

# Sinkhorn Divergences for Unbalanced Optimal Transport

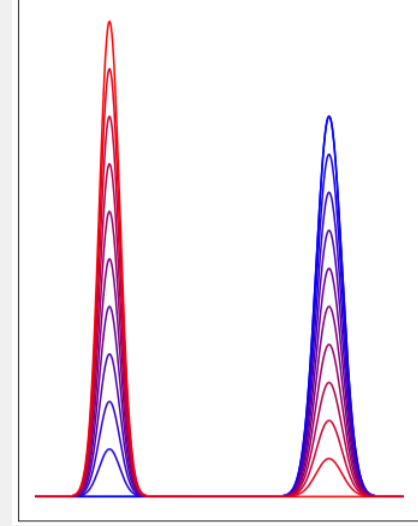
Providing a loss between arbitrary positive measures, fastly computable, with no bias.

Thibault Séjourné<sup>1</sup> Jean Feydy<sup>1,2</sup> François-Xavier Vialard<sup>3</sup> Alain Trounev<sup>2</sup> Gabriel Peyré<sup>1</sup>

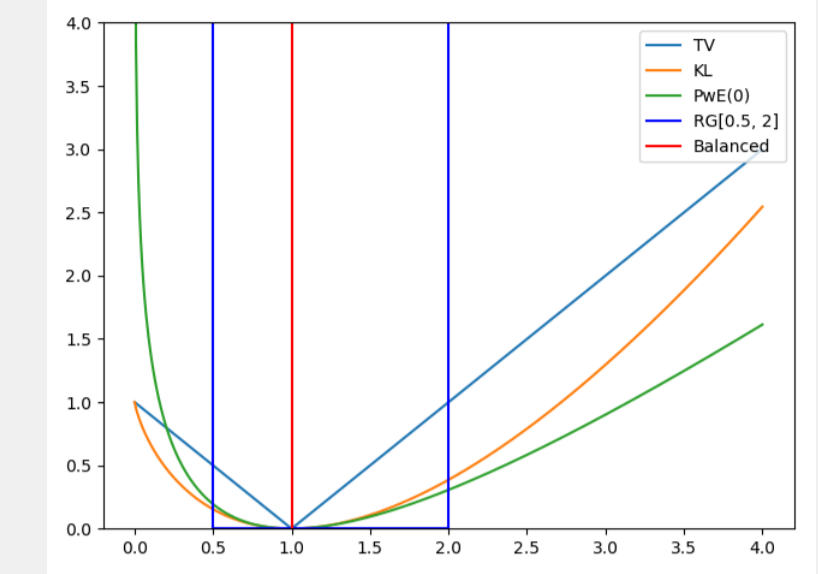
<sup>1</sup>DMA, École Normale Supérieure <sup>2</sup>CMLA, ENS Paris-Saclay <sup>3</sup>LIGM, UPEM

## 1 Csiszar divergences (Csi67) $\approx$ Vertical Geometry

- **Entropy**  $\phi$ : positive, l.s.c., convex on  $\mathbb{R}_+$  s.t.  $\phi(1) = 0$
- **Recession constant**:  $\phi^\infty = \lim_{x \rightarrow \infty} \phi(x)/x$
- **Lebesgue decomposition**:  $\forall(\alpha, \beta), \alpha = \frac{d\alpha}{d\beta} \beta + \alpha^\top$
- **$\phi$ -divergence**:  $D_\phi(\alpha, \beta) \stackrel{\text{def}}{=} \int_{\mathcal{X}} \phi\left(\frac{d\alpha}{d\beta}\right) d\beta + \phi^\infty \int_{\mathcal{X}} d\alpha^\top$



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

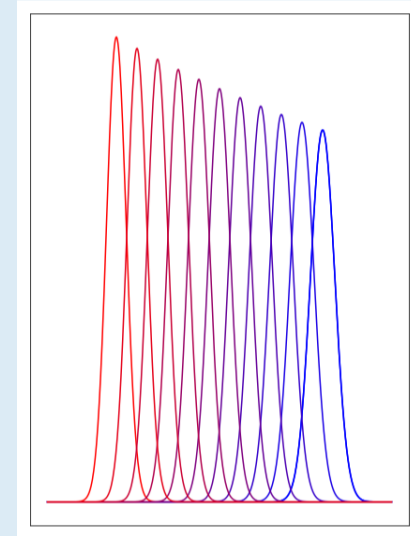


## 2 OT (Kan42) $\approx$ Vertical geometry

Consider a **Compact domain**  $\mathcal{X}$ . Take a cost  $C : (x, y) \mapsto C(x, y)$  continuous, symmetric, Lipschitz (e.g.  $\frac{1}{p} \|x - y\|^p$ ). One defines

$$OT_0(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\pi \geq 0} \{ \langle \pi, C \rangle : \pi_1 = \alpha, \pi^\top 1 = \beta \},$$

where  $\langle \pi, C \rangle \stackrel{\text{def}}{=} \int_{\mathcal{X}^2} C(x, y) d\pi(x, y)$ . OT compares measures by accounting for the geometry. It metrizes the **convergence in law**.

