

Optimal transport and dimension reduction: Entropic Wasserstein Component Analysis

Antoine Collas

Postdoc supervised by Alexandre Gramfort and Rémi Flamary

Paper published at MLSP (Machine Learning for Signal Processing) 2023, Rome

Work done with Titouan Vayer and Arnaud Breloy



Inria



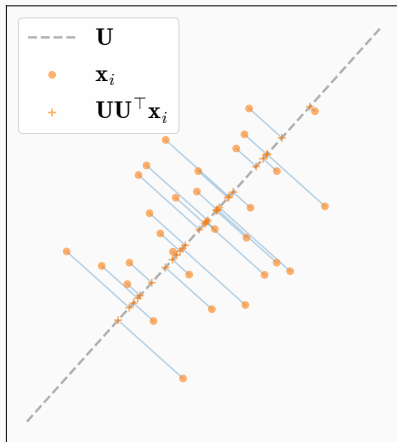
université
PARIS-SACLAY

A few reminders

Principal Component Analysis (PCA)

Subspace learning from data: $(\mathbf{x}_1, \dots, \mathbf{x}_n) \in (\mathbb{R}^d)^n$.

Goal: find a subspace \mathbf{U} such that $\mathbf{x}_i \approx \mathbf{U}\mathbf{U}^\top \mathbf{x}_i$.



Principal Component Analysis (PCA)

PCA: solution to

$$\mathbf{U}^{\text{PCA}} \in \arg \min_{\mathbf{U} \in \text{St}(d,k)} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^{\text{T}}\mathbf{x}_i\|_2^2$$

with $\text{St}(d, k) \triangleq \{ \mathbf{U} \in \mathbb{R}^{d \times k} \mid \mathbf{U}^{\text{T}}\mathbf{U} = \mathbf{I}_k \}$.

Solution computation:

$$\mathbf{X} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_n] \stackrel{\text{SVD}}{=} [\mathbf{U}^{\text{PCA}} \mid \mathbf{U}_{\perp}] \boldsymbol{\Sigma} \mathbf{V}^{\text{T}}$$

Optimal Transport (OT): Wasserstein distance

Given $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ and $(\mathbf{z}_1, \dots, \mathbf{z}_n)$ in \mathbb{R}^d and their empirical measures

$$\mu = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i} \quad \text{and} \quad \nu = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{z}_i}$$

the squared 2-Wasserstein distance with the ℓ^2 metric is

$$\mathcal{W}_2^2(\mu, \nu) = \underset{\pi \in \Pi(\frac{1}{n}\mathbf{1}_n, \frac{1}{n}\mathbf{1}_n)}{\text{minimize}} \sum_{i,j} \pi_{ij} \|\mathbf{x}_i - \mathbf{z}_j\|_2^2$$

with

$$\Pi(\mathbf{a}, \mathbf{b}) \triangleq \{ \pi \in \mathbb{R}^{n \times n} \mid \pi_{ij} \geq 0, \pi \mathbf{1}_n = \mathbf{a}, \pi^\top \mathbf{1}_n = \mathbf{b} \}.$$

[Peyré et al. 2019]

Optimal Transport (OT): entropic regularization

Entropic regularized OT:

$$\text{minimize}_{\pi \in \Pi(\frac{1}{n}\mathbf{1}_n, \frac{1}{n}\mathbf{1}_n)} \sum_{i,j}^{n,n} \pi_{ij} \|\mathbf{x}_i - \mathbf{z}_j\|_2^2 - \varepsilon \mathbf{H}(\pi)$$

with $\mathbf{H}(\pi) \triangleq -\sum_{i,j}^{n,n} \pi_{ij} \log \pi_{ij}$ and $\varepsilon > 0$.

