


Félix Laplante   Université de Paris Saclay  
Christophe Ambroise <sup>1</sup>   Laboratoire de Mathématiques et Modélisation d'Evry, Université  
Paris-Saclay, CNRS, Univ Evry,

Date published: 2025-07-06   Last modified: 2024-06-28

### Abstract

In this paper, Spectral-Bridges, a novel clustering algorithm, is presented. This algorithm builds upon the traditional k-means and spectral clustering frameworks by subdividing data into small Voronoi regions, which are subsequently assessed for their connectivity. Drawing inspiration from Support vector machine, a non-parametric clustering approach is proposed. This approach is characterized by minimal hyperparameters and intuitive usability, thereby augmenting adaptability and enabling the delineation of intricate, non-convex cluster structures.

Both global and local data arrangements are aimed to be discerned by Spectral-Bridges in a scale-invariant manner.

The empirical results underscore Spectral-Bridges as a fast, robust, and versatile tool for sophisticated clustering tasks spanning diverse domains. Its efficacy is observed to extend seamlessly to large-scale scenarios encompassing both real-world and synthetic datasets.

*Keywords:* spectral clustering, vector quantization, scalable, non-parametric

## Contents

1	<b>1 Introduction</b>	2
3	<b>2 Background</b>	2
4	<b>3 Spectral bridges</b>	2
5	3.1 Bridge gain affinity . . . . .	3
6	3.2 Algorithm . . . . .	4
7	<b>4 Numerical experiments</b>	5
8	4.1 Real-world Data . . . . .	5
9	4.2 Synthetic Data . . . . .	5
10	4.2.1 Datasets Summary & Class Balance . . . . .	5
11	4.3 Metrics . . . . .	6
12	4.4 Platform . . . . .	6
13	4.5 Accuracy . . . . .	6
14	<b>5 Conclusive remarks</b>	9
15	<b>6 Appendix</b>	10
16	6.1 Derivation of the bridge gain . . . . .	10

---

<sup>1</sup>Corresponding author: [christophe.ambroise@univ-evry.fr](mailto:christophe.ambroise@univ-evry.fr)

17	<b>References</b>	<b>10</b>
18	<b>Session information</b>	<b>11</b>

## 19 **1 Introduction**

20 Clustering is a fundamental technique for exploratory data analysis. It partition a set of objects into  
 21 a certain number of homogeneous groups, each referred to as a cluster. It is extensively utilized  
 22 across diverse fields such biology, social sciences, and psychology. Clustering is frequently employed  
 23 in conjunction with supervised learning as a pre-processing step, where it helps to structure and  
 24 simplify data, thereby enhancing the performance and interpretability of subsequent predictive  
 25 models.

26 There are numerous approaches to clustering, each defined by how similarity between objects is  
 27 measured, either through a similarity measure or more strictly through a distance metric.

28 Density-based methods identify regions within the data that have a high concentration of points,  
 29 corresponding to the modes of the joint density. A notable non-parametric example of this approach  
 30 is DBSCAN (Ester et al. 1996). In contrast, model-based clustering, such as Gaussian mixture models,  
 31 represents a parametric approach to density-based methods.

32 Geometric approaches, such as kmeans (MacQueen et al. 1967) are distance-based and aim to partition  
 33 the data in a way that optimizes a criterion reflecting group homogeneity.

34 Graph-based methods treat data as a graph, with vertices representing data points and edges weighted  
 35 to reflect the affinity between these points.

36 The algorithm proposed in this paper draws from numerous clustering techniques. The initial  
 37 intuition is to detect high-density areas. To this end, a vector quantization is used to divide the  
 38 space into a Voronoi tessellation. An original geometric criterion is then employed to detect pairs  
 39 of Voronoi regions that are either distant from each other or separated by a low-density boundary.  
 40 Finally, this affinity measure is considered as the weight of an edge in a complete graph connecting  
 41 the centroids of the tessellation, and a spectral clustering algorithm is used to find a partition of this  
 42 graph.