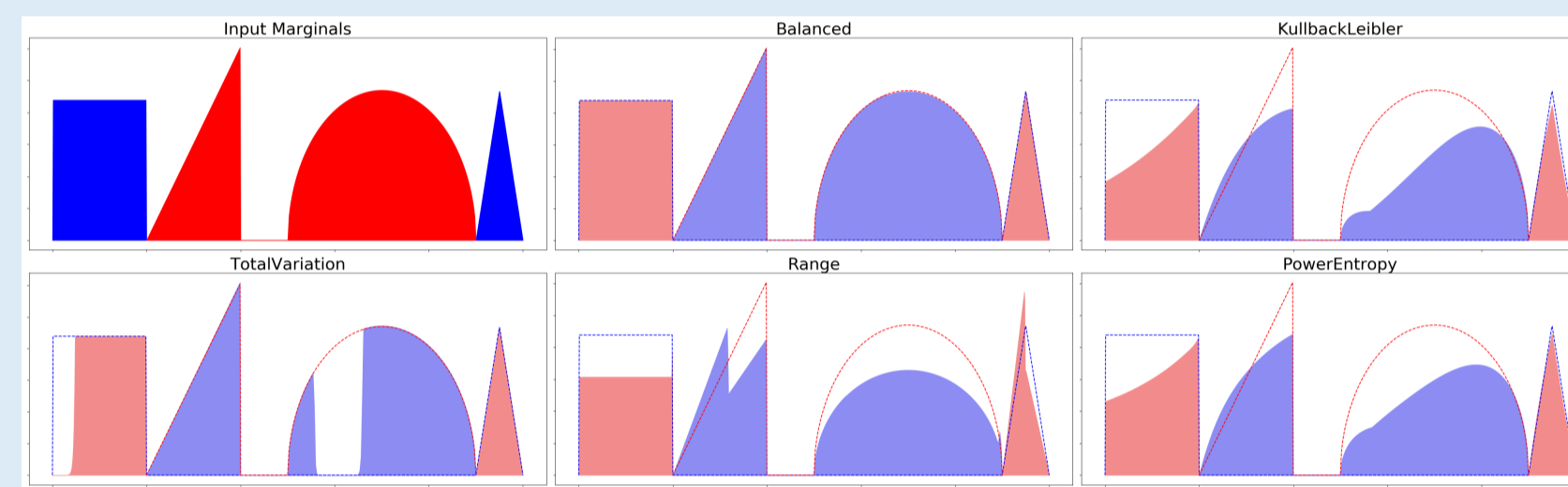


## 3 Unbalanced Optimal Transport (LMS15)

OT only compares normalized measures. Its generalization reads

$$OT_0(\alpha, \beta) \stackrel{\text{def}}{=} \inf_{\pi \geq 0} \langle \pi, C \rangle + \rho D_\phi(\pi_1, \alpha) + \rho D_\phi(\pi_2, \beta).$$

It hybridizes vertical and horizontal geometries.



## 4 Entropic UOT

(U)OT has a complexity  $O(n^3 \log n)$  and is non differentiable. Adding entropy improves both aspects (Sch31; KY94; Cut13; CPSV18).

$$OT_\epsilon(\alpha, \beta) \stackrel{\text{def}}{=} \inf_{\pi \geq 0} \langle \pi, C \rangle + \rho D_\phi(\pi_1, \alpha) + \rho D_\phi(\pi_2, \beta) + \epsilon \text{KL}(\pi, \alpha \otimes \beta)$$

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

$$OT_\epsilon(\alpha, \beta) = \sup_{f, g \in \mathcal{C}(\mathcal{X})} \langle \alpha, -(\rho\phi)^*(-f) \rangle + \langle \beta, -(\rho\phi)^*(-g) \rangle - \epsilon \langle \alpha \otimes \beta, e^{\frac{f(x)+g(y)-C(x,y)}{\epsilon}} - 1 \rangle.$$

