



Spectral Bridges

Scalable Spectral Clustering free from hyperparameters

ISSN 2824-7795

Félix Laplante Université de Paris Saclay

Christophe Ambroise ¹ Laboratoire de Mathématiques et Modélisation d'Evry, Université Paris-Saclay, CNRS, Univ Evry,

Date published: 2025-07-06 Last modified: 2024-07-09

Abstract

In this paper, Spectral Bridges, a novel clustering algorithm, is introduced. This algorithm builds upon the traditional k-means and spectral clustering frameworks by subdividing data into small Voronoï regions, which are subsequently merged according to a connectivity measure. Drawing inspiration from Support Vector Machine's margin concept, a non-parametric clustering approach is proposed, building an affinity margin between each pair of Voronoï regions. This approach is characterized by minimal hyperparameters and delineation of intricate, non-convex cluster structures.

The numerical experiments underscore Spectral Bridges as a fast, robust, and versatile tool for sophisticated clustering tasks spanning diverse domains. Its efficacy extends to large-scale scenarios encompassing both real-world and synthetic datasets.

The Spectral Bridge algorithm is implemented both in Python (<https://pypi.org/project/spectral-bridges>) and R (<https://github.com/cambroise/spectral-bridges-Rpackage>).

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

1 Contents

2	1 Introduction	2
3	2 Related Work	3
4	3 Spectral Bridges	3
5	3.1 Bridge affinity	4
6	3.2 Algorithm	6

¹Corresponding author: christophe.ambroise@univ-evry.fr

7	4 Numerical experiments	7
8	4.1 Datasets	7
9	4.1.1 Real-world data	7
10	4.1.2 Synthetic data	7
11	4.1.3 Datasets Summary & Class Balance	8
12	4.2 Metrics	8
13	4.3 Platform	8
14	4.4 Hyperparameter settings	8
15	4.5 Time complexity	9
16	4.6 Accuracy	9
17	4.7 Noise robustness	13
18	4.8 Hyperparameter values effect on accuracy	13
19	5 Conclusive remarks	14
20	6 Appendix	14
21	6.1 Derivation of the bridge affinity	14
22	6.2 Code	15
23	6.2.1 Implementation	15
24	6.2.2 Affinity matrix computation	15
25	References	16
26	Session information	17

27 1 Introduction

28 Clustering is a fundamental technique for exploratory data analysis, organizing a set of objects into
29 distinct homogeneous groups known as clusters. It is extensively utilized across various fields, such
30 as biology for gene expression analysis (Eisen et al. 1998), social sciences for community detection in
31 social networks (Latouche, Birmelé, and Ambroise 2011), and psychology for identifying behavioral
32 patterns. Clustering is often employed alongside supervised learning as a pre-processing step, helping
33 to structure and simplify data, thus enhancing the performance and interpretability of subsequent
34 predictive models (Verhaak et al. 2010). Additionally, clustering can be integrated into supervised
35 learning algorithms, such as mixture of experts (Jacobs et al. 1991), as part of a multi-objective
36 strategy.

37 There are various approaches to clustering, and the quality of the results is largely determined by
38 how the similarity between objects is defined, either through a similarity measure or a distance
39 metric. Clustering techniques originate from diverse fields of research, such as genetics, psychometry,
40 statistics, and computer science. Some methods are entirely heuristic, while others aim to optimize
41 specific criteria and can be related to statistical models.

42 Density-based methods identify regions within the data with a high concentration of points, corre-
43 sponding to the modes of the joint density. A notable non-parametric example of this approach is
44 DBSCAN (Ester et al. 1996). In contrast, model-based clustering, such as Gaussian mixture models,
45 represents a parametric approach to density-based methods. Model-based clustering assumes that
46 the data is generated from a mixture of underlying probability distributions, typically Gaussian
47 distributions. Each cluster is viewed as a component of this mixture model, and the Expectation-
48 Maximization (EM) algorithm is often used to estimate the parameters. This approach provides a
49 probabilistic framework for clustering, allowing for the incorporation of prior knowledge and the

50 ability to handle more complex cluster shapes and distributions (McLachlan and Peel 2000).

51 Geometric approaches, such as k-means (MacQueen et al. 1967), are distance-based methods that aim
52 to partition data by optimizing a criterion reflecting group homogeneity. The k-means++ algorithm
53 (Arthur and Vassilvitskii 2006) enhances this approach by providing faster and more reliable results.
54 However, a key limitation of these methods is the assumption of linear boundaries between clusters,
55 implying that clusters are convex. To address non-convex clusters, the kernel trick can be applied,
56 allowing for a more flexible k-means algorithm. This approach is comparable to spectral clustering in
57 handling complex cluster boundaries (Dhillon, Guan, and Kulis 2004). The k-means algorithm can also
58 be interpreted within the framework of model-based clustering under specific assumptions (Govaert
59 and Nadif 2003), revealing that it is essentially a special case of the more general Gaussian mixture
60 models, where clusters are assumed to be spherical Gaussian distributions with equal variance.

61 Graph-based methods represent data as a graph, with vertices symbolizing data points and edges
62 weighted to indicate the affinity between these points. Spectral clustering can be seen as a relaxed
63 version of the graph cut algorithm (Shi and Malik 2000). However, traditional spectral clustering faces
64 significant limitations due to its high time and space complexity, greatly hindering its applicability
65 to large-scale problems (Von Luxburg 2007).

66 The method we propose aims to find non-convex clusters in large datasets, without relying on a
67 parametric model, by using spectral clustering based on an affinity that characterizes the local density
68 of the data. The algorithm described in this paper draws from numerous clustering approaches. The
69 initial intuition is to detect high-density areas. To this end, vector quantization is used to divide the
70 space into a Voronoï tessellation. An original geometric criterion is then employed to detect pairs
71 of Voronoï regions that are either distant from each other or separated by a low-density boundary.
72 Finally, this affinity measure is considered as the weight of an edge in a complete graph connecting
73 the centroids of the tessellation, and a spectral clustering algorithm is used to find a partition of this
74 graph. The only parameters of the algorithm are the number of Voronoï Cells and the number of
75 clusters.

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

79 **2 Related Work**

80 Spectral clustering is a graph-based approach that computes the eigen-vectors of the graph’s Laplacian
81 matrix. This technique transforms the data into a lower-dimensional space, making the clusters
82 more discernible. A standard algorithm like k-means is then applied to these transformed features
83 to identify the clusters (Von Luxburg 2007). Spectral clustering enables capturing complex data
84 structures and discerning clusters based on the connectivity of data points in a transformed space,
85 effectively treating it as a relaxed graph cut problem.

