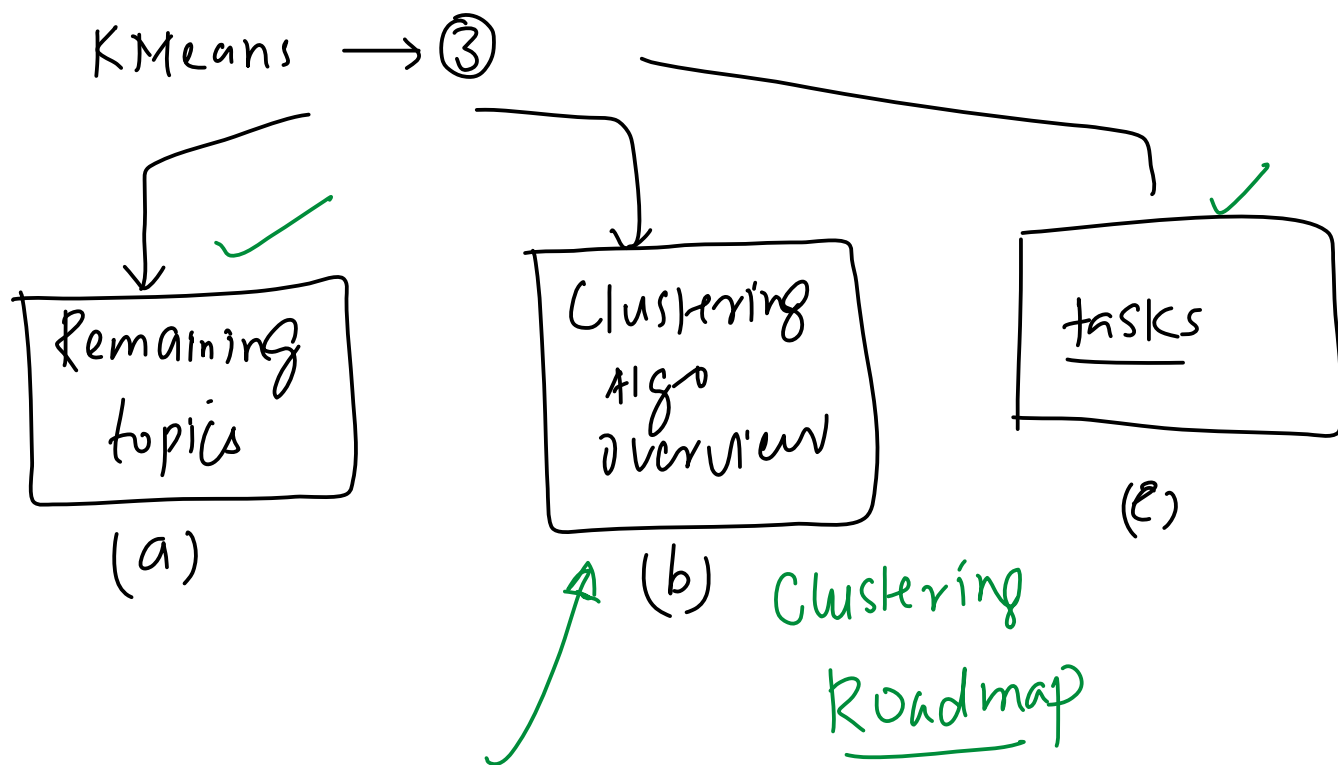


Recap

05 December 2023

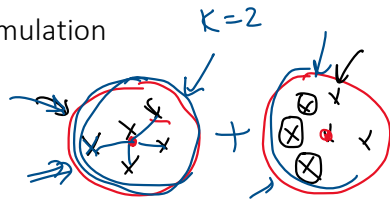
17:15



Kmeans Mathematical Formulation

05 December 2023 17:16

Kmeans



[Lloyd's algorithm] Juggad
↓
approximation

$$J = \underset{M_1, M_2, \dots, M_K}{\operatorname{argmin}} \left[\sum_{i=1}^K \left[\sum_{x \in S_i} \|x - \mu_i\|^2 \right] \right] \quad \text{such that} \quad [S_i \cap S_j = \emptyset]$$