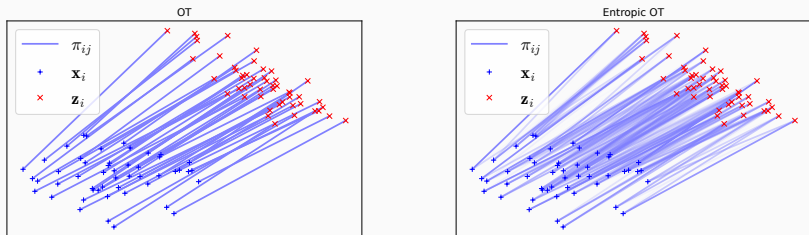


Figure adapted from POT library [Flamary et al. 2021]

Optimal Transport (OT): Sinkhorn-Knopp algorithm

Solution to the entropic regularized OT problem:

$$\boldsymbol{\pi} = \text{diag}(\mathbf{u})\mathbf{K} \text{diag}(\mathbf{v})$$

with

$$K_{ij} \triangleq \exp(-\|\mathbf{x}_i - \mathbf{z}_j\|_2^2/\varepsilon)$$

and \mathbf{u} and \mathbf{v} obtained by iterating

$$\begin{aligned}\mathbf{u} &\leftarrow \frac{1}{n}\mathbf{1}_n \oslash \mathbf{K}\mathbf{v} \\ \mathbf{v} &\leftarrow \frac{1}{n}\mathbf{1}_n \oslash \mathbf{K}^\top \mathbf{u}.\end{aligned}$$

[Cuturi 2013]

Entropic Wasserstein Component Analysis (EWCA)

Motivation

Given the empirical measures

$$\mu = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i} \quad \text{and} \quad \nu = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{U}\mathbf{U}^\top \mathbf{x}_i}$$

we have

$$\mathbf{U}^{\text{PCA}} \in \arg \min_{\mathbf{U} \in \text{St}(d,k)} \left\{ \mathcal{W}_2^2(\mu, \nu) = \underset{\pi \in \Pi(\frac{1}{n}\mathbf{1}_n, \frac{1}{n}\mathbf{1}_n)}{\text{minimize}} \sum_{i,j} \pi_{ij} \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^\top \mathbf{x}_j\|_2^2 \right\}$$

and optimal coupling

$$\pi^* = \frac{1}{n} \mathbf{I}_n.$$

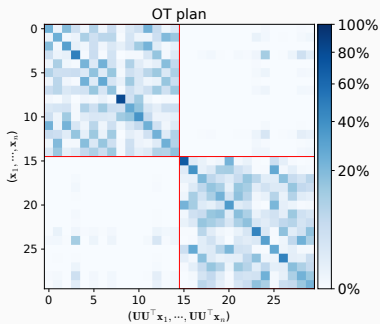
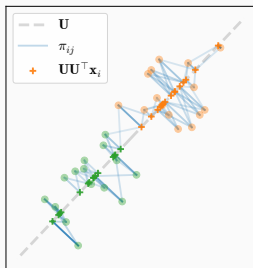
Motivation: add entropy to have $\pi^* \neq \frac{1}{n} \mathbf{I}_n$ and thus minimize reconstruction error of clusters (and not points).

Entropic Wasserstein Component Analysis (EWCA) problem

Entropic Wasserstein Component Analysis (EWCA):

$$\begin{aligned} & \underset{\substack{\pi \in \Pi(\frac{1}{n}\mathbf{1}_n, \frac{1}{n}\mathbf{1}_n) \\ \mathbf{U} \in \text{St}(d, k)}}}{\text{minimize}} \sum_{i,j=1}^{n,n} \pi_{ij} \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^\top \mathbf{x}_j\|_2^2 - \varepsilon \text{H}(\pi). \end{aligned}$$

where H is the entropy function.

