

# Business Intelligence

## SS 2018

### Cross-sectional Analysis 3

### Clustering Methods

W. Grossmann

# Content

- Methodology Cluster analysis
- Definition of distances
- Hierarchical Clustering
- Partitioning Methods
- Evaluation of Cluster solutions
- Other Methods

# Methodology Cluster Analysis

## Template: Cluster Analysis

- **Relevant Business and Data:** Customer behavior represented as cross-sectional data for process instances with a matrix  $X$  of explanatory variables
- **Analytical Goals:**
  - Find a segmentation of the data into clusters which allows an interpretation from domain point of view
  - Determine representatives for the clusters
- **Modeling Tasks:** Definition of a model for data description either based on the distances between the observations or by a mixture model for the distribution
- **Analysis Tasks:**
  - *Splitting Data:* If necessary split the data randomly into one set for training and another for testing the model
  - *Model Estimation:* Estimate the cluster solutions
  - *Model Assessment:* Assess the quality of models with respect to homogeneity of the clusters, separation between clusters, validity and reliability
  - *Model Selection:* Select a model by specifying the number of clusters
- **Evaluation and Reporting Task:** Evaluate the selected model using test data

# Methodology Cluster Analysis

- In case of clustering the data contain no output variable and we want to find a group structure in the data, such that the observations in the groups are rather homogeneous with respect to the variables

- Data  $N$  observations in  $p$  variables

Variables:  $\vec{X} = (X_1, X_2, \dots, X_p)$

Observations:  $\vec{x} = (x_1, x_2, \dots, x_p)$

- The “Relevant Business and Data” task should be done in the same way as for classification

# Definition of Distances

- Main prerequisite for clustering is a distance between observations
- Most important distance for quantitative variables is the Euclidean distance:

$$d^2(\vec{x}, \vec{z}) = \sum_{i=1}^p (x_i - z_i)^2$$

- Alternatives:

- Absolute deviation  $d_1(\vec{x}, \vec{z}) = \sum_{i=1}^p |x_i - z_i|$

- Maximum distance:  $d_s(\vec{x}, \vec{z}) = \max_{i=1}^p |x_i - z_i|$

# Definition of Distances

- In case of binary variables the Hamming distance is frequently used, which is defined by the number of different values in the variables
  - The Hamming distance is equivalent to the Euclidean distance

# Definition of Distances

- Another measure frequently used is the Cosine similarity defined by the angle between two vectors
  - The cosine defines a similarity of two vectors; if they have the same direction, the similarity is 1, if the vectors are orthogonal the similarity is 0, if the vectors have opposite direction the similarity is -1
  - From statistical point of view the cosine similarity measures the correlation between two observation vectors
  - Cosine similarity is frequently used if the vectors are binary vectors

# Definition of Distances

- In case of categorical (qualitative) variables one can reduce the problem to a problem for binary variables by defining indicator variables for each attribute value
  - $q$  different values lead to  $\ln_2(q)$  binary variables
- Combination of qualitative and quantitative variables: procedure **daisy** in R



# Definition of Distances

- For more complex structures like string variables (text) or graphs the distance calculation can be based on kernels
  - String kernels: based on counting the simultaneous occurrence of substrings of certain length

# Definition of Distances

- An important issue is many times standardization of the variables
  - All variables are standardized with mean zero and variance 1
  - All variables are standardized such that the values are in the interval  $[0, 1]$ , or  $[-1, 1]$

# Hierarchical Clustering

- Most frequently used is agglomerative clustering
- Basic outline of an agglomerative cluster algorithm:

---

## Agglomerative clustering

---

```
1 Define clusters  $C_k, 1 \leq k \leq N$  by the observations,  $N_{cl} = N$ ;  
2 for  $k = 1$  to  $N - 1$  do  
3   | Merge clusters  $C_r$  and  $C_s$  for which  $d(C_r, C_s) = \min_{(l,k)} D((C_l, C_k))$ ;  
4   |  $N_{cl} = N_{cl} - 1$ ;  
5 end
```

