Large Scale computation of Means and Clusters for Persistence Diagrams using Optimal Transport

Theo Lacombe $^{(1)}$, Marco Cuturi $^{(2)}$, Steve Oudot $^{(1)}$ (1)Inria Saclay, datashape. (2)CREST, ENSAE & Google Brain

Overview

Topological Data Analysis:

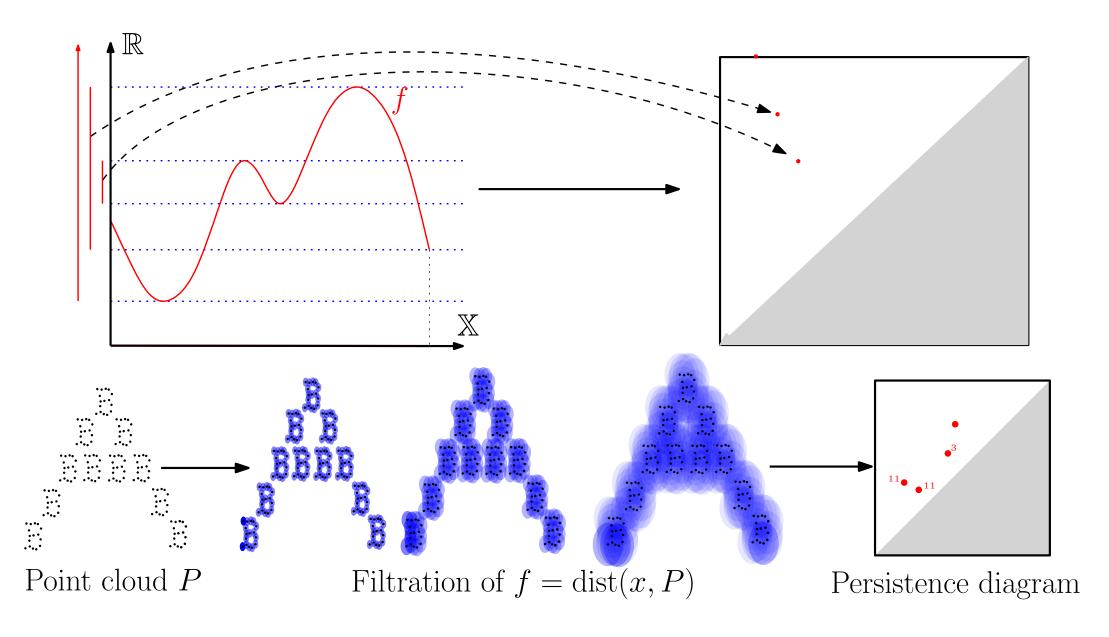
- Provide descriptors, called **persistence diagrams** (PDs), of the topology of an object at all scales.
- Compare PDs with partial matching metrics.

Problem motivation:

- Hard to compute statistical tools for PDs, even as elementary as barycenters.
- Current algorithm [1] to estimate PD barycenters is non-convex and intractable on large data.

Our contributions:

- Reformulate PD metrics as exact OT problems.
- Adapt the OT entropic smoothing [2] for PD metrics, in particular convolution on regular grids [3] allowing parallelization and GPU computations.
- Propose a convex formulation and scalable algorithm for PD barycenter estimation.



I. Persistence diagrams and metrics

Persistence diagrams (PDs) are finite point measures, i.e. $\mu = \sum \delta_{x_i}$, with $x_i \in \{(t_1, t_2) \in \mathbb{R}^2, t_2 > t_1\}$. The distance between two diagrams μ and ν is $(p \ge 1)$:

$$d_p(\boldsymbol{\mu}, \boldsymbol{\nu}) := \left(\min_{\zeta \in \Gamma(\boldsymbol{\mu}, \boldsymbol{\nu})} \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in \zeta} \|\boldsymbol{x} - \boldsymbol{y}\|^p + \sum_{s \notin \zeta} \|s - \pi_{\Delta}(s)\|^p \right)^{\frac{1}{p}},$$

with $\Gamma(\mu, \nu)$: **partial** matchings between μ and ν , and $\pi_{\Delta}(s)$ the orthogonal projection of s onto the diagonal.