86 Classical spectral clustering involves two phases: construction of the affinity matrix and eigen-
87 decomposition. Constructing the affinity matrix requires $O(n^2d)$ time and $O(n^2)$ memory, while
88 eigen-decomposition demands $O(n^3)$ time and $O(n^2)$ memory, where n is the data size and d is the
89 dimension. As n increases, the computational load escalates significantly (Von Luxburg 2007).

90 To mitigate this computational burden, one common approach is to sparsify the affinity matrix and
91 use sparse eigen-solvers, reducing memory costs but still requiring computation of all original matrix
92 entries (Von Luxburg 2007). Another strategy is sub-matrix construction. The Nyström method
93 randomly selects m representatives from the dataset to form an $n \times m$ affinity sub-matrix (Chen et
94 al. 2010). Cai et al. extended this with the landmark-based spectral clustering method, which uses

95 k-means to determine m cluster centers as representatives (Cai and Chen 2014). Ultra-scalable spectral
 96 clustering (U-SPEC) employs a hybrid representative selection strategy and a fast approximation
 97 method for constructing a sparse affinity sub-matrix (Huang et al. 2019).

98 Other approaches use the properties of the small initial clusters for the affinity computation. Cluster-
 99 ing Based on Graph of Intensity Topology (GIT) estimates for example a global topological graph
 100 (topo-graph) between local clusters (Gao et al. 2021). It then uses the Wasserstein Distance between
 101 predicted and prior class proportions to automatically cut noisy edges in the topo-graph and merge
 102 connected local clusters into final clusters.

103 The issue of characterizing the affinity between two clusters to create an edge weight is central to
 104 the efficiency of a spectral clustering algorithm operating from a submatrix.

105 Notice that the clustering robustness of many Spectral clustering algorithms heavily relies on the
 106 proper selection of kernel parameter, which is difficult to find without prior knowledge (Ng, Jordan,
 107 and Weiss 2001).

108 3 Spectral Bridges

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

115 3.1 Bridge affinity

116 The basic idea involves calculating the difference in inertia achieved by projecting onto a segment
 117 connecting two centroids, rather than using the two centroids separately (see Figure 1). If the
 118 difference is small, it suggests a low density between the classes. Conversely, if this difference is large,
 119 it indicates that the two classes may reside within the same densely populated region.

120 Let us consider a sample $X = (\mathbf{x}_i)_{i \in \{1, \dots, m\}}$ of vectors $\mathbf{x}_i \in \mathbb{R}^d$ and a set of m coding vectors $(\boldsymbol{\mu}_k)_{k \in \{1, \dots, m\}}$
 121 defining a partition $P = \{\mathcal{V}_1, \dots, \mathcal{V}_m\}$ of \mathbb{R}^d into m Voronoï regions:

$$\mathcal{V}_k = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x} - \boldsymbol{\mu}_k\| \leq \|\mathbf{x} - \boldsymbol{\mu}_j\| \text{ for all } j \neq k\}.$$

122 In the following a ball denotes the subset of X in a Voronoï region. The inertia of two balls \mathcal{V}_k and
 123 \mathcal{V}_l is

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

124 We define a bridge as a structure defined by a segment connecting two centroids $\boldsymbol{\mu}_k$ and $\boldsymbol{\mu}_l$. The
 125 inertia of a bridge between \mathcal{V}_k and \mathcal{V}_l is defined as

$$B_{kl} = \sum_{\mathbf{x}_i \in \mathcal{V}_k \cup \mathcal{V}_l} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2,$$

126 where

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

127 with

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

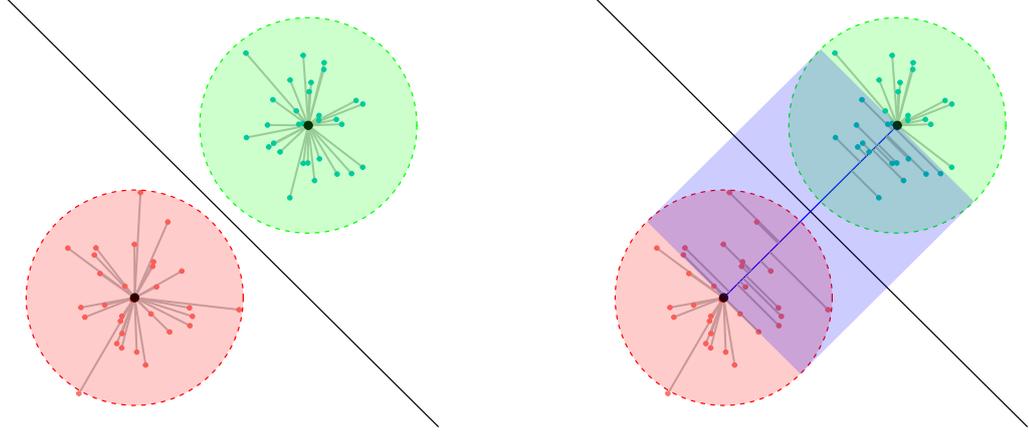


Figure 1: Balls (left) versus Bridge (right). The inertia of each structure is the sum of the squared distances represented by grey lines.

128 Considering two centroids, the normalized average of the difference between
 129 inertia (See [Appendix](#)) constitutes the basis of our affinity measure between two regions:

$$\begin{aligned} \frac{B_{kl} - I_{kl}}{(n_k + n_l)\|\mu_k - \mu_l\|^2} &= \frac{\sum_{\mathbf{x}_i \in \mathcal{V}_k} \langle \mathbf{x}_i - \mu_k | \mu_l - \mu_k \rangle_+^2 + \sum_{\mathbf{x}_i \in \mathcal{V}_l} \langle \mathbf{x}_i - \mu_l | \mu_k - \mu_l \rangle_+^2}{(n_k + n_l)\|\mu_k - \mu_l\|^4}, \\ &= \frac{\sum_{\mathbf{x}_i \in \mathcal{V}_k \cup \mathcal{V}_l} \alpha_i^2}{n_k + n_l}, \end{aligned}$$

130 where

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

131 The basic intuition behind this affinity is that t_i represents the relative position of the projection of \mathbf{x}_i
 132 on the segment $[\mu_k, \mu_l]$. α_i represents the relative position on the segment, with the centroid of the
 133 class to which \mathbf{x}_i belongs as the reference point.