column inertia

$[S_i \cap S_j = \emptyset]$

$\mu_1, \mu_2, \dots, \mu_K \rightarrow$ centroid
K centroids

$[S_i] \rightarrow$ current cluster

$$[S_i \cap S_j = \emptyset]$$

n 1000, K 2 cluster
↓
 K
↓
 n points K cluster

$$J = \underset{M_1, M_2, \dots, M_K}{\operatorname{argmin}} \left[\sum_{i=1}^K \sum_{x \in S_i} \|x - \mu_i\|^2 \right] \rightarrow \text{minimize}$$

such that $S_i \cap S_j = \emptyset$

NP hard

non-convex \rightarrow local minimal

approximation best sol

very hard to solve

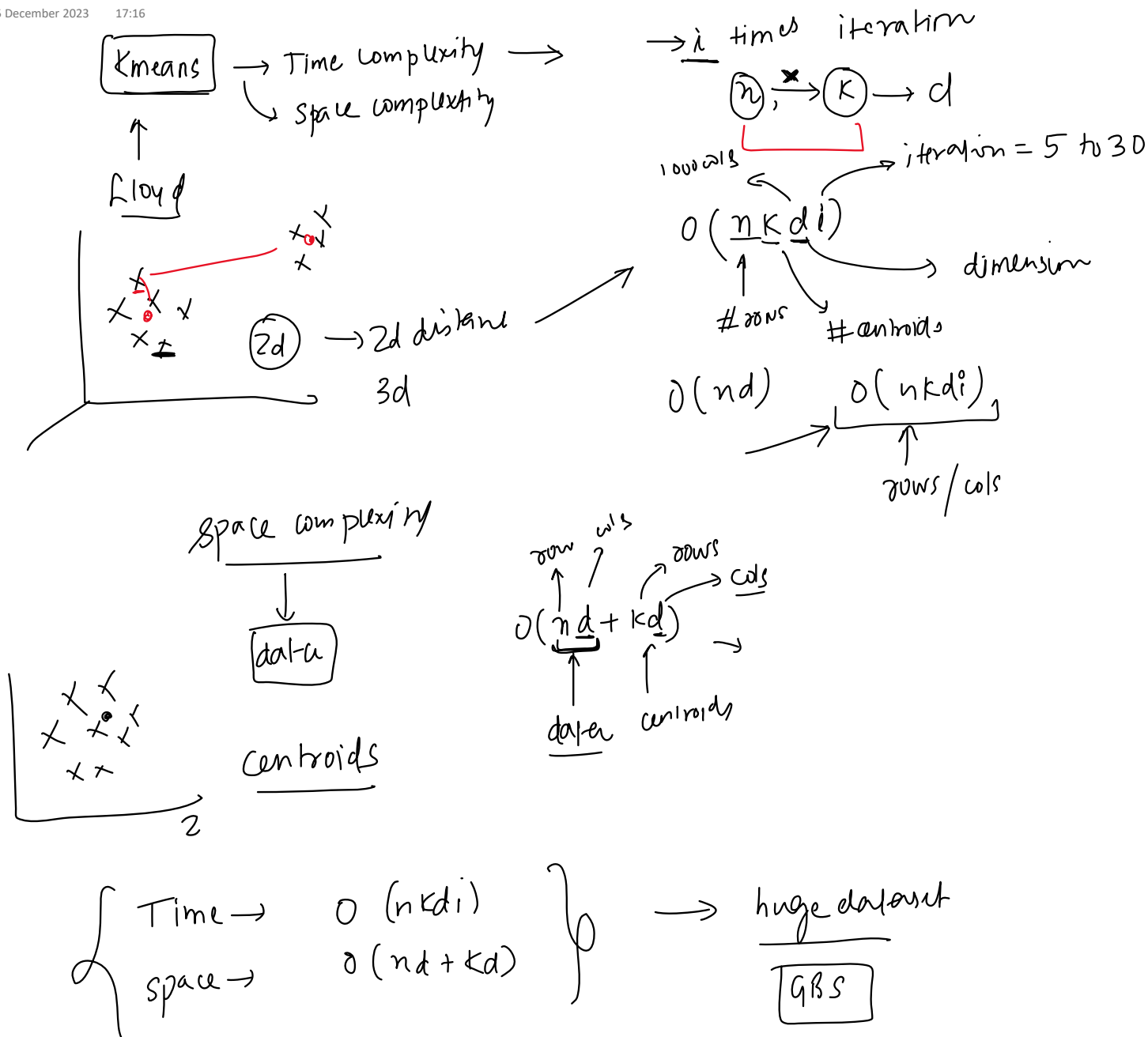
Juggad

Lloyd

non-optimal local minima

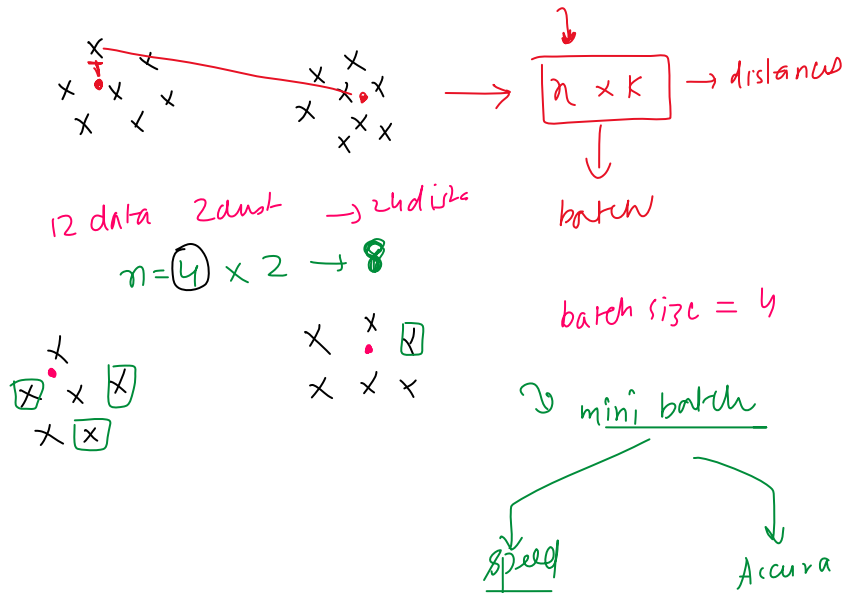
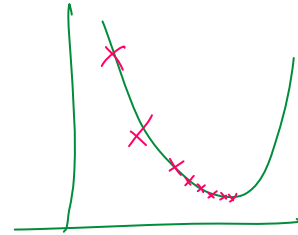
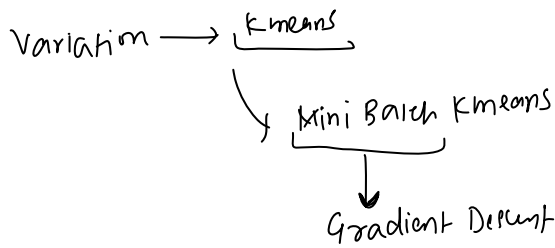
Kmeans Time and Space Complexity