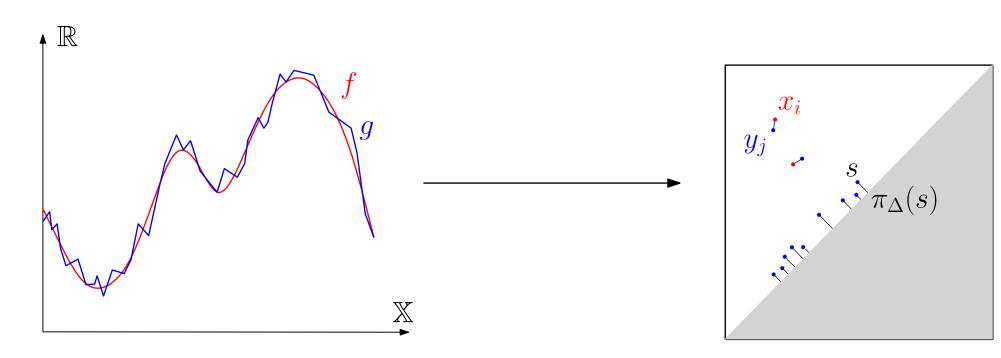
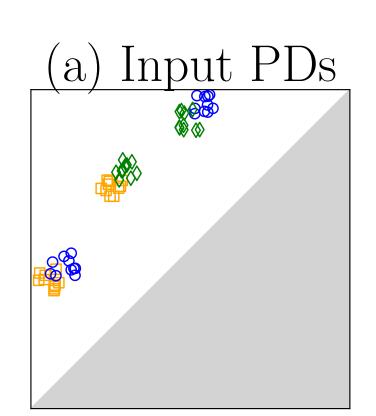
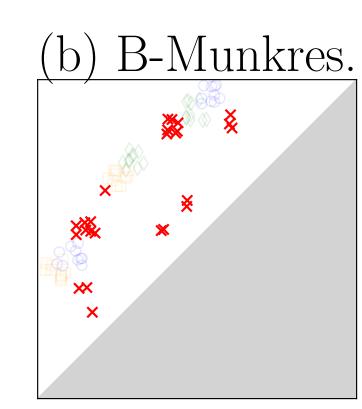
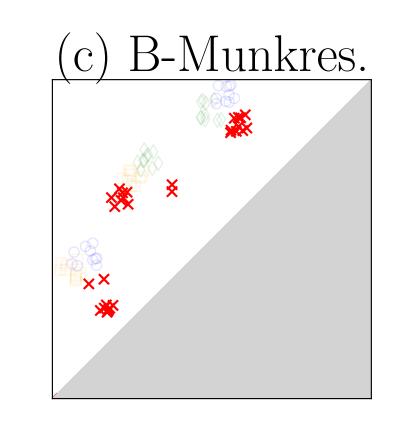
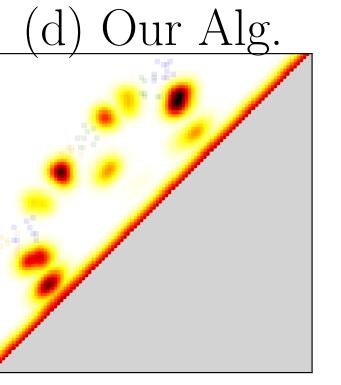


Figure 1: (left) Two functions $f,g:\mathbb{X}\to\mathbb{R}$. (right) Corresponding PDs and an optimal partial matching ζ (edges).





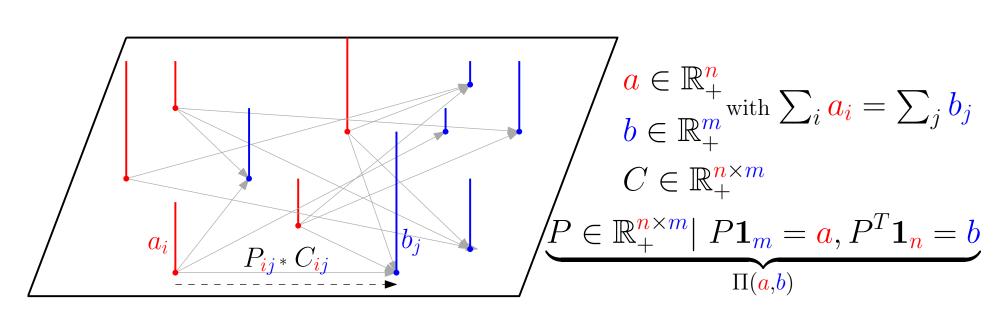




Our Alg. (d)	0.542
B-Munkres (c)	0.555
B-Munkres (b)	0.589
	Final energy

Figure 2:Illustration of our approach on a simple example. (a) 3 PDs for which we want to estimate a barycenter. (b,c) Outputs of B-Munkres algorithm [1] for two different initializations. Variability is due to non-convexity. (d) The output of our convex formulation. It performs better (lower energy).

