

# Data Mining: Lecture 2

## Cluster Analysis

S. Nõmm

<sup>1</sup>Department of Software Science, Tallinn University of Technology

08.09.2020

# Introduction

Given a set of data points, partition them into groups with respect to chosen similarity criteria.

- Data summarization.
- Discover the structure of the set.
- Part of preprocessing.

# Feature selection

Given a set of data points, partition them into groups with respect to chosen similarity criteria.

- Filter Models

- ▶ Predictive Attribute Dependence
- ▶ Entropy

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

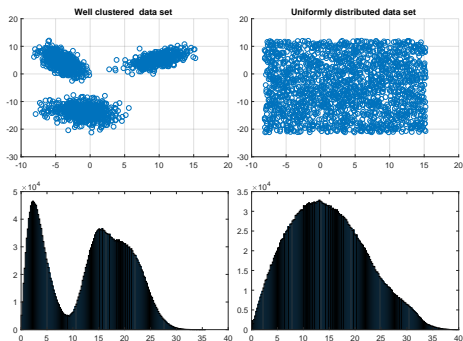
- ▶ Hopkins Statistic

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

- Wrapper models

# 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  $\beta_1, \dots, \beta_r$  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.

# 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 provide different application specific insides.
- **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.

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 element belongs. NB! Similarity and distances are synonyms.

## $K$ - 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

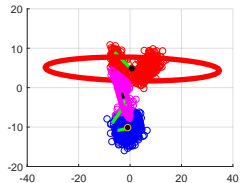
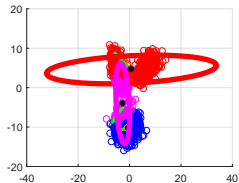
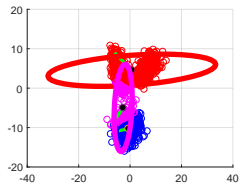
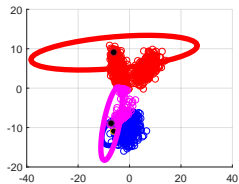
## $K$ -means:

- 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.

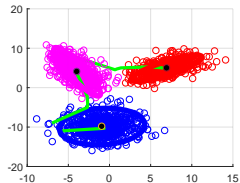
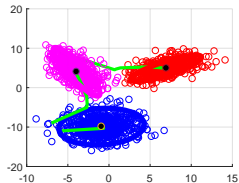
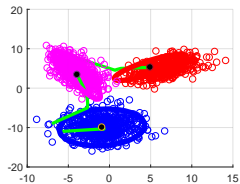
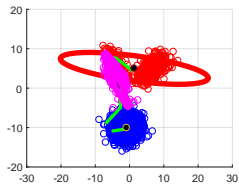
# $K$ -means clustering example

Steps 1 - 4



# $K$ -means clustering example

Steps 5 - 8



## $K$ - 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.

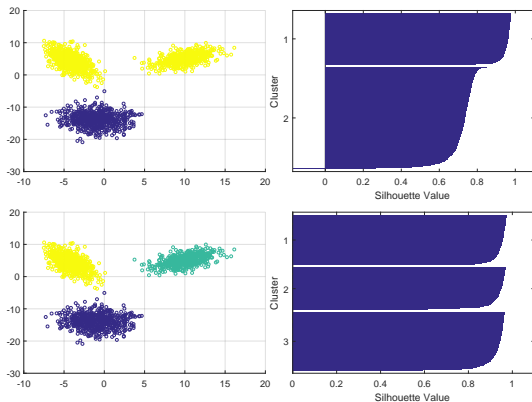
- **Silhouette coefficient**

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

where  $D_{avg_i}^{in}$  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  $D_{min_i}^{out}$  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.

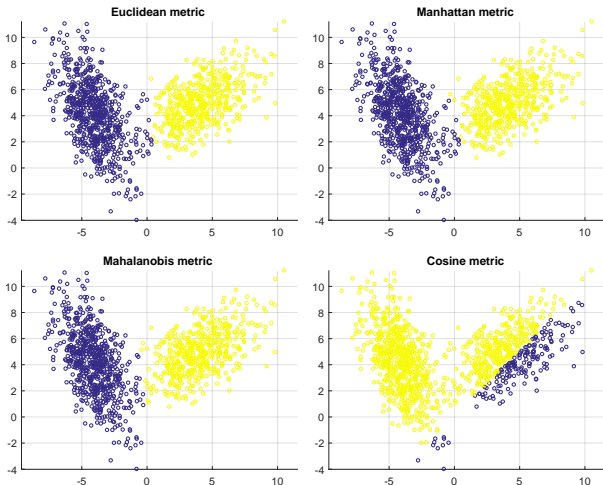
# Silhouette coefficient

- 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.

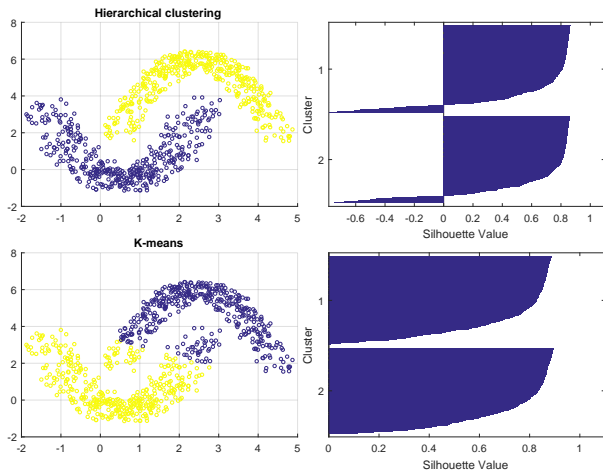


# Impact of distance functions

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



# Limitations



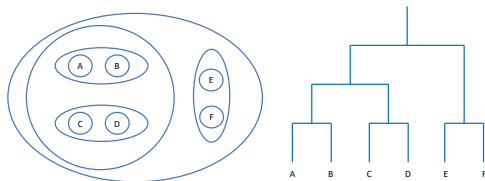


# Hierarchical clustering: Agglomerative clustering

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*.