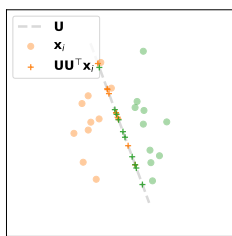


Entropic Wasserstein Component Analysis (EWCA) problem

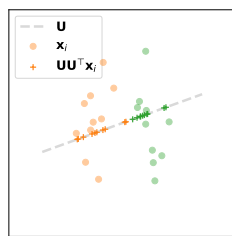
Limit cases:

$$(\pi_\varepsilon, \mathbf{U}_\varepsilon) = \arg \min_{\substack{\pi \in \Pi(\frac{1}{n}\mathbf{1}_n, \frac{1}{n}\mathbf{1}_n) \\ \mathbf{U} \in \text{St}(d, k)}} \sum_{i,j=1}^{n,n} \pi_{ij} \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^\top \mathbf{x}_j\|_2^2 - \varepsilon \mathbf{H}(\pi)$$

- $\varepsilon \rightarrow 0 \implies \pi_\varepsilon \rightarrow \frac{1}{n}I_n$ and $\mathbf{U}_\varepsilon \rightarrow$ top k eigenvectors of $\frac{1}{n}\mathbf{X}\mathbf{X}^\top$; we recover PCA !
- $\varepsilon \rightarrow +\infty \implies \pi_\varepsilon \rightarrow \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top$ and $\mathbf{U}_\varepsilon \rightarrow$ last k eigenvectors of $\frac{1}{n}\mathbf{X}\mathbf{X}^\top$.



(a) $\varepsilon = 0.1$



(b) $\varepsilon = 50$

Block coordinate descent algorithm

$$\underset{\substack{\boldsymbol{\pi} \in \Pi(\frac{1}{n}\mathbf{1}_n, \frac{1}{n}\mathbf{1}_n) \\ \mathbf{U} \in \text{St}(d, k)}}}{\text{minimize}} \sum_{i,j=1}^{n,n} \pi_{ij} \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^\top \mathbf{x}_j\|_2^2 - \varepsilon H(\boldsymbol{\pi})$$

Given the current estimate $(\boldsymbol{\pi}^{(t)}, \mathbf{U}^{(t)})$,

- **$\boldsymbol{\pi}$ -step**: compute $\boldsymbol{\pi}^{(t+1)}$ using Sinkhorn-Knopp algorithm,
- **\mathbf{U} -step**: compute $\mathbf{U}^{(t+1)}$ as the k first eigenvectors of

$$\mathbf{X} \left(2 \text{sym}(\boldsymbol{\pi}^{(t+1)}) - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \mathbf{X}^\top.$$

Problem: \mathbf{U} -step requires SVD of a $d \times d$ matrix.

Majorization-minimization (MM) algorithm

Majorization-minimization over the Stiefel manifold

$$\underset{\mathbf{U} \in \text{St}(d, k)}{\text{minimize}} f(\mathbf{U})$$

Given iterate $\mathbf{U}^{(t)}$,

- **Majorization:**

$$f(\mathbf{U}) \leq g(\mathbf{U} | \mathbf{U}^{(t)}) \text{ for all } \mathbf{U} \in \text{St}(d, k)$$

such that

$$g(\mathbf{U} | \mathbf{U}^{(t)}) = 2 \text{Tr}(\mathbf{U}^\top \mathbf{M} \mathbf{U}^{(t)}) + \text{const. (linearity)}$$

for some $\mathbf{M} \in \mathbb{R}^{d \times d}$.

- **Minimization:**

$$\mathbf{U}^{(t+1)} = \text{pf}(-\mathbf{M} \mathbf{U}^{(t)}) = \underset{\mathbf{U} \in \text{St}(d, k)}{\text{arg min}} g(\mathbf{U} | \mathbf{U}^{(t)})$$

where pf returns the orthogonal factor of the polar decomposition.

[Breloy et al. 2021]

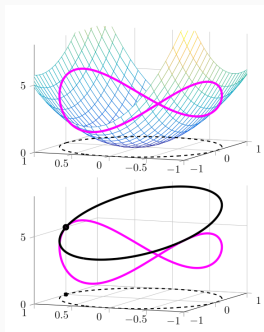


Figure 2: A quadratic form over $\text{St}(2, 1)$ (pink) and its surrogate (black). Figure from [Breloy et al. 2021].

