad\qquad \longrightarrow \text{FedEx}$$
$$\text{s.t. } \qquad f(x_i) + g(y_j) \leqslant C(x_i, y_j),$$

$\Longrightarrow$ **Combinatorial** problem on the simplex

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

For $\varepsilon > 0$ :

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

$$+ \; \varepsilon \underbrace{\sum_{i,j} \pi_{i,j} \cdot \log \frac{\pi_{i,j}}{\alpha_i \, \beta_j}}_{\text{entropic barrier}}$$

subject to

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

19

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

$$\text{s.t.} \qquad \pi\, \mathbf{1} \;=\; \alpha, \qquad \pi^\mathsf{T} \mathbf{1} \;=\; \beta$$

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

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

$$= \max_{f, g} \langle \alpha, f \rangle + \langle \beta, g \rangle \qquad \longrightarrow \text{Cheeky FedEx}$$

$$- \underbrace{\varepsilon \langle \alpha \otimes \beta, e^{(f \oplus g - \mathrm{C})/\varepsilon} - 1 \rangle}_{\text{soft constraint } f \oplus g \leqslant \mathrm{C}}$$

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

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

$$= \max_{f, g} \langle \alpha, f \rangle + \langle \beta, g \rangle \qquad \longrightarrow \text{Cheeky FedEx}$$

$$- \underbrace{\varepsilon \langle \alpha \otimes \beta, e^{(f \oplus g - \mathsf{C})/\varepsilon} - 1 \rangle}_{\text{soft constraint } f \oplus g \leqslant \mathsf{C}}$$

$\implies$ **Strictly convex** problem on the simplex

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

$$\text{s.t.} \qquad \pi\, \mathbf{1} \;=\; \alpha, \qquad \pi^{\mathsf{T}} \mathbf{1} \;=\; \beta$$

$$=\; \max_{f,g} \langle\, \alpha\,,\, f\,\rangle \;+\; \langle\, \beta\,,\, g\,\rangle \qquad \longrightarrow \text{ Cheeky FedEx}$$

$$-\; \underbrace{\varepsilon \langle\, \alpha \otimes \beta\,,\, e^{(f \oplus g - \mathsf{C})/\varepsilon} - 1\,\rangle}_{\text{soft constraint } f \oplus g \leqslant \mathsf{C}}$$

$\Longrightarrow$ **Strictly convex** problem on the simplex

$$\text{At the optimum,} \quad \pi \;=\; e^{(f \oplus g - \mathsf{C})/\varepsilon} \,\cdot\, \alpha \otimes \beta$$

$$\text{i.e.} \qquad \pi_{i,j} \;=\; \alpha_i\, e^{f_i/\varepsilon}\, e^{-\mathsf{C}(x_i, y_j)/\varepsilon}\, e^{g_j/\varepsilon}\, \beta_j.$$

$$\pi_{i,j} = \Delta(U\alpha) \cdot K_{x,y} \cdot \Delta(V\beta)$$

with

- a kernel function $k$

$$k(x_i - y_j) = e^{-C(x_i, y_j)/\varepsilon}.$$

- $U = e^{f/\varepsilon}$ and $V = e^{g/\varepsilon}$, positive weights on $\{x_i\}$ and $\{y_j\}$.

$\rightarrow$ Enforce the **constraints**

$$\pi \, \mathbf{1} = \alpha, \qquad \pi^{\mathsf{T}} \mathbf{1} = \beta$$

Source and target.

