



Subspace learning

Given data $(\mathbf{x}_1, \dots, \mathbf{x}_n) \in (\mathbb{R}^d)^n$, the goal is to find a subspace \mathbf{U} such that $\mathbf{x}_i \approx \mathbf{U}\mathbf{U}^\top \mathbf{x}_i$.

Principal Component Analysis (PCA):

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

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

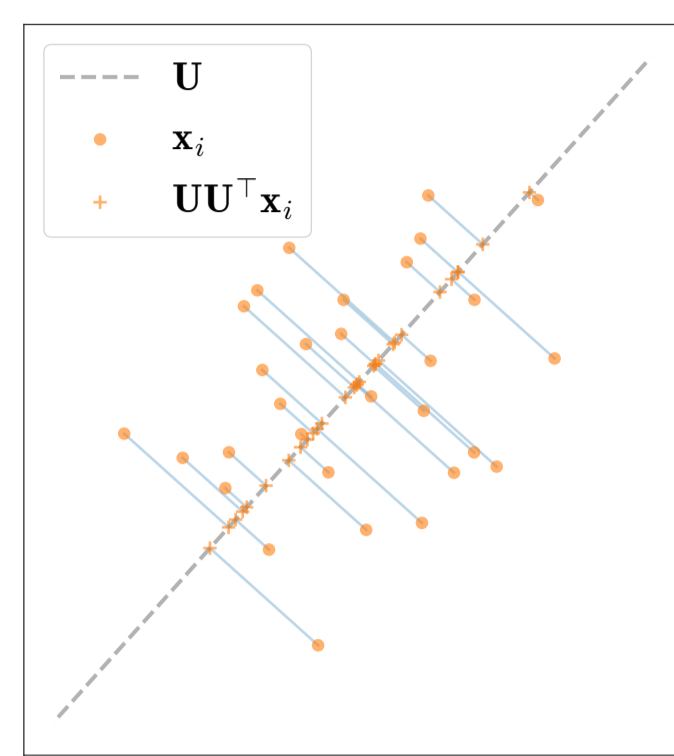


Figure 1. PCA illustration

Optimal transport (OT)

Given $(\mathbf{x}_1, \dots, \mathbf{x}_n), (\mathbf{z}_1, \dots, \mathbf{z}_n) \in (\mathbb{R}^d)^n$ 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 [1]

$$\mathcal{W}_2^2(\mu, \nu) = \min_{\pi \in \Pi} \sum_{i,j} \pi_{ij} \|\mathbf{x}_i - \mathbf{z}_j\|_2^2$$

with

$$\Pi \triangleq \left\{ \pi \in \mathbb{R}^{n \times n} \mid \pi_{ij} \geq 0, \pi \mathbf{1}_n = \frac{1}{n} \mathbf{1}_n, \pi^\top \mathbf{1}_n = \frac{1}{n} \mathbf{1}_n \right\}.$$

Entropic regularized OT

Given the entropy $H(\pi) \triangleq -\sum_{i,j} \pi_{ij} \log \pi_{ij}$ and $\varepsilon > 0$,

$$\min_{\pi \in \Pi} \sum_{i,j} \pi_{ij} \|\mathbf{x}_i - \mathbf{z}_j\|_2^2 - \varepsilon H(\pi).$$

Solved with the Sinkhorn-Knopp algorithm.