43 The paper begins with a section dedicated to presenting the context and related algorithms, followed  
 44 by a detailed description of the proposed algorithm. Experiments and comparisons with reference  
 45 algorithms are then conducted on both real and synthetic data.

## 46 **2 Background**

47 Spectral clustering is a graph-based approach that computes the eigenvectors of the graph's Laplacian  
 48 matrix. This technique transforms the data into a lower-dimensional space, making the clusters  
 49 more discernible. A standard algorithm like k-means is then applied to these transformed features to  
 50 identify the clusters(Von Luxburg 2007).

51 It enables to capture complex data structures and discern clusters based on the connectivity of data  
 52 points in a transformed space. Notice that spectral clustering can be seen as a relaxed graph cut  
 53 problem.

### 3 Spectral bridges

The proposed algorithm uses K-means centroids for vector quantization defining Voronoi region, and a strategy is proposed to link these regions, with the “affinity” gauged in terms of minimal margin between pairs of classes. These affinities are considered as weight of edges defining a completely connected graph whose vertices are the regions. Spectral clustering on the region provide a partition of the input space. The sole parameters of the algorithm are the number of Voronoi region and the number of final cluster.

#### 3.1 Bridge gain affinity

The basic idea involves calculating the difference in inertia achieved by projecting onto a segment connecting two centroids, rather than using the two centroids separately. If the difference is small, it suggests a low density between the classes. Conversely, if this difference is large, it indicates that the two classes may reside within the same densely populated region.

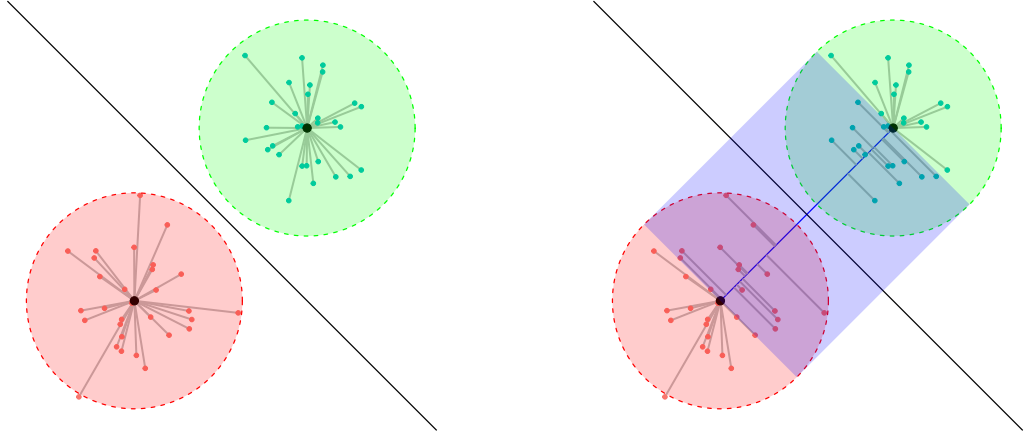


Figure 1: Balls (left) versus Bridge (right). The inertia of each structure is the sum of the squared distances represented by grey lines.

The inertia of two balls  $k$  and  $l$  is

$$I_{kl} = \sum_{i \in k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 + \sum_{i \in l} \|\mathbf{x}_i - \boldsymbol{\mu}_l\|^2.$$

The inertia of a bridge between  $k$  and  $l$  is defined as

$$B_{kl} = \sum_{i \in kl} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2,$$

where

$$\mathbf{p}_{kl}(\mathbf{x}_i) = \boldsymbol{\mu}_k + t_i(\boldsymbol{\mu}_l - \boldsymbol{\mu}_k),$$

with

$$t_i = \min \left( 1, \max \left( 0, \frac{\langle \mathbf{x}_i - \boldsymbol{\mu}_k | \boldsymbol{\mu}_l - \boldsymbol{\mu}_k \rangle}{\|\boldsymbol{\mu}_l - \boldsymbol{\mu}_k\|^2} \right) \right).$$

