

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: 2024-11-10

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 hyperparameter sensitivity and delineation of intricate, non-convex cluster structures.

The numerical experiments underscore Spectral Bridges as a fast, robust, and versatile tool for 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

Contents

1	1 Introduction	2
2	2 Related Work	3
3	3 Spectral Bridges	3
4	3.1 Bridge affinity	4
5	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	13
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	17
26	Session information	18

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

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

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

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

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

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

2 Related Work

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

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

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

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

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

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

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

3 Spectral Bridges

The proposed algorithm uses k-means centroids for vector quantization defining Voronoï region, and a strategy is proposed to link these regions, with an “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 Voronoï region and the number of final cluster.

3.1 Bridge 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 (see Figure 1). 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.

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

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

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

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

We define a bridge as a structure defined by 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

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

where

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

with

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

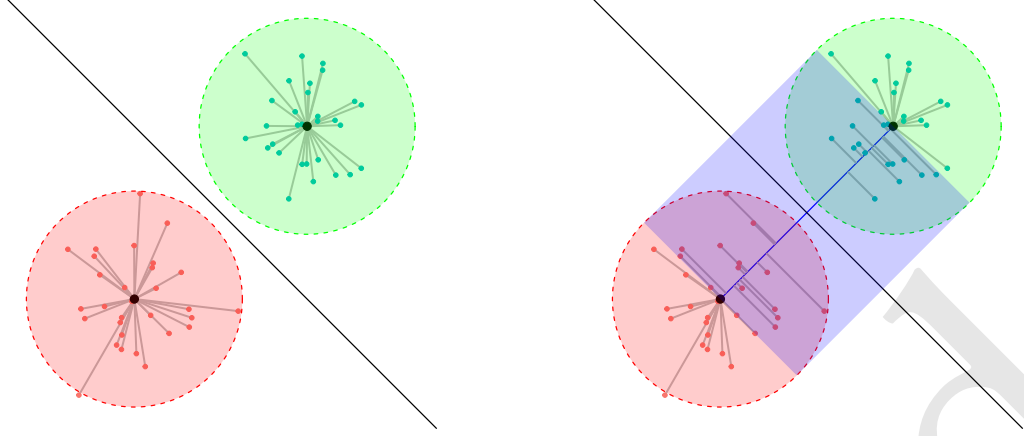


Figure 1: Balls (left) versus Bridge (right). The inertia of each structure is the sum of the squared distances represented by grey lines.

Considering two centroids, the normalized average of the difference between Bridge and balls inertia (see [Appendix](#)) constitutes the basis of our affinity measure between two regions:

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

where

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

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

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

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 segment, the distance from a point to the hyperplane is

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

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

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

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

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

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

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

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

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

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

The inclusion of the square root can be understood as redefining the affinity measure. Instead of considering the variance and the squared Euclidean norm, we interpret the affinity as the ratio between the standard deviation and the length of the segment connecting two centroids.

This regularization is crucial: with a bounded affinity metric, exponentiation enhances the separation between low and high-density regions, controlled by a scaling parameter, as in traditional spectral clustering. Redefining the metric with a square root helps mitigate a converse issue, where machine error can cause numerical instability when solving the Laplacian eigenproblem, especially if values become too large. Since the range of affinity values can become too wide when the initial ratio between the largest and smallest non-zero unscaled bridge affinities is high, applying the square root reduces the maximum values in the affinity matrix while preserving the metric's interpretability and distance-like properties; importantly, this adjustment is not intended for outlier detection.

