

INTERPOLATING BETWEEN OPTIMAL TRANSPORT & KL REGULARIZED OPTIMAL TRANSPORT WITH RÉNYI DIVERGENCES

joint work with



Jonas Bresch, TU Berlin

University of South Carolina, Columbia, 12.09.2024.

Graduate Colloquium (Alec Helm, Jonah Klein).

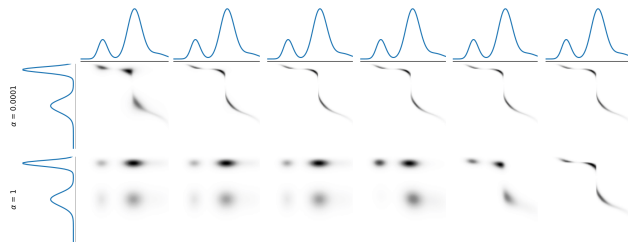
Optimal transport (OT) **distance on probability measures** via **transport plan**.

Problem: $O(N^3)$ for N samples.

Solution: Entropic OT (Cuturi, NeurIPS'13): add ε times KL-regularizer to OT problem for $\varepsilon > 0$.

Sinkhorn algorithm $\rightsquigarrow O(N^{1+\frac{1}{d}} \ln(N))$.

Problem in practice: need ε very small to get accurate plan, but \rightsquigarrow numerical instabilities.



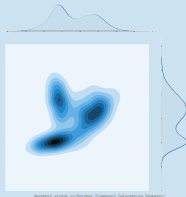
Our solution: Add instead ε times different ($=\alpha$ -Rényi) regularizer and let $\alpha \searrow 0$ instead of $\varepsilon \searrow 0$.

1. Tsallis divergence and α -Rényi divergence



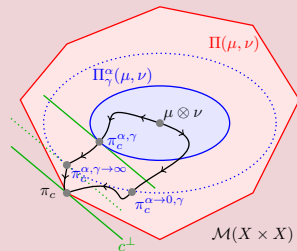
<https://arxiv.org/abs/1405.0487v2>
https://www.wikipedia.org/wiki/Convergence_of_information_theory

2. Optimal transport and its regularization



Source: J. G. L. (2019). Optimal Transport, Information Geometry

3. Rényi-regularized OT



4. Dual formulation

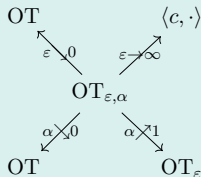


Polyan, Oshin, 2020

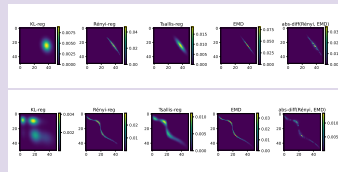
$$\min_{\pi \in \mathcal{P}(X^2)} \langle c, \pi \rangle + \varepsilon R_\alpha(\pi)$$

$$\max_{h \in \mathcal{C}(X^2)} \langle h, \pi \rangle - \varepsilon \ln(\gamma_h^\alpha)$$

5. Interpolation properties



6. Numerical results



DEFINITION (α -RÉNYI DIVERGENCE)

The α -Rényi divergence of order $\alpha \in (0, 1)$ is

$$R_\alpha: \mathcal{P}(X) \times \mathcal{P}(X) \rightarrow [0, \infty], \quad (\mu \mid \nu) \mapsto \frac{1}{\alpha - 1} \ln \left(\int_X \left(\frac{\rho_\mu(x)}{\rho_\nu(x)} \right)^\alpha d\nu(x) \right).$$

where for $\sigma \in \mathcal{P}(X)$, ρ_σ is the density w.r.t. $\frac{1}{2}(\mu + \nu)$, and $\ln(0) := -\infty$.

