

Machine Learning

Cluster Analysis I

S. Nõmm

¹Department of Software Science, Tallinn University of Technology

06.02.2024

Clustering

Clustering belongs to the non supervised machine learning techniques. It is used to discover the structure of a given data set.

Definition

Clustering is a process of grouping elements of the given data set into groups with respect to chosen similarity criteria.

This definition requires one to determine the following "parameters" either in process of learning or before it.

- Similarity criteria (distance or metric function).
- Algorithm and its goodness criteria.
- Validation.

Hyperparameter is the parameter which value is not determined during the learning.

As a result of clustering each element is assigned label describing which cluster it belongs to. NB! Similarity and distances are synonyms in the area of machine learning.

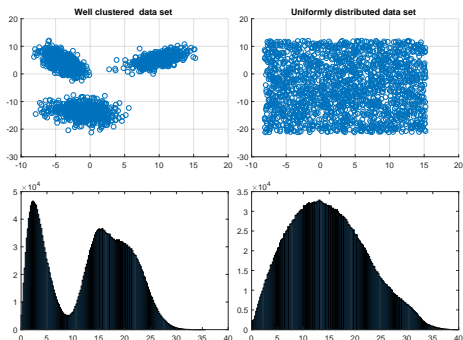
Classification of clustering techniques

Most common clustering techniques may be classified as follows:

- **Representative based techniques:** k-means, k-medians, k-medoids, etc. Each cluster has a representative which is either the element of the data set or an element from the same space as all other elements of the dataset. Shape of the clusters is affected by the choice of distance function. Number of clusters is usually a hyperparameter.
- **Hierarchical clustering techniques:** Agglomerative and Divisive techniques. Not always relies on the distance function. Different levels of clustering granularity provide different application specific insights.
- **Grid and Density based techniques:** Relies on the local density of the data points. Well suited for the clusters of irregular shapes.
- **Probabilistic algorithms:** EM and EM-like algorithms.

Feature selection

- Underlying idea is that features with uniformly distributed values carry less information compared to those distributed non uniformly.
- Distance distributions of well-clustered sets should be different from those uniformly distributed.



Measures

- *Entropy*

$$E = -\sum_{i=1}^m [p_i \log(p_i) + (1 - p_i) \log(1 - p_i)].$$

where p_i is the proportion of the points in the region i , m - total number of regions. Large values of E indicate poor clustering behaviour.

- *Hopkins statistics*. Let \mathcal{D} be the data set to investigate and \mathcal{R} is a representative sample of \mathcal{D} , of power r . \mathcal{S} is a synthetic data set of r data points randomly generated from the same domain. Let $\alpha_1, \dots, \alpha_r$ be the distances of each point of \mathcal{R} to the nearest neighbour in \mathcal{D} and β_1, \dots, β_r are the distances of each point of \mathcal{S} to the nearest neighbour in \mathcal{D} . The Hopkins statistic is defined as follows:

$$H = \frac{\sum_{i=1}^r \beta_i}{\sum_{i=1}^r (\alpha_i + \beta_i)}.$$

Higher values of H indicate highly clustered data.

Feature selection

- Filter Methods: Use Entropy or Hopkins Statistics to decide set of features leads best clustering behaviour. Filter methods may be applied on the stage of preprocessing.
- Wrapper models: clustering algorithm is used to evaluate the quality of subset of features.

K - means

K - means is one of the most popular algorithms belongs to the class of iterative descent methods.

- It is intended for the quantitative variables.
- Squared Euclidean distance as dissimilarity measure.
- The idea is to assign close points to the same cluster. Minimize natural loss ("energy") function.

$$W(C) = \frac{1}{2} \sum_{k=1}^K N_k \sum_{C(i)=k} |x_i - \bar{x}_k|^2.$$

where \bar{x}_k is the mean vector associated with the k th cluster (*centroid*). $N_k = \sum_{i=1}^N I(C(i) = k)$.

- Iterative descent algorithm is used to achieve this goal.

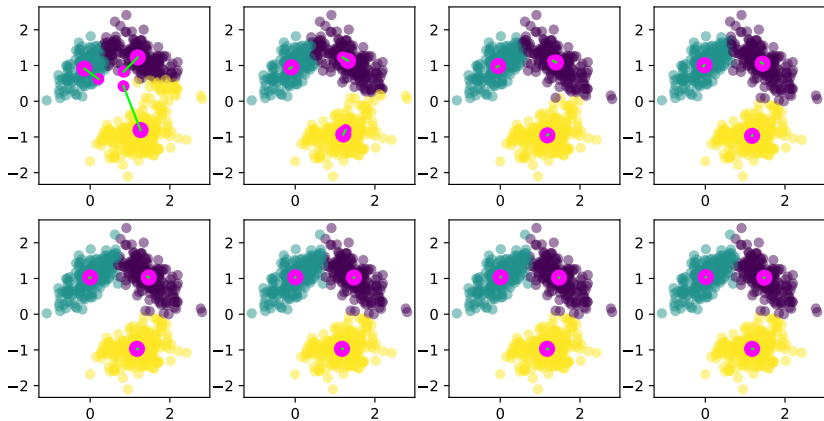
Representative based clustering

K -means:

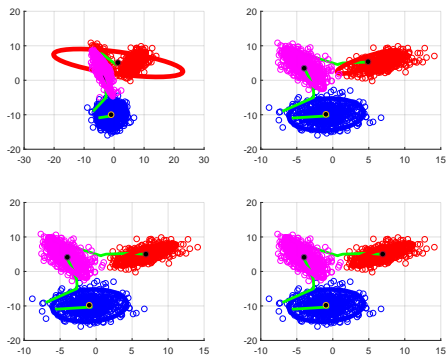
- Hyperparameters: K - number of desired clusters, distance function.
- Initialize: generate K random points from the same limits as initial dataset. These points are referred as *centroid*.
- Repeat:
 - ▶ For each point assign the label of closest centroid.
 - ▶ For each label recompute centroid as the mean of all points with given label.
- Until converge.
- Report labels of each point.

Other representative based techniques differ only by the way representative is find.

K -means clustering example



K -means clustering example



K - means, example discussion

- Convergence criteria?
 - ▶ Assignments do not change?
 - ▶ Minimum of a loss function?
- Relations to the EM-algorithm? Instead of maximizing likelihood K - means minimizes loss function.
- K - means best perform when clustered dataset composed of spherical or similar subsets.
- How to validate quality of clustering?

Validation

- **Sum of square distances to centroids.** (SSQ) This criteria is suited for K -means since it minimizes the loss function. (With reservations)
- **Intracluster to intercluster distance ratio.** Sample r points from the data set. Let P be the set of pairs that belong to the same cluster and Q the set of remaining pairs.

$$II = \frac{\sum_{(x_i, x_j) \in P} S(x_i, x_j) / |P|}{\sum_{(x_i, x_j) \in Q} S(x_i, x_j) / |Q|}$$

Small values of the ratio indicate better clustering behaviour.

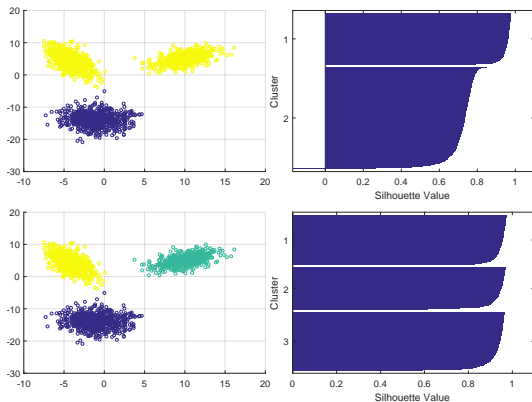
- **Silhouette coefficient**

$$s(i) = \frac{D_{\min_i}^{\text{out}} - D_{\text{avg}_i}^{\text{in}}}{\max\{D_{\min_i}^{\text{out}}, D_{\text{avg}_i}^{\text{in}}\}}$$

where $D_{\text{avg}_i}^{\text{in}}$ is the average distance of point x_i to points within the cluster it belong to. Compute average distance of point x_i to the points of each cluster. Let $D_{\min_i}^{\text{out}}$ is the minimum of these average distances. $s(i) \in (-1, 1)$. Overall coefficient is the average of the individual points coefficients. Large positive values indicate highly separated clusters.

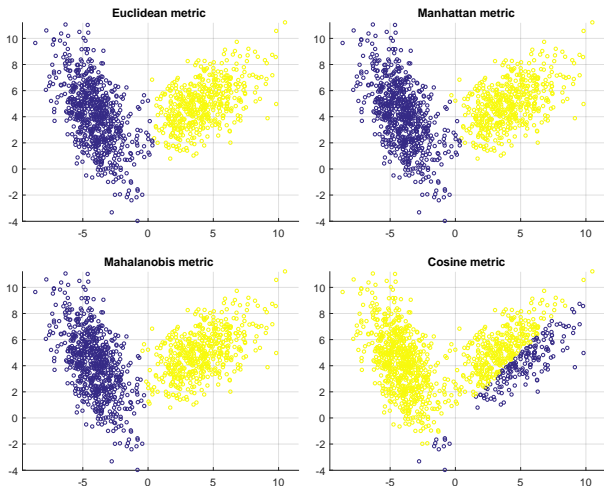
Silhouette coefficient

- Considered to be most popular criteria for clustering validation.
- Silhouette plot is the graphic representation of the silhouette coefficient.
- Overall silhouette coefficient may be used to determine number of clusters.