05 December 2023 17:16



Mini Batch Kmeans

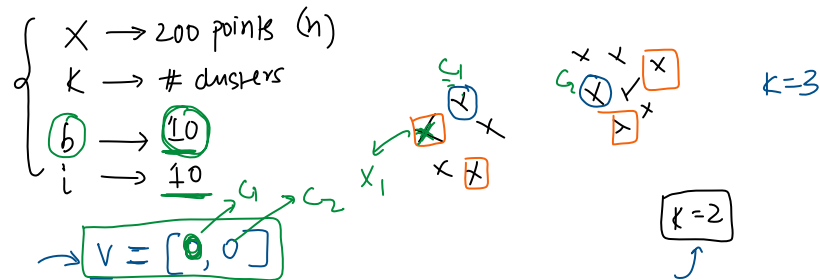
05 December 2023 17:15



Algorithm 1 Mini-batch k -Means.

```

1: Given:  $k$ , mini-batch size  $b$ , iterations  $t$ , data set  $X$ 
2: Initialize each  $c \in C$  with an  $x$  picked randomly from  $X$ 
3:  $v \leftarrow 0$ 
4: for  $i = 1$  to  $t$  do
5:    $M \leftarrow b$  examples picked randomly from  $X$  (mini batch)
6:   for  $x \in M$  do
7:      $d[x] \leftarrow f(C, x)$  // Cache the center nearest to  $x$ 
8:   end for
9:   for  $x \in M$  do
10:     $c \leftarrow d[x]$  // Get cached center for this  $x$ 
11:     $v[c] \leftarrow v[c] + 1$  // Update per-center counts
12:     $\eta \leftarrow \frac{1}{v[c]}$  // Get per-center learning rate
13:     $c \leftarrow (1 - \eta)c + \eta x$  // Take gradient step
14:   end for
15: end for
    
```

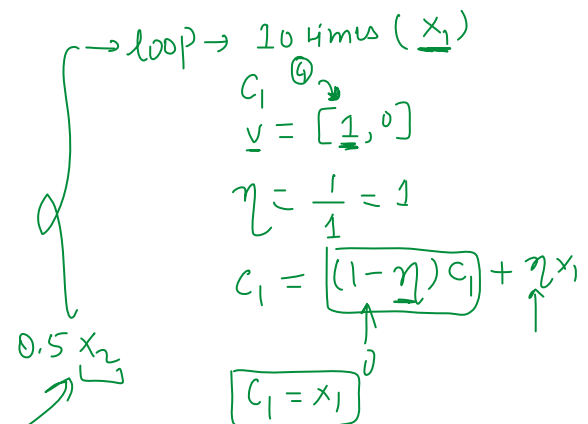
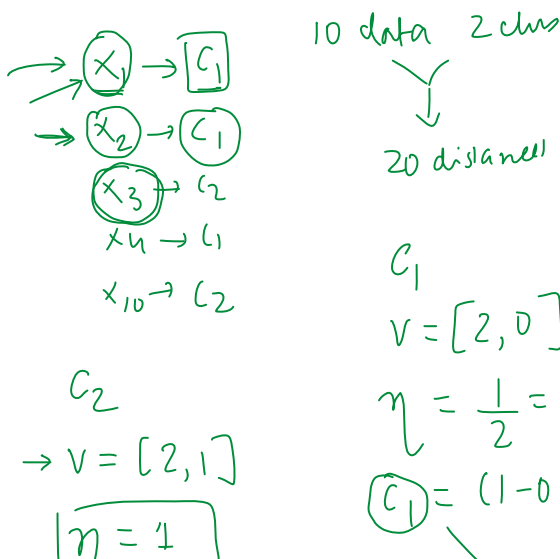


loop $\rightarrow 1-t / 10$ times

[randomly 10 \rightarrow 200 points] mini batch

loop $\rightarrow 10$ times

\rightarrow [dict] (x, c)