---

### Sinkhorn Iterative Algorithm

**Input** : source $\alpha = \sum_i \alpha_i \delta_{x_i}$
    target $\beta = \sum_j \beta_j \delta_{y_j}$

**Parameter** : $k : x \mapsto e^{-|x|^2/\varepsilon}$

1: $U \leftarrow ones(size(\alpha))$
2: $V \leftarrow ones(size(\beta))$
3: **while** updates $>$ tol **do**
4:     $U \leftarrow 1 \; ./ \; \mathrm{K} \; \cdot (V\beta)$
5:     $V \leftarrow 1 \; ./ \; \mathrm{K}^{\mathsf{T}} \cdot (U\alpha)$
6: **return** $\varepsilon \left( \langle \alpha, \log(U) \rangle + \langle \beta, \log(V) \rangle \right)$

**Output** : fidelity $\mathrm{OT}_\varepsilon(\alpha, \beta)$

---

Seen by the kernel $k$.

| **Sinkhorn** Iterative Algorithm |
|---|
| **Input** : source $\alpha = \sum_i \alpha_i \delta_{x_i}$ |
| target $\beta = \sum_j \beta_j \delta_{y_j}$ |
| **Parameter :** $k : x \mapsto e^{-\lvert x\rvert^2/\varepsilon}$ |
| 1: $U \leftarrow \mathit{ones}(\mathit{size}(\alpha))$ |
| 2: $V \leftarrow \mathit{ones}(\mathit{size}(\beta))$ |
| 3: **while** updates $>$ tol **do** |
| 4: $\quad U \leftarrow 1 \;./\; \mathrm{K} \cdot (V\beta)$ |
| 5: $\quad V \leftarrow 1 \;./\; \mathrm{K}^{\mathsf{T}} \cdot (U\alpha)$ |
| 6: **return** $\varepsilon \left( \langle \alpha, \log(U) \rangle + \langle \beta, \log(V) \rangle \right)$ |
| **Output :** fidelity $\mathrm{OT}_\varepsilon(\alpha, \beta)$ |

Sinkhorn Iteration 000

Starting estimate.

### **Sinkhorn** Iterative Algorithm

**Input** : source $\alpha = \sum_i \alpha_i \delta_{x_i}$

            target $\beta = \sum_j \beta_j \delta_{y_j}$

**Parameter** : $k : x \mapsto e^{-|x|^2/\varepsilon}$

1: $U \leftarrow ones(size(\alpha))$

2: $V \leftarrow ones(size(\beta))$
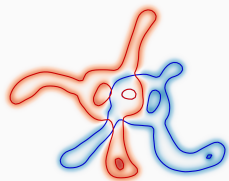
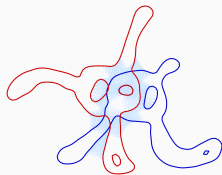3: **while** updates > tol **do**

4:     $U \leftarrow 1 ./ K \cdot (V\beta)$

5:     $V \leftarrow 1 ./ K^\mathsf{T} \cdot (U\alpha)$

6: **return** $\varepsilon \left( \langle \alpha, \log(U) \rangle + \langle \beta, \log(V) \rangle \right)$

    **Output** : fidelity $\mathrm{OT}_\varepsilon(\alpha, \beta)$

Sinkhorn Iteration 250

Computing the OT plan.

---

**Sinkhorn** Iterative Algorithm

---

**Input** : source $\alpha = \sum_i \alpha_i \delta_{x_i}$

target $\beta = \sum_j \beta_j \delta_{y_j}$

**Parameter** : $k : x \mapsto e^{-|x|^2/\varepsilon}$

1: $U \leftarrow ones(size(\alpha))$

2: $V \leftarrow ones(size(\beta))$

3: **while** updates > tol **do**

4: $\quad U \leftarrow 1 \ ./ \ \mathsf{K} \ \cdot (V\beta)$

5: $\quad V \leftarrow 1 \ ./ \ \mathsf{K}^\mathsf{T} \cdot (U\alpha)$

6: **return** $\varepsilon \left( \langle \alpha, \log(U) \rangle + \langle \beta, \log(V) \rangle \right)$

**Output** : fidelity $\mathrm{OT}_\varepsilon(\alpha, \beta)$

---

Sinkhorn Iteration 250

Computing the OT plan.