70 The normalized average of the difference between Bridge and balls inertia is (See [Appendix](#))

$$\frac{B_{kl} - I_{kl}}{(n_k + n_l)\|\mu_k - \mu_l\|^2} = \frac{\sum_{i \in k} \langle \mathbf{x}_i - \mu_k | \mu_l - \mu_k \rangle_+^2 + \sum_{i \in l} \langle \mathbf{x}_i - \mu_l | \mu_k - \mu_l \rangle_+^2}{(n_k + n_l)\|\mu_k - \mu_l\|^4}.$$

71 From this reduction, we define the bridge affinity between centroids  $k$  and  $l$  as:

$$a_{kl} = \begin{cases} 0, & \text{if } k = l, \\ \sqrt{\frac{B_{kl} - I_{kl}}{(n_k + n_l)\|\mu_k - \mu_l\|^2}}, & \text{otherwise.} \end{cases}$$

72 The basic intuition behind this affinity is that  $t_i$  represents the relative position of the projection of  $\mathbf{x}_i$   
73 on the segment  $[\mu_k, \mu_l]$ . For each  $\mathbf{x}_i$  an affinity value  $\alpha_i$  is defined as

$$\alpha_i = \begin{cases} t_i, & \text{if } t_i \in [0, 1/2] \\ 1 - t_i, & \text{if } t_i \in ]1/2, 1], \end{cases}$$

74 This value represents the relative position on the segment, with the centroid of the class to which  $\mathbf{x}_i$   
75 belongs as the starting point.

76 The boundary that separates the two clusters defined by centroids  $\mu_k$  and  $\mu_l$  is a hyperplane. This  
77 hyperplane is orthogonal to the line segment connecting the centroids and intersects this segment at  
78 its midpoint.

79 If we consider all points  $\mathbf{x}_i \in kl$  which are not projected on centroids but somewhere on the segments,  
80 The distance from a point to the hyperplane is large,

$$\|p_{kl}(\mathbf{x}_i) - \mu_{kl}\| = (1/2 - \alpha_i)\|\mu_k - \mu_l\|.$$

81 This distance is similar to the concept of margin in Support Vector Machine (Cortes and Vapnik  
82 1995).

83 When the  $\alpha_i$  values are small (close to zero since  $\alpha_i \in [0, 1/2]$ ), the margins to the hyperplane are  
84 large, indicating a low density between the classes. Conversely, if the margins are small, it suggests  
85 that the two classes may reside within the same densely populated region. Consequently, the sum of  
86 the  $\alpha_i$  or  $\alpha_i^2$  increases with the density of the region between the classes.

87 Note that the criterion is local and indicates the relative difference in densities between the balls and  
88 the bridge, rather than evaluating a global score for the densities of the structures.

### 89 3.2 Algorithm

## 90 4 Numerical experiments

91 In this section, we present the results obtained from testing our algorithm on various datasets, both  
92 small and large scale, including real-world and well-known synthetic datasets. These experiments  
93 assess the accuracy, time and space complexity, ease of use, robustness, and adaptability of our  
94 algorithm. We compare **Spectral-Bridges (SB)** against several state-of-the-art methods, including  
95 **k-means++ (KM)** (MacQueen et al. 1967; **arthur2007k?**), **Expectation-Maximization (EM)**  
96 (**dempster1977maximum?**), **Ward Clustering (WC)** (**ward1963hierarchical?**), and **DBSCAN**  
97 (**DB**) (Ester et al. 1996). This comparison establishes baselines across centroid-based clustering  
98 algorithms, hierarchical methods, and density-based methods. We evaluate the algorithms on both  
99 raw and PCA-processed data with varying dimensionality. For synthetic datasets, we introduce  
100 Gaussian and/or uniform noise to evaluate the robustness of our algorithm.