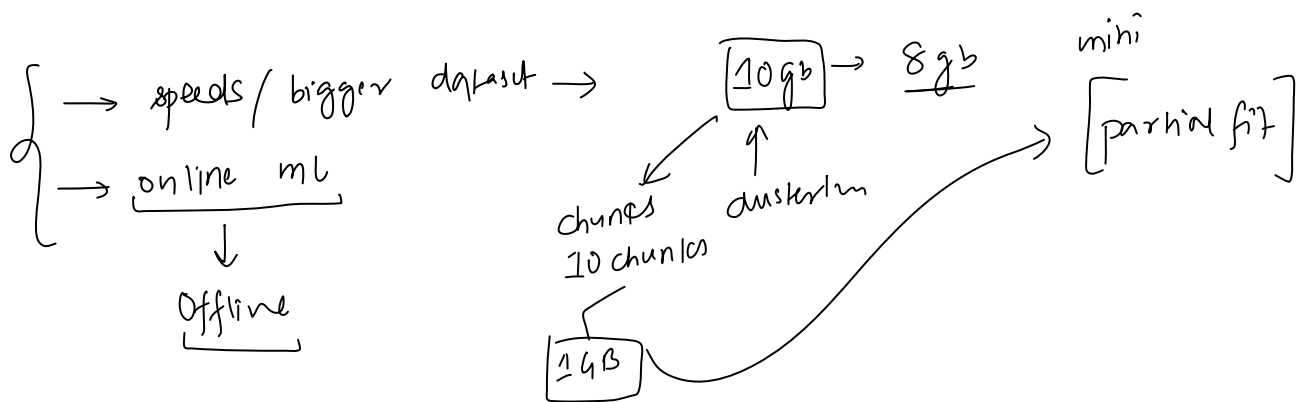


$\eta = 1$
 $\rightarrow C_2 = (1-1)C_2 + 1 \times x_3$
 $C_2 = x_3$

$C_1 = (1-0.5)C_1 + 0.5 \times x_2$

$C_1 = x_1$

When should you use Minibatch



downside
accuracy

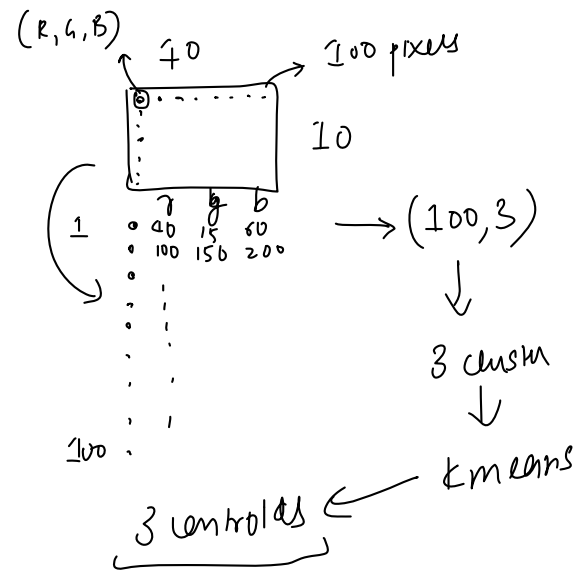
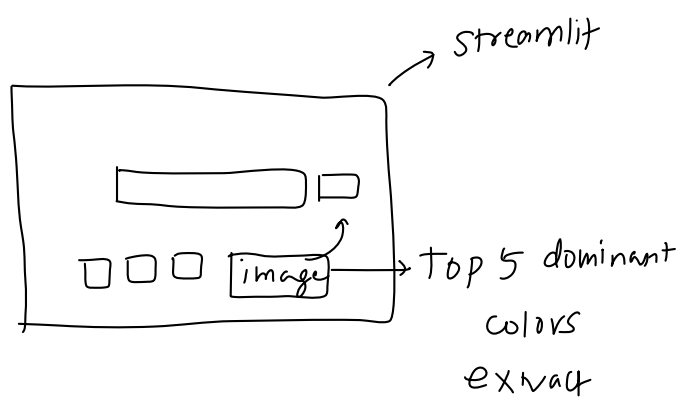
Task 1