**Sinkhorn** Iterative Algorithm

**Input** : source $\alpha = \sum_i \alpha_i \delta_{x_i}$
target $\beta = \sum_j \beta_j \delta_{y_j}$
**Parameter** : $k : x \mapsto e^{-|x|^2/\varepsilon}$

1: $U \leftarrow ones(size(\alpha))$
2: $V \leftarrow ones(size(\beta))$
3: **while** updates > tol **do**
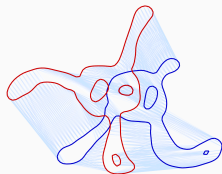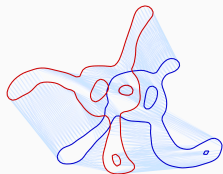4: $\quad U \leftarrow 1 \ ./ \ \mathsf{K} \cdot (V\beta)$
5: $\quad V \leftarrow 1 \ ./ \ \mathsf{K}^\mathsf{T} \cdot (U\alpha)$
6: **return** $\varepsilon \left( \langle \alpha, \log(U) \rangle + \langle \beta, \log(V) \rangle \right)$

**Output :** fidelity $\mathrm{OT}_\varepsilon(\alpha, \beta)$

$\simeq$

*TPS-RPM algorithm*,
Chui and Rangarajan, CVPR **2000**

*Optimal Transport for diffeomorphic registration*, Feydy et al., MICCAI **2017**

*TPS-RPM algorithm,*
Chui and Rangarajan, CVPR **2000**

*Optimal Transport for diffeomorphic registration*, Feydy et al., MICCAI **2017**

$\implies$ We've added weights, orientations, convergence analysis...
But shouldn't we go a bit **further**?

# It's 2019 now:
# What's new?

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

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

## Fact 1 : if $\varepsilon > 0$, $\mathrm{OT}_\varepsilon$ is *not* a valid divergence

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

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

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

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

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



**Bad news:** for $0 < \varepsilon \leqslant +\infty$, we converge towards $\alpha$ such that

$$\mathrm{OT}_\varepsilon(\alpha, \beta) \;\; < \;\; \mathrm{OT}_\varepsilon(\beta, \beta).$$

*TPS-RPM algorithm*, Chui and Rangarajan, CVPR 2000

*TPS-RPM algorithm*, Chui and Rangarajan, CVPR 2000

$\implies$ **Cumbersome** and brittle workaround,
with parameters to tune.

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

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

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

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

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

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

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

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

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

$$\mathsf{S}_\varepsilon(\alpha, \beta) = \text{OT}_\varepsilon(\alpha, \beta) - \tfrac{1}{2}\text{OT}_\varepsilon(\alpha, \alpha) - \tfrac{1}{2}\text{OT}_\varepsilon(\beta, \beta)$$

26

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

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

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

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

$$\mathsf{S}_\varepsilon(\alpha, \beta) \ = \ \mathsf{OT}_\varepsilon(\alpha, \beta) \ - \ \tfrac{1}{2}\mathsf{OT}_\varepsilon(\alpha, \alpha) \ - \ \tfrac{1}{2}\mathsf{OT}_\varepsilon(\beta, \beta)$$

$$\mathsf{Wasserstein}_{+\mathsf{C}}(\alpha, \beta) \ \xleftarrow{\ \varepsilon \to 0\ } \ \mathsf{S}_\varepsilon(\alpha, \beta) \ \xrightarrow{\ \varepsilon \to +\infty\ } \ \mathsf{Kernel}_{-\mathsf{C}}(\alpha, \beta)$$

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

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

$$\mathrm{OT}_\varepsilon(\alpha, \beta) \xrightarrow{\quad \varepsilon \to +\infty \quad} \langle \alpha \otimes \beta, \mathsf{C} \rangle = \langle \alpha, \mathsf{C} \star \beta \rangle$$

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

$$\mathrm{S}_\varepsilon(\alpha, \beta) = \mathrm{OT}_\varepsilon(\alpha, \beta) - \tfrac{1}{2}\mathrm{OT}_\varepsilon(\alpha, \alpha) - \tfrac{1}{2}\mathrm{OT}_\varepsilon(\beta, \beta)$$

$$\text{Wasserstein}_{+\mathsf{C}}(\alpha, \beta) \xleftarrow{\varepsilon \to 0} \mathrm{S}_\varepsilon(\alpha, \beta) \xrightarrow{\varepsilon \to +\infty} \text{Kernel}_{-\mathsf{C}}(\alpha, \beta)$$