Majorization-minimization over the Stiefel manifold

U -step:

$$\text{minimize}_{U \in \text{St}(d,k)} \left\{ \sum_{i,j=1}^{n,n} \pi_{ij} \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^\top \mathbf{x}_j\|_2^2 \propto \text{Tr}(\mathbf{U}^\top \mathbf{M}\mathbf{U}) \right\}$$

for some $\mathbf{M}^\top = \mathbf{M}$ and $\mathbf{M} \preceq \mathbf{0}$ (negative semi-definite).

Given the current estimate $\mathbf{U}^{(t)}$,

- **Majorization** (by concavity):

$$\text{Tr}(\mathbf{U}^\top \mathbf{M}\mathbf{U}) \leq 2 \text{Tr}(\mathbf{U}^\top \mathbf{M}\mathbf{U}^{(t)}) + \text{const.},$$

- **Minimization:**

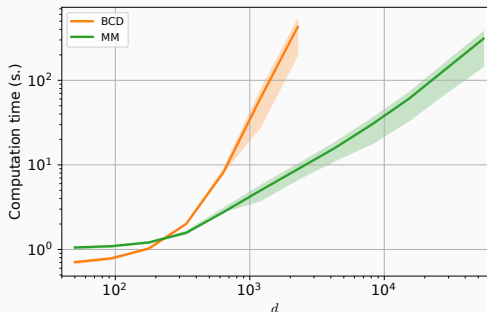
$$\mathbf{U}^{(t+1)} = \text{pf}(-\mathbf{M}\mathbf{U}^{(t)})$$

BCD vs block-MM: computational complexity

Overall computational complexity per iteration:

- **BCD**: $\mathcal{O}(n^2d + nd^2 + d^3)$,
- **Block-MM**: $\mathcal{O}(n^2d)$.

Lower complexity of the block-MM but requires more iterations...



Numerical experiments

Numerical experiments: classification

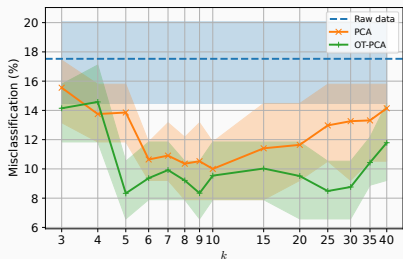
Datasets of gene expressions:

- Breast: $d = 54675$, $n = 151$, and 6 classes [Feldes et al. 2019],
- Khan2001: $d = 2308$, $n = 63$, and 4 classes [Khan et al. 2001].

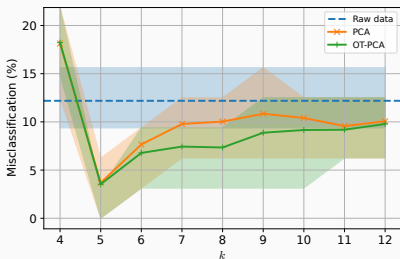
Classification pipeline:

- 1-Nearest neighbor classifier on the projected data $\mathbf{U}^T \mathbf{x}_i$,
- two algorithms: PCA and EWCA,
- evaluation over 100 random splits of the data (50% training, 50% testing),
- hyperparameter ε tuned by cross-validation on the training set.

Numerical experiments: classification



(a) *Breast* dataset



(b) *kxan2001* dataset

Figure 3: Misclassification rate (%) versus subspace dimension k (the lower the better). Mean, 1st and 3rd quartiles are reported.