3.2 Algorithm

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

In spectral clustering, the time complexity is usually dominated by the eigen-decomposition step, which is $O(n^3)$. However, in the case of Spectral Bridges, the k-means algorithm has a time complexity of $O(n \times m \times d)$. For datasets with large n , this can be more significant than the $O(m^3)$ time complexity of the Spectral Bridges eigen-decomposition. As for the affinity matrix construction, there are m^2 coefficients to be calculated. Each a_{kl} coefficient requires the computation of $n_k + n_l$ dot products as well as the norm $\|\mu_k - \mu_l\|$, the latter often being negligible. Assuming that the Voronoï regions are roughly balanced in cardinality, we have $n_k \approx \frac{n}{m}$. Since m should always be less than n , therefore $\frac{n}{m} > 1$ 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 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( $X$ , centroids, 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( $X$ , labels, voronoiRegions)  
10:  return clusters  $\triangleright$  Return cluster labels for data points in  $X$   
11: end procedure
```

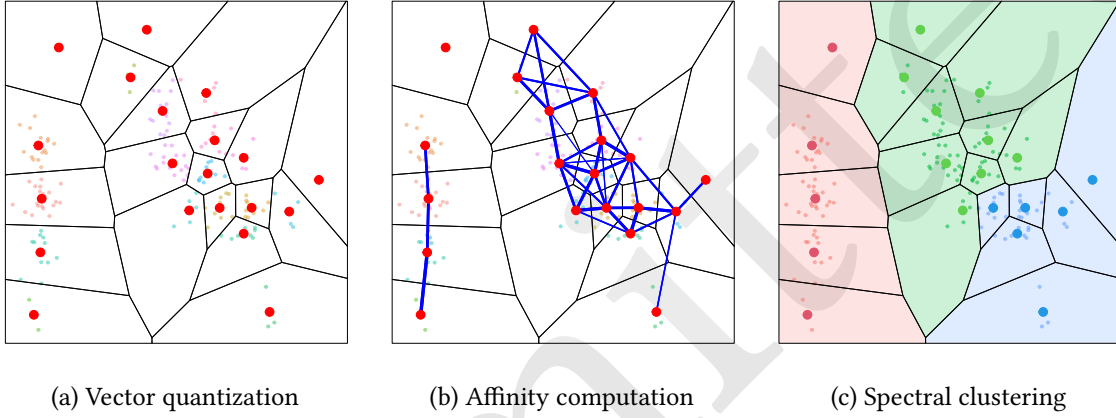


Figure 2: Illustration of the Spectral bridges algorithm with the Iris dataset (first principal plane). Vector quantization (Step 1 of Algorithm 1), Affinity computation (Step 2 of Algorithm 1), Spectral clustering and spreading (Step 3-4 of Algorithm 1).

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.

4.1 Datasets

4.1.1 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.

4.1.2 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.

4.1.3 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%

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

4.2 Metrics

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.

4.3 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.

4.4 Hyperparameter settings