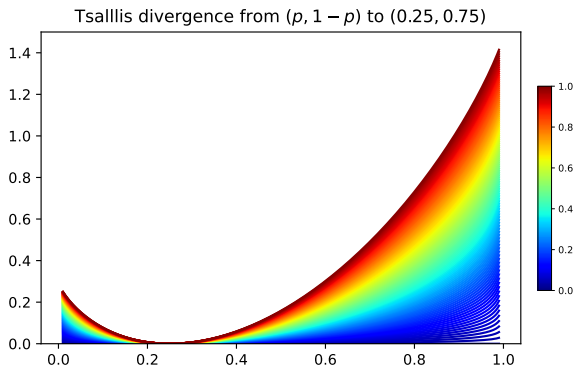
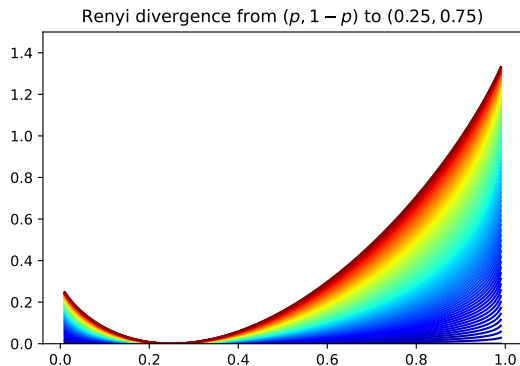
Muzellec et. al (AAAI 2017) examine Tsallis-regularized OT.

DEFINITION (q -TSALLIS DIVERGENCE)

The q -Tsallis divergence of order $q > 0$, $q \neq 1$, is

$$T_q = \frac{1}{q - 1} [\exp((q - 1)R_q) - 1]: \mathcal{P}(X) \times \mathcal{P}(X) \rightarrow [0, \infty], \quad (\mu \mid \nu) \mapsto \frac{1}{q - 1} \left[\int_X \left(\frac{\rho_\mu(x)}{\rho_\nu(x)} \right)^q d\nu(x) - 1 \right]$$

Tsallis = 1st order approximation of Rényi since $\ln(y) \approx y - 1$ (1st order Taylor).



THEOREM (PROPERTIES OF THE RÉNYI DIVERGENCE)

- **Divergence property:** $R_\alpha(\mu \mid \nu) \geq 0$ and $R_\alpha(\mu \mid \nu) = 0$ if and only if $\mu = \nu$.
- R_α is **nondecreasing** and **continuous** in $\alpha \in [0, 1]$ with $\lim_{\alpha \nearrow 1} R_\alpha = \text{KL}$ pointwise.
- R_α **jointly convex**, **jointly weakly lower semicontinuous** for $\alpha \in (0, 1]$.

Let (X, d) metric space, with d lower semicontinuous.

Let $p \in [1, \infty)$, $\mathcal{P}(X)$ the set of probability measures.

$$\mathcal{P}_p(X) := \left\{ \mu \in \mathcal{P}(X) : \int_X d(x, x_0)^p d\mu(x) < \infty \right\}, \quad x_0 \in X.$$

On $\mathcal{P}_p(X)$, the **Wasserstein- p metric** is

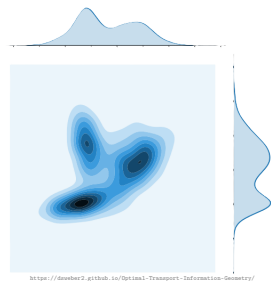
$$\text{OT}(\mu, \nu)^p = \min_{\pi \in \Pi(\mu, \nu)} \int_{X \times X} d(x, y)^p d\pi(x, y), \quad \mu, \nu \in \mathcal{P}_p(X),$$

where the **transport polytope** is

$$\Pi(\mu, \nu) := \{ \pi \in \mathcal{P}(X \times X) : \pi(A \times X) = \mu(A), \pi(X \times A) = \nu(A) \forall A \}$$

The product measure $\mu \otimes \nu \in \Pi(\mu, \nu)$.

Notation: $\langle f, \mu \rangle := \int_X f(x) d\mu(x)$, so we can write $\text{OT}(\mu, \nu)^p = \min \{ \langle d^p, \pi \rangle : \pi \in \Pi(\mu, \nu) \}$.



