

# Cluster analysis

Petr Ocelík

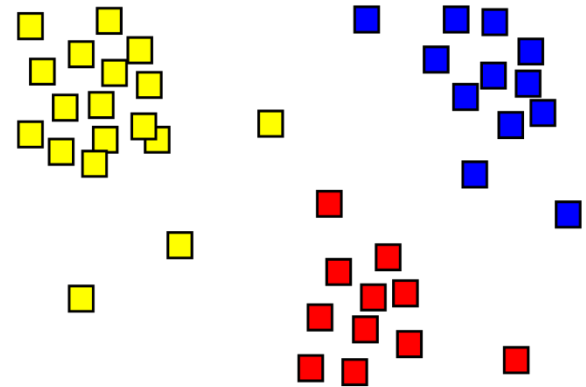
MVPd002 Quantitative Research in International and European Politics

# Plan for today

- Intuition
- Cluster analysis step-by-step
- Exercise

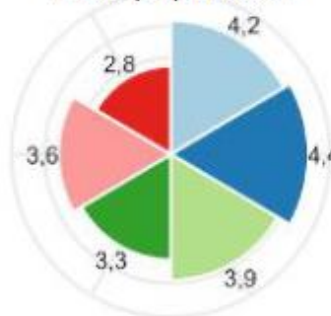
# Cluster analysis

- **Data reduction** technique
- **Cluster:** a grouping of similar objects
- Basic idea: identifying groups of mutually **similar objects** based on particular variable(s)
- **Unsupervised** technique

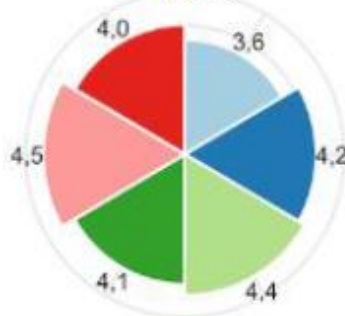


secured middle class   emerging cosmopolitan class   traditional working class

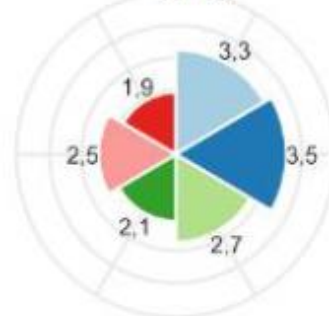
22 % population



12 %



14 %



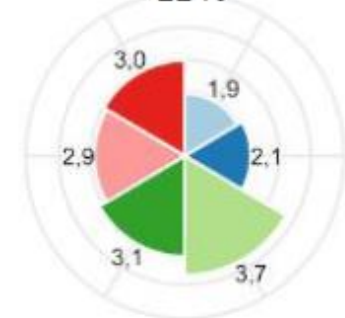
class of local ties

12 %



endangered class

22 %

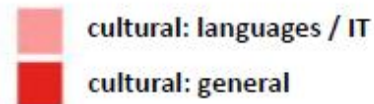
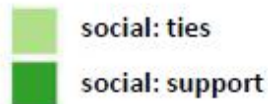
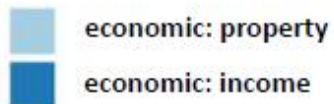