**In practice,** $\mathrm{S}_\varepsilon$ 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 $\alpha$, $\beta$ and regularization $\varepsilon > 0$ :*

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

*For all probability measures $\alpha$, $\beta$ and regularization $\varepsilon > 0$:*

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

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

*For all probability measures $\alpha$, $\beta$ and regularization $\varepsilon > 0$ :*

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

$$\alpha \mapsto S_\varepsilon(\alpha, \beta) \text{ is convex and differentiable}$$

27

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

*For all probability measures $\alpha$, $\beta$ and regularization $\varepsilon > 0$ :*

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

$$\alpha \mapsto S_\varepsilon(\alpha, \beta) \ \text{is convex and differentiable}$$

*These results can be generalized to arbitrary **feature** spaces*
*– e.g. (position,orientation,curvature).*

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

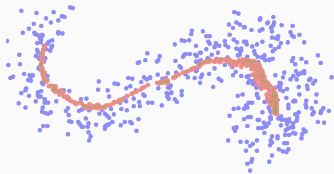*For all probability measures $\alpha$, $\beta$ and regularization $\varepsilon > 0$ :*

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

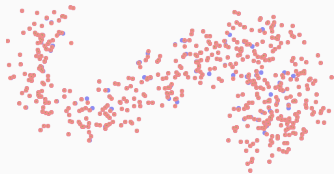$$\alpha \mapsto S_\varepsilon(\alpha, \beta) \quad \text{is convex and differentiable}$$

*These results can be generalized to arbitrary **feature** spaces*
*– e.g. (position,orientation,curvature).*



Loss = $OT_\varepsilon$            Loss = $S_\varepsilon$

Unfortunately,

$$k(x_i, y_j) \simeq 0 \qquad \text{if } \varepsilon \text{ is too small.}$$

Unfortunately,

$$k(x_i, y_j) \simeq 0 \qquad \text{if } \varepsilon \text{ is too small.}$$

$$\text{OT}_\varepsilon(\alpha, \beta) = \max_{f,g} \langle \alpha, f \rangle + \langle \beta, g \rangle \qquad \longrightarrow \text{Cheeky FedEx}$$
$$- \underbrace{\varepsilon \langle \alpha \otimes \beta, e^{(f \oplus g - C)/\varepsilon} - 1 \rangle}_{\text{soft constraint } f \oplus g \leqslant C}$$

Unfortunately,

$$k(x_i, y_j) \simeq 0 \qquad \text{if } \varepsilon \text{ is too small.}$$

$$\mathrm{OT}_\varepsilon(\alpha, \beta) = \max_{f,g} \langle \alpha, f \rangle + \langle \beta, g \rangle \qquad \longrightarrow \text{Cheeky FedEx}$$

$$- \underbrace{\varepsilon \langle \alpha \otimes \beta, e^{(f \oplus g - C)/\varepsilon} - 1 \rangle}_{\text{soft constraint } f \oplus g \leqslant C}$$

Equivalent to the constraints on $\pi$, the optimality conditions read:

$$f(x_i) = -\varepsilon \log \sum_j \beta_j \exp \tfrac{1}{\varepsilon}(g(y_j) - C(x_i, y_j)),$$

$$g(y_j) = -\varepsilon \log \sum_i \alpha_i \exp \tfrac{1}{\varepsilon}(f(x_i) - C(x_i, y_j)).$$