---

# Hierarchical Clustering

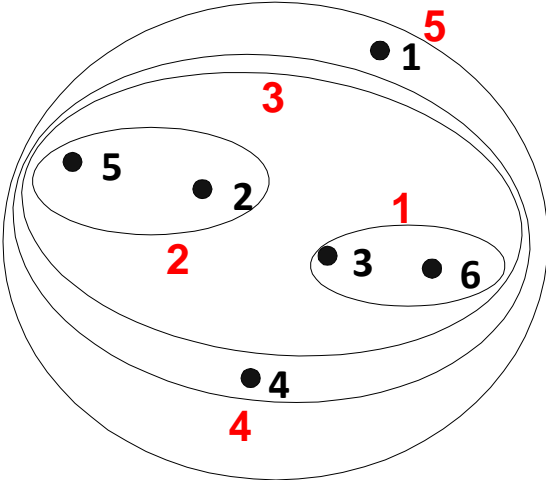
- Main problem is the definition of the distance between the clusters based on the distance between the objects, and determination of number of clusters
- The distance is called the linkage of the clusters
- Different specifications are possible
  - Single linkage: Distance between clusters is the distance of the closest points (minimum spanning tree)
  - Complete linkage: Distance between clusters is the distance of the farthest points

# Hierarchical Clustering

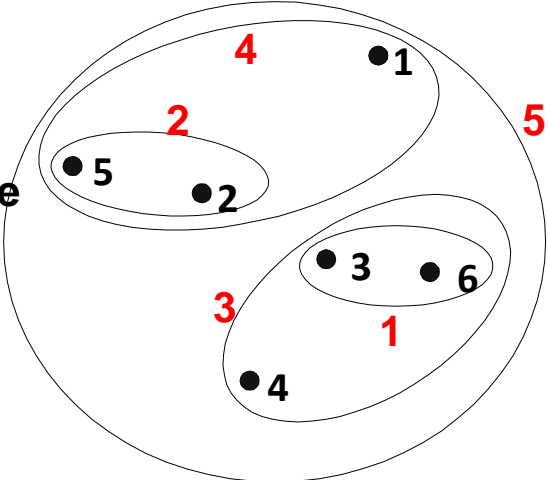
- Average linkage: mean distance between all the point in the two clusters
- Ward distance: difference between the total within cluster sum of squares for the two clusters separately, and the within cluster sum of squares resulting from merging the two clusters in cluster
- In general average linkage and Ward's method are recommended
- Single linkage usually not recommended because the clusters are often a chain-shaped and not ball-shaped

# Hierarchical Clustering, Comparison of Linkages

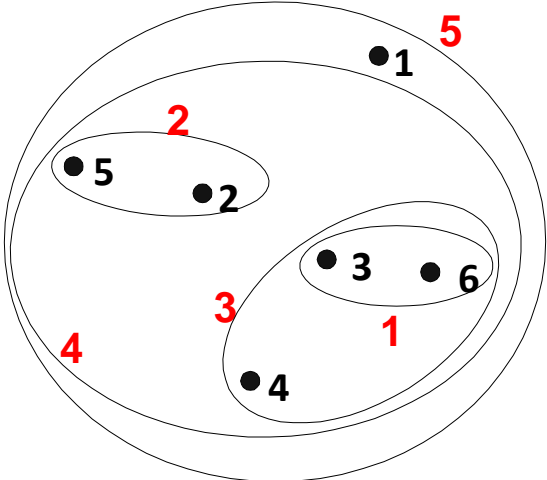
(Clusters shown in red)