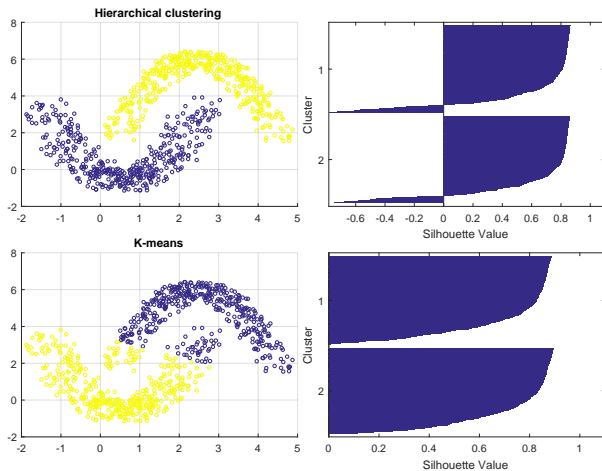


Impact of distance functions

NB! Always observe if distance function is defined for the given dataset and if using it makes sense from the viewpoint of interpretation.



Limitations

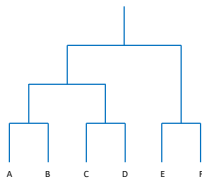
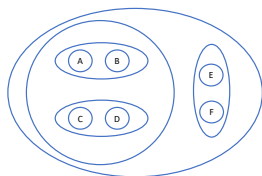


Hierarchical clustering: Agglomerative clustering

Some times referred ad *bottom-up*

Algorithm

- Initialize $n \times n$ distance matrix \mathcal{M}
- **Repeat**
 - ▶ Choose closest pair of clusters (i, j) based on \mathcal{M} .
 - ▶ Merge clusters i and j and update matrix \mathcal{M} .
- **Until** termination criterion.
- Return cluster labels for each point.



Group-based statistics

Also referred as *linkage*.

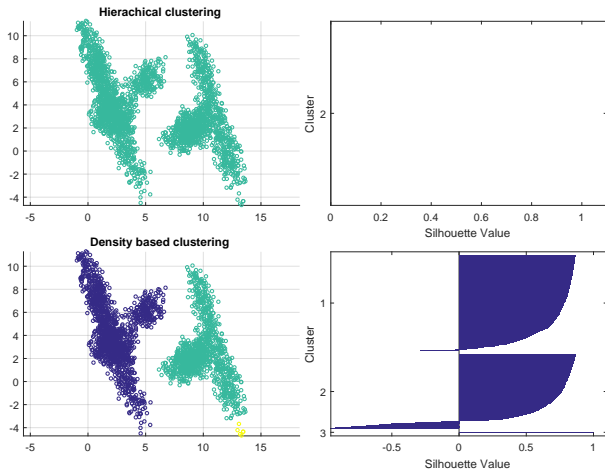
- Best (single) linkage. Distance is equal to the minimum distance between all pairs of elements (from two groups). Suitable to discover clusters of arbitrary shape. Drawback noise points may merge distant clusters.
- Worst (complete) linkage. (Complete linkage method) Distance is equal to the maximum distance between all pairs of elements (from two groups). Attempts to minimize maximal diameter of the cluster.
- Group average linkage. Distance between two groups is equal to the average of the distances between all pairs of elements (from two groups).
- Closest centroid. Clusters with closest centroid are merged.
- Variance based criterion. Minimizes the change in the objective function a result of merging.
- Ward's method, like previous but instead of variance observes changes in some squared error.

Top-down divisive methods

Algorithm

- Initialize tree \mathcal{T} to root containing dataset \mathcal{D}
- **Repeat**
 - ▶ Select a lead node \mathcal{L} in \mathcal{T} based on predefined criterion.
 - ▶ Use splitting algorithm \mathcal{A} to split \mathcal{L} into $\mathcal{L}_1, \dots, \mathcal{L}_k$.
 - ▶ add $\mathcal{L}_1, \dots, \mathcal{L}_k$ as children of \mathcal{L} in \mathcal{T} .
- **Until** termination criteria.

Limitations



Grid- and Density-based clustering

Explores the idea, that clusters are of a different density than space between them. May be seen as the sub class of agglomerative methods.

Generic Grid:

Hyperparameters: Ranges and density threshold τ .

- Discretize each dimension into p ranges.
- Determine *dense* grid cells at level τ .
- Create graph in which dense grids are connected if they are adjacent.
- Determine connected components of the graph.
- Return cluster indexes for each point.

DBSCAN

Let \mathcal{D} denote the data set, τ - density threshold and ϵ - radius of the neighborhood.

Definition

Core point: A data point is defined as the core point, if its ϵ - neighbourhood contains at least τ data points.

Definition

Border point: A data point is defined as the border point, if its ϵ - neighbourhood contains at least one another data point of \mathcal{D} and at least one core point.

Definition

Noise point: Is defined as data point of \mathcal{D} which neither core point nor border point.

DBSCAN

Algorithm:

- Determine Core, border and noise points for given ϵ and τ .
- Create graph in which core points are connected (if they are within ϵ of one another).
- Assign each border point to a connected component.
- Return cluster indexes for each point.