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

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

By adjusting m in this manner, we aim to balance the need for detailed representation with the risk of overfitting, ensuring that each Voronoï region meaningfully captures the underlying data distribution. The sensitivity or lack thereof is illustrated later on by Figure 11.

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

4.5 Time complexity

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

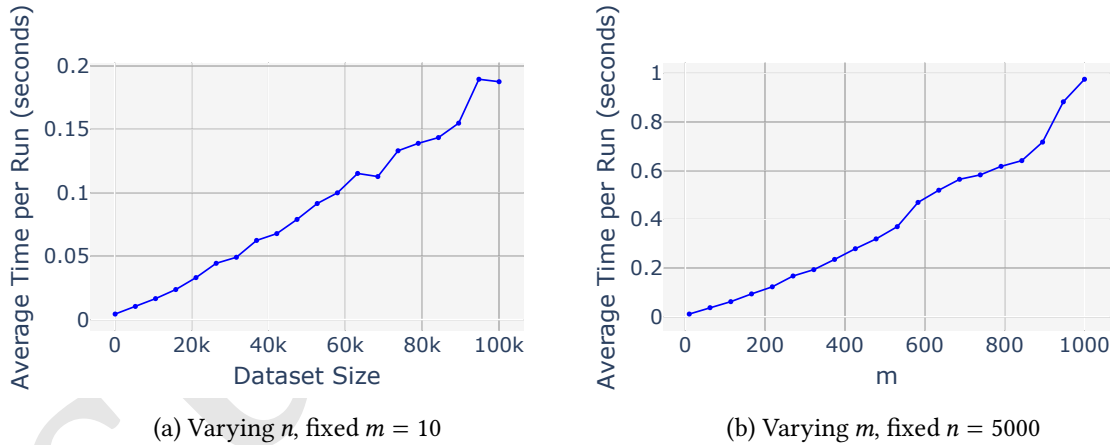


Figure 3: Average time taken per model fit.

4.6 Accuracy

The algorithm’s accuracy was first evaluated 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 10 consecutive runs with a random seed set for reproducibility purposes.

Let h denote the embedding dimension of the dataset. Spectral Bridges was tested both on the raw

MNIST dataset without preprocessing ($h = 784$) and after reducing its dimension using PCA to $h \in \{8, 16, 32, 64\}$ (see Figure 4).

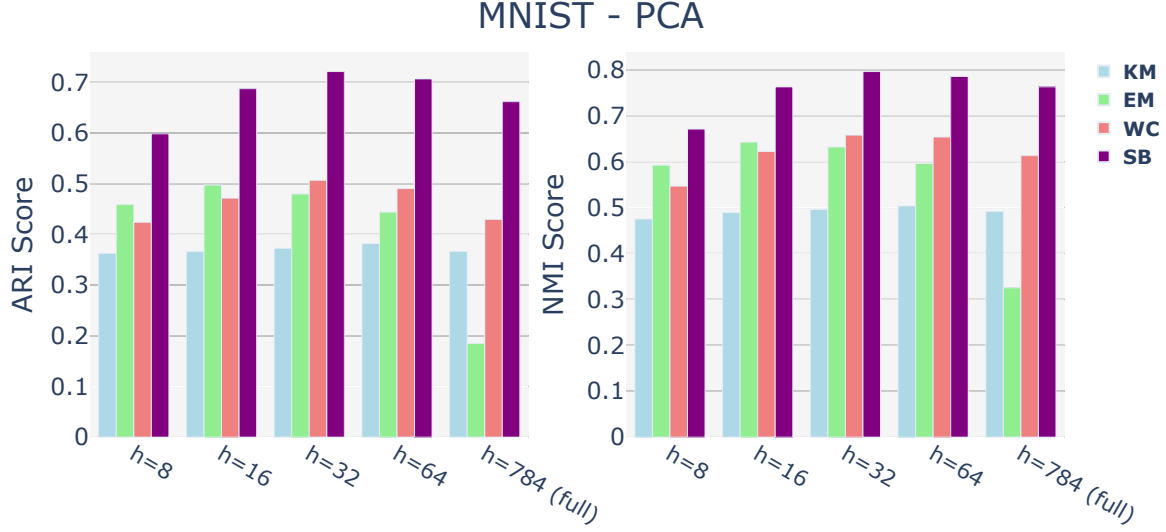


Figure 4: ARI and NMI scores of k-means++ (blue), EM (green), Ward Clustering (red), and Spectral Bridges (purple) on PCA embedding and full MNIST.

Furthermore, the proposed algorithm was also tested on the same MNIST dataset after reducing its dimension to $h \in \{2, 4, 8, 16\}$ with UMAP (McInnes et al. 2018), one of the most advanced non-linear dimension reduction algorithm (see Figure 5).

Note the Spectral Bridges is substantially better than other traditional methods and shines even with quite simple dimension reduction algorithms.

For visualization purposes, the predicted clusters by Spectral Bridges and k-means++ were projected using UMAP to compare them against the ground truth labels and to better understand the cluster shapes (see Figure 6). Note this projection was not used in the 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, $h = 32$ and 250 Voronoï regions. 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 significance on clustering quality.