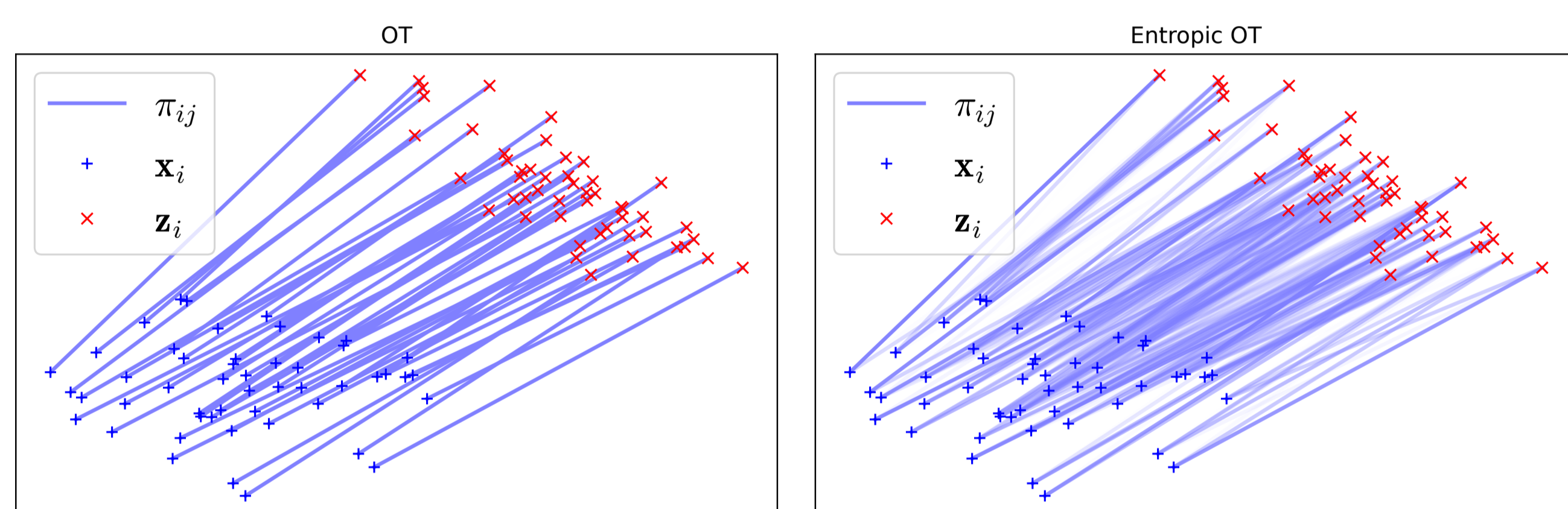


Figure 2. OT creates a one-to-one correspondence between the two datasets (left). Entropic regularization allows for non one-to-one correspondences (right). Figure adapted from POT library [2].

Motivation of the paper

Given the empirical measures

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

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

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).

Problem formulation

Entropic Wasserstein Component Analysis (EWCA):