# The SoftMin interpolates between a minimum and a sum

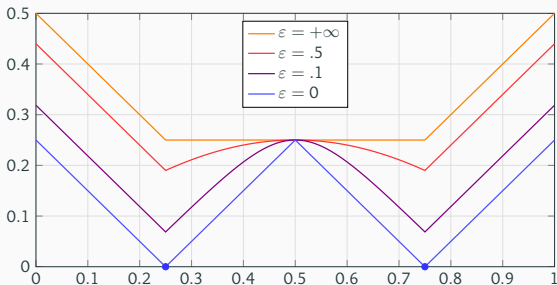$$\log\left(e^c + e^d\right) = \max(c,d) + \log\left(\underbrace{e^{c-\max(c,d)} + e^{d-\max(c,d)}}_{\in[1,2]}\right)$$

# The SoftMin interpolates between a minimum and a sum

$$\log\left( e^c + e^d \right) \;=\; \max(c,d) \;+\; \log\big( \underbrace{e^{c-\max(c,d)} + e^{d-\max(c,d)}}_{\in [1,2]} \big)$$

Building on this, for a regularization parameter $\varepsilon > 0$, we define

$$b^\varepsilon(x) \;=\; \min_{\substack{y \sim \beta}}{}_\varepsilon \; \|x-y\| \;=\; -\varepsilon \, \log \sum_{j=1}^{M} \beta_j \, \exp\big( -\tfrac{1}{\varepsilon}\|x-y_j\| \big)$$



$b^\varepsilon(x)$, with $\beta = \tfrac{1}{2}\delta_{.25} + \tfrac{1}{2}\delta_{.75}$

The optimality conditions read:

$$f(x_i) \;=\; b(x) \;=\; -\varepsilon \, \log \sum_j \beta_j \, \exp \tfrac{1}{\varepsilon} \big[\, g(y_j) - \mathsf{C}(x_i, y_j) \,\big],$$

$$g(y_j) \;=\; a(y) \;=\; -\varepsilon \, \log \sum_i \alpha_i \, \exp \tfrac{1}{\varepsilon} \big[\, f(x_i) - \mathsf{C}(x_i, y_j) \,\big].$$

The optimality conditions read:

$$f(x_i) = b(x) = \min_{y \sim \beta}{}^{\varepsilon} \big[\, C(x,y) - a(y) \,\big] \qquad ,$$

$$g(y_j) = a(y) = \min_{x \sim \alpha}{}^{\varepsilon} \big[\, C(x,y) - b(x) \,\big] \qquad .$$

The optimality conditions read:

$$f(x_i) = b(x) = \min_{\substack{y \sim \beta}}{}^{\varepsilon} \left[ C(x,y) - a(y) \right] \quad ,$$

$$g(y_j) = a(y) = \min_{\substack{x \sim \alpha}}{}^{\varepsilon} \left[ C(x,y) - b(x) \right] \quad .$$

Final cost:

$$\mathrm{OT}_{\varepsilon}(\alpha, \beta) = \langle \alpha, f \rangle + \langle \beta, g \rangle = \langle \alpha, b \rangle + \langle \beta, a \rangle,$$

$$\mathrm{S}_{\varepsilon}(\alpha, \beta) = \langle \alpha, b^{\beta \to \alpha} - a^{\alpha \leftrightarrow \alpha} \rangle + \langle \beta, a^{\alpha \to \beta} - b^{\beta \leftrightarrow \beta} \rangle.$$
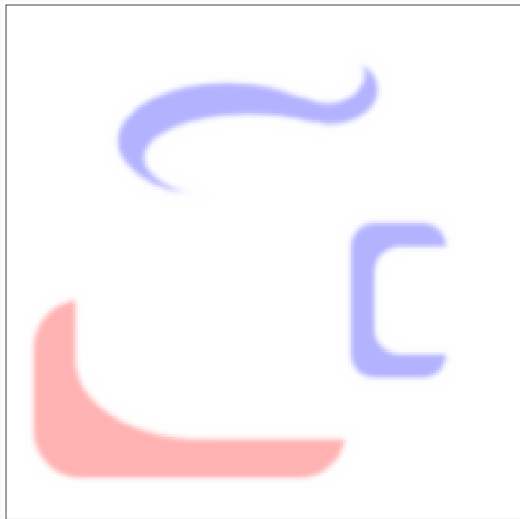
The optimality conditions read:

$$f(x_i) = b(x) = \min_{\substack{\varepsilon \\ y \sim \beta}} \big[\, C(x,y) - a(y) \,\big] \qquad ,$$

$$g(y_j) = a(y) = \min_{\substack{\varepsilon \\ x \sim \alpha}} \big[\, C(x,y) - b(x) \,\big] \qquad .$$