## 5 Sinkhorn algorithm

Define the following operators

- **Softmin / LogSumExp**  
 $\text{Smin}_\alpha^\epsilon(f) \stackrel{\text{def}}{=} -\epsilon \log \langle \alpha, e^{-f/\epsilon} \rangle$
- **Anisotropic Prox (CR13)**  
 $\text{aprox}(p) = \arg \min_{q \in \mathbb{R}} \epsilon e^{(p-q)/\epsilon} + \phi^*(q)$

The dual optimality condition defines the Sinkhorn algorithm

$$g_{t+1}(y) = -\text{aprox}(-\text{Smin}_\alpha^\epsilon(C(\cdot, y) - f_t))$$

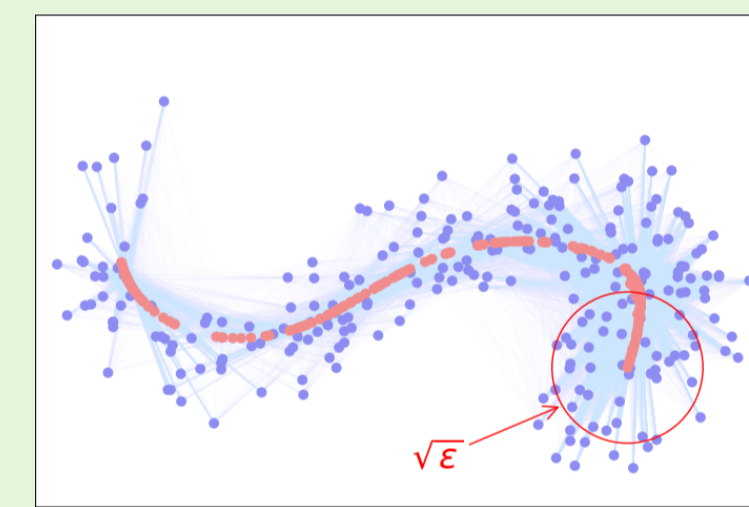
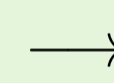
$$f_{t+1}(x) = -\text{aprox}(-\text{Smin}_\beta^\epsilon(C(x, \cdot) - g_{t+1})).$$

**Ex:** (Balanced)  $\text{aprox}(p) = p$ ,  $(\rho \text{KL}) \text{aprox}(p) = \frac{\rho}{\rho+1} p$

## 6 Removing the entropic bias

When  $\epsilon > 0$ , **fuzzy transport plans** induce shrinking artifacts (CR03):

Minimize  $OT_\epsilon(\alpha, \beta)$  with respect to  $\alpha$



$\implies$  Use the **unbiased** Sinkhorn divergence (RTC17; GPC18; SZRM18):

$$S_\epsilon(\alpha, \beta) = OT_\epsilon(\alpha, \beta) - \frac{1}{2} OT_\epsilon(\alpha, \alpha) - \frac{1}{2} OT_\epsilon(\beta, \beta) + \frac{\epsilon}{2} (m(\alpha) - m(\beta))^2,$$

$$\underbrace{OT(\alpha, \beta)}_{\text{Wasserstein}} \xleftarrow{\epsilon \rightarrow 0} \underbrace{S_\epsilon(\alpha, \beta)}_{\text{Easy to compute}} \xrightarrow{\epsilon \rightarrow \infty} \underbrace{MMD-C(\alpha, \beta)}_{\text{Kernel MMD}}$$

## 7 Contributions 1 & 2

**Theorem:** If  $e^{-C(x,y)/\epsilon}$  is a positive definite kernel, then for any strictly convex  $\phi^*$

$$S_\epsilon(\beta, \beta) = 0 \leq S_\epsilon(\alpha, \beta)$$

$$S_\epsilon(\alpha, \beta) = 0 \iff \alpha = \beta$$

$$S_\epsilon(\alpha_n, \beta) \rightarrow 0 \iff \alpha_n \rightarrow \beta$$

**Loss**  $\beta : \alpha \mapsto S_\epsilon(\alpha, \beta)$  is **convex**.

**Theorem:** For all entropies mentioned above, the Sinkhorn algorithm converges linearly towards the optimal dual potentials  $(f, g)$  with a rate independent of the number of sample points.