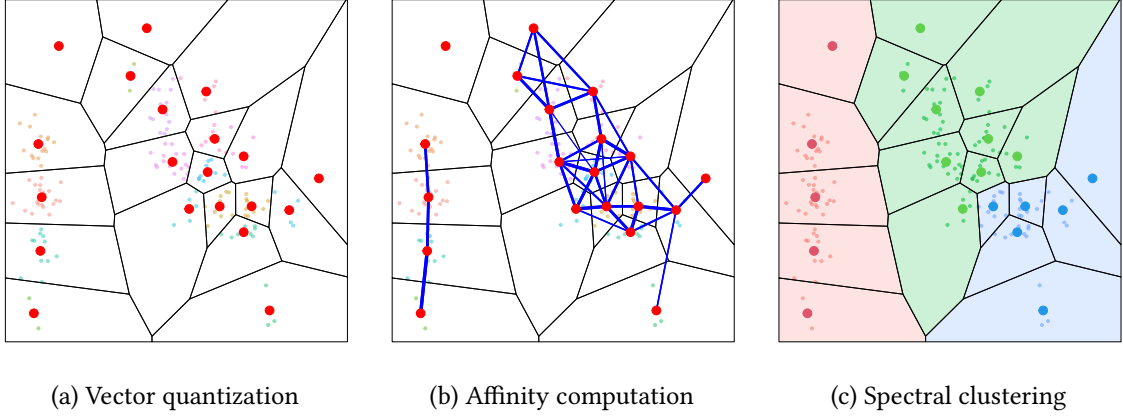


Figure 2: Illustration of the Spectral bridges algorithm with the Iris dataset (first principal plane). Vector quantization (left), Affinity computation (center), Spectral clustering and spreading (right).

---

**Algorithm 1** Spectral Bridges

---

```

1: procedure SPECTRALBRIDGES( $X, k, m$ )  $\triangleright X$ : input dataset,  $k$ : number of clusters,  $m$ : number of
   Voronoi regions
2:   Step 1: Vector Quantization
3:    $centroids \leftarrow \text{KMEANS}(X, m)$   $\triangleright$  Initial centroids using k-means++
4:    $voronoiRegions \leftarrow \text{SUBDIVIDE}(X, centroids)$   $\triangleright$  Subdivide data into Voronoi regions
5:   Step 2: Affinity Matrix Computation  $A = \{a_{kl}\}$ 
6:   Step 3: Spectral Clustering  $\triangleright$  Affect each region to a cluster
7:    $labels \leftarrow \text{SPECTRALCLUSTERING}(X, k)$ 
8:   Step 4: Propagate  $\triangleright$  Each data point is affected to the cluster of its region
9:    $clusters \leftarrow \text{PROPAGATE}(X, labels, centroids)$ 
10:  return  $clusters$   $\triangleright$  Cluster labels for data points in  $X$ 
11: end procedure

```

---

## 4.1 Real-world Data

- **MNIST**: A large dataset containing 60,000 handwritten digit images in ten balanced classes, commonly used for image processing benchmarks. Each image consists of  $28 \times 28 = 784$  pixels.
- **UCI ML Breast Cancer Wisconsin**: A dataset featuring computed attributes from digitized images of fine needle aspirates (FNA) of breast masses, used to predict whether a tumor is malignant or benign.

## 4.2 Synthetic Data

- **Impossible**: A synthetic dataset designed to challenge clustering algorithms with complex patterns.
- **Moons**: A two-dimensional dataset with two interleaving half circles.
- **Circles**: A synthetic dataset of points arranged in two non-linearly separable circles.
- **Smile**: A synthetic dataset with points arranged in the shape of a smiling face, used to test the separation of non-linearly separable data.

### 4.2.1 Datasets Summary & Class Balance

Table 1: Datasets Summary &amp; Class Balance

Dataset	#Dims	#Samples	#Classes	Class Proportions
MNIST	784	60000	10	9.9%, 11.2%, 9.9%, 10.3%, 9.7%, 9%, 9.9%, 10.4%, 9.7%, 9.9%
Breast Cancer	30	569	2	37.3%, 62.7%
Impossible	2	3594	7	24.8%, 18.8%, 11.3%, 7.5%, 12.5%, 12.5%, 12.5%
Moons	2	1000	2	50%, 50%
Circles	2	1000	2	50%, 50%
Smile	2	1000	4	25%, 25%, 25%, 25%

Class proportions are presented in ascending order starting from label 0.

