

# Assignment - 16

## Clustering Algorithms

---

HARSH SIDDHAPURA  
1230169813

04/05/2024

1. Clustering is an important task in machine learning and data analysis that aims to divide a dataset into groups where data points in each group are more similar compared to the content in other categories. This technique is essential for uncovering hidden patterns and patterns in data and aids in many applications such as image segmentation, customer segmentation, and suspect detection.

As an unsupervised learning activity, clustering stands in contrast to supervised learning methods by operating on datasets devoid of labeled outputs. Clustering algorithms do not rely on a predefined list but automatically identify patterns in data, making them useful for data analysis. Integration helps understand complex data by exploring relationships and relationships between data points and helps make informed decisions in different areas, thus increasing the efficiency and effectiveness of the data-driven process.

2. Hierarchical clustering and partition clustering are the two main methods of clustering algorithms, and each offers different methods and applications. Hierarchical clustering creates a hierarchical cluster structure, represented as a tree diagram, by combining or separating clusters based on their similarities. In agglomerative hierarchical clustering, individual data points start from separate clusters and gradually merge; distributed hierarchical clustering, on the other hand, starts with all data points in the cluster and vice versa. This approach provides an overview of the organization of the dataset and helps understand the relationships between data elements. Although the hierarchical representation is explicit, hierarchical integration is computed in the time complexity of  $O(n^2)$ .

In contrast, partition clustering divides data into non overlapping clusters, as in algorithms such as K-means, where the number of clusters (K) must be specified. The k-means uniformity method reassigns data points to the nearest centroid until convergence occurs and adjusts centroids to the mean of the points in each group. Although clustering appears to be more efficient with  $O(n)$  time complexity, it does not have the hierarchical structure that hierarchical clustering provides and cannot provide fine-grained but sparse information.

3. Complete clustering and partial clustering are two different approaches in clustering algorithms. Complete clustering, also known as complete clustering, aims to split the dataset at a fixed location such that all data points are in a single cluster. This approach provides

clear and concise solutions, making it suitable for situations that require different and non-overlapping categories. An example of successful clustering is the K-means algorithm.

A partial group, on the other hand, allows data points to go into multiple groups at once. This is done by assigning each group a possible degree of membership or participation. This method is especially useful when the boundaries of categories are not clear or where data points overlap between different categories. Examples of partial clustering algorithms include fuzzy C-means clustering and Gaussian mixture models. In summary, the entire group provides all the information that points to a group, while the partial group allows for overlapping memberships, allowing complex information to be represented flexibly.

4. The K-Means algorithm is a popular method used in data science and machine learning to partition datasets into K distinct, non-overlapping bins. The main idea behind K-Means is to divide points into fixed clusters, maximizing the similarity within each cluster and minimizing the similarity between clusters.

Here is a step-by-step explanation of how the K-Means algorithm works:

- **Initialization:** The algorithm first randomly selects K points from the dataset. This point is the first center of the K group.
- **Assignment:** All data points in the file are assigned to the nearest center. The distance between a data point and the center of mass is usually calculated using the Euclidean distance, but other distance measurements can also be used. This step will create K groups of data points, with each group corresponding to a center.
- **Update:** Recalculate the clusters' centroids. This is usually done by averaging all data points in each category. This step moves the center of gravity to the center of the cluster.
- **Iteration:** Repeat steps 2 and 3 until the center of mass does not move or the maximum number of iterations is reached. This iterative process increases the stack at each step, thus improving the quality of the stack.
- **Termination:** When the centroid is stable, the algorithm stops and indicates that the data points are divided into K groups.

The K-Means algorithm is simple and effective, suitable for big data. However, it has some limitations. For example, how many groups of K need to be specified in advance, which may

not always be known. It is also important for the initial selection of the center of gravity and may fall into local optima.

5. Aggregate clustering, also known as hierarchical clustering, is an algorithm that creates a group hierarchy by merging or splitting existing groups. The agglomerative approach is a "bottom-up" concept in which each data point starts in its own cluster, and when one cluster moves, two sets of clusters merge.