II. Smoothed optimal transport (OT)



Smoothed OT problem ($\gamma > 0$):

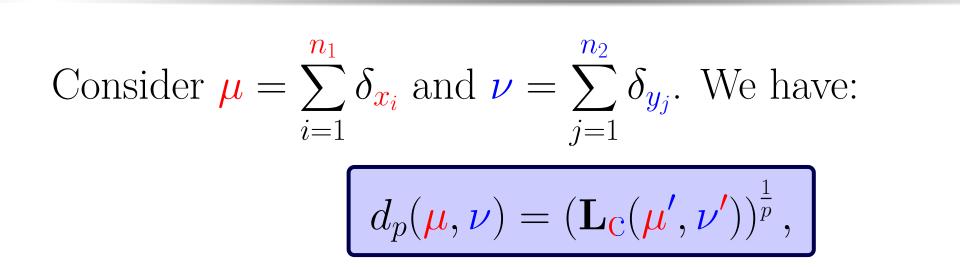
$$\mathbf{L}_{C}^{\gamma}(\mathbf{a}, \mathbf{b}) := \min_{P \in \Pi(\mathbf{a}, \mathbf{b})} \langle P, C \rangle - \gamma h(P)$$

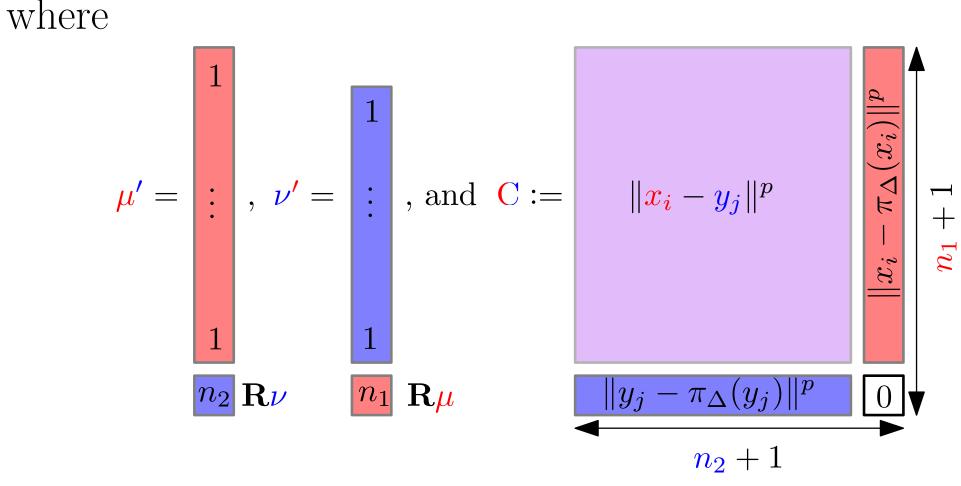
where $h(P) := -\sum_{ij} P_{ij} (\log P_{ij} - 1)$.

Advantages:

- Solved by iterating $(\mathbf{u}, \mathbf{v}) \mapsto \left(\frac{\mathbf{a}}{K\mathbf{v}}, \frac{\mathbf{b}}{K^T\mathbf{u}}\right)$, with $K := e^{-\frac{C}{\gamma}}$.
- Converges to $\mathbf{L}_C(\boldsymbol{a}, \boldsymbol{b}) := \min\{\langle P, C \rangle; P \in \Pi(\boldsymbol{a}, \boldsymbol{b})\}$ when **Idea:** Approximate d_p with L_C^{γ} . $\gamma \to 0$, with controllable error (upper and lower bounds).
- Numerically efficient to solve: GPU + Parallelism.
- Differentiable, with tractable gradient.