134 The boundary that separates the two clusters defined by centroids μ_k and μ_l is a hyperplane. This
 135 hyperplane is orthogonal to the line segment connecting the centroids and intersects this segment at
 136 its midpoint.

137 If we consider all points $\mathbf{x}_i \in \mathcal{V}_k \cup \mathcal{V}_l$ which are not projected on centroids but somewhere on the
 138 segment, the distance from a point to the hyperplane is

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

139 This distance is similar to the concept of margin in Support Vector Machine (Cortes and Vapnik 1995).
 140 When the α_i values are small (close to zero since $\alpha_i \in [0, 1/2]$), the margins to the hyperplane are
 141 large, indicating a low density between the classes. Conversely, if the margins are small, it suggests
 142 that the two classes may reside within the same densely populated region. Consequently, the sum of
 143 the α_i or α_i^2 increases with the density of the region between the classes.

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

146 Eventually, we define the bridge affinity between centroids k and l as:

$$a_{kl} = \begin{cases} 0, & \text{if } k = l, \\ \frac{\sum_{\mathbf{x}_i \in \mathcal{V}_k \cup \mathcal{V}_l} \alpha_i^2}{n_k + n_l}, & \text{otherwise.} \end{cases}$$

147 To allow points with large margin to dominate and make the algorithm more robust to noise and
 148 outliers we consider the following exponential transformation:

$$\tilde{a}_{kl} = g(a_{kl}) = \exp(\gamma \sqrt{a_{kl}}).$$

149 where γ is a scaling factor. This factor is set to ensure a large enough separation between the final
 150 coefficients. This factor is determined by the equation:

$$\gamma = \frac{\log(M)}{\sqrt{q_{90}} - \sqrt{q_{10}}}$$

151 where q_{10} and q_{90} are respectively the 10th and 90th percentiles of the original affinity matrix
 152 and $M > 1$. Thus, since the transformation is order-preserving, the 90th percentile of the newly
 153 constructed matrix is M times greater than the 10th percentile. By default, M is arbitrarily set to a
 154 large value of 10^4 .

155 The inclusion of the square root can be understood as redefining the affinity measure. Instead of
 156 considering the variance and the squared Euclidean norm, we interpret the affinity as the ratio
 157 between the standard deviation and the length of the segment connecting two centroids. This
 158 reinterpretation greatly enhances numerical stability, contributing to more reliable clustering results.

159 3.2 Algorithm

160 The Spectral Bridges algorithm first identifies local clusters to define Voronoï regions, computes
 161 edges with affinity weights between these regions, and ultimately cuts edges between regions with
 162 low inter-region density to determine the final clusters (See Algorithm 1 and Figure 2).

163 In spectral clustering, the time complexity is usually dominated by the eigen-decomposition step,
 164 which is $O(n^3)$. However, in the case of Spectral Bridges, the k-means algorithm has a time complexity
 165 of $O(n \times m \times d)$. For datasets with large n , this can be more significant than the $O(m^3)$ time complexity
 166 of the Spectral Bridges eigen-decomposition. As for the affinity matrix construction, there are m^2
 167 coefficients to be calculated. Each a_{kl} coefficient requires the computation of $n_k + n_l$ dot products as
 168 well as the norm $\|\boldsymbol{\mu}_k - \boldsymbol{\mu}_l\|$, the latter often being negligible. Assuming that the Voronoï regions are
 169 roughly balanced in cardinality, we have $n_k \approx \frac{n}{m}$. Since m should always be less than n , therefore
 170 $\frac{n}{m} > 1$ and the time complexity of the affinity matrix is $O(\frac{n}{m} \times m^2 \times d) = O(n \times m \times d)$ given the
 171 acceptable range of values for m . Nonetheless, this is rarely the bottleneck.