Here is a detailed explanation of how agglomerative hierarchical clustering works:

- **Initialization:** The algorithm first treats each data point as a separate group. This means that if there are N elements, we start with N groups.
- **Find the closest pair:** This algorithm calculates the distance between each cluster. The distance between two groups is often measured by the Euclidean distance between the two closest points (also called connectivity), the two farthest points (fully connected), or the two midpoints (mid-connection).
- **Merge:** Merge two groups that are close to each other. Reduces the total stack by one.
- **Update Distance Matrix:** After merging, the distance between the new group and all old groups will be updated.
- **Iteration:** Repeat steps 2 through 4 until all items are placed in a cluster.
- **Resulting Dendrogram:** The results of hierarchical clustering can be represented by a tree-like structure called a dendrogram. The root of the tree is the group that contains all the samples, and the leaves are the group that contains only one sample.

The advantage of aggregate hierarchical clustering is that it does not require the user to specify the number of clusters in advance. However, for large data the computational cost can be high and the results are sensitive to the choice of connection method.

6. DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a noise-based clustering algorithm that is different from classification and hierarchical clustering methods.

Here is a detailed explanation of how DBSCAN works:

- **Parameter settings:** DBSCAN requires two parameters: 'eps' and 'minPts'. 'eps' is the maximum radius of the neighborhood around data points and 'minPts' is the minimum number of data points required to make the area thick.
- **Region Query:** For each data point in the dataset, the algorithm first searches for all points up to "eps". This leads to the field of data points.
- **Main Points, Breakpoints and Unsatisfied Points:** A data point is classified as a main point if it has at least "minPts" in neighborhood strength "eps". The number of boundary points in the "eps" neighborhood is less than "minPts" but is in the "eps" area of the main points. Noise points are random points that are not core points or boundary points.
- **Cluster formation:** The algorithm then proceeds to combine key points that are close together (i.e. within "eps" distance). Each group of links points to each other as a separate group. While noise is not assigned, boundary points are assigned to their groups.
- **Iteration:** This algorithm continues to process all parameters in the dataset until all points are entered and either grouped or labeled as noise.

DBSCAN has many advantages over other integration programs. It can detect groups of suspicious images, does not require the user to predetermine the number of groups, and does not understand the initialization of the algorithm. However, it will not work well when the group has different sizes or the data has high attributes.

## 7. Summarizing various clustering methods in different conditions:

	KMeans	Single-Link	Complete-Link	DBSCAN
Noise	Sensitive to noise and outliers. Since it minimizes the variance within clusters, a few noisy points can significantly shift the cluster centroids and	Sensitive to noise and outliers. A single noisy point can cause two clusters to merge, leading to less accurate clusters.	Less sensitive to noise and outliers than Single-Link, as it considers the farthest points in clusters for merging. Noise can still affect the shape and	Handles noise well by classifying them as outliers. It only merges points that are densely packed together, leaving out the sparse points as

	distort the clusters.		compactness of the clusters.	noise.
<b>Clusters with Different Sizes</b>	Assumes clusters of similar sizes, and hence, struggles with clusters of different sizes. Tends to partition space into regions of roughly equal sizes.	Can handle clusters of different sizes, as it gradually merges the closest points or clusters, allowing clusters to grow unevenly.	Can also handle clusters of different sizes, as it merges the clusters based on the farthest points, which allows for more flexibility in cluster sizes.	Can handle clusters of different sizes, as it forms clusters based on the density of points, allowing clusters to grow as long as the density of points remains high.
<b>Clusters with Different Densities</b>	Struggles with clusters of different densities. It tends to create clusters that cover equal-sized regions, regardless of the density of points.	Can lead to a chaining effect in the presence of varying densities, where a series of close points can cause distant clusters to merge.	Less prone to the chaining effect than Single-Link, as it considers the farthest points in clusters for merging. May still struggle with clusters of different densities.	Can handle clusters of different densities, as it forms clusters based on the density of points. It can separate dense clusters from less dense areas.
<b>Non-globularly-shaped Clusters</b>	Struggles with non-globular clusters. It tends to create spherical clusters due to its use of Euclidean distance.	Can handle non-globular shapes, as it gradually merges the closest points or clusters, allowing clusters to take on any shape.	Can handle non-globular shapes to some extent, as it merges the clusters based on the farthest points.	Can handle non-globular shapes, as it forms clusters based on the density of points, allowing clusters to take on any shape.