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

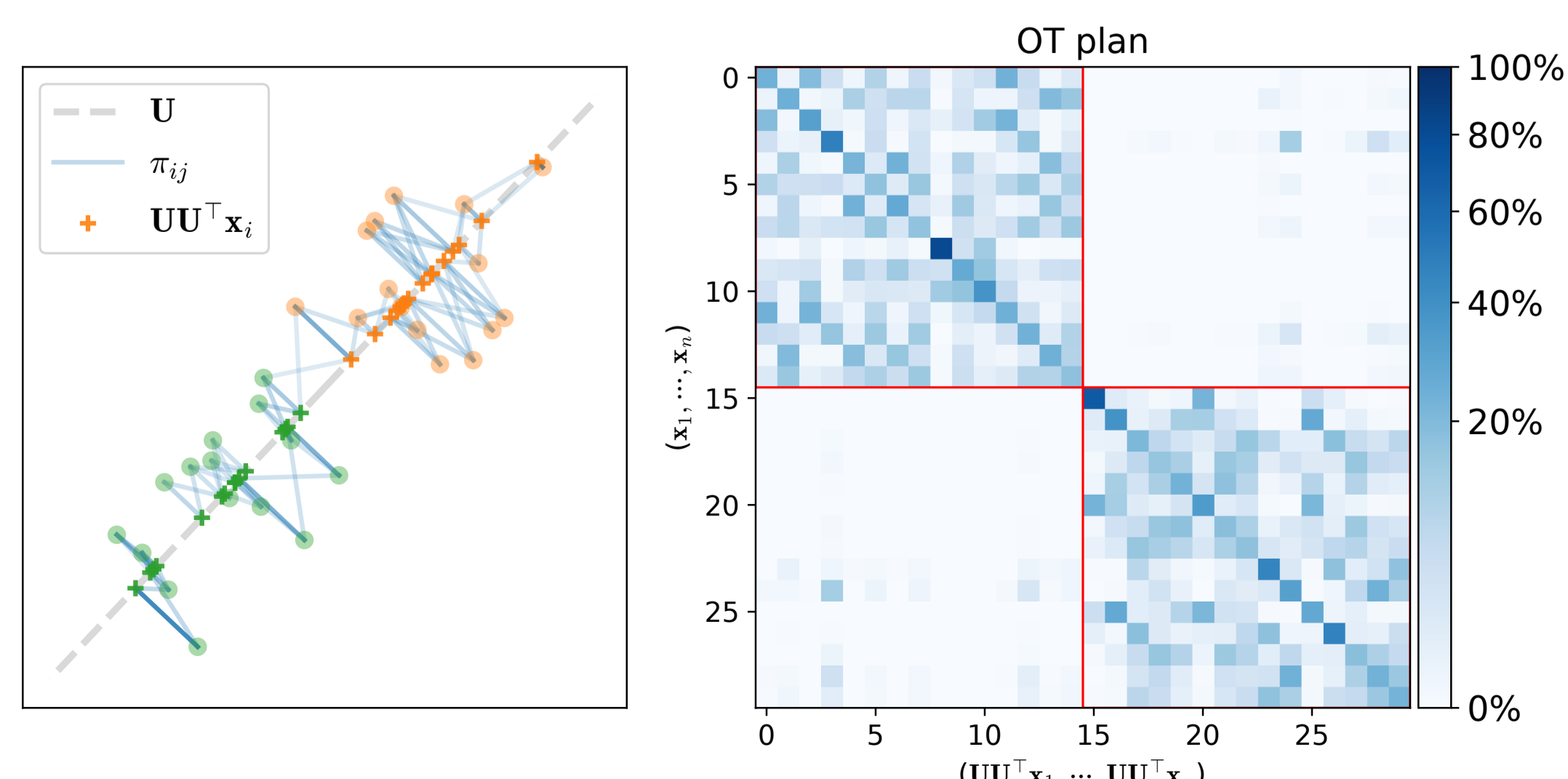


Figure 3. Illustration of EWCA on a two-class dataset. On the left are the samples and their 1D projections, and on the right is the corresponding OT plan.

Limit cases & interpretation

Denoting

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

- $\varepsilon \rightarrow 0 \implies \pi_\varepsilon \rightarrow \frac{1}{n} \mathbf{I}_n$ and $\mathbf{U}_\varepsilon \rightarrow$ top k eigenvectors $\frac{1}{n} \mathbf{X} \mathbf{X}^\top$; 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 $\frac{1}{n} \mathbf{X} \mathbf{X}^\top$.