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, voronoiRegions  $\leftarrow$  KMEANS( $X, m$ )  $\triangleright$  Initial centroids and Voronoi regions using  
   k-means++  
4:   Step 2: Affinity Computation  
5:    $A = \{g(a_{kl})\}_{kl} \leftarrow$  AFFINITY( $X, \text{centroids}, \text{voronoiRegions}$ )  $\triangleright$  Compute affinity matrix  $A$   
6:   Step 3: Spectral Clustering  $\triangleright$  Assign each region to a cluster  
7:   labels  $\leftarrow$  SPECTRALCLUSTERING( $A, k$ )  
8:   Step 4: Propagate  $\triangleright$  Assign each data point to the cluster of its region  
9:   clusters  $\leftarrow$  PROPAGATE( $X, \text{labels}, \text{voronoiRegions}$ )  
10:  return clusters  $\triangleright$  Return cluster labels for data points in  $X$   
11: end procedure
```

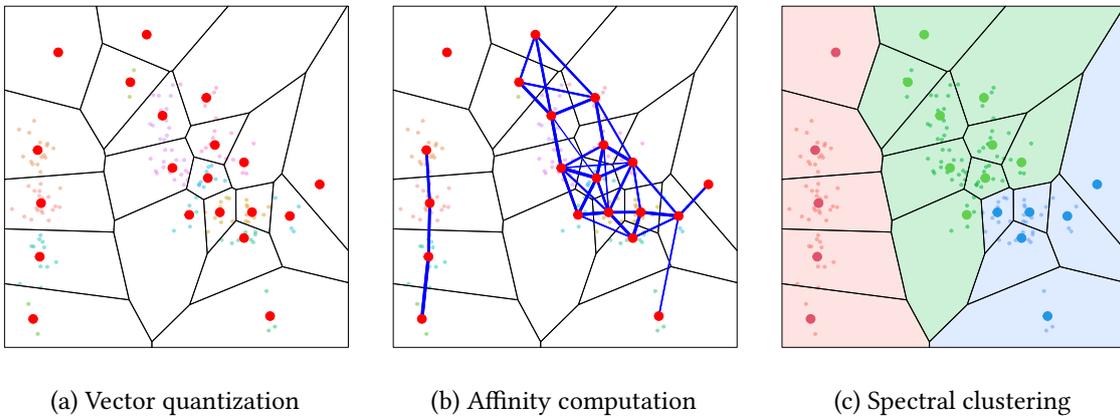


Figure 2: Illustration of the Spectral bridges algorithm with the Iris dataset (first principal plane). Vector quantization (Step 1 of Algorithm 1), Affinity computation (Step 2 of Algorithm 1), Spectral clustering and spreading (Step 3-4 of Algorithm 1).

4 Numerical experiments

In this section, the results obtained from testing the Spectral Bridges algorithm on various datasets, both small and large scale, including real-world and well-known synthetic datasets, are presented. These experiments assess the accuracy, time and space complexity, ease of use, robustness, and adaptability of our algorithm. We compare Spectral Bridges (SB) against several state-of-the-art methods, including k-means++ (KM) (MacQueen et al. 1967; Arthur and Vassilvitskii 2006), Expectation-Maximization (EM) (Dempster, Laird, and Rubin 1977), Ward Clustering (WC) (Ward Jr 1963), and DBSCAN (DB) (Ester et al. 1996). This comparison establishes baselines across centroid-based clustering algorithms, hierarchical methods, and density-based methods.

The algorithms are evaluated on both raw and PCA-processed data with varying dimensionality. For synthetic datasets, Gaussian and/or uniform noise is introduced to assess the robustness of the algorithm.

184 4.1 Datasets

185 4.1.1 Real-world data

- 186 • **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.
- 187 • **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.

191 4.1.2 Synthetic data

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

198 4.1.3 Datasets Summary & Class Balance

Table 1: Datasets Summary & 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%

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

200 4.2 Metrics

201 To evaluate the performance of the clustering algorithm, the Adjusted Rand Index (ARI) (Halkidi, Batistakis, and Vazirgiannis 2002) and Normalized Mutual Information (NMI) (Cover and Thomas 1991) are used. 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, the variability of scores across multiple runs is also reported due to the random initialization in k-means, though k-means++ generally provides stable and reproducible results.

208 4.3 Platform

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

211 4.4 Hyperparameter settings

212 The hyperparameters of the Spectral Bridges algorithm were based on the size of each dataset, n ,
213 and the number of clusters, K . A larger number of clusters typically suggests that a higher value
214 for the number of Voronoï regions is optimal. Conversely, using a high number of Voronoï regions
215 for a small dataset might result in nearly empty regions that do not adequately represent any local
216 structure.

217 A good yet not very precise way of setting the number of Voronoï regions m is to observe the Within
218 Cluster Sum of Squares (WCSS) or inertia in a way akin to the elbow method. Since m should be set