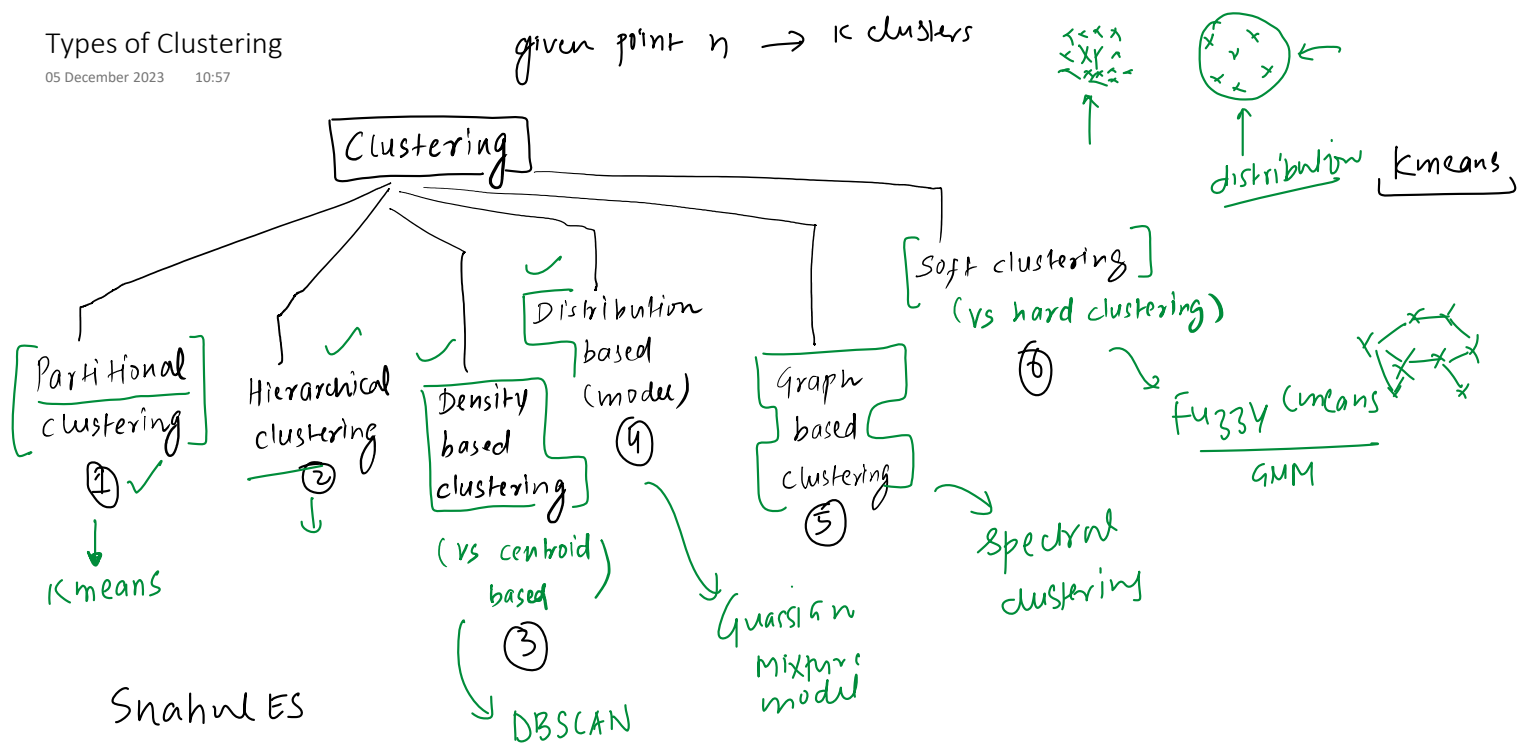
05 December 2023

17:16

Task 2

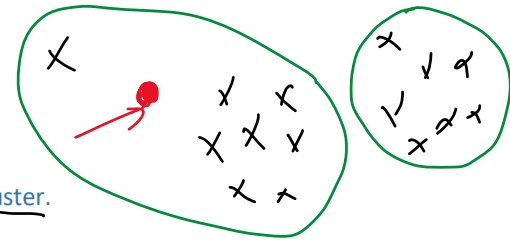
05 December 2023 17:16



given point $n \rightarrow k$ clusters

Partitional Clustering

05 December 2023 10:57



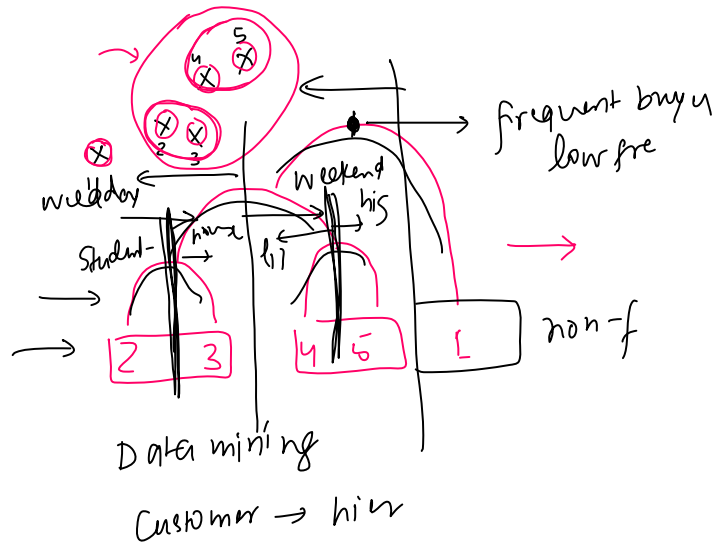
1. **Basic Concept:** Partitioning clustering algorithms divide a dataset into a set of non-overlapping subgroups or clusters, where each data point belongs to exactly one cluster.
2. **Examples:** The most famous partitioning clustering algorithm is k-means. It assigns data points to clusters in such a way that each point belongs to the cluster with the closest mean, which serves as a prototype of the cluster.
3. **Defining Number of Clusters:** A key requirement is pre-specifying the number of clusters (k). The selection of k significantly affects the outcome and quality of the clustering. elbow
silhouette
4. **Iterative Process:** These algorithms typically use an iterative refinement technique. For instance, in k-means, the process involves repeatedly assigning points to the nearest cluster centroid and then recalculating the centroids. iteration
5. **Objective Function Optimization:** They aim to optimize an objective function, such as minimizing the total within-cluster variance or the sum of squared distances between data points and their respective cluster centroids. inertia
6. **Suitability for Certain Data Shapes:** Partitioning methods are most effective when clusters are spherical or globular in shape. They assume homogeneity in cluster shapes and sizes.
7. **Sensitivity to Initial Conditions:** These algorithms can be sensitive to the initial starting conditions (like initial cluster centroids in k-means). Different initializations can lead to different clustering results. kmeans++
8. **Handling of Outliers:** Partitioning algorithms can be influenced by outliers, as these can significantly skew the mean or centroid of a cluster.
9. **Scalability and Efficiency:** They are generally more scalable and efficient for larger datasets compared to hierarchical clustering, making them suitable for many practical applications.