Regularizer: **Kullback-Leibler divergence**

$$\text{KL}(\cdot \mid \mu \otimes \nu) : \Pi(\mu, \nu) \rightarrow [0, \infty),$$

$$\pi \mapsto \int_{X \times X} \ln \left(\frac{d\pi}{d\mu \otimes \nu}(x, y) \right) d\mu(x) d\nu(y)$$

KL-regularized OT:

$$\text{OT}_\varepsilon(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \langle d^P, \pi \rangle + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu)$$

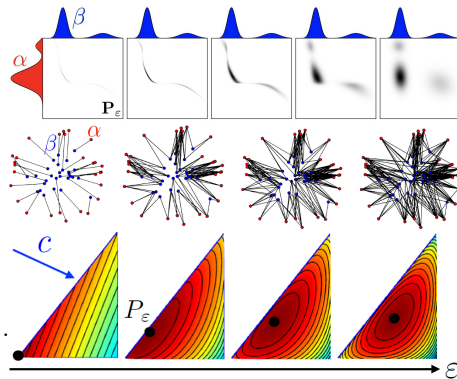
$$= \max_{f, g \in \mathcal{C}(X)} \left\langle f \oplus g - \varepsilon \exp \left(-\frac{1}{\varepsilon} (f \oplus g - d^P) \right), \mu \otimes \nu \right\rangle.$$

$$(f \oplus g)(x, y) := f(x) + g(y) \text{ for } f, g \in \mathcal{C}(X).$$

Primal-dual relation: $\hat{\pi}^\varepsilon = \exp \left(\frac{\hat{f} \oplus \hat{g} - d^P}{\varepsilon} \right) \cdot \mu \otimes \nu$

$$\text{argmin} \{ \text{KL}(\pi \mid \mu \otimes \nu) : \langle d^P, \pi \rangle = \text{OT}(\mu, \nu) \} \xleftarrow{\varepsilon \searrow 0} \hat{\pi}_\varepsilon \xrightarrow{\varepsilon \rightarrow \infty} \mu \otimes \nu$$

$$\text{OT}(\mu, \nu) \xleftarrow{\varepsilon \searrow 0} \text{OT}_\varepsilon(\mu, \nu) \xrightarrow{\varepsilon \rightarrow \infty} \langle d^P, \mu \otimes \nu \rangle$$



Here, $c = d^P$. ©G. Pérye, M. Cuturi, 2019

Discretize $X \approx (x_i)_{i=1}^N$

$\mu, \nu \in \mathcal{P}(X)$ become vectors $\mathbf{r} := (\mu(x_i))_{i=1}^N, \mathbf{c} := (\nu(x_i))_{i=1}^N \in \Sigma_N$,

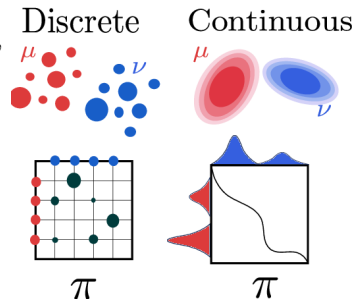
where

$$\Sigma_N := \{x \in [0, 1]^N : \sum_{i=1}^N x_i = 1\}.$$

cost matrix: $\mathbf{M} := (d(x_i, x_j)^p)_{i,j=1}^N$.

transport polytope

$$\Pi(\mathbf{r}, \mathbf{c}) := \{\mathbf{P} \in \Sigma_{N \times N} : \mathbf{P} \mathbf{1}_N = \mathbf{r}, \mathbf{P}^T \mathbf{1}_N = \mathbf{c}\}$$



©T. Vayer.

Optimal transport plan as KL-projection of Gibbs kernel

$$\hat{\mathbf{P}}^\varepsilon = \operatorname{argmin}_{\mathbf{P} \in \Pi(\mathbf{r}, \mathbf{c})} \operatorname{KL} \left(\mathbf{P} \mid \exp \left(\frac{-\mathbf{M}}{\varepsilon} \right) \right)$$

Sinkhorn algorithm finds this projection via matrix scaling.