### 4.3 Metrics

To evaluate the performance of our clustering algorithm, we use the Adjusted Rand Index (**ARI**) (**halkidi2002cluster?**) and Normalized Mutual Information (**NMI**) (**cover1991information?**). ARI measures the similarity between two clustering results, ranging from  $-0.5$  to  $1$ , with  $1$  indicating perfect agreement. NMI ranges from  $0$  to  $1$ , with higher values indicating better clustering quality. In some tests, we also report the variability of scores across multiple runs due to the random initialization in k-means, though k-means++ generally provides stable and reproducible results.

### 4.4 Platform

All experiments were conducted on an Archlinux machine with Linux 6.9.3 Kernel, 8GB of RAM, and an AMD Ryzen 3 7320U processor.

### 4.5 Accuracy

We first evaluated our algorithm’s accuracy on the MNIST dataset. Metrics were collected to compare our method with k-means++, EM, and Ward clustering. Metric were estimated by taking the empirical average over 100 consecutive runs with the same random seed for each method.

Let  $h$  denote the embedding dimension of the dataset. We tested our method on the raw MNIST dataset without preprocessing ( $h = 784$ ) and after reducing its dimension using PCA to  $h \in \{8, 16, 32, 64\}$  (see fig.1).

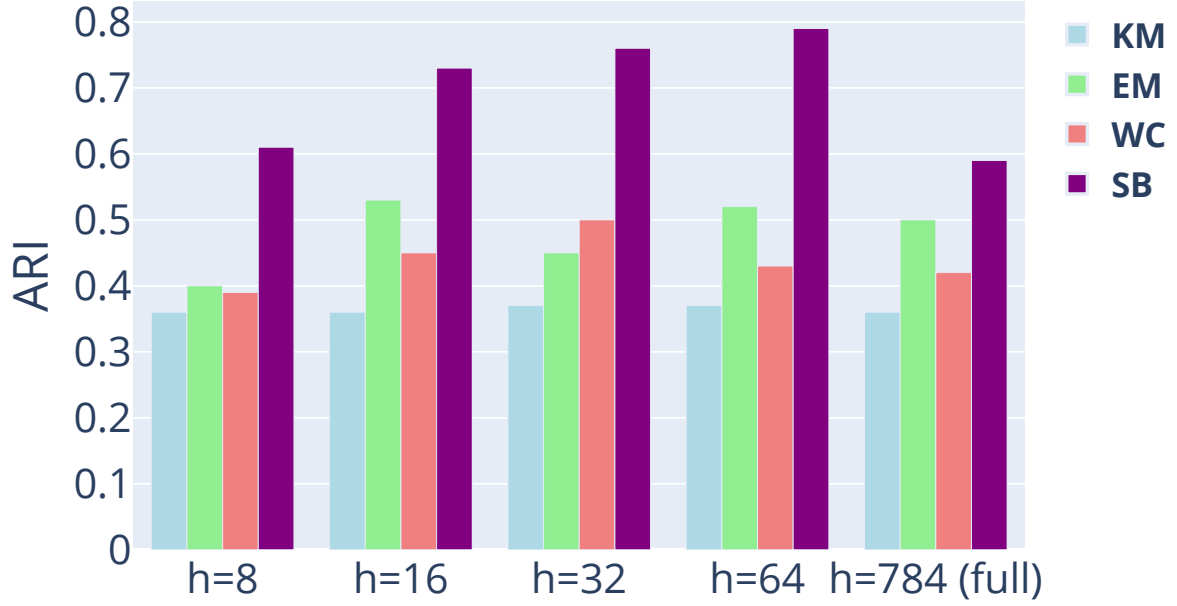


Figure 3: Comparison of **k-means++** (blue), **EM** (green), **Ward Clustering** (red), and **Spectral-Bridges** (purple) on PCA embedding and full MNIST

For visualization purposes, we projected with UMAP the predicted clusters from our algorithm and other methods to compare them against the ground truth labels to better understand the cluster shapes (see table 2). Note that the projection was not used in our experiments as an embedding, and thus does not play any role in the clustering process itself. As a matter of fact, the embedding used was obtained with PCA,  $h = 64$ . Note that the label colors match the legend only in the case of the ground truth data. Indeed, the ordering of the labels have no impact on clustering quality.