219 to a value strictly greater than K , we plot the WCSS for varying values of m , and find a value such
220 that the WCSS- m relationship becomes quasi-linear.

221 By adjusting m in this manner, we aim to balance the need for detailed representation with the
222 risk of overfitting, ensuring that each Voronoï region meaningfully captures the underlying data
223 distribution. The sensitivity or lack thereof is illustrated later on by Figure 10.

224 For other algorithms, such as DBSCAN, labels were used to determine the best hyperparameter
225 values to compare our method against the “best case scenario”, thus putting the Spectral Bridges
226 algorithm at a voluntary disadvantage.

227 4.5 Time complexity

228 To assess the algorithm’s time complexity, the average execution times over 50 runs were computed
229 for varying numbers of Voronoï regions m as well as dataset sizes. With a constant number of clusters
230 $K = 5$ and an embedding dimension of $d = 10$, the results (see Figure 3) highlight Spectral Bridges
231 algorithm’s efficacy. As discussed previously, we observe a linear relationship between m and the
232 execution time because the matrix construction is highly optimized and the time taken is almost
233 negligible compared to that of the initial k-means++ centroids initialization.

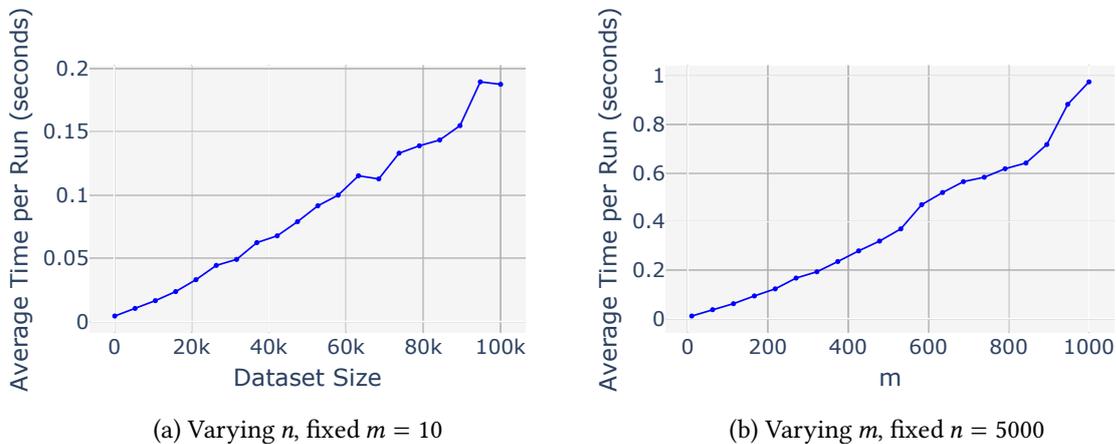


Figure 3: Average time taken per model fit.

234 4.6 Accuracy

235 The algorithm’s accuracy was first evaluated on the MNIST dataset. Metrics were collected to
236 compare our method with k-means++, EM, and Ward clustering. Metric were estimated by taking
237 the empirical average over 10 consecutive runs with the same random seed for each method. Since
238 our computational capabilities were too limited, a sample of 20,000 (one third) data points was chosen
239 at random for each iteration.

240 Let h denote the embedding dimension of the dataset. Spectral Bridges was tested both on the raw
 241 MNIST dataset without preprocessing ($h = 784$) and after reducing its dimension using PCA to
 242 $h \in \{8, 16, 32, 64\}$ (see Figure 4).

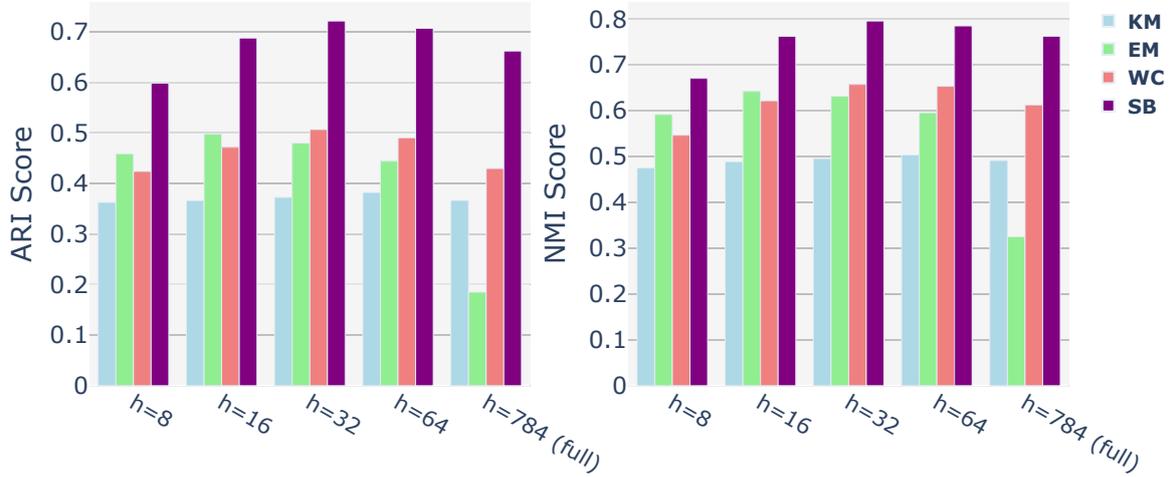


Figure 4: ARI and NMI scores of k-means++ (blue), EM (green), Ward Clustering (red), and Spectral Bridges (purple) on PCA embedding and full MNIST.

243 For visualization purposes, the predicted clusters by Spectral Bridges and k-means++ were projected
 244 using UMAP to compare them against the ground truth labels and to better understand the cluster
 245 shapes (see Figure 5). Note that the projection was not used in the experiments as an embedding, and
 246 thus does not play any role in the clustering process itself. As a matter of fact, the embedding used
 247 was obtained with PCA, $h = 32$ and 250 Voronoi regions. Note that the label colors match the legend
 248 only in the case of the ground truth data. Indeed, the ordering of the labels have no significance on
 249 clustering quality.

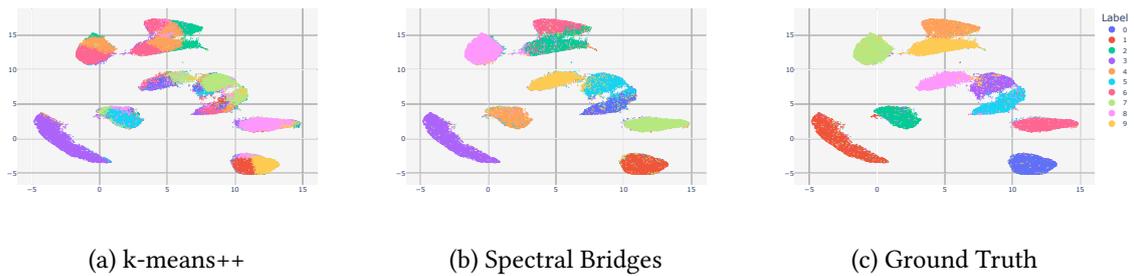


Figure 5: UMAP projection of predicted clusters against the ground truth labels.

250 The Spectral Bridges algorithm was also put to the test against the same competitors using scikit-
 251 learn’s UCI Breast Cancer data. Once again, this new method performs well although the advantage
 252 is not as obvious in this case (see Figure 6). However, in none of our tests has it ranked worse than
 253 k-means++. The results are displayed as a boxplot generated from 200 iterations of each algorithm
 254 using a different seed, in order to better grasp the variability lying in the seed dependent nature of
 255 the k-means++, Expectation Maximization and Spectral Bridges algorithms.