8. Bisection K-means is a modification of the K-means algorithm designed to overcome some of its limitations. The main idea behind splitting K-means into two is to use the

divide-and-conquer strategy to improve cluster quality and reduce sensitivity to the initial choice of centroids.

The process of splitting the word K into two is as follows:

- **Initialization:** Initialize with all data points in a group.
- **Split:** Select the group to split. Selection can be based on various criteria, such as the cluster with the highest number of squared errors (SSE) or the largest cluster. Use the K-Means algorithm ( $K = 2$ ) on this cluster to create two clusters.
- **Iteration:** Repeat the bisection step until the desired group is obtained. After each iteration, the next group to split is selected according to the selection process.

Advantages of bisecting K-means over standard K-means are:

- **Clustering Quality:** Dividing K-means into two generally produces clusters, better than the K-means model. This is because the algorithm has a higher chance of finding the optimal world by searching different solution spaces.
- **Robustness of initialization:** Unlike K-Means, which are sensitive to the choice of initial center of gravity, bisection K-means are inherently more sensitive to initialization.
- **Flexibility:** Dividing K-means in half can create a number of groups, providing greater flexibility and visibility

However, it is worth noting that bisecting K-means can be more computationally accurate. Due to the overhead of repeating the K-means algorithm, it is more expensive than the K-means model, especially for large data sets.

9. Although the K-Means algorithm is powerful and efficient, it has some limitations that may affect its ability to detect distinct clusters. The reasons for this are:

- **Equal clustering assumption:** K-means works with the assumption that all clusters are approximately equal and have a spherical shape. This assumption can result in decreased performance when dealing with clusters of different sizes or densities. If one cluster is larger or denser than another cluster, K-Means can be split into smaller clusters while merging smaller clusters.
- **Distance measurement:** K-Means uses Euclidean distance to assign cluster points. This measure treats all instructions equally, which can be problematic when groups

are different. Larger clusters can turn into many smaller clusters because points in the larger cluster may be closer to the center of the smaller cluster than their centroids.

- **Sensitivity to initial centroids:** The final result of K-Means can be greatly affected by the initial position of the centroids. If the first median is placed near smaller clusters, K-Means will not correctly identify larger clusters.
- **Number of clusters:** K-Means must first specify the number of clusters. If the number chosen does not reflect the distribution of cluster size in the data, K-Means will not correctly identify distinct clusters.

In summary, K-Means' assumptions and feature algorithms will limit its effectiveness in analyzing groups of different variables. Other algorithms, such as DBSCAN or hierarchical clustering, do not make the same assumptions about cluster size or shape and may be better for such tasks.

10. Although the K-Means algorithm is efficient and widely used, it has some shortcomings in identifying non-spherical or spherical clusters. Here is an explanation:

- **Euclidean distance:** K-Means uses Euclidean distance to measure the similarity between data points. This measure is isotropic, meaning it treats all directions equally, favoring spherical clusters. Therefore, K-Means may have trouble processing data sets whose natural clusters are stretched or irregularly shaped.
- **Centrum Calculation:** This algorithm calculates the centroid of the cluster based on the average of all data points in the cluster. This method inherently supports sphericity because the definition does not represent the shape of the data points in the cluster.
- **Data point assignment:** In K-Means, each data point is assigned to the cluster with the closest centroid. This can lead to erroneous data points that belong to non-spherical clusters but are close to the centroid of a different cluster.
- **Flexibility Limitation:** K-Means assumes that all clusters are convex and isotropic; this can lead to poor performance on data with complex patterns.

In summary, K-Means' inherent assumptions and feature algorithms hinder its ability to effectively identify non-spherical clusters. Other algorithms, such as DBSCAN or spectral clustering, do not imply clustering and may be better for this task.