Final cost:

$$\mathrm{OT}_\varepsilon(\alpha, \beta) = \langle\, \alpha, f \,\rangle + \langle\, \beta, g \,\rangle = \langle\, \alpha, b \,\rangle + \langle\, \beta, a \,\rangle,$$

$$\mathrm{S}_\varepsilon(\alpha, \beta) = \langle\, \alpha, b^{\beta \to \alpha} - a^{\alpha \leftrightarrow \alpha} \,\rangle + \langle\, \beta, a^{\alpha \to \beta} - b^{\beta \leftrightarrow \beta} \,\rangle.$$

Discrete, computational OT [Cuturi, 2013, Peyré and Cuturi, 2018]:

       Start from an $\varepsilon$-smoothed **Hausdorff** distance, but let

       the influence fields $a$ and $b$ **interact** with each other.

       Enforce a **mass spreading** constraint on the spring system:

             all of $\alpha$ should be linked to all of $\beta$.

30

Iteration 0

Iteration 1

Iteration 2

Iteration 3

Iteration 4

Iteration 5

Iteration 6

Iteration 7

Iteration 8

Leverages the KeOps library [Charlier et al., 2018]:

$$\Longrightarrow \texttt{pip install pykeops} \Longleftarrow$$

**Loss + gradient** in 3D, on a cheap laptop's GPU (GTX960M)

$t = .00$

$t = .25$

$t = .50$

$t = 1.00$

$t = 5.00$

$t = 10.00$

Iteration 0

Iteration 1

Iteration 2

Iteration 10

Iteration 0

Iteration 1

Iteration 2

Iteration 10

Iteration 0

Iteration 1

Iteration 2

Iteration 10

Iteration 0

Iteration 1

Iteration 2

Iteration 10

# Conclusion

The **true** $OT_0$ problem is **hard.**

The **true** $OT_0$ problem is **hard**.
Approximating it with **subsampling** or **smoothing** is easy:
this is what **SoftAssign** is all about.

The **true** $OT_0$ problem is **hard**.
Approximating it with **subsampling** or **smoothing** is easy:
this is what **SoftAssign** is all about.

Remarkably, $S_\varepsilon(\alpha, \beta)$ is a cheap approximation of $OT_0(\alpha, \beta)$
that defines a **positive definite** cost between the **discrete samples**.
It is the first known way of doing so.

$$\text{Loss}(\alpha, \beta) \,=\, \max_{f \in B} \langle\, \alpha - \beta \,,\, f \,\rangle,$$

$$\text{look for } \theta^* \,=\, \arg\min_{\theta} \max_{f \in B} \langle\, \alpha(\theta) - \beta \,,\, f \,\rangle$$

$$\text{Loss}(\alpha, \beta) \; = \; \max_{f \in B} \langle \, \alpha - \beta \, , \, f \, \rangle,$$

$$\text{look for } \; \theta^* \; = \; \arg\min_{\theta} \max_{f \in B} \langle \, \alpha(\theta) - \beta \, , \, f \, \rangle$$

- $B = \{ \, \|f\|_\infty \leqslant 1 \, \} \implies \text{Loss} = \text{TV norm}:$
  - zero geometry
  - **too many** test functions

$$\text{Loss}(\alpha, \beta) \, = \, \max_{f \in B} \langle \, \alpha - \beta \, , \, f \, \rangle,$$

$$\text{look for} \ \theta^* \, = \, \arg \min_\theta \max_{f \in B} \langle \, \alpha(\theta) - \beta \, , \, f \, \rangle$$

