# Sinkhorn Divergences for Unbalanced Optimal Transport

---

Thibault Séjourné

Mokameeting – 6th April, 2020

Joint work with Jean Feydy, Francois-Xavier Vialard, Alain Trouvé and Gabriel Peyré

# Outline

# Introduction

# Machine Learning setting

- Given an empirical measure $\beta$,
- And a model $\alpha_\theta$ parametrized by $\theta$.





Shape registration    Supervised Learning    Unsupervised Learning

- Then we optimize via GD & backpropagation a loss $\mathcal{L}$

$$\theta^* \in \arg\min_\theta \mathcal{L}(\alpha_\theta, \beta).$$

Which loss $\mathcal{L}$ should we use to compare probability measures ?

Desired properties of $\mathcal{L}$:

- Positive, definite and convex
- Metrizes the weak* convergence
  $\alpha_n \rightharpoonup \alpha \Leftrightarrow \forall f \in \mathcal{C}(\mathcal{X}), \int f d\alpha_n \to \int f d\alpha.$
- Differentiable

Possible losses between measures:

- Csiszar divergences (KL, TV, Hellinger, etc...)
- Maximum mean discrepancies / kernel norms
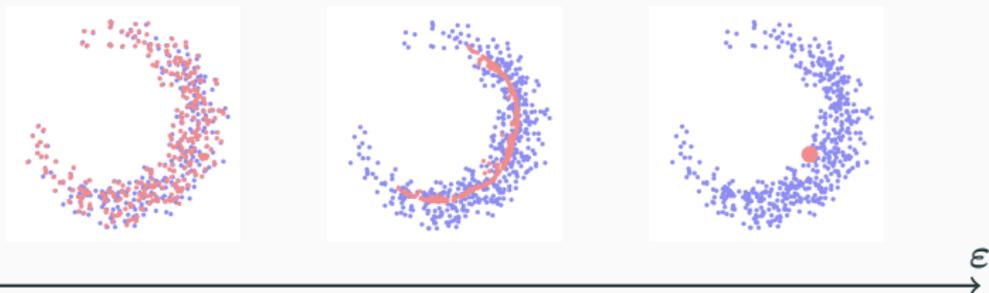- Optimal transport distances

OT issues: non-smooth + complexity + curse of dimensionality

$$\mathrm{OT}_\varepsilon(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\pi \in \mathcal{U}(\alpha,\beta)} \langle \pi, \, \mathrm{C} \rangle + \varepsilon \mathrm{KL}(\pi, \alpha \otimes \beta)$$

Problem: $\mathrm{OT}_\varepsilon$ does not metrize weak* convergence for $\varepsilon > 0$. ☹

$$\exists \alpha \in \mathcal{M}_1^+(\mathcal{X}), \mathrm{OT}_\varepsilon(\alpha, \beta) < \mathrm{OT}_\varepsilon(\beta, \beta).$$



$\varepsilon$

$\Rightarrow$ Debias: $\mathrm{S}_\varepsilon(\alpha, \beta) = \mathrm{OT}_\varepsilon(\alpha, \beta) - \frac{1}{2}\mathrm{OT}_\varepsilon(\alpha, \alpha) - \frac{1}{2}\mathrm{OT}_\varepsilon(\beta, \beta)$
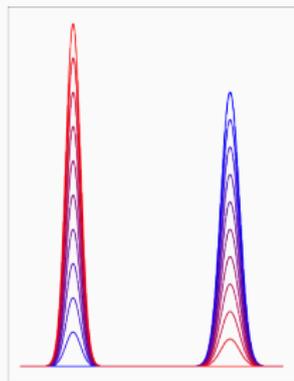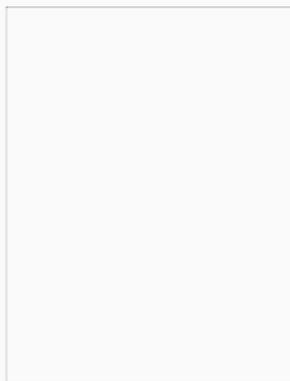
# Unbalanced OT

# Goal of Unbalanced OT

Mitigate between vertical and horizontal geometries on $\mathcal{X}$.
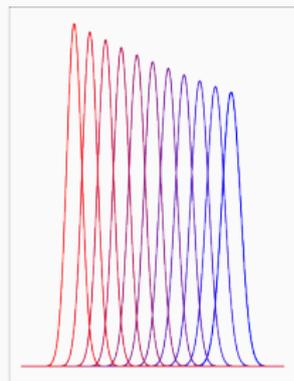
+ avoids normalizing data and geometric outliers.
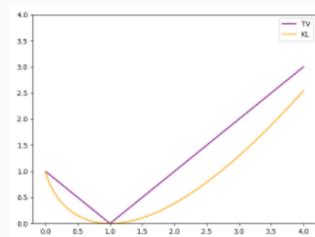