The Spectral Bridges algorithm was also put to the test against the same competitors using scikit-learn’s UCI Breast Cancer data. Once again, this new method performs well although the advantage is not as obvious in this case (see Figure 7). However, in none of our tests has it ranked worse than k-means++. The results are displayed as a boxplot generated from 200 iterations of each algorithm using a different seed, in order to better grasp the variability lying in the seed dependent nature of the k-means++, Expectation Maximization and Spectral Bridges algorithms.

Since the Spectral Bridges algorithm is expected to excel at discerning complex and intricate cluster structures, an array of four toy datasets was collected, as illustrated in Figure 8.

Multiple algorithms, including the proposed one, were benchmarked in the exact same manner as for the UCI Breast Cancer data. The results show that the proposed method outperforms all tested algorithms (DBSCAN, k-means++, Expectation Maximization, and Ward Clustering) while requiring few hyperparameters. As previously discussed, DBSCAN’s parameters were optimized using the ground truth labels to represent a best-case scenario; however, in practical applications,

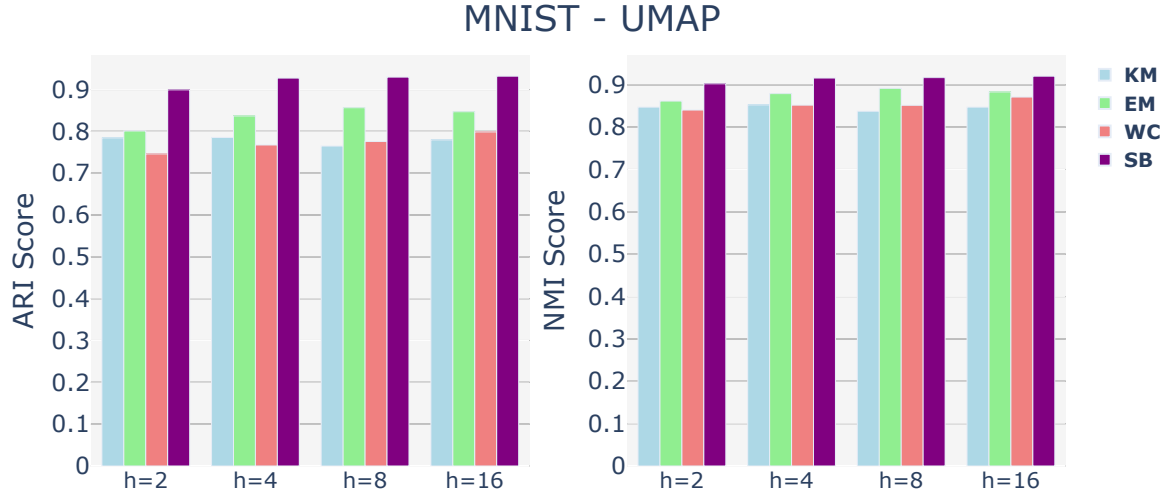


Figure 5: ARI and NMI scores of k-means++ (blue), EM (green), Ward Clustering (red), and Spectral Bridges (purple) on UMAP embedding

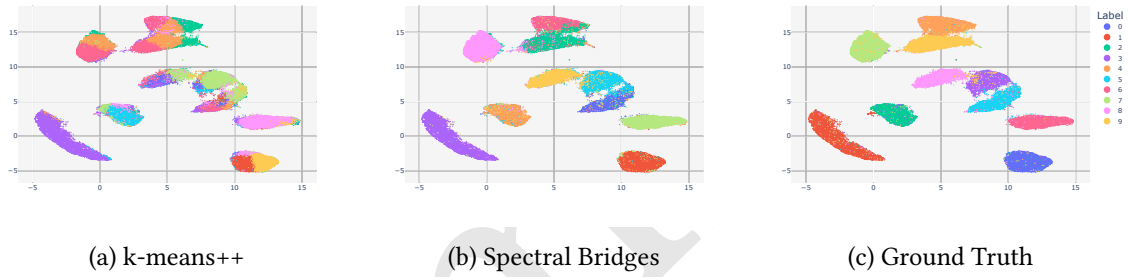


Figure 6: UMAP projection of predicted clusters against the ground truth labels.

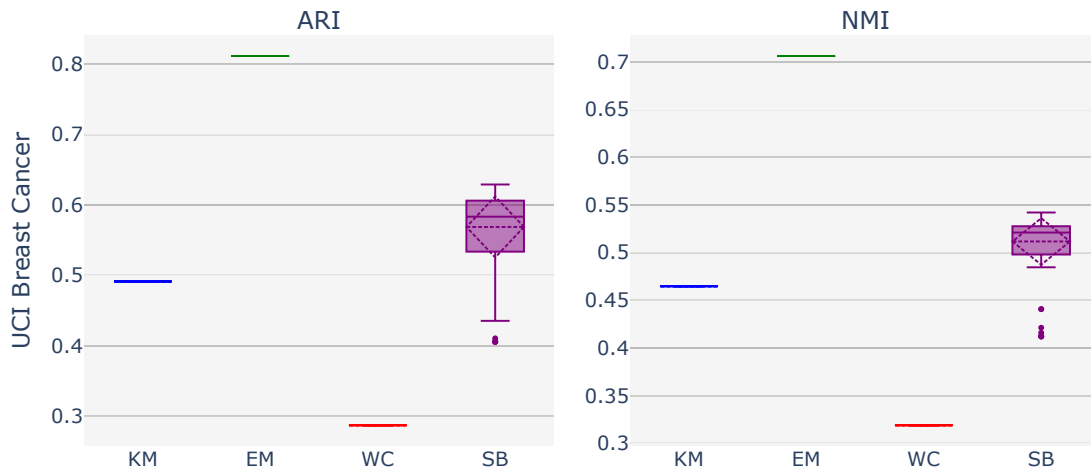


Figure 7: ARI and NMI scores of k-means++ (blue), EM (green), Ward Clustering (red), and Spectral Bridges (purple) on the UCI Breast Cancer dataset.

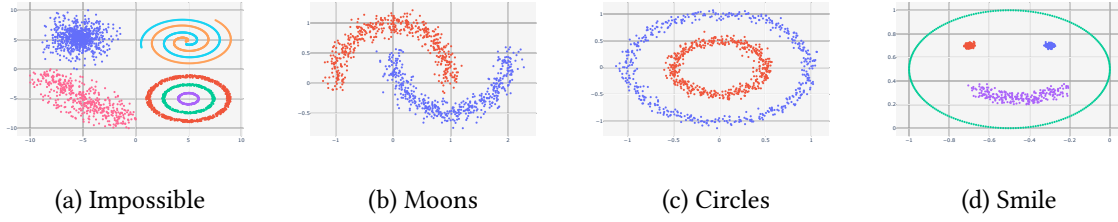


Figure 8: Four toy datasets.