10. **Use Cases and Limitations:** While widely used in various fields like market research, pattern recognition, and image processing, these algorithms have limitations in handling non-spherical clusters, varying cluster sizes, and noisy datasets. Advanced versions and variations of partitioning algorithms have been developed to address some of these limitations.

Hierarchical Clustering

05 December 2023 10:58

- Nature of Clustering:** Hierarchical clustering builds a hierarchy of clusters either by successively merging smaller clusters into larger ones (Agglomerative approach) or by successively splitting larger clusters into smaller ones (divisive approach).
- No Need to Specify Number of Clusters:** Unlike partitioning algorithms like k-means, hierarchical clustering does not require pre-specifying the number of clusters. The number of clusters can be determined by analysing the dendrogram.
- Dendrogram Visualization:** It provides a tree-like diagram called a dendrogram, which is a visual representation of the clustering process showing the order of cluster combination and the distance at which clusters are merged.
- Distance Metrics and Linkage Criteria:** Hierarchical clustering uses various distance metrics (like Euclidean or Manhattan distance) and linkage criteria (like single linkage, complete linkage, average linkage, and Ward's method) to decide which clusters to merge or split.
- Flexibility in Identifying Cluster Shapes:** Hierarchical clustering can identify clusters with various shapes and sizes, unlike partitioning methods that generally assume spherical clusters.
- Computational Complexity:** It is generally more computationally intensive than partitioning methods, especially for large datasets, due to the need to compute and store distances between all pairs of points.
- Sensitivity to Noise and Outliers:** The method can be sensitive to noise and outliers, as these can influence the formation of clusters and the structure of the dendrogram.
- Applications:** It is widely used in fields like biology (for gene and protein sequencing), social science, and linguistics, and is particularly useful for exploratory data analysis where understanding the hierarchical relationship between objects is important.