For regularization parameter $\gamma \in [0, \infty]$ and $\alpha \in (0, 1)$, the **restricted transport polytope**,

$$\Pi_\gamma^\alpha(\mu, \nu) := \{\pi \in \Pi(\mu, \nu) : R_\alpha(\pi \mid \mu \otimes \nu) \leq \gamma\},$$

is weakly compact, since $R_\alpha(\cdot \mid \mu \otimes \nu)$ is weakly lsc. and $\Pi(\mu, \nu)$ is weakly compact.

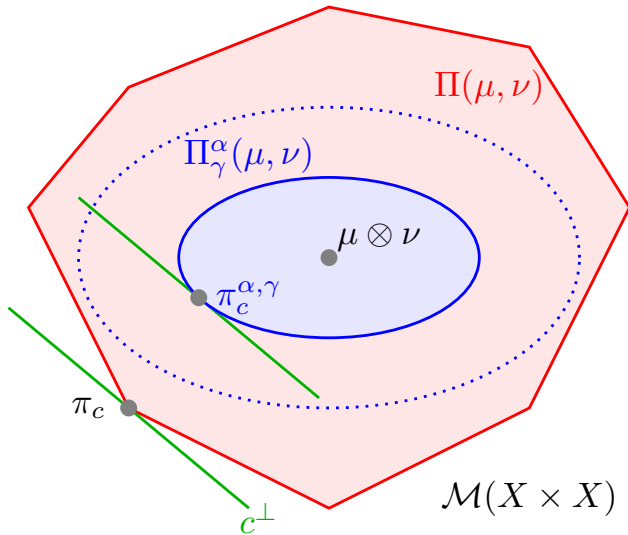
DEFINITION (RÉNYI-SINKHORN DISTANCE)

The **Rényi-Sinkhorn distance** between $\mu, \nu \in \mathcal{P}_p(X)$ is

$$d_{\gamma, \alpha} : \mathcal{P}_p(X) \times \mathcal{P}_p(X) \rightarrow \mathbb{R}, \quad (\mu, \nu) \mapsto \min \left\{ \langle d^p, \pi \rangle^{\frac{1}{p}} : \pi \in \Pi_\gamma^\alpha(\mu, \nu) \right\}. \quad (1)$$

THEOREM (BRESCH, S. '24)

- For $(\mu, \nu) \in \mathcal{P}_p(X)$, the optimization problem (1) is **convex** and has a **unique** minimizer.
- $\mathcal{P}_p(X)^2 \ni (\mu, \nu) \mapsto \mathbb{1}_{[\mu \neq \nu]}(\mu, \nu) d_{\gamma, \alpha}(\mu, \nu)$ is a **metric** for $\alpha \in (0, 1)$, $\gamma \in [0, \infty]$.



Transport polytope $\Pi(\mu, \nu)$, restricted transport polytope $\Pi_\gamma^\alpha(\mu, \nu)$ for $c = d^p$.
 (Plot inspired by (Cuturi, 2013).)

Instead of restricting the problems domain, penalize the Rényi divergence constraint in (1).

DEFINITION (DUAL RÉNYI-DIVERGENCE-SINKHORN DISTANCE)

The **dual Rényi-Divergence-Sinkhorn distance** for $\alpha \in (0, 1)$, $\varepsilon \in [0, \infty)$ is

$$d^{\alpha, \varepsilon} : \mathcal{P}_p(X) \times \mathcal{P}_p(X) \rightarrow \mathbb{R}, \quad (\mu, \nu) \mapsto \langle d^p, \pi^{\alpha, \varepsilon}(\mu, \nu) \rangle^{\frac{1}{p}},$$

where $\pi^{\alpha, \varepsilon}(\mu, \nu) \in \operatorname{argmin} \{ \langle d^p, \pi \rangle + \varepsilon R_\alpha(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \}.$ (2)

THEOREM (LAGRANGIAN POINT OF VIEW AND PRE-METRIC [BRESCH, S. '24])