suboptimal performance is more likely. Despite this optimization, the Spectral-Bridge algorithm still demonstrates superior ability to capture and represent the underlying cluster structures.

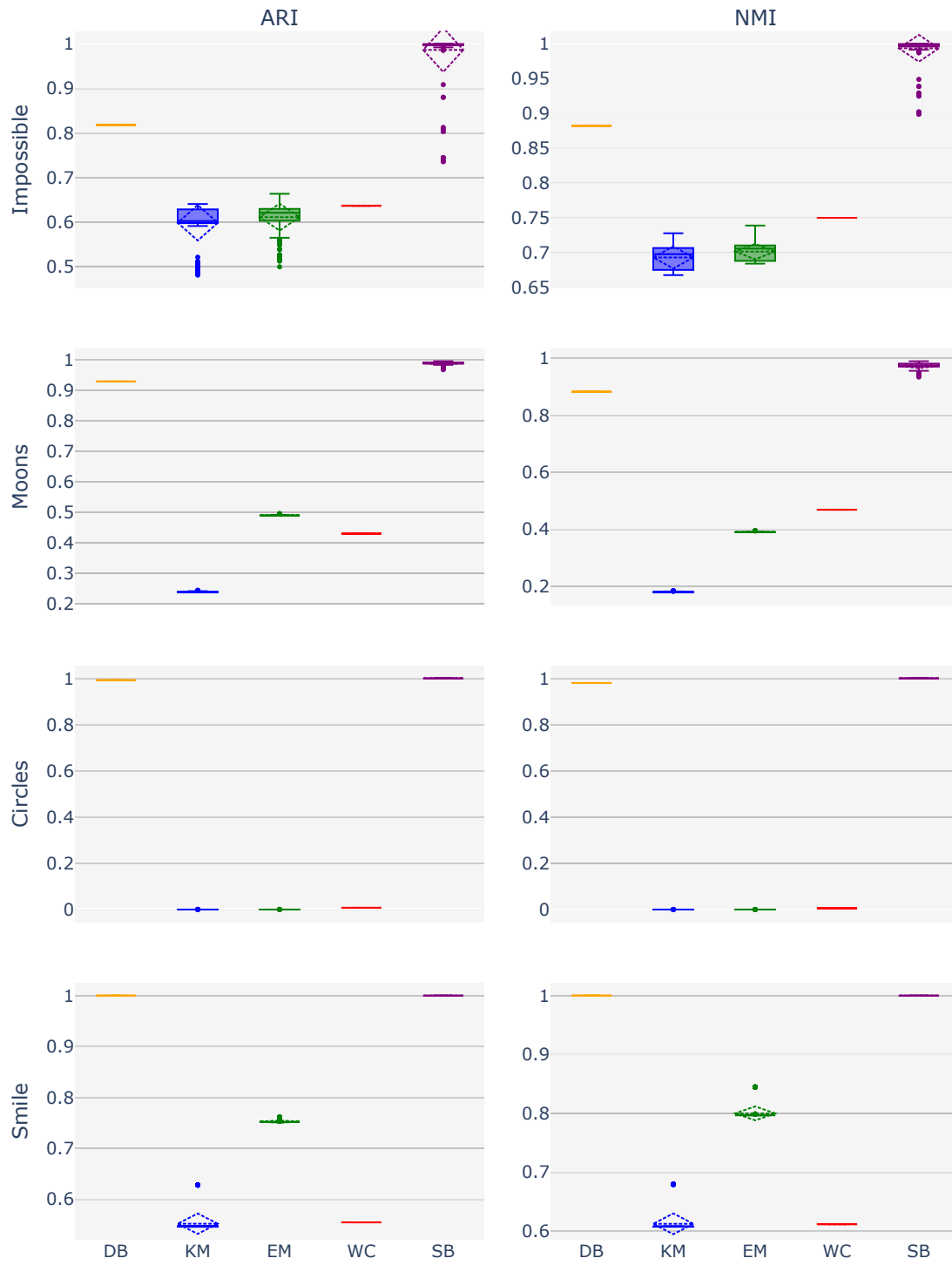


Figure 9: ARI and NMI scores of Spectral Bridges and competitors on standard synthetic toy datasets.

4.7 Noise robustness

To evaluate the noise robustness of the algorithm, two experimental setups were devised: one involved introducing Gaussian-distributed perturbations to the data, and the other involved concatenating uniformly distributed points within a predefined rectangular region (determined by the span of the dataset) to the existing dataset. As illustrated in Figure 10, the tests demonstrate that in both scenarios, the algorithm exhibits a high degree of insensitivity to noise.

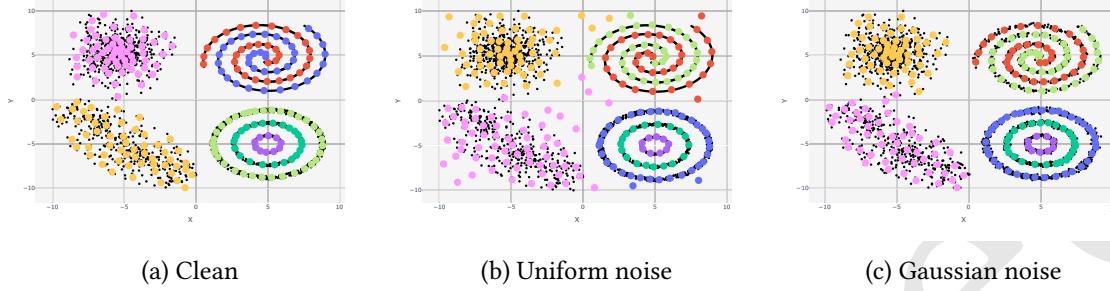


Figure 10: Three representations of the algorithm’s predicted cluster centers are displayed as colored dots, with each point of the Impossible dataset shown as a small black dot. In the left graph, the dataset is unmodified. In the center graph, 250 uniformly distributed samples were added. In the right graph, Gaussian noise perturbations with $\sigma = 0.1$ were applied.

4.8 Hyperparameter values effect on accuracy

The proposed algorithm inputs three parameters: the number of clusters K , the number of Voronoi regions m , as well as a scaling parameter used in the spectral clustering phase.

Metrics such as the Gap Statistic could be a possible approach to determine the number of clusters. Comparing the total within-cluster variation for different numbers of clusters with that expected under a null reference distribution of the data, the optimal value of K is where the observed clustering yields the largest gap statistic. Another approach could leverage the Laplacian’s eingengap, as in traditional spectral clustering (Von Luxburg 2007). K would here be chosen to maximize the difference between the $(K + 1)$ -th and K -th eigenvalues sorted in ascending order.