Vertical (Csiszar)        In between ?        Horizontal (OT)

**Definitions [Csiszàr'67]**

- Entropy $\varphi$: nonnegative, l.s.c., convex on $\mathbb{R}_+$ s.t. $\varphi(1) = 0$
- Recession constant: $\varphi'^\infty = \lim_{x \to \infty} \varphi(x)/x$
- Lebesgue decomposition: $\forall (\alpha, \beta), \, \alpha = \frac{\mathrm{d}\alpha}{\mathrm{d}\beta}\beta + \alpha^\top$
- $\varphi$-divergence: $\mathrm{D}_\varphi(\alpha, \beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi(\frac{\mathrm{d}\alpha}{\mathrm{d}\beta})\mathrm{d}\beta + \varphi'^\infty \int_{\mathcal{X}} \mathrm{d}\alpha^\top$
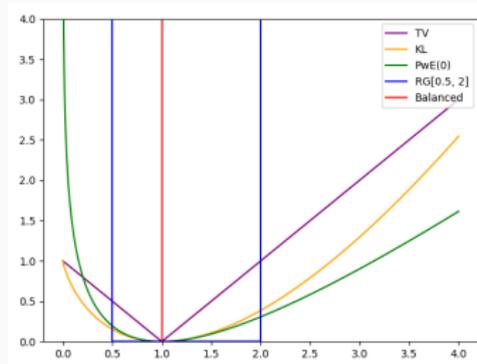
Examples:

- KL: $\varphi(x) = x \log x - x + 1$, $\varphi'^\infty = +\infty$,
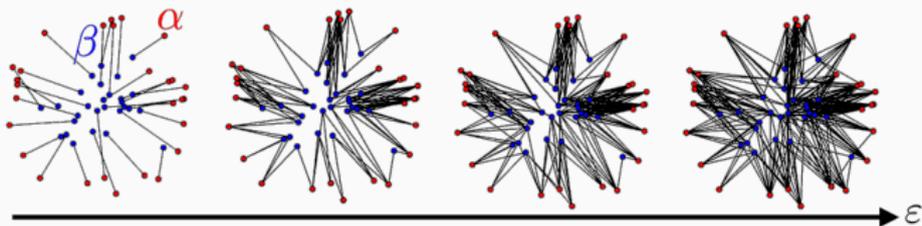- TV: $\varphi(x) = |x - 1|$ and $\varphi'^\infty = 1$.



7

- **Balanced**: $\varphi(x) = \iota_{\{1\}}(x)$ with $D_\varphi(\pi_1, \alpha) = \iota_{(=)}(\pi_1, \alpha)$.
- **TV**: $\varphi(x) = |x - 1|$
- **KL**: $\varphi(x) = x \log x - x + 1$
- **Power entropy**: $\varphi(x) = \frac{1}{p(p-1)}(x^p - p(x - 1) - 1)$, $p \in \mathbb{R}$.
- **Range**: $\varphi(x) = \iota_{[a,b]}(x)$ $(a \leq 1 \leq b)$, i.e $a\alpha \leq \pi_1 \leq b\alpha$.

Entropic Unbalanced OT [Chizat'18]

$$\mathrm{OT}_{\varepsilon,\rho}(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\pi \geq 0} \langle \pi, \mathrm{C} \rangle + \rho \mathrm{D}_\varphi(\pi_1, \alpha) + \rho \mathrm{D}_\varphi(\pi_2, \beta)$$
$$+ \varepsilon \mathrm{KL}(\pi, \alpha \otimes \beta)$$



choice of C and $\mathrm{D}_\varphi$ = priors on geometry and mass dynamics

Reminder: Local mass creation and destruction is allowed

- Shows how $\alpha$ is matched onto $\beta$ and vice versa through $\pi$.
- Plots $\pi_1 \approx \alpha$ and $\pi_2 \approx \beta$
- Input marginals are dashed.