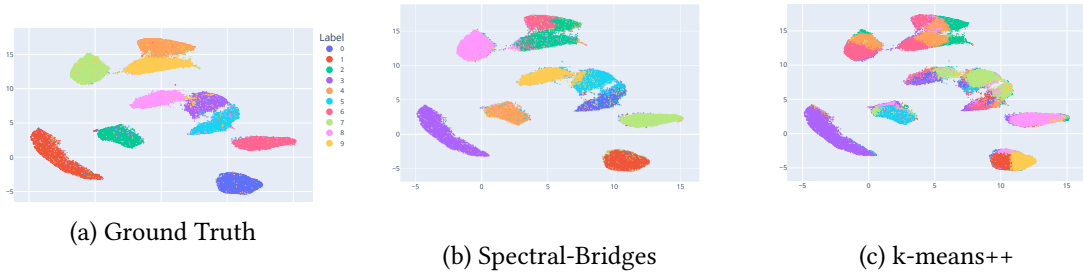


Figure 4: UMAP projection of predicted clusters: **Ground Truth (left)**, **Spectral-Bridges (center)**, **k-means++ (right)** UMAP projection of predicted clusters : **Ground Truth (top)**, **Spectral-Bridges (middle)**, **k-means++ (bottom)**

TODO CANCER

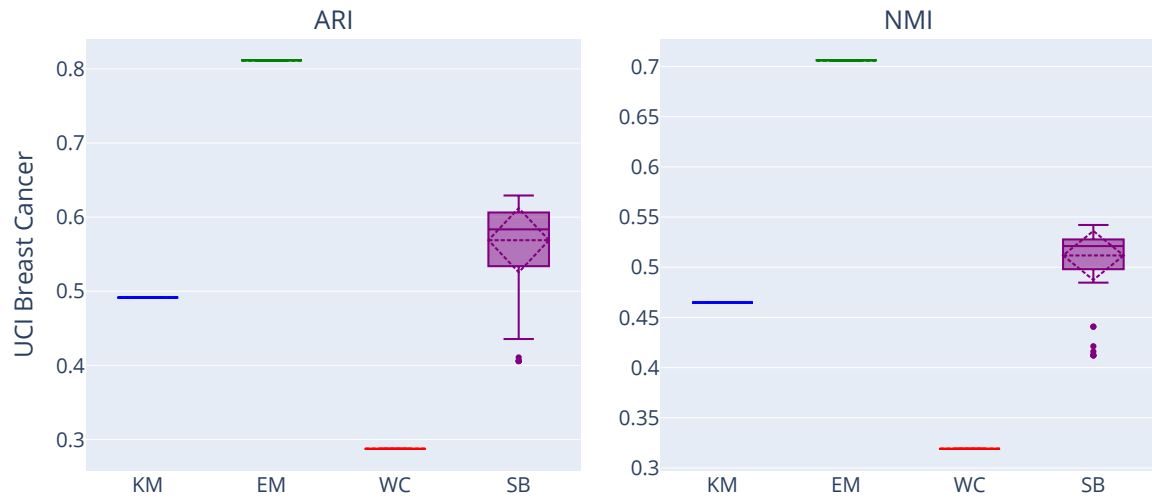


Figure 5: ARI and NMI scores of **k-means++** (blue), **EM** (green), **Ward Clustering** (red), and **Spectral-Bridges** (purple) on the UCI Breast Cancer dataset



