**The Learning Machine**

Statistics    Machine Learning    Applied ML    Contribute

# Clustering
## Differences between clustering algorithms

**Clustering** is the task of dividing data points into a number of groups such that data points in the same groups are more similar to other data points in the same group and dissimilar to the data points in other groups. It is basically a collection of objects on the basis of similarity and dissimilarity between them.

**Clustering Methods:**

1. **Density-Based Methods:** These methods consider the clusters as the dense region having some similarity and different from the lower dense region of the space. These methods have good accuracy and the ability to merge two clusters. Example DBSCAN (Density-Based Spatial Clustering of Applications with Noise), OPTICS (Ordering Points to Identify Clustering Structure) etc.
2. **Hierarchical Based Methods:** The clusters formed in this method forms a tree-type structure based on the hierarchy. New clusters are formed using the previously formed one. It is divided into two category
   --> *Agglomerative* (bottom-up approach)
   --> *Divisive* (top-down approach).
   Examples CURE (Clustering Using Representatives), BIRCH (Balanced Iterative Reducing Clustering and using Hierarchies) etc.
3. **Centroid (Partitioning) Methods:** These methods partition the objects into k clusters and each partition forms one cluster. This method is used to optimize an objective criterion similarity function such as when the distance is a major parameter example K-means, CLARANS (Clustering Large Applications based upon randomized Search) etc.
4. **Distribution Methods:** These clustering models are based on the notion of how probable is it that all data points in the cluster belong to the same distribution (For example: Normal, Gaussian). These models often suffer from overfitting. A popular example of these models is Expectation-maximization algorithm which uses multivariate normal distributions.

## Density-Based Models

| Clustering Model | Pros | Cons |
| --- | --- | --- |
| **DBSCAN** (Density-Based Spatial Clustering of Applications with Noise) | • Clusters of arbitrary shape and size<br>• Robust to noise<br>• Does not need a pre-set number of clusters<br>• Deterministic | • Requires connected regions of sufficiently high density<br>• Data sets with varying densities are problematic |
| **Mean-Shift Clustering** | • Does not need a pre-set number of clusters<br>• Simple to understand | • Selection of the window size/radius "r" can be non-trivial |

## Hierarchical Models

| Clustering Model | Pros | Cons |
| --- | --- | --- |
| **Hierarchical Clustering** | • The optimal number of clusters can be obtained by the model itself<br>• Practical visualisation with the dendrogram | • Not appropriate for large datasets |
| **BIRCH** (Balanced Iterative Reducing Clustering and using Hierarchies) | • Designed for clustering a large amount of numerical data<br>• Works well only for spherical clusters | • Can handle only numeric data<br>• Sensitive to the order of the data records |

## Centroid Methods

| Clustering Model | Pros | Cons |
| --- | --- | --- |
| **K-Means** (MiniBatch K-Means, K-Means ++) | • Simple to understand<br>• Easily adaptable<br>• Works well on small or large datasets<br>• Fast, efficient and performant | • Need to choose the number of clusters<br>• Assumes the clusters as spherical, so does not work efficiently with complex geometrical shaped data(Mostly Non-Linear)<br>• Hard Assignment might lead to mis grouping. |
| **Affinity Propagation** | • Does not need a pre-set number of clusters<br>• Works well on small or large datasets<br>• Clusters of arbitrary shape and size<br>• Does better clusters than K-Means | • Much slower than K-Means |

• Does better clusters than K-Means

| | | |
|---|---|---|
| **Spectral Clustering** | • Elegant, and well-founded mathematically<br>• Works quite well when relations are approximately transitive<br>• Excellent quality under many different data forms | • Not appropriate for very noisy datasets<br>• Much slower than KMeans |

# Distribution Methods

| Clustering Model | Pros | Cons |
|---|---|---|
| **GMM**<br>Expectation-Maximization using Gaussian Mixture Models | • A lot more flexible in terms of cluster covariance than K-Means<br>• Have multiple clusters per data points, I.e GMMs support mixed membership<br>• Does not assume clusters to be of any geometry. Works well with non-linear geometric distributions as well.<br>• Does not bias the cluster sizes to have specific structures as does by K-Means (Circular). | • Slow convergence<br>• Inability to provide estimation to the asymptotic variance-covariance matrix of the maximum likelihood estimator (MLE)<br>• Difficult to interpret. |

# Overall Model Performance