The same technique could also be used to estimate the optimal value for the scaling parameter; other methods involve local scaling which empirically does not translate well to the bridge affinity (Zelnik-Manor and Perona 2004).

At last, the choice of the the number of Voronoi regions m has been discussed previously. To better grasp the sensitivity regarding this parameter, Spectral Bridges was run on the PCA $h = 32$ embedded MNIST dataset with varying values of $m \in \{10, 120, 230, 340, 450, 560, 670, 780, 890, 1000\}$. The case $m = 10$ is equivalent to the k-means++ algorithm. ARI and NMI scores are recorded over 20 consecutive iterations and subsequently plotted. As shown by Figure 11, the accuracy seems to be consistently increasing with values of m , with the largest observed gap occurring between values of $m = 10$ and $m = 120$, and flattening thereafter, indicating a tremendous improvement over the classical k-means++ framework even for empirically suboptimal hyperparameter values.

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

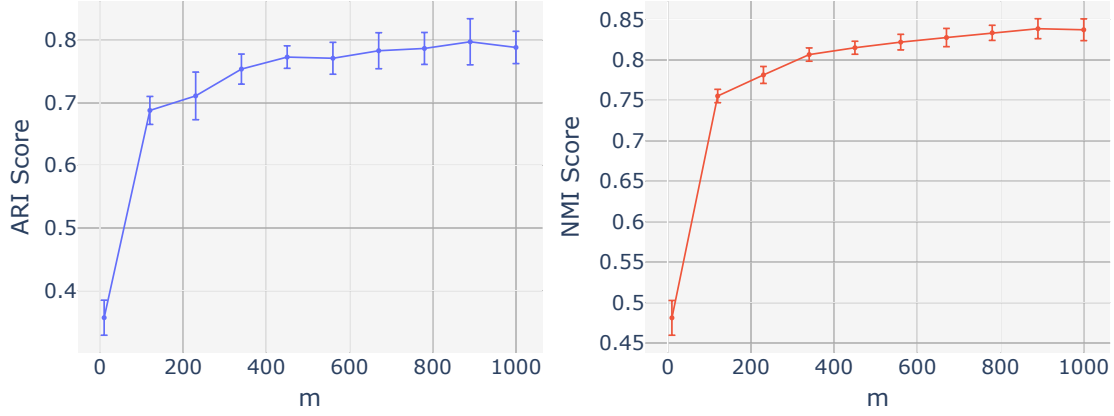


Figure 11: ARI and NMI scores of Spectral Bridges with varying values of m .

simple affinity measure for spectral clustering, which is derived from the minimal margin between pairs of Voronoï 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 Voronoï 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:

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

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

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

329 and each part further decomposed using Pythagore

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

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

330 Thus

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

331 6.2 Code

332 6.2.1 Implementation

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

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