- $B = \{ \, \|f\|_\infty \leqslant 1 \, \} \implies \text{Loss} = \text{TV norm:}$
  - zero geometry
  - **too many** test functions

- $B = \{ \, \|f\|_2^2 + \|\nabla f\|_2^2 + \cdots \leqslant 1 \, \} \implies \text{Loss} = \text{kernel norm:}$
  - may saturate at infinity
  - **screening** artifacts

$$\text{Loss}(\alpha, \beta) \; = \; \max_{f \in B} \langle \, \alpha - \beta \, , \, f \, \rangle,$$

$$\text{look for } \; \theta^* \; = \; \arg \min_{\theta} \max_{f \in B} \langle \, \alpha(\theta) - \beta \, , \, f \, \rangle$$

- $B = \{ \, f \text{ is 1-Lipschitz} \, \} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0\text{)}:$

$$\text{Loss}(\alpha, \beta) \;=\; \max_{f \in B} \langle\, \alpha - \beta\,,\, f\,\rangle,$$

$$\text{look for } \theta^* \;=\; \arg\min_{\theta}\, \max_{f \in B} \langle\, \alpha(\theta) - \beta\,,\, f\,\rangle$$

- $B = \{\, f \text{ is 1-Lipschitz}\,\} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0\text{)}:$
  - $S_\varepsilon$ is nearly as efficient as a **closed formula**

$$\text{Loss}(\alpha, \beta) \ = \ \max_{f \in B} \langle \, \alpha - \beta \, , f \, \rangle,$$

$$\text{look for} \ \ \theta^* \ = \ \arg\min_{\theta} \max_{f \in B} \langle \, \alpha(\theta) - \beta \, , f \, \rangle$$

- $B = \{\, f \text{ is 1-Lipschitz} \,\} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0\text{)}$:
  - $S_\varepsilon$ is nearly as efficient as a **closed formula**
  - relevant in **low dimensions**
  - **useless** in $(\mathbb{R}^{512 \times 512}, \| \cdot \|_2)$: the ground cost makes no sense

$$\text{Loss}(\alpha, \beta) \; = \; \max_{f \in B} \langle \, \alpha - \beta \, , \, f \, \rangle,$$

$$\text{look for} \;\; \theta^* \; = \; \arg \min_{\theta} \max_{f \in B} \langle \, \alpha(\theta) - \beta \, , \, f \, \rangle$$

- $B = \{ \, f \text{ is 1-Lipschitz} \, \} \Longrightarrow \text{Loss} = \text{Wasserstein-1 (OT}_0\text{)}$:
  - $S_\varepsilon$ is nearly as efficient as a **closed formula**
  - relevant in **low dimensions**
  - **useless** in $(\mathbb{R}^{512 \times 512}, \| \cdot \|_2)$: the ground cost makes no sense

- $B \simeq \{ \, f \text{ is 1-Lipschitz} \, \} \bigcap \{ \, f \text{ is a } \textbf{CNN} \, \}$
  $\Longrightarrow \text{Loss} = \text{Wasserstein GAN}$ :

$$\text{Loss}(\alpha, \beta) = \max_{f \in B} \langle \alpha - \beta, f \rangle,$$

$$\text{look for } \theta^* = \arg\min_{\theta} \max_{f \in B} \langle \alpha(\theta) - \beta, f \rangle$$

- $B = \{ f \text{ is } 1\text{-Lipschitz} \} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0\text{)}:$
  - $S_\varepsilon$ is nearly as efficient as a **closed formula**
  - relevant in **low dimensions**
  - **useless** in $(\mathbb{R}^{512 \times 512}, \|\cdot\|_2)$: the ground cost makes no sense

- $B \simeq \{ f \text{ is } 1\text{-Lipschitz} \} \bigcap \{ f \text{ is a } \textbf{CNN} \}$
  $\implies \text{Loss} = \text{Wasserstein GAN}:$
  - use **perceptually sensible** test functions

$$\text{Loss}(\alpha, \beta) \ = \ \max_{f \in B} \langle \, \alpha - \beta \, , f \, \rangle,$$

$$\text{look for } \theta^* \ = \ \arg\min_\theta \max_{f \in B} \langle \, \alpha(\theta) - \beta \, , f \, \rangle$$