256 Since the Spectral Bridges algorithm is expected to excel at discerning complex and intricate cluster

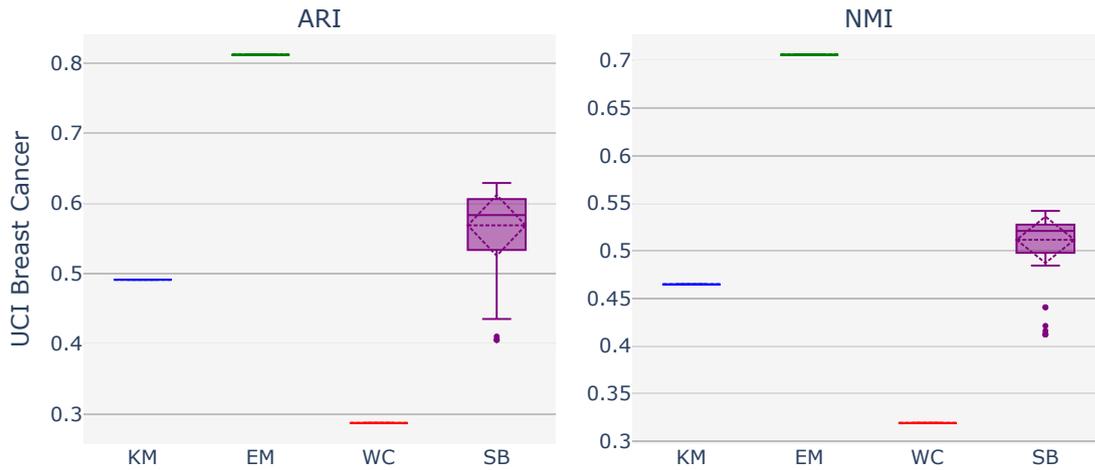


Figure 6: ARI and NMI scores of k-means++ (blue), EM (green), Ward Clustering (red), and Spectral Bridges (purple) on the UCI Breast Cancer dataset.

257 structures, an array of four toy datasets was collected, as illustrated in Figure 7.

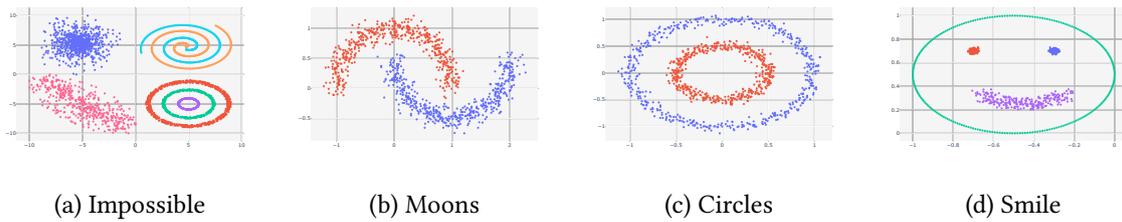


Figure 7: Four toy datasets.

258 Multiple algorithms, including the proposed one, were benchmarked in the exact same manner
 259 as for the UCI Breast Cancer data. The results show that the proposed method outperforms all
 260 tested algorithms (DBSCAN, k-means++, Expectation Maximization, and Ward Clustering) while
 261 requiring few hyperparameters. As previously discussed, DBSCAN’s parameters were optimized
 262 using the ground truth labels to represent a best-case scenario; however, in practical applications,
 263 suboptimal performance is more likely. Despite this optimization, the Spectral-Bridge algorithm still
 264 demonstrates superior ability to capture and represent the underlying cluster structures.

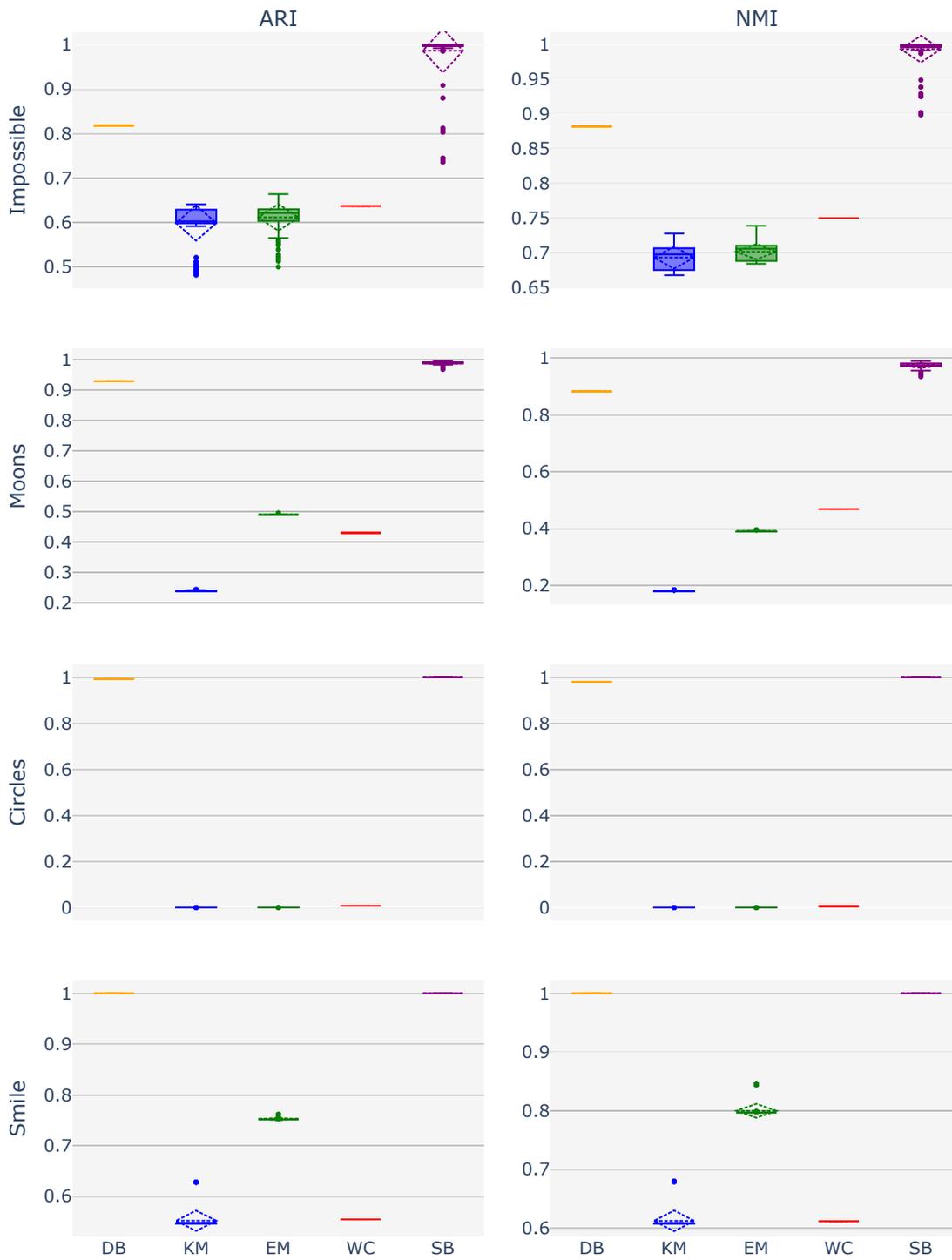


Figure 8: ARI and NMI scores of Spectral Bridges and competitors on standard synthetic toy datasets.

265 **4.7 Noise robustness**

266 To evaluate the noise robustness of the algorithm, two experimental setups were devised: one involved
 267 introducing Gaussian-distributed perturbations to the data, and the other involved concatenating
 268 uniformly distributed points within a predefined rectangular region (determined by the span of the
 269 dataset) to the existing dataset. As illustrated in Figure 9, the tests demonstrate that in both scenarios,
 270 the algorithm exhibits a high degree of insensitivity to noise.

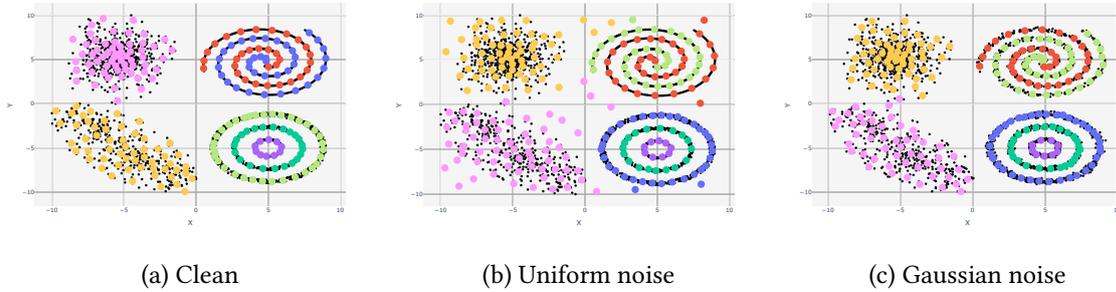


Figure 9: Three representations of the algorithm’s predicted cluster centers are displayed as colored dots, with each point of the Impossible dataset shown as a small black dot. In the left graph, the dataset is unmodified. In the center graph, 250 uniformly distributed samples were added. In the right graph, Gaussian noise perturbations with $\sigma = 0.1$ were applied.

271 **4.8 Hyperparameter values effect on accuracy**

272 To better understand and measure the significance of choosing the right values for the hyper-
 273 parameters of the proposed algorithm, that it to say the number of Voronoï regions m , Spectral
 274 Bridges was run on the PCA $h = 32$ embedded MNIST dataset with varying values of $m \in$
 275 $\{10, 120, 230, 340, 450, 560, 670, 780, 890, 1000\}$. The case $m = 10$ is equivalent to the k-means++ algo-
 276 rithm. ARI and NMI scores are recorded over 20 consecutive iterations and subsequently plotted. As
 277 shown by Figure 10, the accuracy seems to be consistently increasing with values of m , although
 278 the largest observed gap occurs between values of $m = 10$ and $m = 120$, indicating a tremendous im-
 279 provement over the classical k-means++ framework even for empirically suboptimal hyperparameter
 280 values.

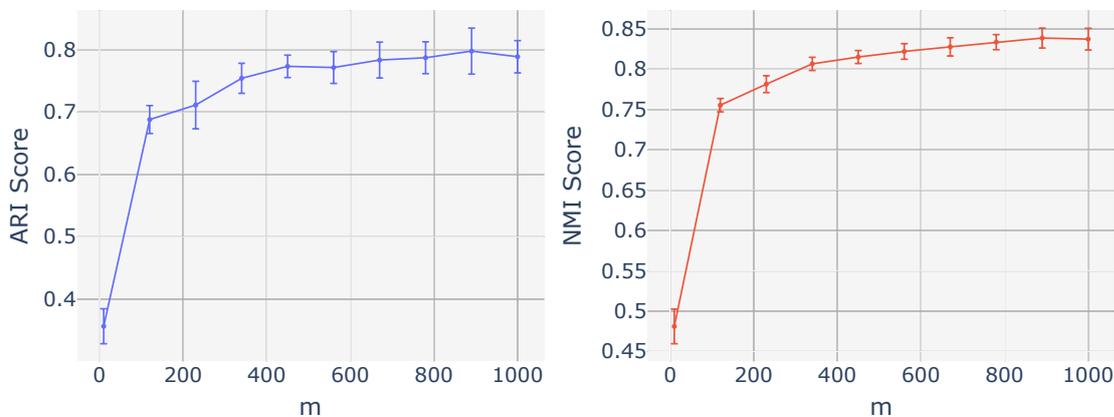


Figure 10: ARI and NMI scores of Spectral Bridges with varying values of m .

5 Conclusive remarks

Spectral Bridges is an original clustering algorithm which presents a novel approach by integrating the strengths of traditional k-means and spectral clustering frameworks. This algorithm utilizes a simple affinity measure for spectral clustering, which is derived from the minimal margin between pairs of Voronoi regions.