suffering class

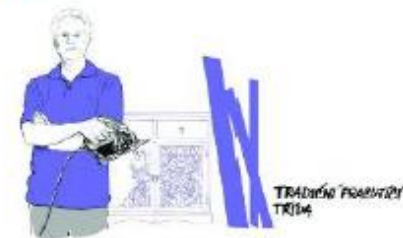
18 %



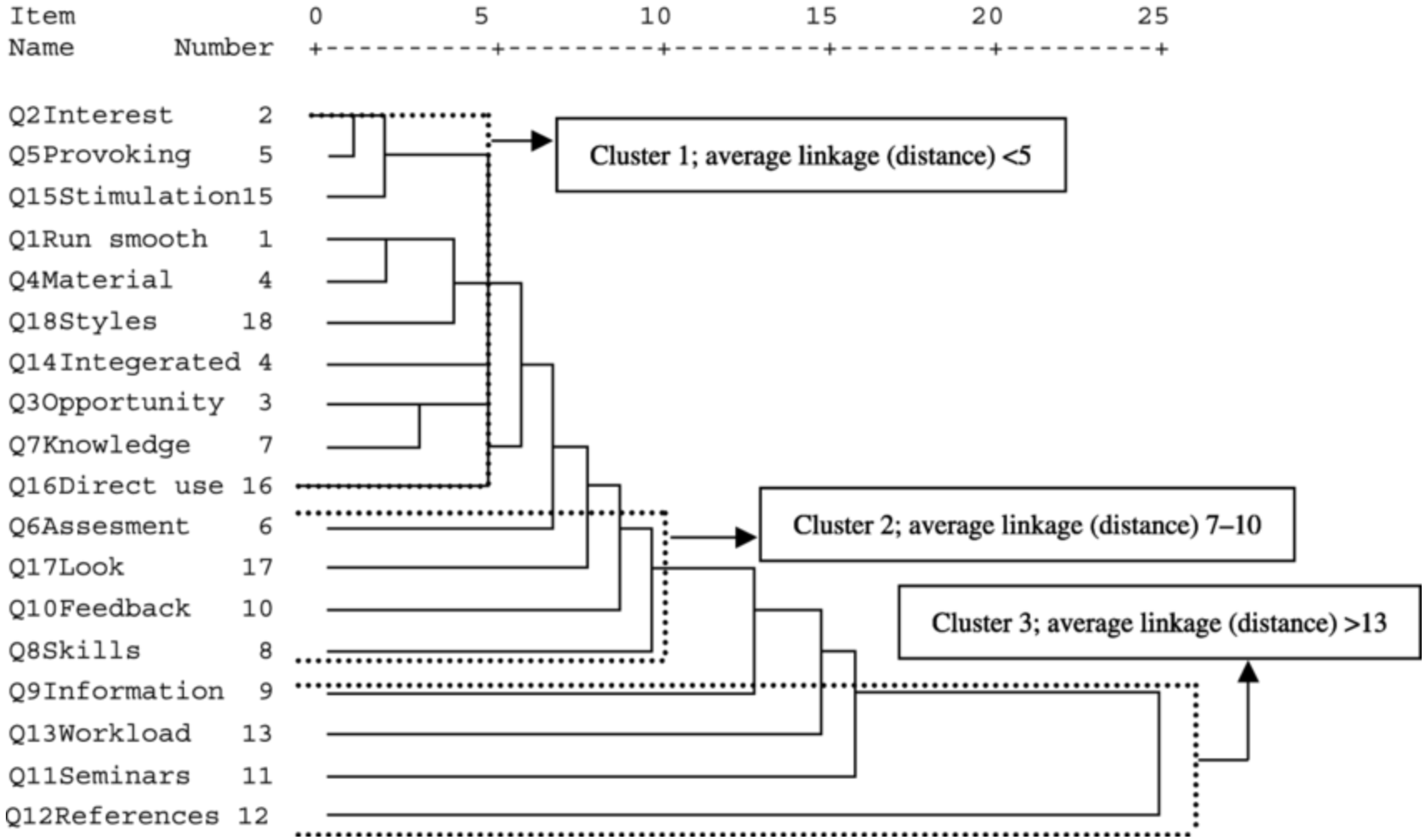
capital forms



Prokop et al. 2019  