- $B = \{\, f \text{ is 1-Lipschitz} \,\} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0)$:
    - $S_\varepsilon$ is nearly as efficient as a **closed formula**
    - relevant in **low dimensions**
    - **useless** in $(\mathbb{R}^{512 \times 512}, \|\cdot\|_2)$: the ground cost makes no sense

- $B \simeq \{\, f \text{ is 1-Lipschitz} \,\} \bigcap \{\, f \text{ is a } \textbf{CNN} \,\}$
    $\implies \text{Loss} = \text{Wasserstein GAN}$ :
    - use **perceptually sensible** test functions
    - no simple formula: use **gradient ascent**

$$\text{Loss}(\alpha, \beta) \; = \; \max_{f \in B} \langle \, \alpha - \beta \, , \, f \rangle,$$

$$\text{look for } \theta^* \; = \; \arg \min_{\theta} \max_{f \in B} \langle \, \alpha(\theta) - \beta \, , \, f \rangle$$

- $B = \{ f \text{ is 1-Lipschitz} \} \implies \text{Loss} = \text{Wasserstein-1 (OT}_0\text{)}$:
  - $S_\varepsilon$ is nearly as efficient as a **closed formula**
  - relevant in **low dimensions**
  - **useless** in $(\mathbb{R}^{512 \times 512}, \| \cdot \|_2)$: the ground cost makes no sense

- $B \simeq \{ f \text{ is 1-Lipschitz} \} \bigcap \{ f \text{ is a } \textbf{CNN} \}$
  $\implies \text{Loss} = \text{Wasserstein GAN}$ :
  - use **perceptually sensible** test functions
  - no simple formula: use **gradient ascent**
  - can we provide relevant **insights** to the ML community?

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
  $\simeq$ a handful of convolutions through the data
  $\rightarrow$ Is it worth it?

Our work:

- Miccai2017 : proof of concept

## Conclusion

Our work:

- Miccai2017 : proof of concept
- ShapeMi2018/AiStats2019 :
  - link with statistics and **computer graphics**
  - reference **implementation** on sparse data
  - theoretical **guarantees**

## Conclusion

Our work:

- Miccai2017 : proof of concept
- ShapeMi2018/AiStats2019 :
  - link with statistics and **computer graphics**
  - reference **implementation** on sparse data
  - theoretical **guarantees**
- 2019 - available soon :
  - unbalanced formulation, to handle **outliers**
  - **evaluation** in varied settings
  - **octree**-like code on the GPU
  - separable **volumetric** implementation

Open questions:

- Can we find a **Brenier-like** formulation for $S_\varepsilon$?

Open questions:

- Can we find a **Brenier-like** formulation for $S_\varepsilon$?
- Link between $S_\varepsilon$ and Sobolev **distances**?

Open questions:

- Can we find a **Brenier-like** formulation for $S_\varepsilon$?
- Link between $S_\varepsilon$ and Sobolev **distances**?
- Proof of convergence for the **multiscale scheme**?

Open questions:

- Can we find a **Brenier-like** formulation for $S_\varepsilon$?
- Link between $S_\varepsilon$ and Sobolev **distances**?
- Proof of convergence for the **multiscale scheme**?
- Interest in the **CVPR/SIGGRAPH** communities?

Thank you for your attention.

Any questions ?

# References

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

📄 Charlier, B., Feydy, J., and Glaunès, J. (2018).
**Kernel operations on the gpu, with autodiff, without memory overflows.**
http://www.kernel-operations.io.
Accessed: 2019-01-20.

📄 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.*

📄 Schmitzer, B. (2016).
**Stabilized sparse scaling algorithms for entropy regularized transport problems.**
*arXiv preprint arXiv:1610.06519.*