The algorithm demonstrates scalability, handling large datasets efficiently through a balanced computational complexity between the k-means clustering and eigen-decomposition steps. As a non-parametric method, Spectral Bridges does not rely on strong assumptions about data distribution, enhancing its versatility across various data types. It performs exceptionally well with both synthetic and real-world data and consistently outperforms conventional clustering algorithms such as k-means, DBSCAN, and mixture models.

The design of Spectral Bridges ensures robustness to noise, a significant advantage in real-world applications. Additionally, the algorithm requires minimal hyperparameters, primarily the number of Voronoi regions, making it straightforward to tune and deploy.

Furthermore, Spectral Bridges can be kernelized, allowing it to handle data in similarity space directly, which enhances its flexibility and applicability. Overall, Spectral Bridges is a powerful, robust, and scalable clustering algorithm that offers significant improvements over traditional methods, making it an excellent tool for advanced clustering tasks across numerous domains.

6 Appendix

6.1 Derivation of the bridge affinity

We denote a bridge as a segment connecting two centroids μ_k and μ_l . The inertia of a bridge between \mathcal{V}_k and \mathcal{V}_l is defined as

$$B_{kl} = \sum_{\mathbf{x}_i \in \mathcal{V}_k \cup \mathcal{V}_l} \|\mathbf{x}_i - \mathbf{p}_{kl}(\mathbf{x}_i)\|^2,$$

where

$$\mathbf{p}_{kl}(\mathbf{x}_i) = \mu_k + t_i(\mu_l - \mu_k),$$

with

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

B_{kl} , the bridge inertia between centroids k and l , can be expressed as the sum of three terms, which represents the projection onto each centroids and onto the segment:

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

The last term may be decomposed in two parts corresponding to the points of the two Voronoi regions which are projected on the segment:

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

309 and each part further decomposed using Pythagore

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

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

310 Thus

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

311 6.2 Code

312 6.2.1 Implementation

313 Numerical experiments have been conducted in Python. The python scripts to reproduce the
314 simulations and figures are available at <https://github.com/flheight/Spectral-Bridges>. The Spectral
315 Bridge algorithm is implemented both in

- 316 • Python: <https://pypi.org/project/spectral-bridges>, and
- 317 • R: <https://github.com/cambroise/spectral-bridges-Rpackage>.

318 6.2.2 Affinity matrix computation

319 Taking a closer look at the second step of Algorithm 1, that is the affinity matrix calculation
320 with a $O(n \times m \times d)$ time complexity, most operations can be parallelized leaving a single loop,
321 bundling together m^2 dot products into only m matrix multiplications, thus allowing for an efficient
322 construction in both high and low level programming languages. Though the complexity of the
323 algorithm remains unchanged, libraries such as Basic Linear Algebra Subprograms can render the
324 calculations orders of magnitude faster. Moreover, the symmetrical nature of the bridge affinity can
325 be used to effectively halve the computation time.

326 The calculation of the affinity matrix is highlighted by the Python code Listing 1. Though it could
327 be even more optimized, the following code snippet is approximately 200 times faster than a naive
328 implementation on a small dataset comprised of $n = 3594$, $d = 2$ points, and a value of $m = 250$.

329 Notice that the Python code is significantly faster than the R code.

Listing 1 Python code for affinity matrix computation