Let $(\mu, \nu) \in \mathcal{P}_p(X).$

- The optimization problem (2) is **convex** and has a **unique** minimizer.
- Rényi-Sinkhorn $d_{\gamma, \alpha}(\mu, \nu)$ and dual Rényi-Sinkhorn $d^{\alpha, \lambda}(\mu, \nu)$ are **equivalent**:

for $\gamma > 0$, there exists $\varepsilon \in [0, \infty]$, such that $\langle d^p, \pi^{\alpha, \varepsilon}(\mu, \nu) \rangle = d_{\gamma, \alpha}(\mu, \nu)^p.$

DEFINITION (RÉNYI-REGULARIZED OT [BRESCH, S. '24])

The **Rényi-regularized OT** problem is

$$\text{OT}_{\varepsilon, \alpha}: \mathcal{P}_p(X) \times \mathcal{P}_p(X) \rightarrow [0, \infty), (\mu, \nu) \mapsto \min_{\pi \in \Pi(\mu, \nu)} \langle c, \pi \rangle + \varepsilon R_{\alpha}(\pi \mid \mu \otimes \nu).$$

THEOREM ($\text{OT}_{\varepsilon, \alpha}$ IS A PRE-METRIC [BRESCH, S. '24])

$\mathcal{P}_p(X)^2 \ni (\mu, \nu) \mapsto \mathbb{1}_{[\mu \neq \nu]} \text{OT}_{\varepsilon, \alpha}(\mu, \nu)$ is a metric for $\alpha \in (0, 1), \varepsilon \in [0, \infty)$.

LEMMA (MONOTONICITY OF RÉNYI REGULARIZED OT [BRESCH, S. '24])

Let $\mu, \nu \in \mathcal{P}_p(X)$, $\alpha, \alpha' \in (0, 1)$ and $\varepsilon, \varepsilon' \geq 0$ with $\alpha > \alpha'$ and $\varepsilon < \varepsilon'$. Then, we have

$$\text{OT}_{\varepsilon', \alpha}(\mu, \nu) \geq \text{OT}_{\varepsilon, \alpha}(\mu, \nu) \geq \text{OT}_{\varepsilon, \alpha'}(\mu, \nu).$$

From now on: **X compact**. The dual space of all finite signed Borel measures on X , $\mathcal{M}(X)$, is $\mathcal{C}(X)$, the space of real-valued continuous functions on X .

Recall $(f \oplus g)(x, y) := f(x) + g(y)$.

THEOREM (DUAL PROBLEM, DUAL REPRESENTATION [BRESCH, S. '24])

We have the strong duality

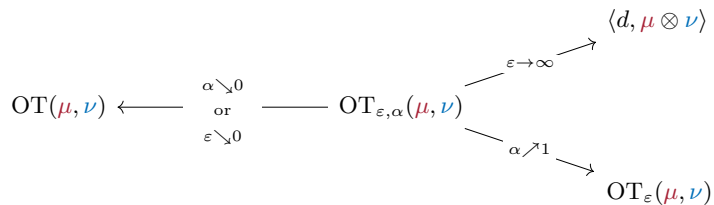
$$\text{OT}_{\varepsilon,\alpha}(\mu, \nu) = \max_{\substack{f, g \in \mathcal{C}(X) \\ f \oplus g \leq d}} \langle f \oplus g, \mu \otimes \nu \rangle - \varepsilon \ln \left(\left\langle (d - f \oplus g)^{\frac{\alpha}{\alpha-1}}, \mu \otimes \nu \right\rangle \right) + C_{\alpha,\lambda}. \quad (3)$$

The optimal dual potentials $\hat{f}, \hat{g} \in \mathcal{C}(X)$ from (3) are unique $\text{supp}(\mu \otimes \nu)$ -a.e. up to additive constants and the unique optimal plan is

$$\pi^{\alpha,\varepsilon} \propto (d - \hat{f} \oplus \hat{g})^{\frac{1}{\alpha-1}} \cdot (\mu \otimes \nu).$$