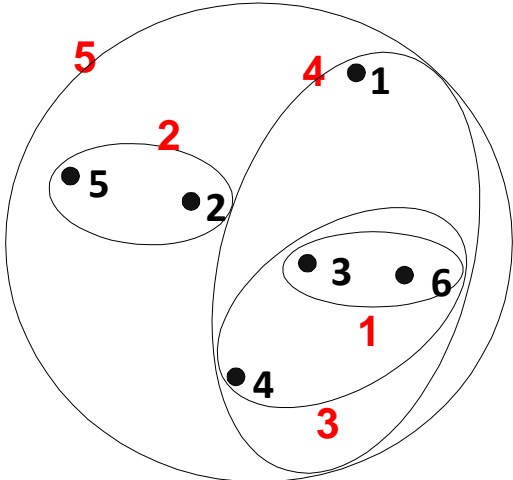
Single



Complete



Average



Ward

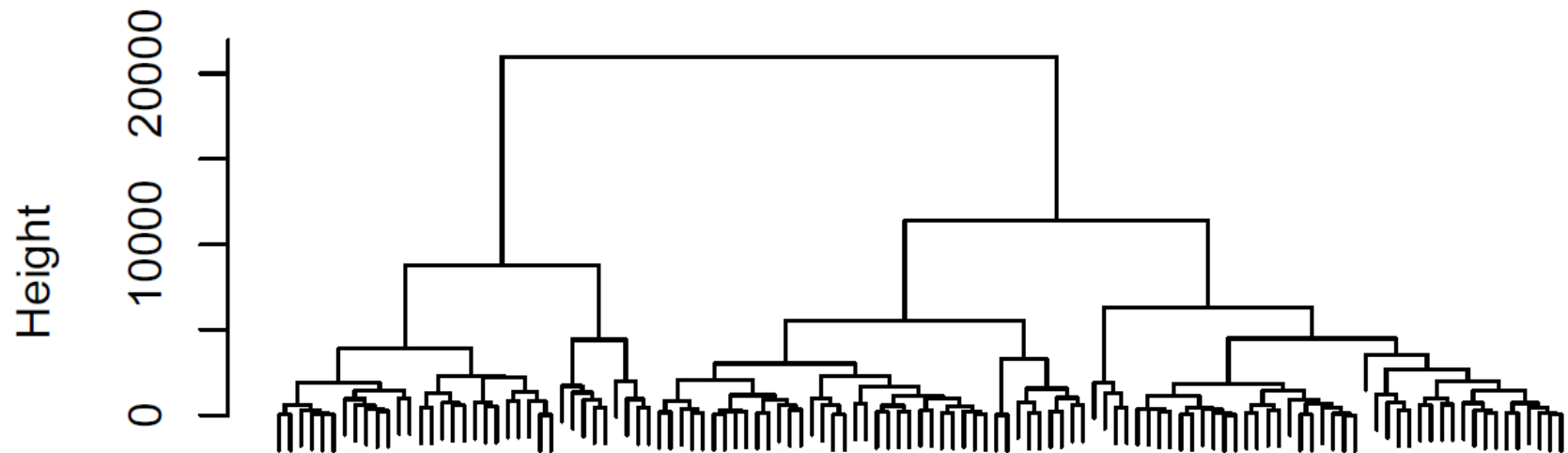
# Hierarchical Clustering

- Determination of the number of clusters is done by visual inspection of the distance of the clusters which are merged
- Most popular method is using the dendrogram
  - A dendrogram is a visual representation of the aggregation process as a tree
  - The leaves of the tree are defined by the objects
  - Other nodes are formed according to the aggregation process
  - The heights of branches is given by the distance