```
# Initialize the affinity matrix
affinity = np.empty((self.n_nodes, self.n_nodes))

# Center each Voronoi region around its centroid
X_centered = [
    X[kmeans.labels_ == i] - kmeans.cluster_centers_[i] for i in range(self.n_nodes)
]

# Count the total number of points in each pair of regions
counts = np.array([X_centered[i].shape[0] for i in range(self.n_nodes)])
counts = counts[np.newaxis, :] + counts[:, np.newaxis]

# Compute the segments between each pair of centroids and their squared Euclidean norm
segments = (
    kmeans.cluster_centers_[np.newaxis, :] - kmeans.cluster_centers_[:, np.newaxis]
)
dists = np.einsum("ijk,ijk->ij", segments, segments)
np.fill_diagonal(dists, 1) # Avoid dividing by zero

# Assign each row of the affinity matrix
for i in range(self.n_nodes):
    projs = np.maximum(np.dot(X_centered[i], segments[i].T), 0)
    affinity[i] = np.einsum("ij,ij->j", projs, projs)

# Symmetrize the matrix and normalize, as well as taking the element-wise square root
affinity = np.sqrt(affinity + affinity.T) / (np.sqrt(counts) * dists)
affinity -= 0.5 * affinity.max() # For numerical stability

# Apply the exponential transformation
q10, q90 = np.quantile(affinity, [0.1, 0.9])

gamma = np.log(self.M) / (q90 - q10)
affinity = np.exp(gamma * affinity)
```

330 **References**

- 331 Arthur, David, and Sergei Vassilvitskii. 2006. "K-Means++: The Advantages of Careful Seeding."
332 Technical Report 2006-13. Stanford InfoLab; Stanford. <http://ilpubs.stanford.edu:8090/778/>.
- 333 Cai, Deng, and Xinlei Chen. 2014. "Large Scale Spectral Clustering via Landmark-Based Sparse
334 Representation." *IEEE Transactions on Cybernetics* 45 (8): 1669–80.
- 335 Chen, Wen-Yen, Yangqiu Song, Hongjie Bai, Chih-Jen Lin, and Edward Y Chang. 2010. "Parallel
336 Spectral Clustering in Distributed Systems." *IEEE Transactions on Pattern Analysis and Machine
337 Intelligence* 33 (3): 568–86.
- 338 Cortes, Corinna, and Vladimir Vapnik. 1995. "Support-Vector Networks." *Machine Learning* 20 (3):
339 273–97.
- 340 Cover, Thomas M, and Joy A Thomas. 1991. "Information Theory and the Stock Market." *Elements of
341 Information Theory*. Wiley Inc., New York, 543–56.
- 342 Dempster, Arthur P, Nan M Laird, and Donald B Rubin. 1977. "Maximum Likelihood from Incomplete
343 Data via the EM Algorithm." *Journal of the Royal Statistical Society: Series B (Methodological)* 39

- (1): 1–22.
- Dhillon, Inderjit S, Yuqiang Guan, and Brian Kulis. 2004. “Kernel k-Means, Spectral Clustering and Normalized Cuts.” In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 551–56. ACM.
- Eisen, Michael B., Paul T. Spellman, Patrick O. Brown, and David Botstein. 1998. “Cluster Analysis and Display of Genome-Wide Expression Patterns.” *Proceedings of the National Academy of Sciences* 95 (25): 14863–68.
- 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.
- Gao, Zhangyang, Haitao Lin, Cheng Tan, Lirong Wu, Stan Li, et al. 2021. “Git: Clustering Based on Graph of Intensity Topology.” *arXiv Preprint arXiv:2110.01274*.
- Govaert, Gérard, and Mohamed Nadif. 2003. “Clustering with Block Mixture Models.” *Pattern Recognition* 36 (2): 463–73.
- Halkidi, Maria, Yannis Batistakis, and Michalis Vazirgiannis. 2002. “Cluster Validity Methods: Part i.” *ACM SIGMOD Record* 31 (2): 40–45.
- Huang, Dong, Chang-Dong Wang, Jian-Sheng Wu, Jian-Huang Lai, and Chee-Keong Kwoh. 2019. “Ultra-Scalable Spectral Clustering and Ensemble Clustering.” *IEEE Transactions on Knowledge and Data Engineering* 32 (6): 1212–26.
- Jacobs, Robert A, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. 1991. “Adaptive Mixtures of Local Experts.” *Neural Computation* 3 (1): 79–87.
- Latouche, Pierre, Etienne Birmelé, and Christophe Ambroise. 2011. “Overlapping stochastic block models with application to the French political blogosphere.” *The Annals of Applied Statistics* 5 (1): 309–36. <https://doi.org/10.1214/10-AOAS382>.
- 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.
- McLachlan, Geoffrey J., and David Peel. 2000. *Finite Mixture Models*. New York: Wiley-Interscience.
- Ng, Andrew, Michael Jordan, and Yair Weiss. 2001. “On Spectral Clustering: Analysis and an Algorithm.” *Advances in Neural Information Processing Systems* 14.
- Shi, Jianbo, and Jitendra Malik. 2000. “Normalized Cuts and Image Segmentation.” *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (8): 888–905.
- Verhaak, Roel G. W., Katherine A. Hoadley, Elizabeth Purdom, Victoria Wang, Yuexin Qi, Matthew D. Wilkerson, Charlie R. Miller, et al. 2010. “Integrated Genomic Analysis Identifies Clinically Relevant Subtypes of Glioblastoma Characterized by Abnormalities in PDGFRA, IDH1, EGFR, and NF1.” *Cancer Cell* 17 (1): 98–110.
- Von Luxburg, Ulrike. 2007. “A Tutorial on Spectral Clustering.” *Statistics and Computing* 17: 395–416.
- Ward Jr, Joe H. 1963. “Hierarchical Grouping to Optimize an Objective Function.” *Journal of the American Statistical Association* 58 (301): 236–44.

382 Session information

383 R version 4.3.2 (2023-10-31)
384 Platform: x86_64-apple-darwin20 (64-bit)
385 Running under: macOS Sonoma 14.3.1
386
387 Matrix products: default
388 BLAS: /Library/Frameworks/R.framework/Versions/4.3-x86_64/Resources/lib/libRblas.0.dylib
389 LAPACK: /Library/Frameworks/R.framework/Versions/4.3-x86_64/Resources/lib/libRlapack.dylib; LAPACK
390
391 locale:

```
392 [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
393
394 time zone: Europe/Paris
395 tzcode source: internal
396
397 attached base packages:
398 [1] stats      graphics  grDevices  utils      datasets  methods   base
399
400 loaded via a namespace (and not attached):
401 [1] digest_0.6.34      fastmap_1.1.1      xfun_0.42         Matrix_1.6-5
402 [5] lattice_0.22-5     reticulate_1.37.0  knitr_1.45        htmltools_0.5.7
403 [9] png_0.1-8          rmarkdown_2.26     cli_3.6.2         grid_4.3.2
404 [13] compiler_4.3.2     rstudioapi_0.15.0  tools_4.3.2       evaluate_0.23
405 [17] Rcpp_1.0.12        yaml_2.3.8         rlang_1.1.3       jsonlite_1.8.8
```