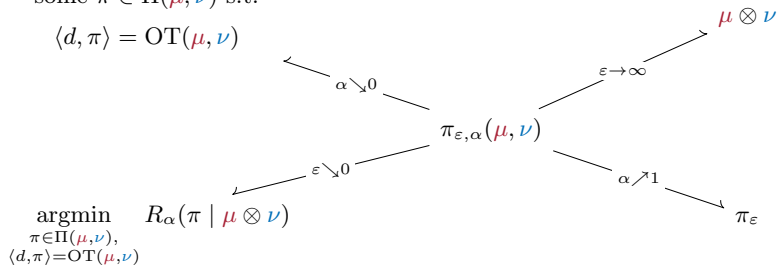
Proof idea. Use Fenchel-Rockafellar theorem, extend objective to $\mathcal{M}(X) \times \mathcal{M}(X)$ by ∞ .

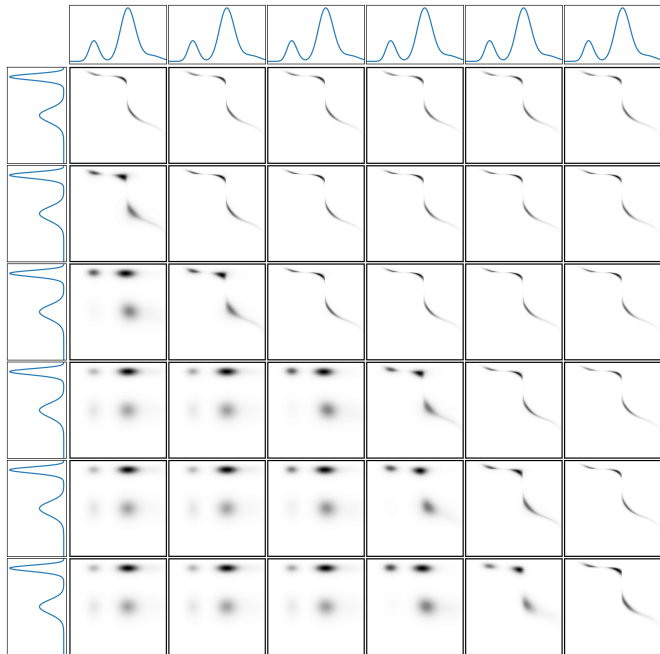
THEOREM (INTERPOLATION PROPERTIES).



some $\pi \in \Pi(\mu, \nu)$ s.t.

$$\langle d, \pi \rangle = \text{OT}(\mu, \nu)$$





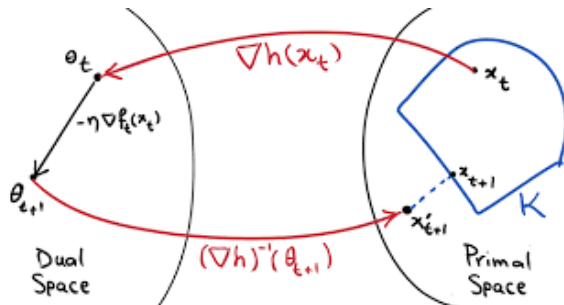
Solve

$$\min_{x \in K} f(x), \quad \text{where } K \subset \mathbb{R}^n \text{ compact.}$$

via the updates

$$x^{(k+1)} = \operatorname{argmin}_{y \in K} D_h \left(y \mid (\nabla h)^{-1} \left(\nabla h(x^{(k)}) - \eta_k \nabla f(x^{(k)}) \right) \right), \quad x^{(0)} \in K, \eta_k > 0, \quad (4)$$

for a convex function $h: \mathbb{R}^n \rightarrow \mathbb{R}$ with special properties.



Choose $K = \Sigma_N$ (probability simplex), $-h = \text{Shannon entropy} \implies D_h = \text{KL}$.