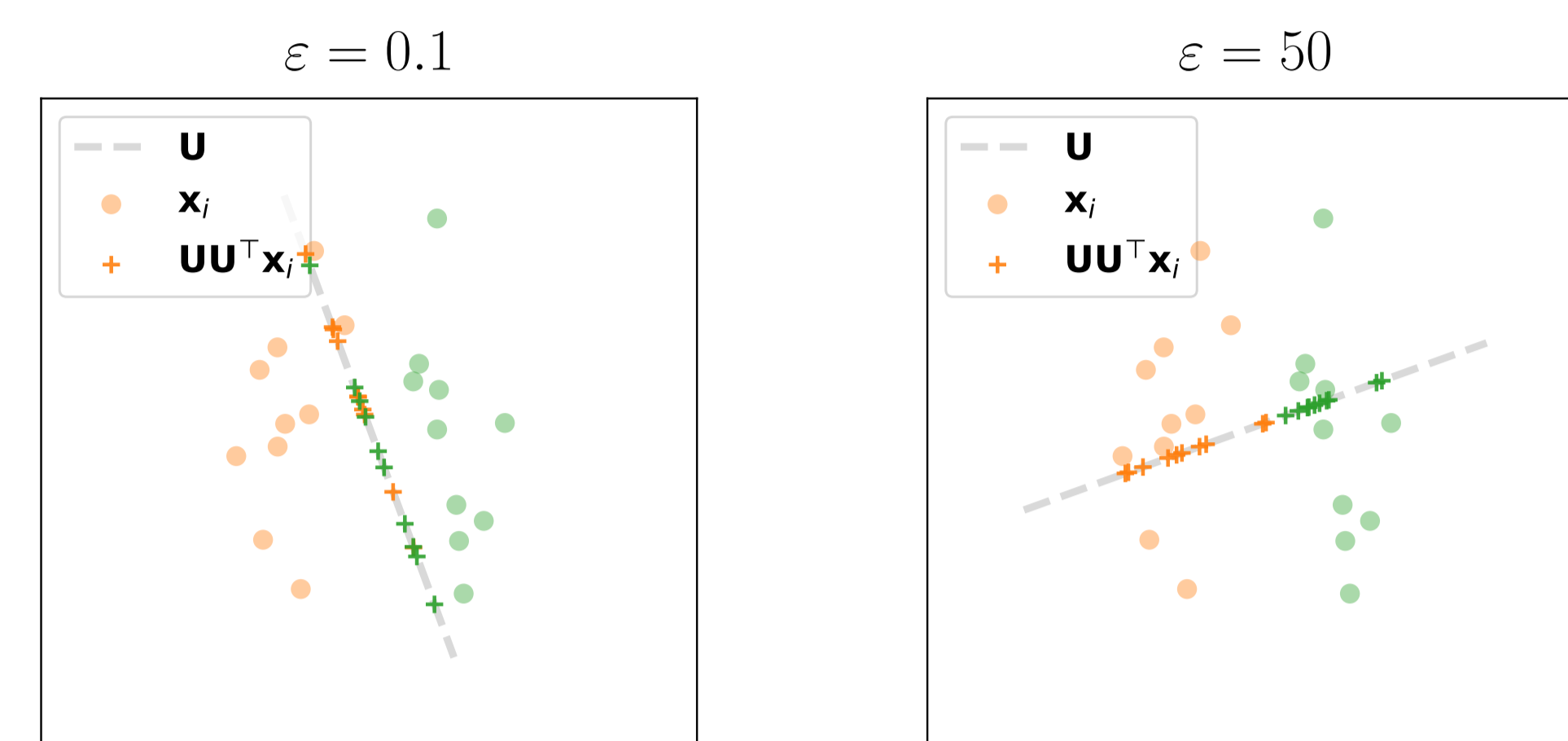


Figure 4. Illustration of EWCA for different values of ε on a two-class dataset.

Resolution: block coordinate descent

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

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

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

Block-majorization-minimization algorithm to avoid SVD of a $d \times d$ matrix in the paper.

Application

Datasets of gene expressions:

- Breast: $d = 54675$, $n = 151$, and 6 classes [3],
- Khan2001: $d = 2308$, $n = 63$, and 4 classes [4].

Classification setup:

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

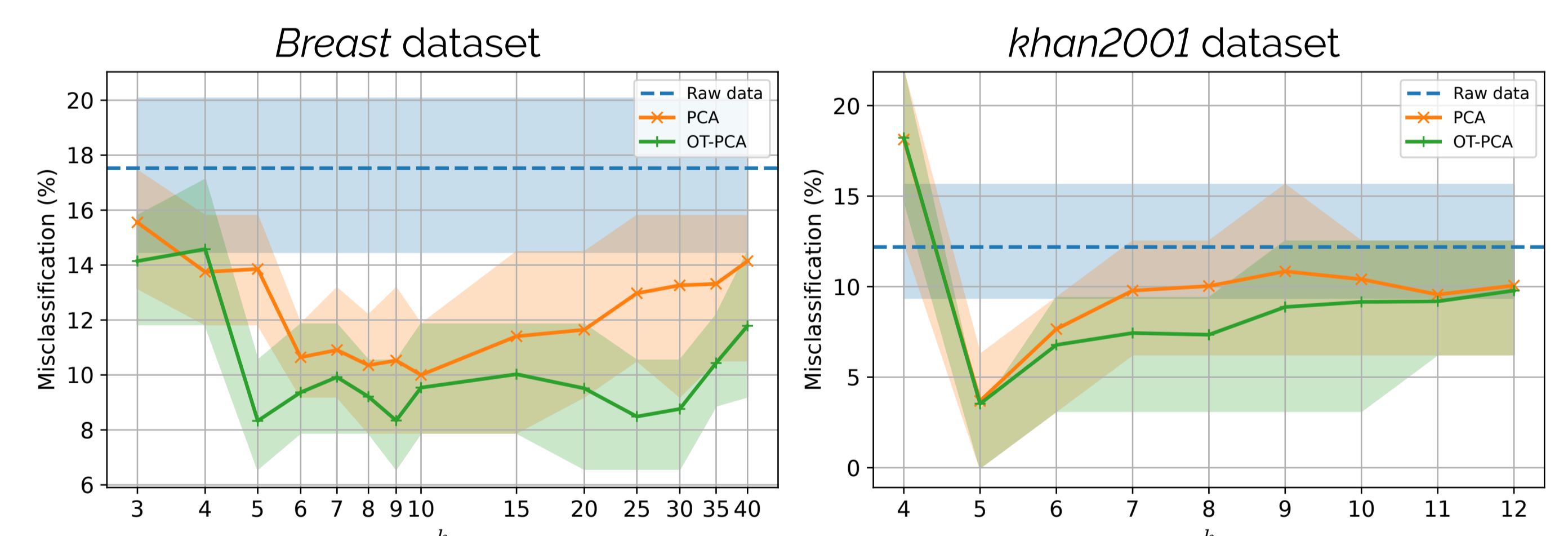


Figure 5. Misclassification rate (%) versus subspace dimension k (the lower the better). Mean, 1st and 3rd quartiles are reported.