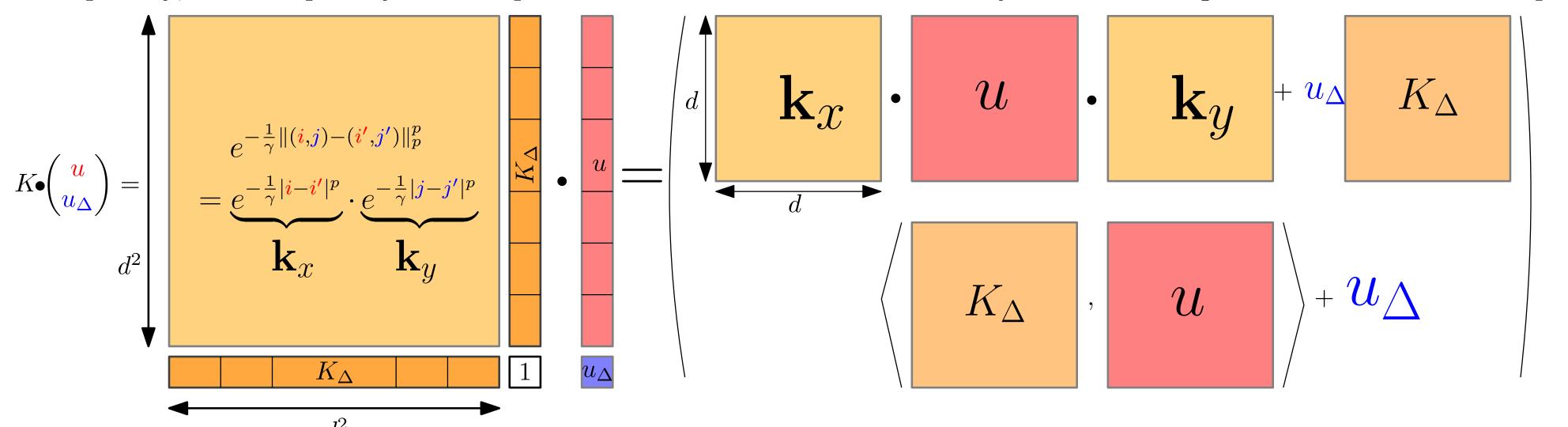
III. OT formulation of d_p



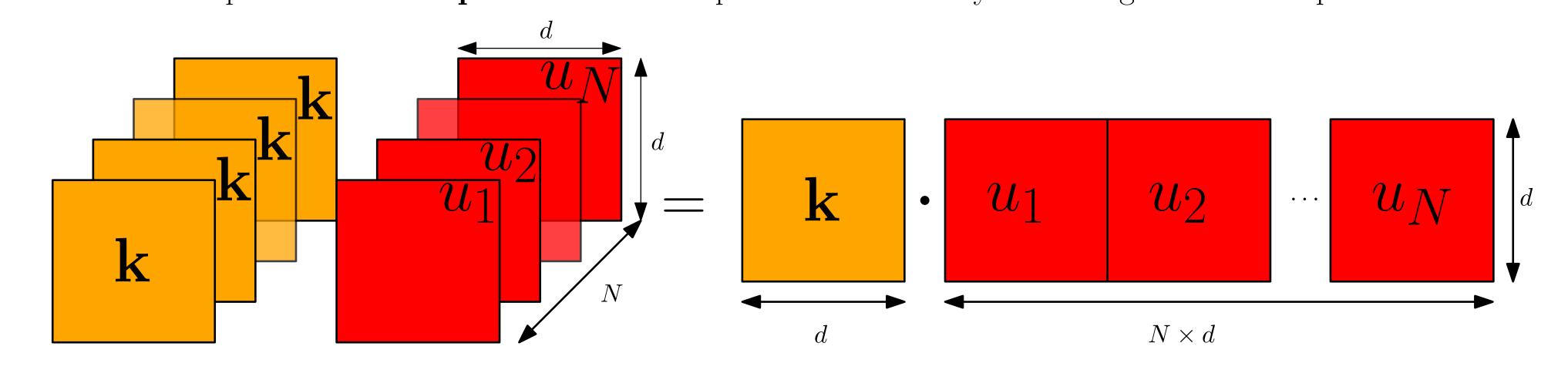


IV. Fast convolutions in the PD space

Discretize PDs on a $d \times d$ grid (+1 for the diagonal), $\Rightarrow (d^2+1)$ histograms (Eulerian approach). C, K are $(d^2+1) \times (d^2+1)$ shaped. Hopefully, the complexity of the operation $u \mapsto Ku$ can be drastically reduced using **convolutions** in the plane.