Numerical experiments: transport plan

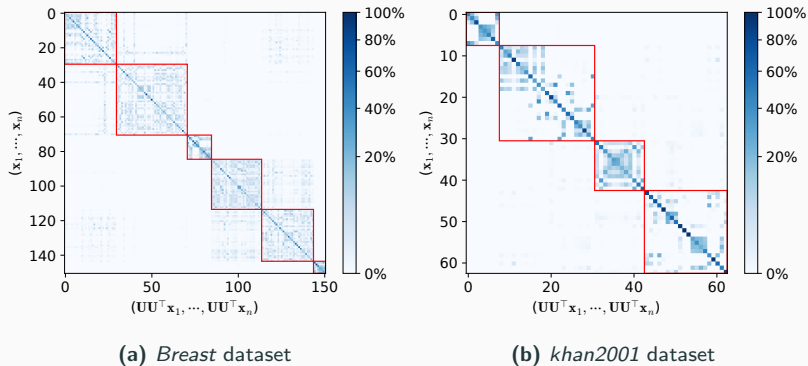


Figure 4: Transport plan π (%) computed with EWCA ($k = 5$). The red squares enclose the data belonging to the same class.

Numerical experiments: TSNE

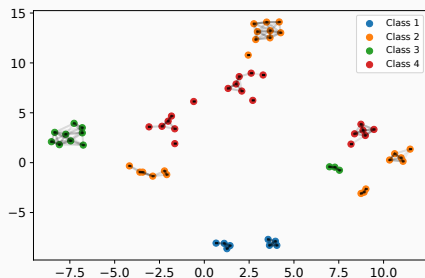


Figure 5: TSNE of the projected data $(\mathbf{U}^\top \mathbf{x}_1, \dots, \mathbf{U}^\top \mathbf{x}_n)$ ($k = 5$) computed with EWCA on the *Khan2001* dataset. The grey links represent the intensity of the values of the transport plan.

Numerical experiments: sensitivity to entropy regularization

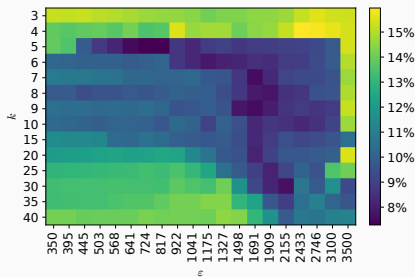


Figure 6: Misclassification rate (%) versus subspace dimension k and entropy intensity ε on the *Breast* dataset (the lower the better).

Conclusions

- Generalization of PCA that takes into account the neighbourhood of data,
- optimization with a BCD and a block-MM,
- use in place of PCA in a classification pipeline on two gene expressions datasets.

Preprint available at

<https://arxiv.org/abs/2303.05119>

Code available at

github.com/antoinecollas/Entropic_Wasserstein_Component_Analysis

Implemented in POT library

https://pythonot.github.io/master/auto_examples/others/plot_EWCA.html

References



Breloy, A. et al. (2021). "Majorization-minimization on the Stiefel manifold with application to robust sparse PCA". In: *IEEE Transactions on Signal Processing* 69, pp. 1507–1520.



Cuturi, M. (2013). "Sinkhorn distances: Lightspeed computation of optimal transport". In: *Advances in neural information processing systems* 26.



Feltes, B. C. et al. (2019). "Cumida: an extensively curated microarray database for benchmarking and testing of machine learning approaches in cancer research". In: *Journal of Computational Biology* 26.4, pp. 376–386.



Flamary, R. et al. (2021). "POT: Python Optimal Transport". In: *Journal of Machine Learning Research* 22.78, pp. 1–8. URL: <http://jmlr.org/papers/v22/20-451.html>.



Khan, J. et al. (2001). "Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks". In: *Nat Med* 7.6, pp. 673–679. ISSN: 10788956. URL: <http://dx.doi.org/10.1038/89044>.



Peyré, G. et al. (2019). *Computational Optimal Transport: With Applications to Data Science*.

Optimal transport and dimension reduction: Entropic Wasserstein Component Analysis

Antoine Collas

Postdoc supervised by Alexandre Gramfort and Rémi Flamary

Paper published at MLSP (Machine Learning for Signal Processing) 2023, Rome

Work done with Titouan Vayer and Arnaud Breloy



Inria



université
PARIS-SACLAY