338 6.2.2 Affinity matrix computation

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

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

Listing 1 Python code for affinity matrix computation

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

#Center each region
X_centered = [
    np.array(
        X[kmeans.labels_ == i] - kmeans.cluster_centers_[i],
        dtype=np.float32,
        order="F",
    )
    for i in range(self.n_nodes)
]

#Cardinal calculation
counts = np.array([X_centered[i].shape[0] for i in range(self.n_nodes)])
counts = counts[None, :] + counts[:, None]

#Calculate the affinity
for i in range(self.n_nodes):
    segments = np.asfortranarray(
        kmeans.cluster_centers_ - kmeans.cluster_centers_[i]
    )
    dists = np.einsum("ij,ij->i", segments, segments)
    dists[i] = 1

    projs = sgemm(1.0, X_centered[i], segments, trans_b=True) / dists
    np.clip(projs, 0, None, out=projs)

    affinity[i] = np.einsum("ij,ij->j", projs, projs)

affinity = np.sqrt((affinity + affinity.T) / counts)
affinity -= 0.5 * affinity.max()

#Scale and exponentiate
q10, q90 = np.quantile(affinity, [0.1, 0.9])

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

References

- Arthur, David, and Sergei Vassilvitskii. 2006. "K-Means++: The Advantages of Careful Seeding." Technical Report 2006-13. Stanford InfoLab; Stanford. <http://ilpubs.stanford.edu:8090/778/>.
- Cai, Deng, and Xinlei Chen. 2014. "Large Scale Spectral Clustering via Landmark-Based Sparse Representation." *IEEE Transactions on Cybernetics* 45 (8): 1669–80.
- Chen, Wen-Yen, Yangqiu Song, Hongjie Bai, Chih-Jen Lin, and Edward Y Chang. 2010. "Parallel Spectral Clustering in Distributed Systems." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33 (3): 568–86.
- Cortes, Corinna, and Vladimir Vapnik. 1995. "Support-Vector Networks." *Machine Learning* 20 (3): 273–97.