Rényi-regularized OT objective

$$\Pi(\mathbf{r}, \mathbf{c}) \rightarrow [0, \infty), \quad \mathbf{P} \mapsto \langle \mathbf{M}, \mathbf{P} \rangle + \varepsilon R_\alpha(\mathbf{P} \mid \mathbf{r}\mathbf{c}^\text{T}).$$

is not Lipschitz continuous, but locally Lipschitz on

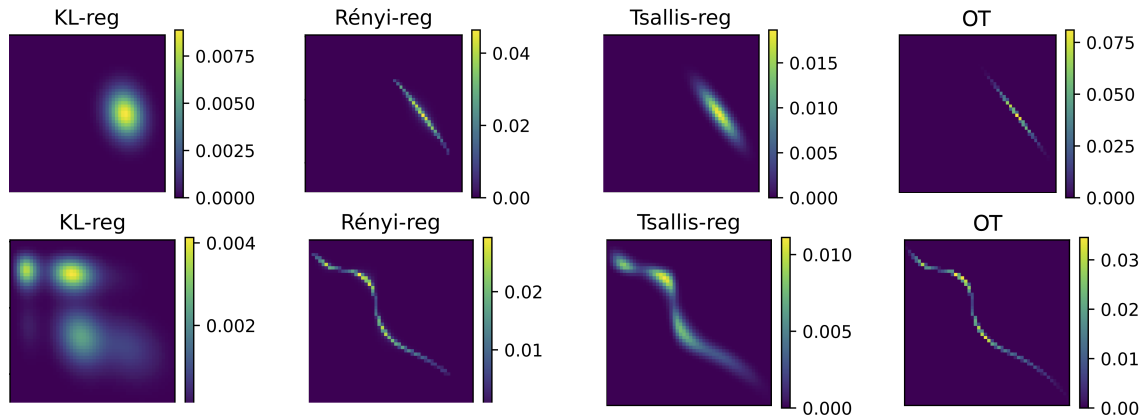
$$\{\mathbf{P} \in \Pi(\mathbf{r}, \mathbf{c}) : \mathbf{P}|_{\text{supp}(\mathbf{r} \otimes \mathbf{c})} > 0\} = \Pi(\mathbf{c}, \mathbf{r}) \cap \mathbb{R}_{>0}^N,$$

which suffices for convergence of a mirror descent with **special step size** $(\eta_k)_{k \in \mathbb{N}}$ (You, Li, 2022).

In each iteration one KL projection onto Σ_N (using Sinkhorn algorithm) is performed:

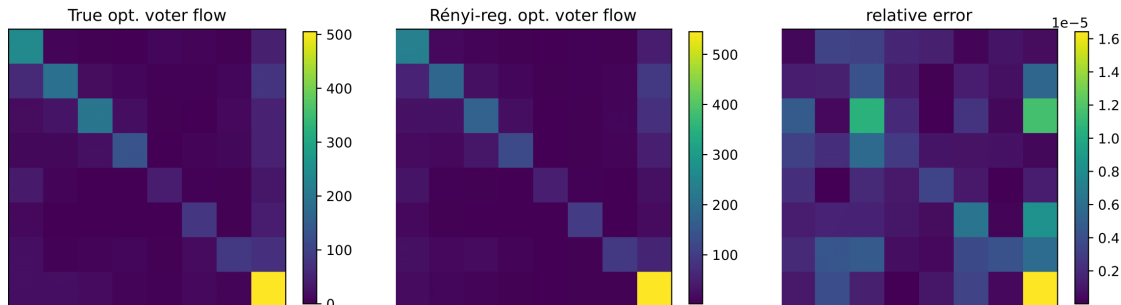
$$\mathbf{P}^{(k)} \leftarrow \text{Sinkhorn} \left(\mathbf{P}^{(k-1)} \odot \exp \left(-\eta_k \mathbf{M} - \frac{\eta_k}{\lambda} \frac{\alpha}{\alpha - 1} \frac{(\mathbf{r}\mathbf{c}^\text{T} \oslash \mathbf{P})^{1-\alpha}}{\langle \mathbf{P}^\alpha, (\mathbf{r}\mathbf{c}^\text{T})^{1-\alpha} \rangle} \right); \mathbf{r}, \mathbf{c} \right), \quad k \in \mathbb{N}$$

RÉNYI REGULARIZATION YIELDS MORE ACCURATE PLANS



Regularized OT plans for Gaussian (*top*) and Poisson (*bottom*) marginals with regularization parameter $\lambda = 10$, Rényi order $\alpha = 0.01$, Tsallis order: $q = 2$.

