

Introduction to Machine Learning: Work 3

Pedro Agúndez*, Bruno Sánchez*, María del Carmen Ramírez*, Antonio Lobo*

December 8, 2024

Abstract

Abstract

*Universitat de Barcelona
pedro.agundez@estudiantat.upc.edu
bruno.sanchez.gomez@estudiantat.upc.edu
maria.del.carmen.ramirez@estudiantat.upc.edu
antonio.lopez@estudiantat.upc.edu

1 Introduction

Introduction.

2 Methodology

Introduction to methodology.

2.1 Datasets

2.2 Ordering Points To Identify the Clustering Structure (OPTICS)

The OPTICS (Ordering Points To Identify the Clustering Structure) method is a density-based clustering algorithm designed to reveal the intrinsic structure of data without requiring explicit cluster assignments or fixed parameter settings. Instead of directly generating clusters, OPTICS produces an ordered list of data points annotated with metrics that reflect their density relationships, such as core distances and reachability distances. This ordering captures the clustering structure across a range of density levels, allowing hierarchical and arbitrary-shaped clusters to be identified. The algorithm is computationally efficient, with performance similar to DBSCAN, and excels in uncovering the natural distribution of complex datasets.

2.3 Spectral

2.4 K-Means

K-Means intro.

2.4.1 Hyperparameters

1. **k:**

- The number of clusters to partition the dataset. Determines the complexity and granularity of the clustering.

2. **Distance Metrics:**

- **Euclidean Distance:** Calculates the root of the sum of squared differences between feature values. Standard metric for continuous data:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- **Manhattan Distance:** Computes the sum of absolute differences between feature values, suitable for both categorical and continuous data:

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

- **Clark Distance:** Accounts for proportional differences between feature values, enhancing interpretability for attributes with varying scales:

$$d(x, y) = \sqrt{\sum_{i=1}^n \left(\frac{|x_i - y_i|}{x_i + y_i + \epsilon} \right)^2}$$

where ϵ is a small constant to avoid division by zero.

3. **Additional Parameters:**

- **Initial Centroids:** Pre-defined initial cluster centers used as the starting point for the clustering algorithm.
- **Maximum Iterations:** Limits the number of iterations to prevent excessive computational time, with a default of 10 iterations.

2.4.2 Clustering Methodology

- **Clustering Process:**
 1. Assign each data point to the nearest centroid using the specified distance metric.
 2. Recalculate centroids by computing the mean of all points in each cluster.
 3. Repeat assignment and recalculation until convergence or maximum iterations are reached.
- **Convergence Criteria:**
 - Clusters are considered stable when centroids no longer significantly change between iterations.
- **Variance Computation:**
 - Total within-cluster variance (E) is calculated by summing squared distances of points to their respective cluster centroids.
 - Provides a measure of clustering compactness and quality.

This methodology allows for flexible clustering configurations, enabling analysis across different datasets and hyperparameter values.

2.5 Global K-Means

Out of the proposed improvements to the K-Means algorithm, the first we chose to implement was the **Global K-Means** algorithm [1], which focuses on following a deterministic and systematic approach to “optimal” centroid initialization and cluster formation. Additionally, we have also implemented the improvements to the Global K-Means algorithm itself, proposed in the original article by Likas et al.: *Fast Global K-Means*, and *Initialization with k-d Trees*. By addressing the limitations of traditional K-Means, this enhanced methodology introduces novel strategies, including PCA-based data partitioning and iterative error-reduction mechanisms, to improve both accuracy and computational efficiency.

This section outlines the hyperparameter configurations and clustering methodology adopted for the Global K-Means algorithm, which was implemented in the `GlobalKMeansAlgorithm` class.

2.5.1 Hyperparameters

We consider the same hyperparameters as for the standard K-Means algorithm (Section 2.4), except for 2 significant modifications:

1. **Initial Centroids:**
 - Global K-Means no longer accepts a collection of initial centroids as a hyperparameter, since the goal of this algorithm is rooted in the deterministic calculation of the “best possible” centroids, which substitutes their random initialization.
2. **Number of Buckets:**
 - Controls initial data partitioning, by defining the number of candidate points that we will consider as possible centroids throughout the algorithm.
 - Its default value is $2 \cdot k$, but we also test values $3 \cdot k$ and $4 \cdot k$.

2.5.2 Clustering Methodology

- **Initialization with k-d Trees:**

1. Use k-d tree partitioning based on Principal Component Analysis (PCA).
2. Recursively partition data samples into buckets.
3. Select candidate points based on bucket centroids.

- **Fast Global K-Means Algorithm:**

1. Initialize first centroid as dataset mean.
2. Iteratively add centroids by:
 - For each $k' = 2, \dots, k$, we already have $k' - 1$ centroids.
 - Compute guaranteed error reduction for candidate points with respect to the $k' - 1$ centroids,

$$b_n = \sum_{j=1}^N \max \left(d_{k'-1}^j - \|x_n - x_j\|^2, 0 \right),$$

where $d_{k'-1}^j$ is the squared distance between x_j and the closest centroid among the $k' - 1$ obtained so far. The pair-wise squared distances between points are precomputed at the start.

- Select point with maximum guaranteed error reduction.
 - Run k' -means with the $k' - 1$ centroids plus the selected point, until convergence.
3. Repeat until k clusters are formed.

This methodology provides a sophisticated approach to centroid initialization and clustering, leveraging PCA-based partitioning and error reduction strategies in order to achieve an improvement in consistency and speed with respect to the base K-Means algorithm.

2.6 Fuzzy C-Means

2.7 Metrics

2.7.1 ARI

2.7.2 F1-Score

2.7.3 DBI

2.7.4 Silhouette

2.7.5 Calinski

3 Results

To systematically evaluate the different configurations of each clustering algorithm, the following procedure is followed to extract results for each of the 3 data sets, in order to later perform an statistical analysis:

1. **Data Preparation:** The dataset is loaded and the data samples are separated from their labels into separate 2 sets. This way, we perform the clustering analysis in a completely non-supervised way, and we then utilize the labels to extract supervised metrics of the clustering results.
2. **Parameter Configuration:** A comprehensive set of values for the algorithm's hyperparameters is defined. These combinations reflect various ways to tune the clustering algorithm.
3. **Model Evaluation:** For each parameter combination, the clustering algorithm is applied on the unlabeled data and then evaluated with different metrics. This step yields the following metrics: total cluster variance (E), Adjusted Rand Score (ARI), F1-score, Davies-Bouldin Index(DBI), Silhouette score, Calinski-Harabasz score, and execution time. Together, these metrics (some supervised, some non-supervised) measure the effectiveness and efficiency of the clustering.

Note: Not all clustering algorithms are based on centroids, hence the total cluster variance is not computed for those which are not.

4. **Results Compilation:** The performance metrics for each parameter combination are recorded in a structured format. These results are saved as a dataset that summarizes the outcomes of all evaluations, forming a basis for analysis.
5. **Statistical Analysis:** After results are compiled across all configurations, statistical analysis is performed to identify the best-performing configurations. This analysis helps determine the most reliable and effective parameter settings for accurate and efficient clustering. We will discuss our results in the following sections.

3.1 K-Means

Results of K-Means.

4 Concluision

Conclusion.

References

- [1] Aristidis Likas, Nikos Vlassis, and Jakob J. Verbeek. The global k-means clustering algorithm. Pattern Recognition, 36(2):451–461, 2003.