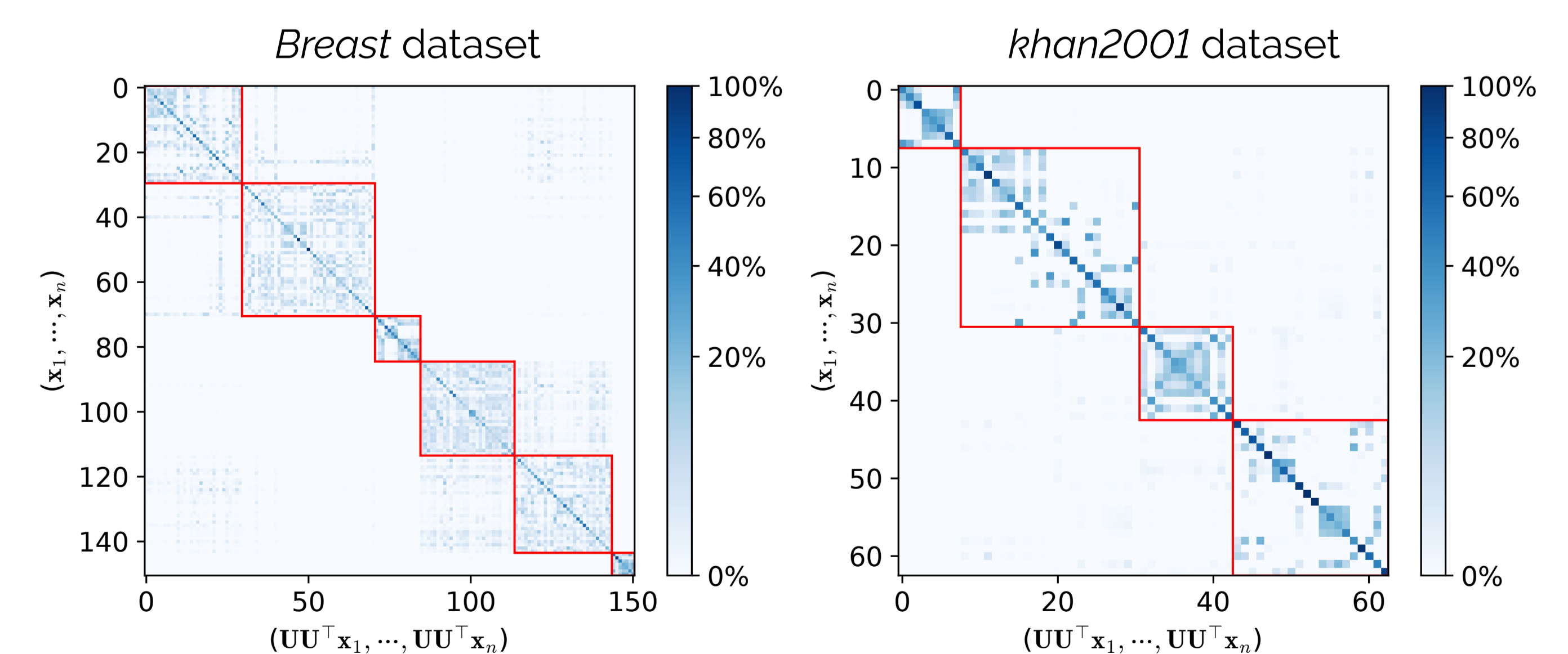


Figure 6. OT plan π (%) computed with EWCA ($k = 5$). The red squares enclose the data belonging to the same class.

References

- [1] G. Peyré and M. Cuturi. *Computational Optimal Transport: With Applications to Data Science*. 2019.
- [2] R. Flamary et al. "POT: Python Optimal Transport". In: *Journal of Machine Learning Research* 22:78 (2021), pp. 1–8. URL: <http://jmlr.org/papers/v22/20-451.html>.
- [3] B. C. Feltes et al. "Cumida: an extensively curated microarray database for benchmarking and testing of machine learning approaches in cancer research". In: *Journal of Computational Biology* 26:4 (2019), pp. 376–386.
- [4] J. Khan et al. "Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks". In: *Nat Med* 7:6 (2001), pp. 673–679. ISSN: 10788956. URL: <http://dx.doi.org/10.1038/89044>.