# Hierarchical Clustering

- Example of a dendrogram

**Dendrogramm, complete linkage**



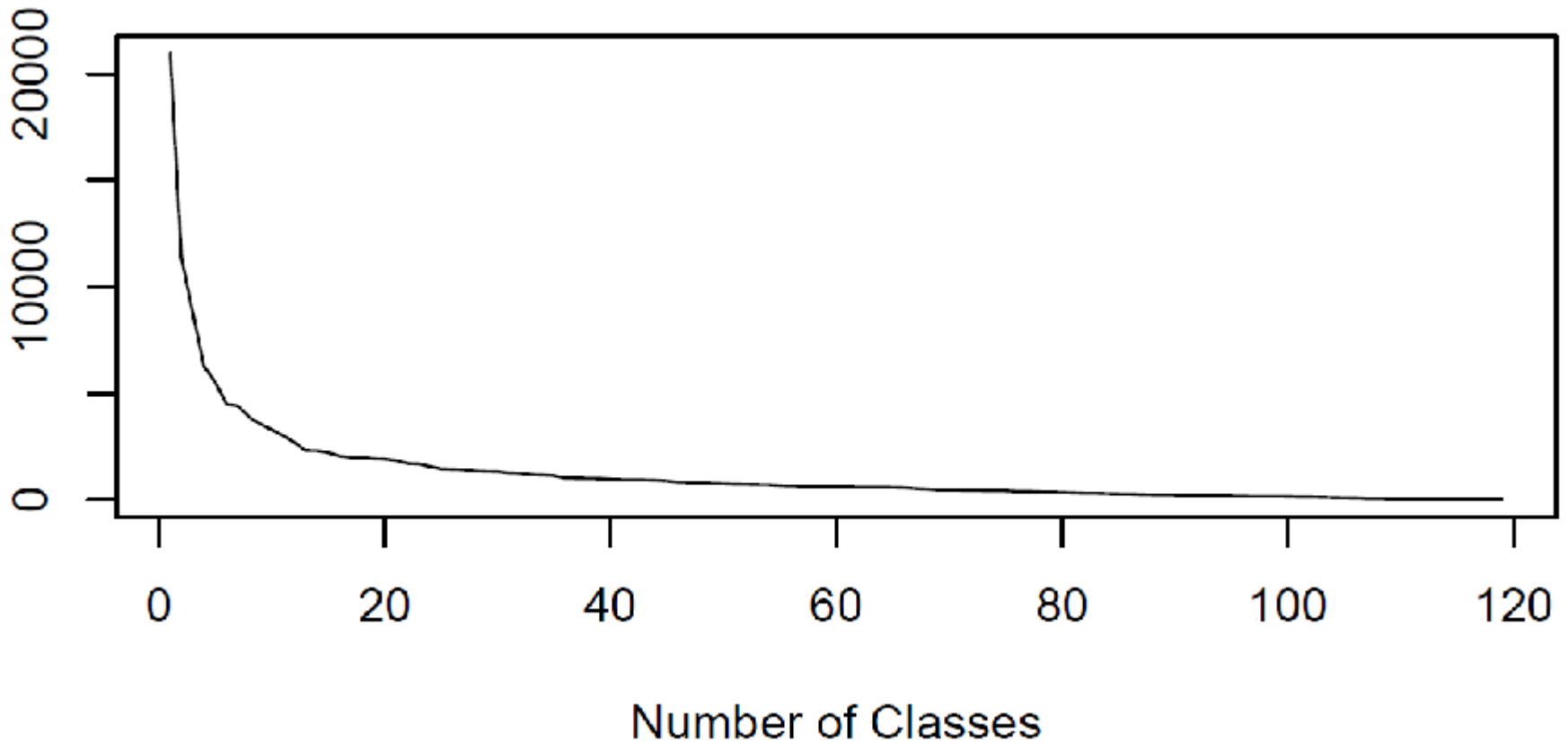


# Hierarchical Clustering

- An alternative to the dendrogram is a scree plot of the distance between the merged classes in dependence of the number of classes
- The decision about the number of classes is defined by the elbow of the scree plot

# Hierarchical Clustering

- Example of a scree plot



# Hierarchical Clustering

- Properties of Hierarchical clustering
  - Good visualization of the solution
  - Decision about the numbers of clusters can be done after analysis
  - Limited to small number of observations
  - The decision about the cluster assignment cannot be changed in the algorithm

# Partitioning Methods

- Partitioning methods define in an iterative way a cluster solution for the observations given the number of clusters in advance
- The most popular method is k-means clustering, where k stands for the number of clusters

# Partitioning Methods

- Basic algorithm

---

## k-Means Algorithm

---

**Data:** Observation matrix  $X$  and distance for the objects; number of clusters  $K$ .

**Result:** Cluster solution for observations

1 **begin**

2     Define an initial solution for the cluster centers  $(c_1, c_2, \dots, c_k)$ ;

3     Assign each observation  $x$  to the cluster which center is closest to the observation;

4     Compute new centers for the clusters as means of the assigned observations;

5     Repeat steps 2 and 3 as long as there is no significant change in the centers;

6 **end**

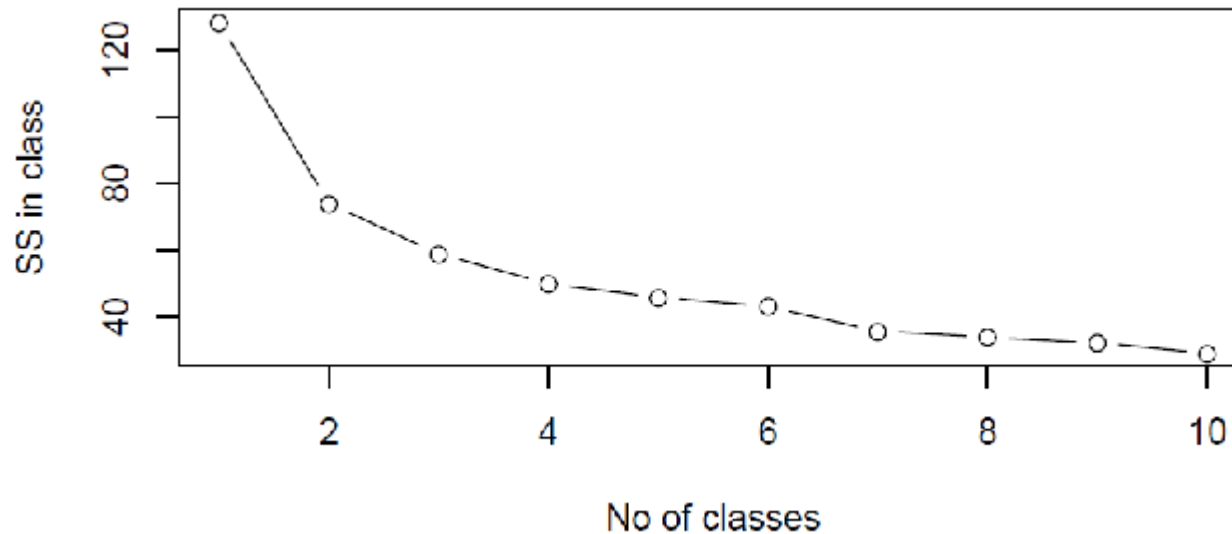
---

# Partitioning Methods

- Main points in application
  - Decision about the number of clusters
  - Finding initial centers of the clusters
- Determination of the number of clusters can be based on a visualization of the sum of squares within the clusters for a solution in dependence of the number of clusters
  - This is similar to the idea of variance decomposition in case of analysis of variance