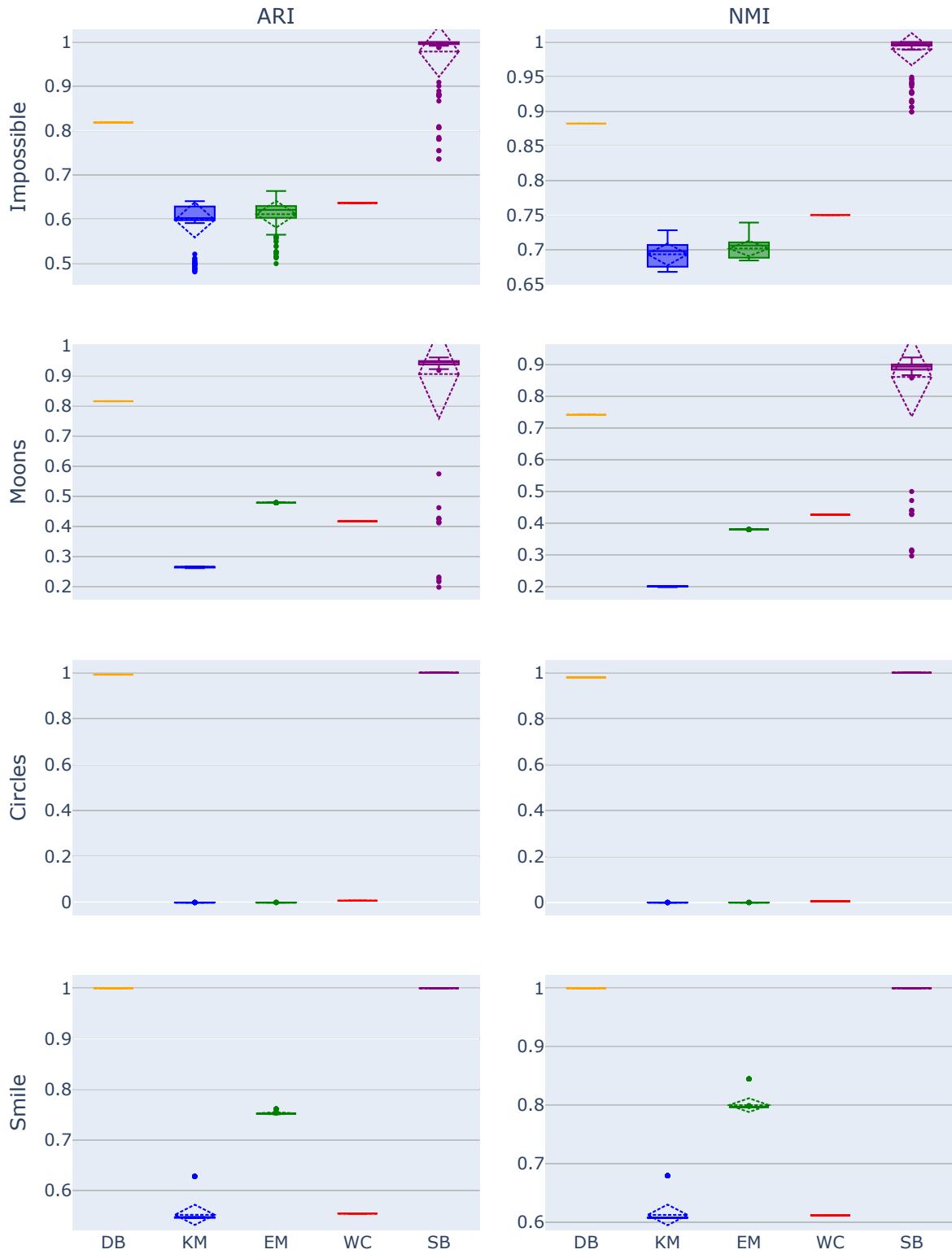


Figure 6: TODO

## 5 Conclusive remarks

Possibility to kernelize

## 6 Appendix

### 6.1 Derivation of the bridge gain

Notice that  $B_{kl}$ , the bridge inertia between centroids  $k$  and  $l$ , can be expressed as the sum of three terms:

$$B_{kl} = \sum_{i|t_i=0} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 + \sum_{i|t_i=1} \|\mathbf{x}_i - \boldsymbol{\mu}_l\|^2 + \sum_{i|t_i \in ]0,1[} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2.$$

