

KMeans Clustering

K-means clustering is a popular unsupervised learning method in machine learning and data analysis. It aims to group unlabeled data points into k distinct clusters based on their similarities. Here's a breakdown of the key aspects:

What it does:

- Partitions a dataset into k clusters, where k is a pre-defined number chosen by the user.
- Points within a cluster are more similar to each other compared to points in other clusters.
- Clusters are represented by their centroids, which are the average of all points within the cluster.

How it works:

1. Initialization: Randomly select k data points as initial cluster centroids.
2. Assignment: Assign each data point to the closest centroid, based on a distance metric (e.g., Euclidean distance).
3. Recalculation: Recompute the centroids by taking the average of all points assigned to each cluster.
4. Iteration: Repeat steps 2 and 3 until the centroids stabilize (no more changes) or a maximum number of iterations is reached.

Strengths:

- Simple and efficient algorithm, making it easy to implement and understand.
- Scalable to large datasets.
- Effective for finding spherical clusters in data.

Weaknesses:

- Sensitive to the choice of k : Choosing the optimal k can be challenging and can significantly impact the clustering results.
- Assumes clusters are spherical: May not work well for complex cluster shapes.
- Outliers can significantly affect the centroids and the clustering results.

Applications:

- Customer segmentation in marketing
- Image segmentation in computer vision
- Anomaly detection in fraud analysis
- Document clustering in information retrieval