# Partitioning Methods

- Usually the plot shows a shape similar to a scree plot and the decision is based on the elbow criterion



# Partitioning Methods

- With respect to the initial solution the standard procedure is choosing the centers randomly and try different solutions
- Properties of k-means clustering
  - Procedure is fast from computational point of view
  - Applicable for large data sets (parallel implementations exist)
  - A found solution can be applied to new observations (cf. nearest neighbor classification)



# Evaluation of the cluster solution

- Evaluation of clusters
  - Homogeneity measures for groups
  - Validation with a test sample
  - Validity of the solution with respect to a subject matter explanation

# Evaluation of the cluster solution

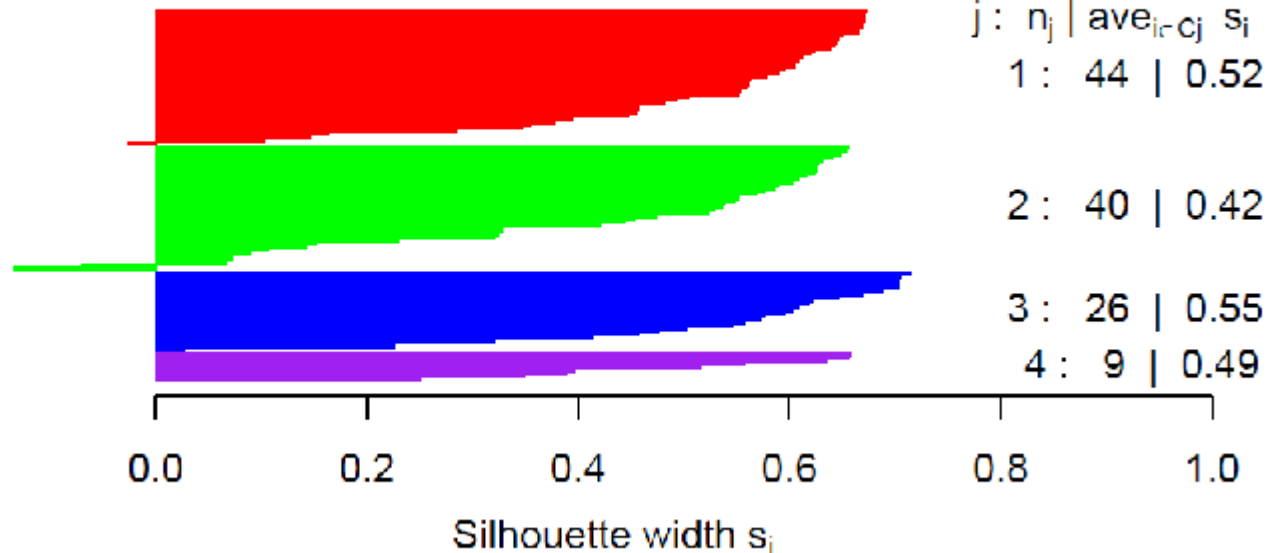
- A method for evaluation of the cluster solution is the silhouette plot
  - The silhouette shows for each point how well the point is located in the cluster
  - A value close to 1 shows that the point is well located in the cluster
  - A negative value indicates that the point is not well assigned
- The silhouette can be used also for comparing different cluster solutions

# Evaluation of the cluster solution

- Example of a silhouette

## Silhouette of Clusters

n = 119



Average silhouette width : 0.49

# Evaluation of the cluster solution

- Another method for visualization useful for a moderate number of observations is using principal component analysis for the data and showing the clusters in the first two components
- An alternative to principal components is using a representation based on multidimensional scaling
- In case of a large number of observations boxplots of the variables grouped by the clusters is often useful
- For visualization see also the demo example

# Other Methods

- There exist numerous clustering algorithms for specific problems
  - Self organizing maps (SOM) are clustering methods based on the idea of a neural net. They can be understood as a k-means clustering defined on a distorted grid
  - Two stage clustering combine the ideas of hierarchical clustering with the ideas of k-means (IBM/SPSS) by using a cluster-feature tree
  - Model based clustering looks at the problem from a more theoretical perspective and defines a model for the data by a mixture of normal distributions