## 8 Sample Complexity

Set  $(\alpha_n, \beta_n)$  the sampled versions of  $(\alpha, \beta)$  with  $n$  points. In  $\mathbb{R}^d$  one has for **Balanced OT**:

Unregularized OT (Dud69; WB17):  
 $\mathbb{E}_{\alpha \otimes \beta} [ |OT(\alpha, \beta) - OT(\alpha_n, \beta_n)| ] = O(n^{-1/d})$

Compact supports (GCB+18):  
 $\mathbb{E}_{\alpha \otimes \beta} [ |OT_\epsilon(\alpha, \beta) - OT_\epsilon(\alpha_n, \beta_n)| ] = O_{\epsilon \rightarrow 0}(\epsilon^{-d/2} n^{-1/2})$

Subgaussian measures (MW19):  
 $\mathbb{E}_{\alpha \otimes \beta} [ |OT_\epsilon(\alpha, \beta) - OT_\epsilon(\alpha_n, \beta_n)| ] = O(\epsilon^{1-[5d/4+3]} n^{-1/2})$

## 9 Contribution 3

Take positive measures  $(\alpha, \beta)$ , set  $\bar{\alpha} = \alpha/m(\alpha)$  and  $\bar{\beta} = \beta/m(\beta)$ . Write

$$\alpha_n = \frac{m(\alpha)}{n} \sum_{i=1}^n \delta_{X_i}, \quad \beta_n = \frac{m(\beta)}{n} \sum_{i=1}^n \delta_{Y_i},$$

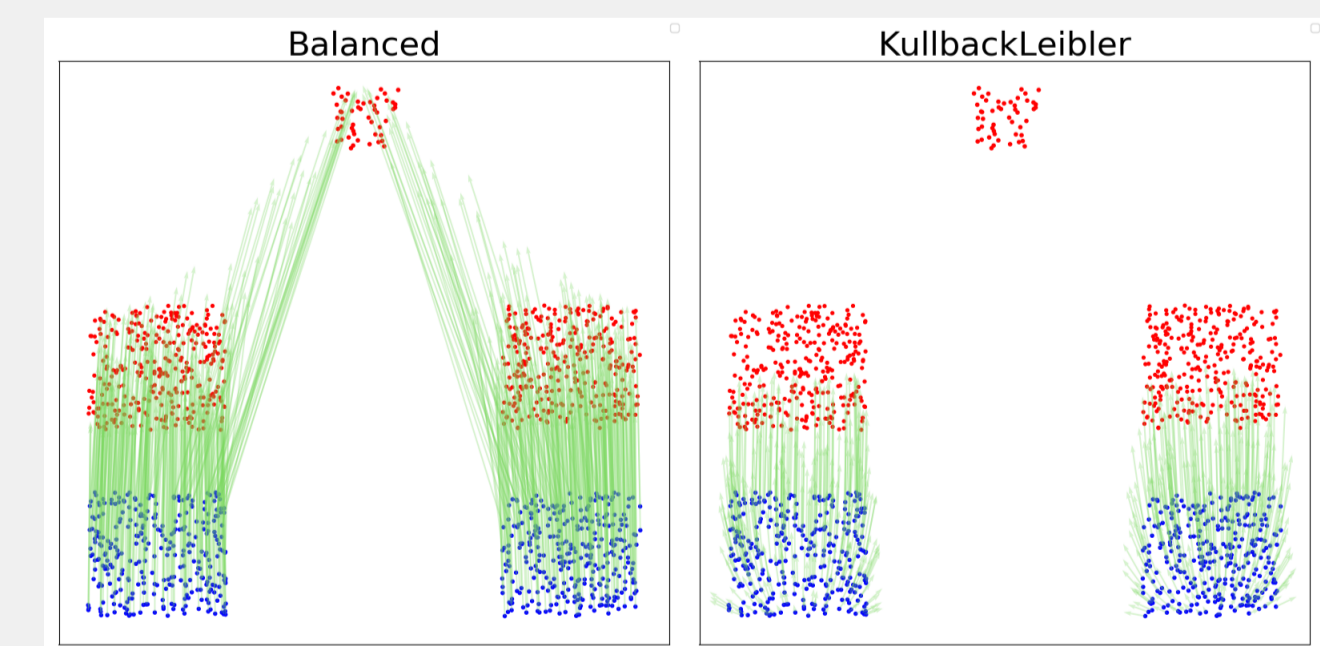
**Theorem:** Assume smooth  $C$  and  $\phi^*$ , assume  $\phi^*$  is strictly convex. Then

$$\mathbb{E}_{\bar{\alpha} \otimes \bar{\beta}} [ |OT_\epsilon(\alpha, \beta) - OT_\epsilon(\alpha_n, \beta_n)| ] = O_{\epsilon \rightarrow 0}(\epsilon^{-d/2} n^{-1/2}) = O_{\epsilon \rightarrow \infty}(n^{-1/2}).$$

Furthermore the rate is linear in  $m(\alpha) + m(\beta)$ .