These matrix manipulations can be **parallelized** and performed efficiently as one big matrix multiplication on a **GPU**.



V. Smoothed barycenters for PDs

For $h_1 \dots h_N$ histograms, a barycenter (Fréchet mean) is a minimizer of the energy:

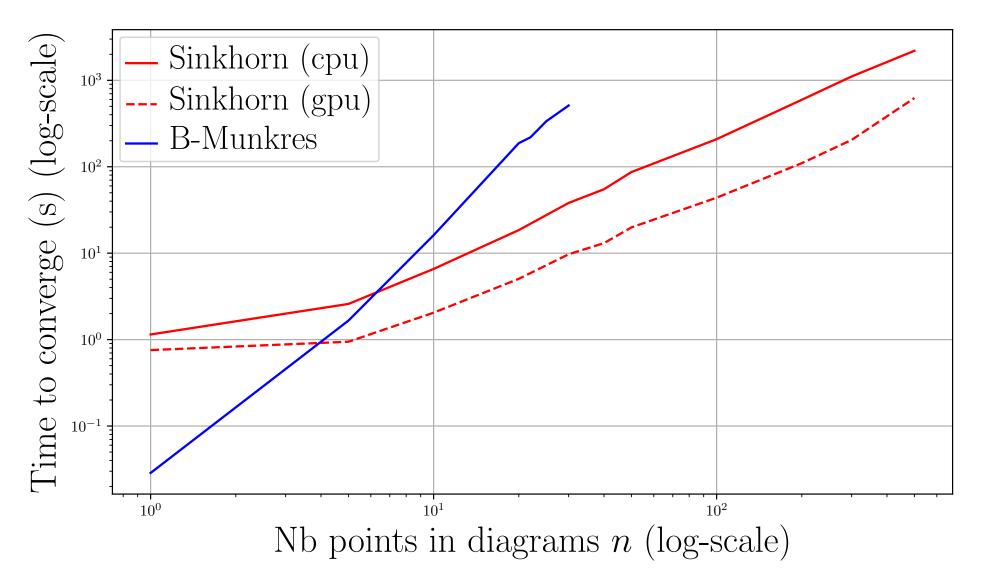
$$\mathcal{E}^{\gamma}: oldsymbol{x} \mapsto \sum_{i=1}^{N} \mathbf{L}_{C}^{\gamma}(oldsymbol{x} + \mathbf{R}oldsymbol{h}_{i}, oldsymbol{h}_{i} + \mathbf{R}oldsymbol{x}),$$

which is **differentiable** with gradient

$$\nabla = \gamma \left(\sum_{i=1}^{N} \log(u_i^{\gamma}) + \mathbf{R}^T \log(v_i^{\gamma}) \right).$$

Advantages:

- Convex formulation: minimize with gradient descent. Gives better estimations in practice.
- GPU + Parallelism: drastically outperform previous algorithm (B-Munkres) developed in [1] on large scales.



Allows for large scale applications, e.g. k-mean clustering on thousands of PDs:

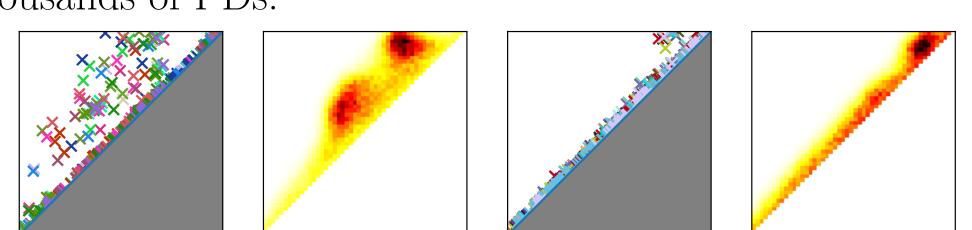


Figure 3:k-means on a real life dataset of 5000 persistence diagrams. Two identified clusters and their centroids.

References

[1] Katharine Turner et al.

Fréchet means for distributions of persistence diagrams. Discrete & Computational Geometry, 52(1):44-70, 2014.

[2] Marco Cuturi.

Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems, pages 2292-2300, 2013.

[3] Solomon et al.

Convolutional Wasserstein distances: Efficient optimal transportation on geometric domains.

In ACM Transactions on Graphics (TOG), 2015.





