

Data Mining

Lab 5: Data Clustering

Summary:

- Aim of data clustering
- Partitional clustering algorithms
- Hierarchical clustering algorithms
- Density-based clustering algorithms

1. Aim of data clustering

Data clustering aims to identify natural groups in data, i.e. subsets of similar data (according to a specific similarity measure) called clusters or classes. The particularity of a clustering task is the fact that the class labels (or even their number) are not known apriori. The goal of a clustering algorithm is to identify the clusters by using the relationships between the data.

2. Partitional clustering algorithms

These algorithms provides a partition of the initial dataset in several clusters (usually the clusters are disjoint but there are also algorithms which lead to overlapping clusters, e.g. fuzzy clustering algorithms). In the case of crisp algorithms (which assign each data to only one cluster) each cluster has a representative or cluster prototype. In the case when the cluster prototype is computed such that it is the average of the data from the cluster then it is called centroid.

The simplest and most popular partitional clustering algorithm is KMeans which generates a set of K centroids and assign each data to the closest centroid.

The general structure of the KMeans algorithm :

Step 1. Initialization: random selection of K centroids from the dataset

Step 2. **Repeat**

- Assign each data to the cluster represented by the nearest centroid
- Recompute the centroids (as averages of data in each cluster)

until the partition has not been changed during the last iteration

Remarks:

- The clustering result is sensitive with respect to the initial values of the centroids
- The KMeans iterative process aim to minimize the intra-cluster variance (the average sum of the distances between data and the centroid of their corresponding cluster)
- KMeans is appropriate for spherical clusters (e.g. data generated by normal distributions) but do not provide a good clustering in the case of arbitrary shaped clusters.

Weka implementations:

- **SimpleKMeans**: standard variant of the algorithm; the user can choose between the Euclidean and Manhattan distances
- **XMeans**: is a variant which estimates the number of clusters (the user provides a minimal and a maximal value for the number of clusters and XMeans applies KMeans for each of

these values and selects the variant leading to the best quality clustering – e.g. smallest intra-variance and largest inter-variance).

Exercise 1:

- a) Open the file “iris.2D.arff” and remove the class attribute
- b) Identify 3 clusters in data by applying KMeans (**Cluster->SimpleKMeans**). Visualize the identified clusters by right clicking on the result (from **Result list**) and select **Visualize Cluster Assignments**
- c) Analyze the values obtained for the intra-cluster variance (SSE = within cluster squared sum of errors) for several values of the number of clusters (the parameter **-N** from **SimpleKMeans**): 2,3,4,5
- d) Apply **XMeans** to the same set of data by using 1 and 10 as minimal and maximal number of clusters, respectively.
- e) Apply **EM** (Expectation-Maximization) to the same set of data using the default values for the parameters.

3. Hierarchical algorithms

These algorithms provide not only one partition but a hierarchy of partitions organized as a tree (dendrogram). The hierarchy can be obtained by one of the following approaches:

- *Agglomerative (bottom-up)*: at the beginning each cluster contains only one data and then at each step the most similar clusters are joined. The similarity between clusters can be measured using different criteria (single-link, complete-link, average-link). The merging process continues until all data belong to one cluster (this corresponds to the root of the dendrogram).
- *Divisive (top-down)*: the process starts with a unique cluster containing all data in the set and apply iteratively a partitioning clustering strategy (which can be KMeans)

Exercise 2:

- a) Open the file “data.arff”
- b) Construct and compare the dendrograms corresponding to the cases when different cluster similarity measures are used: **single-link**, **complete-link**, **average-link**. Hint: select **Cluster->Hierarchical**. To visualize the tree: right click on the result (from **Result List**) and select **VisualizeTree**

4. Density based clustering algorithms

The main idea of these algorithms is that the data are classified as core points, border points or noise based on the data density. For each data the density is estimated by counting the number of other data belonging to a neighborhood of a given radius.

Remark: the density based algorithms are used for spatial data and they allow to identify clusters of arbitrary shapes.

Exercise 3:

- a) Open the file “iris.2D.arff”
- b) Apply the algorithm DBSCAN for different values of the parameter Eps (neighborhood radius) and MinPoints (minimal number of data in the neighborhood) and compare the results. Test values: Eps in {0.2,0.4,0.6,0.8}, MinPoints in {5,10,15,20,30}