## 10 Numerical Highlight

Unbalanced OT allows to **discard geometric outliers!**



### Notation

**Banach duality:**  $f \in \mathcal{C}(\mathcal{X}), \alpha \in \mathcal{M}_+(\mathcal{X})$

**Dual Bracket:**  $\langle \alpha, f \rangle = \int_{\mathcal{X}} f d\alpha = \mathbb{E}_\alpha[f]$

**Discrete Encoding:**  $(\alpha_i)_i \in \mathbb{R}^N, (x_i)_i \in \mathbb{R}^{N \times D}, (f_i)_i \in \mathbb{R}^N$

**Discrete Setting:**  $\alpha = \sum_i \alpha_i \delta_{x_i} \Rightarrow \langle \alpha, f \rangle = \sum_i \alpha_i f_i$

### References

- [CPSV18] L. Chizat, G. Peyré, B. Schmitzer, and F.-X. Vialard. Scaling algorithms for unbalanced transport problems. to appear in Mathematics of computation, 2018.
- [CR03] H. Chui and A. Rangarajan. A new point matching algorithm for non-rigid registration. Computer Vision and Image Understanding, 80(3):141–161, 2003.
- [CR13] P. L. Combettes and N. N. Reyes. Moreau's decomposition in banach spaces. Mathematical Programming, 116(3):151–174, 2013.
- [Cui07] I. Csiszar. Information-type measures of difference of probability distributions and indirect observation. Studia Scientiarum Mathematicarum Hungarica, 22:29–38, 1967.
- [Cut13] M. Cuturi. Sinkhorn distances: Light-speed computation of optimal transport. In Adv. In Neural Information Processing Systems, pages 2292–2300, 2013.
- [Dud69] R. M. Dudley. The speed of mean Glivenko-Cantelli convergence. The Annals of Mathematical Statistics, 40(1):40–50, 1969.
- [Fey17] J. Feydy, F.-X. Vialard, S.-A. Amaral, A. Trounev, and G. Peyré. Interpolating between optimal transport and mmd using Sinkhorn divergences. arXiv preprint arXiv:1810.08276, 2018.
- [GPC18] A. Genevay, L. Chizat, F. Bach, M. Cuturi, and G. Peyré. Sample complexity of Sinkhorn divergences. arXiv preprint arXiv:1810.07213, 2018.
- [Kant02] L. Kantorovich. On the transfer of masses in Russian. Doklady Akademii Nauk, 375:227–229, 1942.
- [Kra94] J. Kratochvíl and K. L. Taylor. The invisible hand algorithm: Solving the assignment problem with statistical physics. Neural networks, 7(3):47–60, 1994.
- [LMS15] M. Lévy, A. Mielke, and G. Savaré. Optimal entropic transport problems and a new Hellinger-Kantorovich distance between positive measures. Inventiones mathematicae, pages 1143–2015.
- [Mén19] G. Ména and J. Weed. Statistical bounds for entropic optimal transport: sample complexity and the central limit theorem. arXiv preprint arXiv:1901.11882, 2019.
- [RT17] A. Ramos, N. G. Trilakis, and M. Cuturi. On Wasserstein two-sample testing and related families of nonparametric tests. Entropy, 19(2), 2017.
- [Sch31] E. Schrödinger. Über die Umkehrung der Naturgesetze. Sitzungsberichte Preuss. Akad. Wiss. Berlin, Phys. Math., 143:144–153, 1931.
- [SZRM18] T. Saitoh, H. Zhang, A. Radford, and D. Metzger. Improving gans using optimal transport. arXiv preprint arXiv:1802.05753, 2018.
- [WB17] J. Weed and F. Bach. Sharp asymptotic and finite-sample rates of convergence of empirical measures in Wasserstein distance. arXiv preprint arXiv:1705.06868, 2017.
- [Wil69] A. G. Wilson. The use of entropy maximizing models, in the theory of trip distribution, mode split and route split. Journal of Transport Economics and Policy, pages 108–136, 1969.

Check the repos at:

[www.github.com/thibsej/unbalanced-ot-functionals](https://www.github.com/thibsej/unbalanced-ot-functionals)  
[www.kernel-operations.io/geomloss](https://www.kernel-operations.io/geomloss)