The last term may be decomposed in two parts

$$\sum_{i|t_i \in ]0,1[} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2 = \sum_{i|t_i \in ]0, \frac{1}{2}[} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2 + \sum_{i|t_i \in [\frac{1}{2}, 1[} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2$$

and each part further decomposed using Pythagore

$$\begin{aligned} \sum_{i|t_i \in ]0, \frac{1}{2}[} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2 &= \sum_{i|t_i \in ]0, \frac{1}{2}[} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 - \sum_{i|t_i \in ]0, \frac{1}{2}[} \|\boldsymbol{\mu}_k - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2 \\ &= \sum_{i|t_i \in ]0, \frac{1}{2}[} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 - \sum_{i|t_i \in ]0, \frac{1}{2}[} t_i (\boldsymbol{\mu}_k - \boldsymbol{\mu}_l)^2, \\ \sum_{i|t_i \in [\frac{1}{2}, 1[} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2 &= \sum_{i|t_i \in [\frac{1}{2}, 1[} \|\mathbf{x}_i - \boldsymbol{\mu}_l\|^2 - \sum_{i|t_i \in [\frac{1}{2}, 1[} \|\boldsymbol{\mu}_l - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2 \\ &= \sum_{i|t_i \in [\frac{1}{2}, 1[} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 - \sum_{i|t_i \in [\frac{1}{2}, 1[} \|(1 - t_i)(\boldsymbol{\mu}_k - \boldsymbol{\mu}_l)\|^2 \end{aligned}$$

Thus

$$\begin{aligned} B_{kl} - I_{kl} &= \sum_{i|t_i \in ]0, \frac{1}{2}[} t_i^2 \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_l\|^2 + \sum_{i|t_i \in [\frac{1}{2}, 1[} (1 - t_i)^2 \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_l\|^2, \\ \frac{B_{kl} - I_{kl}}{\|\boldsymbol{\mu}_k - \boldsymbol{\mu}_l\|^2} &= \sum_{i|t_i \in ]0, \frac{1}{2}[} t_i^2 + \sum_{i|t_i \in [\frac{1}{2}, 1[} (1 - t_i)^2, \\ \frac{B_{kl} - I_{kl}}{(n_k + n_l) \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_l\|^2} &= \frac{\sum_{i \in k} \langle \mathbf{x}_i - \boldsymbol{\mu}_k | \boldsymbol{\mu}_l - \boldsymbol{\mu}_k \rangle_+^2 + \sum_{i \in l} \langle \mathbf{x}_i - \boldsymbol{\mu}_l | \boldsymbol{\mu}_k - \boldsymbol{\mu}_l \rangle_+^2}{(n_k + n_l) \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_l\|^4}. \end{aligned}$$

## References

- Cortes, Corinna, and Vladimir Vapnik. 1995. "Support-Vector Networks." *Machine Learning* 20 (3): 273–97.
- Ester, Martin, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. 1996. "A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise." In *Kdd*, 96:226–31.
- MacQueen, James et al. 1967. "Some Methods for Classification and Analysis of Multivariate Observations." In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1:281–97. Oakland, CA, USA.
- Von Luxburg, Ulrike. 2007. "A Tutorial on Spectral Clustering." *Statistics and Computing* 17: 395–416.

## 159 Session information

```
160 R version 4.3.2 (2023-10-31)
161 Platform: x86_64-apple-darwin20 (64-bit)
162 Running under: macOS Sonoma 14.3.1
163
164 Matrix products: default
165 BLAS: /Library/Frameworks/R.framework/Versions/4.3-x86_64/Resources/lib/libRblas.0.dylib
166 LAPACK: /Library/Frameworks/R.framework/Versions/4.3-x86_64/Resources/lib/libRlapack.dylib; LAPACK
167
168 locale:
169 [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
170
171 time zone: Europe/Paris
172 tzcode source: internal
173
174 attached base packages:
175 [1] stats graphics grDevices utils datasets methods base
176
177 loaded via a namespace (and not attached):
178 [1] compiler_4.3.2 fastmap_1.1.1 cli_3.6.2 tools_4.3.2
179 [5] htmltools_0.5.7 rstudioapi_0.15.0 yaml_2.3.8 rmarkdown_2.26
180 [9] knitr_1.45 jsonlite_1.8.8 xfun_0.42 digest_0.6.34
181 [13] rlang_1.1.3 evaluate_0.23
```