

# DISADVANTAGES OF CLUSTERING

## i) K Means Clustering -

- Choosing  $K$  manually.
- Works only good for Circular shape data
- Clustering data of varying sizes and density
- Distance based model
- Clustering Outliers.
- ~~Do~~ Don't work for Non-Circular shaped data.
- Scaling with number of dimensions (Curse of dimensionality)
- With global clusters, it didn't work well (Global means it can be subdivided into two/three more clusters).
- Lack of flexibility in cluster shape
- Lack of probabilistic cluster assignment

## ii) Hierarchical Clustering -

- If we have a large dataset, it become difficult to determine the correct number of clusters by dendrogram.
- Sensitivity to noise and outliers.
- Breaking large clusters become difficult.

## iii) DBScan Clustering -

- Does not work well when dealing with clusters of varying densities. While DBScan is great at separating high density clusters from low density clusters, DBScan struggles with cluster of similar density.
- Struggle with high dimensionality data. DBScan suffers badly with high dimension.

## iv) Gaussian (EM) Clustering -

- Algorithm is very complex in nature.
- Algorithm simply would not work for datasets where objects do not follow Gaussian distribution.
- Distribution based model.



## CHOOSING RIGHT CLUSTERING ALGORITHM IN DATASET

→ Clustering based on computation of distances between the objects of the whole dataset is called connectivity based or hierarchical. Depending on the direction of algorithm, it can unite or inversely divide the information i.e., agglomerative and divisive.

Most prominent example of connectivity based clusterization is classification of plants. The "tree" dataset starts with particular species and end with a few category of plants.

→ Centroid based clustering, aimed at classifying each object of dataset to a particular cluster.

→ When it is a spherical shaped, well separated data then go for  $k$  means clustering.

When it is not so spherical or non-spherical data and we want probability of each data point to a cluster then use GMM clustering. But GMM will not work if data do not follow Gaussian distribution.

→ Density based clustering, DBScan best used when we have outliers in the data. DBScan can be used when data are in arbitrary shapes, and they are extremely accurate. Beside this algorithm doesn't need number of clusters from outside - it is determined automatically.

→  $k$  means clustering is also known as partition based clustering.