# Sinkhorn divergence

**Definition**

Setting $\mathrm{m}(\mu)$ to be the total mass of the measure $\mu$, we define

$$\mathrm{S}_{\varepsilon,\rho}(\alpha,\beta) \overset{\mathrm{def.}}{=} \mathrm{OT}_{\varepsilon,\rho}(\alpha,\beta) - \tfrac{1}{2}\mathrm{OT}_{\varepsilon,\rho}(\alpha,\alpha) - \tfrac{1}{2}\mathrm{OT}_{\varepsilon,\rho}(\beta,\beta)$$
$$+ \tfrac{\varepsilon}{2}(\mathrm{m}(\alpha) - \mathrm{m}(\beta))^2.$$

It extends the balanced case from [Ramdas '15][Genevay '18].

Impact of KL: entropic bias + mass bias $(\mathrm{m}(\pi) \to \mathrm{m}(\alpha \otimes \beta))$.

**Proposition**

Assuming $\varphi^*$ strictly convex, denoting $\mathrm{k}_\varepsilon \overset{\text{def.}}{=} \mathrm{e}^{-\frac{\mathrm{C}}{\varepsilon}}$ and $\mathrm{f}_\alpha$ and $\mathrm{g}_\beta$ the optimal symmetric potentials of $\mathrm{OT}_\varepsilon(\alpha, \alpha)$ and $\mathrm{OT}_\varepsilon(\beta, \beta)$ respectively, one has

$$\mathrm{S}_{\varepsilon,\rho}(\alpha, \beta) \geq \frac{\varepsilon}{2}\|\alpha \mathrm{e}^{\frac{\mathrm{f}_\alpha}{\varepsilon}} - \beta \mathrm{e}^{\frac{\mathrm{g}_\beta}{\varepsilon}}\|_{\mathrm{k}_\varepsilon}^2.$$

Theorem [S., Feydy, Vialard, Trouve, Peyre '19]

For any Lipschitz cost C s.t. $k_\varepsilon \stackrel{\text{def.}}{=} e^{-\frac{C}{\varepsilon}}$ is a positive universal kernel, for any $\varepsilon > 0$ and strictly convex $\varphi^*$

- $S_{\varepsilon,\rho}$ is convex, positive, definite.
- It is (weakly) differentiable.
- $S_{\varepsilon,\rho}(\alpha, \beta) \to 0 \Leftrightarrow \alpha \rightharpoonup \beta$.

Berg's Theorem: $e^{-C/\varepsilon}$ positive kernel $\Leftrightarrow -C$ positive kernel.

# Sinkhorn algorithm

Writing $\varphi^*(\mathrm{x}) = \sup_{\mathrm{y} \geq 0} \mathrm{xy} - \varphi(\mathrm{y})$, the dual reads

$$\mathrm{OT}_{\varepsilon,\rho}(\alpha, \beta) = \sup_{\mathrm{f,g} \in \mathcal{C}(\mathcal{X})} \langle \alpha, -(\rho\varphi)^*(-\mathrm{f}) \rangle + \langle \beta, -(\rho\varphi)^*(-\mathrm{g}) \rangle$$

$$- \varepsilon \langle \alpha \otimes \beta, \, \mathrm{e}^{\frac{\mathrm{f(x)+g(y)-C(x,y)}}{\varepsilon}} - 1 \rangle$$

**Proposition - Unbalanced Sinkhorn algorithm**

Define the following operators

- (Softmin / LogSumExp) $\mathrm{Smin}_{\alpha}^{\varepsilon} (\mathrm{f}) \overset{\text{def.}}{=} -\varepsilon \log\langle \alpha, \, \mathrm{e}^{-\mathrm{f}/\varepsilon} \rangle$

- (Anisotropic Prox) $\mathsf{aprox}(\mathrm{p}) = \arg\min_{\mathrm{q} \in \mathbb{R}} \varepsilon \mathrm{e}^{(\mathrm{p-q})/\varepsilon} + \varphi^*(\mathrm{q})$

The optimality condition defines the Sinkhorn algorithm

$$\mathrm{g}_{t+1}(\mathrm{y}) = -\mathsf{aprox}( -\mathrm{Smin}_{\alpha}^{\varepsilon} (\mathrm{C}(.,\mathrm{y}) - \mathrm{f}_t) )$$

$$\mathrm{f}_{t+1}(\mathrm{x}) = -\mathsf{aprox}( -\mathrm{Smin}_{\beta}^{\varepsilon} (\mathrm{C}(\mathrm{x},.) - \mathrm{g}_{t+1}) ).$$

## Entropy and Aprox



$$D_\varphi = \iota_{(=)}$$
$$\varphi(\mathrm{x}) = \iota_{\{1\}}(\mathrm{x})$$
$$\mathsf{aprox}(\mathrm{x}) = \mathrm{x}$$