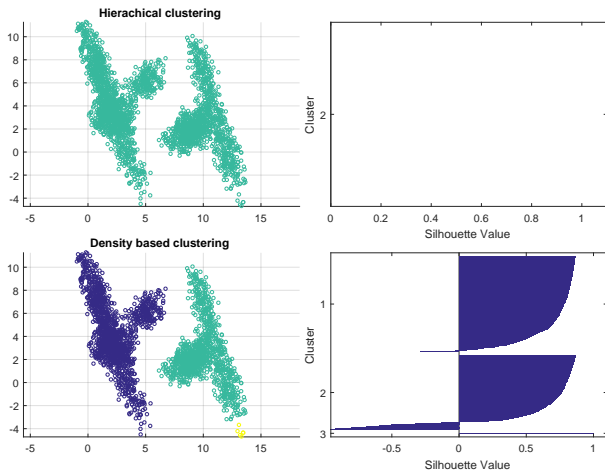
- Best (single) linkage. Distance is equal to the minimum distance between all pairs of elements (from two groups). Suitable to discover clusters of arbitrary shape. Drawback noise points may merge distant clusters.
- Worst (complete) linkage.(Complete linkage method) Distance is equal to the maximum 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 od 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 methods

One of the major problems with distance-based and probabilistic methods is that the shape of the underlying clusters is already defined implicitly by the underlying distance function or probability distribution. Possible solutions:

- Grid- based methods
- Density- based methods
- Graph- based algorithms
- Nonnegative matrix factorization

# 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  $\tau$ .

- Discretize each dimension into  $p$  ranges.
- Determine *dense* grid cells at level  $\tau$ .
- 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,  $\tau$  - density threshold and  $\epsilon$  - radius of the neighborhood.

## Definition

*Core point:* A data point is defined as the core point, if its  $\epsilon$  - neighbourhood contains at least  $\tau$  data points.

## Definition

*Border point:* A data point is defined as the border point, if its  $\epsilon$  - 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  $\epsilon$  and  $\tau$ .
- Create graph in which core points are connected (if they are within  $\epsilon$  of one another ).
- Assign each border point to a connected component.
- Return cluster indexes for each point.



# EM-algorithm

Let us consider K-Means from the probabilistic point of view.

- (E-step) Each data point of the set  $\mathcal{D}$  has a probability belonging to cluster  $j$ , which is proportional to the scaled and exponentiated Euclidean distance to each representative  $Y_j$ . In the k-means algorithm, this is done in a "hard" way, by choosing the smallest Euclidean distance to the representative of  $Y_j$ .
- (M-step) The center  $Y_j$  is the weighted mean over all the data points where the weight is defined by the probability of assignment to cluster  $j$ . The hard version of this is used in k-means, where each data point is either assigned to a cluster or not assigned to a cluster (i.e., 0-1 probabilities).

# EM-algorithm

Assumption: the data was generated from a mixture of  $k$  distributions with probability distributions  $\mathcal{G}_1 \dots \mathcal{G}_k$ . Each distribution  $\mathcal{G}_i$  represents a cluster and is also referred to as a mixture component.

- (E-Step) Given the current value of the parameters in  $\Theta$ , estimate the posterior probability  $P(\mathcal{G}_i|X_j, \Theta)$  of the component  $\mathcal{G}_i$  having been selected in the generative process, given that we have observed data point  $X_j$ . The quantity  $P(\mathcal{G}_i|X_j, \Theta)$  is also the soft cluster assignment probability that we are trying to estimate. This step is executed for each data point  $X_j$  and mixture component  $\mathcal{G}_i$ .
- (M-Step) Given the current probabilities of assignments of data points to clusters, use the maximum likelihood approach to determine the values of all the parameters in  $\Theta$  that maximize the log-likelihood fit on the basis of current assignments.

## Cluster Purity. NB! Not unsupervised any more!!!

- Let  $m_{ij}$  represent the number of data points from class (ground-truth cluster)  $i$  that are mapped to (algorithm determined) cluster  $j$ .
- Denote number of data points in true cluster  $i$  are by  $N_i$ , the number of data points in algorithm-determined cluster  $j$  by  $M_j$ .

$$N_i = \sum_{j=1}^{k_d} m_{ij}; \quad M_j = \sum_{i=1}^{k_t} m_{ij};$$

- For a given algorithm-determined cluster  $j$ , the number of data points  $P_j$  in its dominant class is:  $P_j = \max_i m_{ij}$ .
- Purity index is defined

$$P_a = \frac{\sum_{j=1}^{k_d} P_j}{\sum_{j=1}^{k_d} M_j}.$$

## Gini index

- Gini index for algorithm determined cluster  $j$  is defined:

$$G_j = 1 - \sum_{i=1}^{k_t} \left( \frac{m_{ij}}{M_j} \right)^2.$$

- Average Gini index is defined as follows:

$$G = \frac{\sum_{j=1}^{k_d} G_j M_j}{\sum_{j=1}^{k_d} M_j}.$$