NUMERICAL EXPERIMENTS - PREDICTING VOTER MIGRATION



regularizer, $\varepsilon = 1$	abs error \pm std	KL error	mean squared error
KL	$2.4221 \times 10^1 \pm 2.848 \times 10^1$	8.422×10^2	9.008×10^4
Tsallis	$9.409 \pm 1.529 \times 10^1$	3.173×10^2	2.063×10^4
OT	$1.845 \times 10^1 \pm 2.358 \times 10^1$	7.655×10^2	5.738×10^4
$\frac{3}{10}$ -Rényi	6.611 ± 7.868	2.128×10^2	6.759×10^3

- **Contribution.** Regularize optimal transport problem using the α -Rényi-divergences R_α for $\alpha \in (0, 1)$. Prove dual formulation and interpolation properties.
- **Prior work.** Regularization with $\text{KL} = \lim_{\alpha \nearrow 1} R_\alpha$ and with q -Tsallis divergence
- **Method.** Solve primal problem with mirror descent and dual problem with subgradient descent.
- **Result.** Rényi-regularized OT plans outperform KL / Tsallis regularized OT plans on real and synthetic data.
- **Novelty.** $R_\alpha \notin \{f\text{-divergence, Bregman divergence}\}$ and R_α not “separable” due to the logarithm.

Thank you for your attention!

I am happy to take any questions.

Paper link: <https://arxiv.org/abs/2404.18834>

My website: <https://viktorajstein.github.io>

- [BT03] Amir Beck and Marc Teboulle, *Mirror descent and nonlinear projected subgradient methods for convex optimization*, Oper. Res. Lett. **31** (2003), no. 3, 167–175.
- [Cut13] Marco Cuturi, *Sinkhorn distances: lightspeed computation of optimal transport*, Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2 (Red Hook, NY, USA), NIPS’13, Curran Associates Inc., 2013, p. 2292–2300.
- [MNPN17] Boris Muzellec, Richard Nock, Giorgio Patrini, and Frank Nielsen, *Tsallis regularized optimal transport and ecological inference*, Proceedings of the AAAI conference on Artificial Intelligence (Hilton San Francisco, San Francisco, California, USA), vol. 31, 2017.
- [NS21] Sebastian Neumayer and Gabriele Steidl, *From optimal transport to discrepancy*, Handbook of Mathematical Models and Algorithms in Computer Vision and Imaging: Mathematical Imaging and Vision (2021), 1–36.
- [NY83] Arkadij Semenovič Nemirovskij and David Borisovich Yudin, *Problem complexity and method efficiency in optimization*, Wiley, New York, 1983.

- [PC19] Gabriel Peyré and Marco Cuturi, *Computational optimal transport*, Found. Trends Mach. Learn. **11** (2019), no. 5-6, 355–607.
- [Rén61] Alfréd Rényi, *On measures of entropy and information*, Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics (Statistical Laboratory of the University of California, Berkeley, California, USA), vol. 4, University of California Press, 1961, pp. 547–562.
- [Tsa88] Constantino Tsallis, *Possible generalization of boltzmann-gibbs statistics*, Journal of statistical physics **52** (1988), 479–487.
- [vEH14] Tim van Erven and Peter Harremos, *Rényi divergence and Kullback-Leibler divergence*, IEEE Trans. Inf. Theory **60** (2014), no. 7, 3797–3820.

$$\text{OT}_{\varepsilon,\alpha}(\mu, \mu) \neq 0$$

To obtain valid, differentiable distance:

$$D_{\varepsilon,\alpha}(\mu, \nu) := \text{OT}_{\varepsilon,\alpha}(\mu, \nu) - \frac{1}{2} \text{OT}_{\varepsilon,\alpha}(\mu, \mu) - \frac{1}{2} \text{OT}_{\varepsilon,\alpha}(\nu, \nu).$$

Can be used for gradient flows.