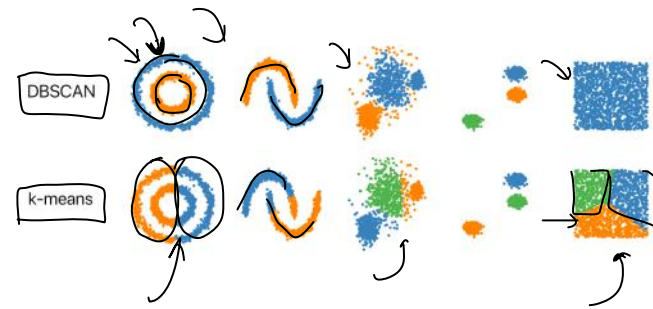
→ [nlp]



Density Based Clustering

05 December 2023 10:58

- ✓ **Principle:** Density-based clustering groups data points based on the density of data points in a region. It defines clusters as areas of high density separated by areas of low density. The algorithm identifies clusters as regions where data points are densely packed together, with areas of low density or noise between them.
- ✓ **Examples:** DBSCAN is one of the most popular density-based clustering algorithms. It is known for its efficiency and ability to find clusters of arbitrary shapes.
- ✓ **No Need to Specify Number of Clusters:** Unlike partitioning methods, density-based clustering doesn't require pre-specifying the number of clusters.
- ✓ **Handling Noise and Outliers:** It is robust to outliers and noise, as these are typically not part of the dense regions that form clusters.
- ✓ **Ability to Find Arbitrary Shapes:** Density-based clustering can discover clusters of arbitrary shapes, unlike methods like k-means which are biased towards spherical clusters.
- ✓ **Parameter Sensitivity:** The performance of these algorithms is sensitive to the input parameters, like the radius of neighborhood (ϵ) and the minimum number of points required to form a dense region (MinPts in DBSCAN).
- ✓ **Scalability Issues:** Some density-based algorithms may struggle with very large datasets due to computational and memory constraints.
- ✓ **Applications:** Widely used in fields such as anomaly detection, geospatial data analysis (like identifying geographic regions of interest), and image processing, especially where the shape of the clusters is not known in advance or the data contains noise.



Distribution/Model based Clustering

05 December 2023 10:58

1. **Statistical Distribution Models:** The central concept of distribution-based clustering is that data points in a cluster follow a certain statistical distribution, most commonly Gaussian or normal distributions.
2. **Parameter Estimation:** These algorithms focus on estimating the parameters (like mean, variance) of the assumed distributions for each cluster. The fit of these parameters to the actual data determines the quality of the clustering.
3. **Expectation-Maximization (EM) Algorithm:** A key algorithm used in distribution-based clustering is EM, which alternates between assigning data points to the most likely distribution (expectation step) and updating the distribution parameters to maximize data fit (maximization step).
4. **Handling of Complex Cluster Shapes:** Unlike methods such as K-Means, distribution-based clustering can identify clusters of various shapes and sizes, making it more flexible in handling real-world data complexities.
5. **Computational Intensity:** The process of estimating distribution parameters and assigning data points can be computationally demanding, especially for large datasets and when the number of features (dimensions) is high.
6. **Handling of Outliers:** These methods can be more robust to outliers, as outliers are less likely to significantly affect the parameters of the overall distribution.
7. **Scalability Issues:** While effective for small to medium-sized datasets, scalability to very large datasets can be challenging due to the computational complexity of the algorithms and the need for more sophisticated optimization techniques.
8. **Soft Clustering:** In soft clustering, each data point is associated with a probability distribution across different clusters, indicating the degree of belonging to each cluster. This is in contrast to hard clustering, where each data point is assigned to exactly one cluster.