- Cover, Thomas M, and Joy A Thomas. 1991. "Information Theory and the Stock Market." *Elements of Information Theory*. Wiley Inc., New York, 543–56.
- Dempster, Arthur P, Nan M Laird, and Donald B Rubin. 1977. "Maximum Likelihood from Incomplete 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.
- McInnes, Leland, John Healy, Nathaniel Saul, and Lukas Großberger. 2018. "UMAP: Uniform Manifold Approximation and Projection." *Journal of Open Source Software* 3 (29): 861. <https://doi.org/10.21105/joss.00861>.
- 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

399 D. Wilkerson, Charlie R. Miller, et al. 2010. "Integrated Genomic Analysis Identifies Clinically
 400 Relevant Subtypes of Glioblastoma Characterized by Abnormalities in PDGFRA, IDH1, EGFR,
 401 and NF1." *Cancer Cell* 17 (1): 98–110.
 402 Von Luxburg, Ulrike. 2007. "A Tutorial on Spectral Clustering." *Statistics and Computing* 17: 395–416.
 403 Ward Jr, Joe H. 1963. "Hierarchical Grouping to Optimize an Objective Function." *Journal of the*
 404 *American Statistical Association* 58 (301): 236–44.
 405 Zelnik-Manor, Lihi, and Pietro Perona. 2004. "Self-Tuning Spectral Clustering." *Advances in Neural*
 406 *Information Processing Systems* 17.

407 Session information

408 R version 4.4.2 (2024-10-31 ucrt)
 409 Platform: x86_64-w64-mingw32/x64
 410 Running under: Windows 11 x64 (build 22631)
 411
 412 Matrix products: default
 413
 414
 415 locale:
 416 [1] LC_COLLATE=French_France.utf8 LC_CTYPE=French_France.utf8
 417 [3] LC_MONETARY=French_France.utf8 LC_NUMERIC=C
 418 [5] LC_TIME=French_France.utf8
 419
 420 time zone: Europe/Paris
 421 tzcode source: internal
 422
 423 attached base packages:
 424 [1] stats graphics grDevices utils datasets methods base
 425
 426 loaded via a namespace (and not attached):
 427 [1] compiler_4.4.2 fastmap_1.2.0 cli_3.6.3 tools_4.4.2
 428 [5] htmltools_0.5.8.1 yaml_2.3.8 rmarkdown_2.27 knitr_1.47
 429 [9] jsonlite_1.8.8 xfun_0.45 digest_0.6.36 rlang_1.1.4
 430 [13] evaluate_0.24.0