Balanced



15

## Entropy and Aprox



$$D_\varphi = \rho\text{KL}$$
$$\varphi(\text{x}) = \rho(\text{x}\log \text{x} - \text{x} + 1)$$
$$\mathsf{aprox}(\text{x}) = \tfrac{\rho}{\rho+\varepsilon}\,\text{x}$$
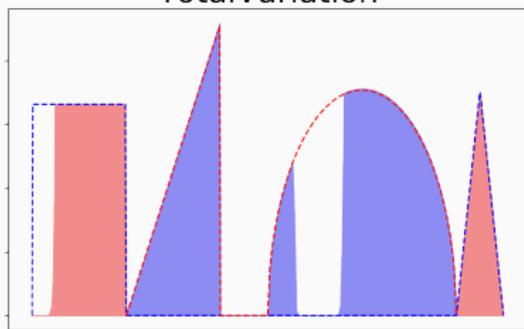
KullbackLeibler

## Entropy and Aprox



$$D_\varphi = \rho TV$$
$$\varphi(x) = \rho|x - 1|$$
$$\text{aprox}(x) = x \text{ if } x \in [-\rho, \rho],\ \rho \text{ if}$$
$$x \geq \rho \text{ and } -\rho \text{ if } x \leq -\rho$$

**TotalVariation**

**Proposition**

One has for any $(f, g) \in \mathcal{C}(\mathcal{X})$

$$|\text{Smin}_{\alpha}^{\varepsilon}(f) - \text{Smin}_{\alpha}^{\varepsilon}(g)| \leq \|f - g\|_{\infty}. \tag{1}$$

$$\left\|\text{aprox}_{\varphi^*}^{\varepsilon}(f) - \text{aprox}_{\varphi^*}^{\varepsilon}(g)\right\|_{\infty} \leq \|f - g\|_{\infty}. \tag{2}$$

$\Rightarrow$ The algorithm is numerically stable. If there exists a fixed point and compactness, the algorithm then converges linearly towards it.

Assume either:

- $\varphi^*$ strictly convex and $\partial\varphi^*$ goes to zero or $+\infty$ as x goes to $0$ or $+\infty$.

- The entropy corresponds to Balanced, TV or Range.

Theorem - Existence and convergence

Assume C is Lipschitz on a compact space $\mathcal{X}$. Then:

1. The space of dual potentials can be restricted to a relatively compact set, thus there is existence of dual maximizers in $\mathcal{C}(\mathcal{X})$.

2. The Sinkhorn algorithm converges towards fixed points which are dual maximizers.

# Numerical illustrations - Gradient flows

Setting adapted from [Chizat '19].

- Position/mass parameterization $x = \{(x_i, r_i)_i\} \in (\mathbb{R}^d \times \mathbb{R})^n$
- Model measure $x \mapsto \alpha(x) = \sum_i^n r_i^2 \delta_{x_i}$
- Cone metric $\langle (x_1, r_1), (x_2, r_2) \rangle_{(x,r)} = \frac{\eta_x}{r^2} \langle x_1, x_2 \rangle_x + \eta_r r_1 r_2$
- Flow $\nabla x(t) = -\nabla_x S_{\varepsilon,\rho}(\alpha(x), \beta)$

**Updates of the coordinates**

$$x_i^{(t+1)} = x_i^{(t)} - \eta_x \nabla_{x_i} S_{\varepsilon,\rho}(\alpha^{(t)}, \beta), \qquad (3)$$

$$r_i^{(t+1)} = r_i^{(t)} . \exp\big( - 2\eta_x \nabla_{r_i} S_{\varepsilon,\rho}(\alpha^{(t)}, \beta)\big) \qquad (4)$$

$\mathrm{OT}_\varepsilon$-KL
$(10^{-3}, 0.3)$

$\mathrm{S}_{\varepsilon,\rho}$-KL
$(10^{-3}, 0.3)$

$\mathrm{S}_{\varepsilon,\rho}$-KL
$(10^{-2}, 0.3)$

## Conclusion

- Family of parametric losses with appealing properties (convexity, differentiability, positivity...)
- Algorithm with linear convergence, compatible with GPU
- Consistent behavior which allows to crossvalidate w.r.t. $\varepsilon$
- Improvement of the statistical complexity, dampening of the curse of dimensionality (Not detailed here)

http://www.kernel-operations.io/geomloss/
https://github.com/thibsej/unbalanced-ot-functionals