Kočí et al. 2019

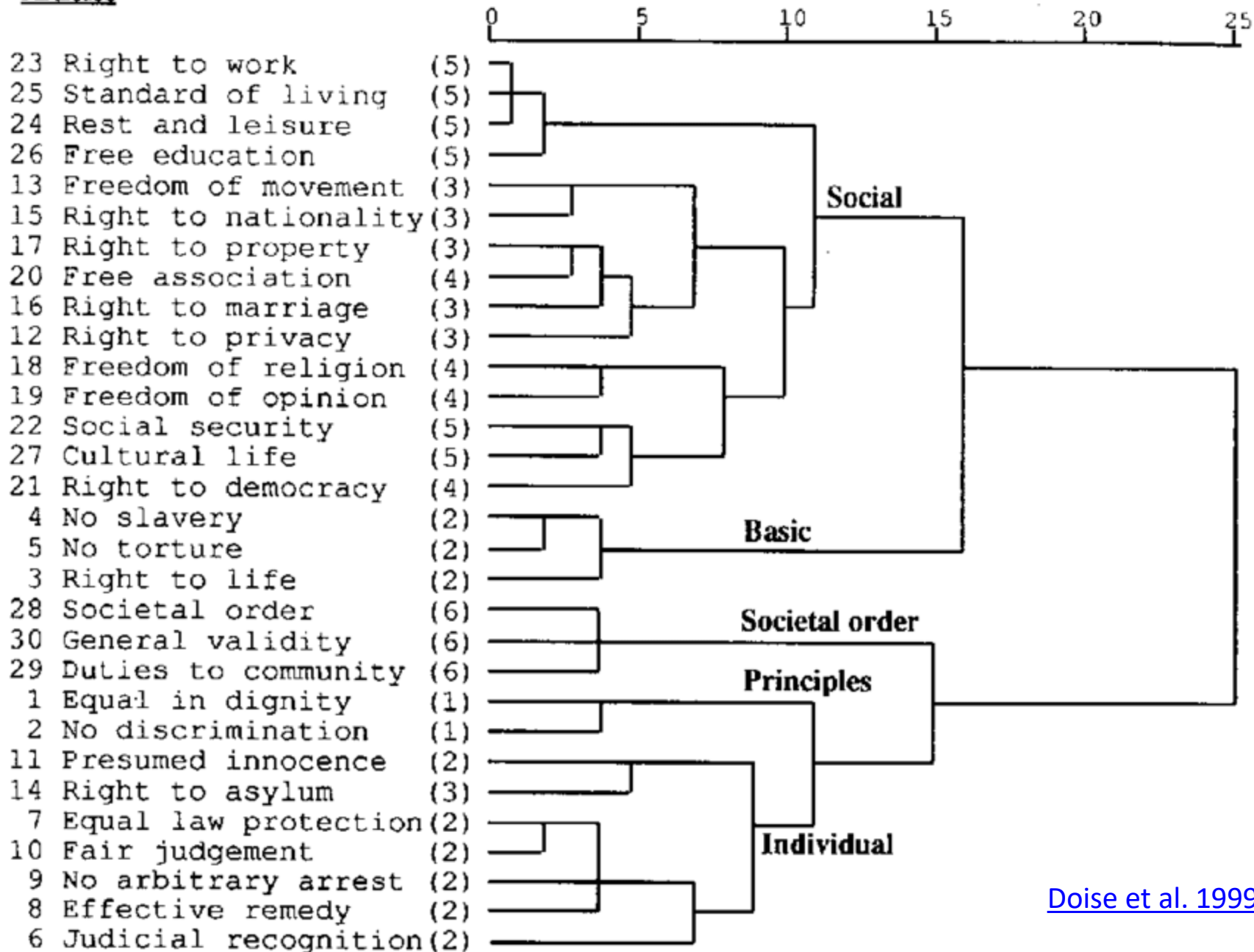


Rescaled Distance Cluster Combine

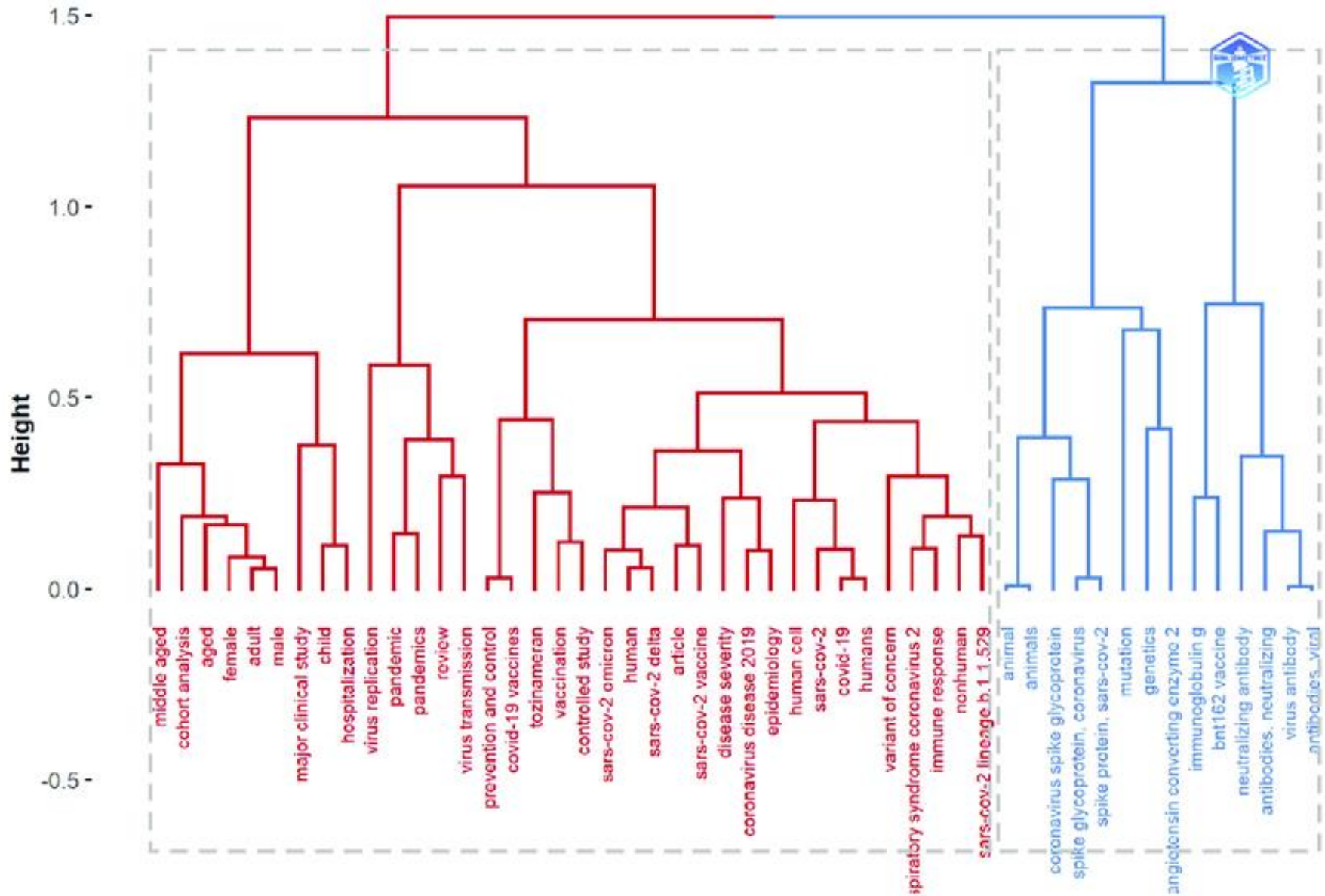


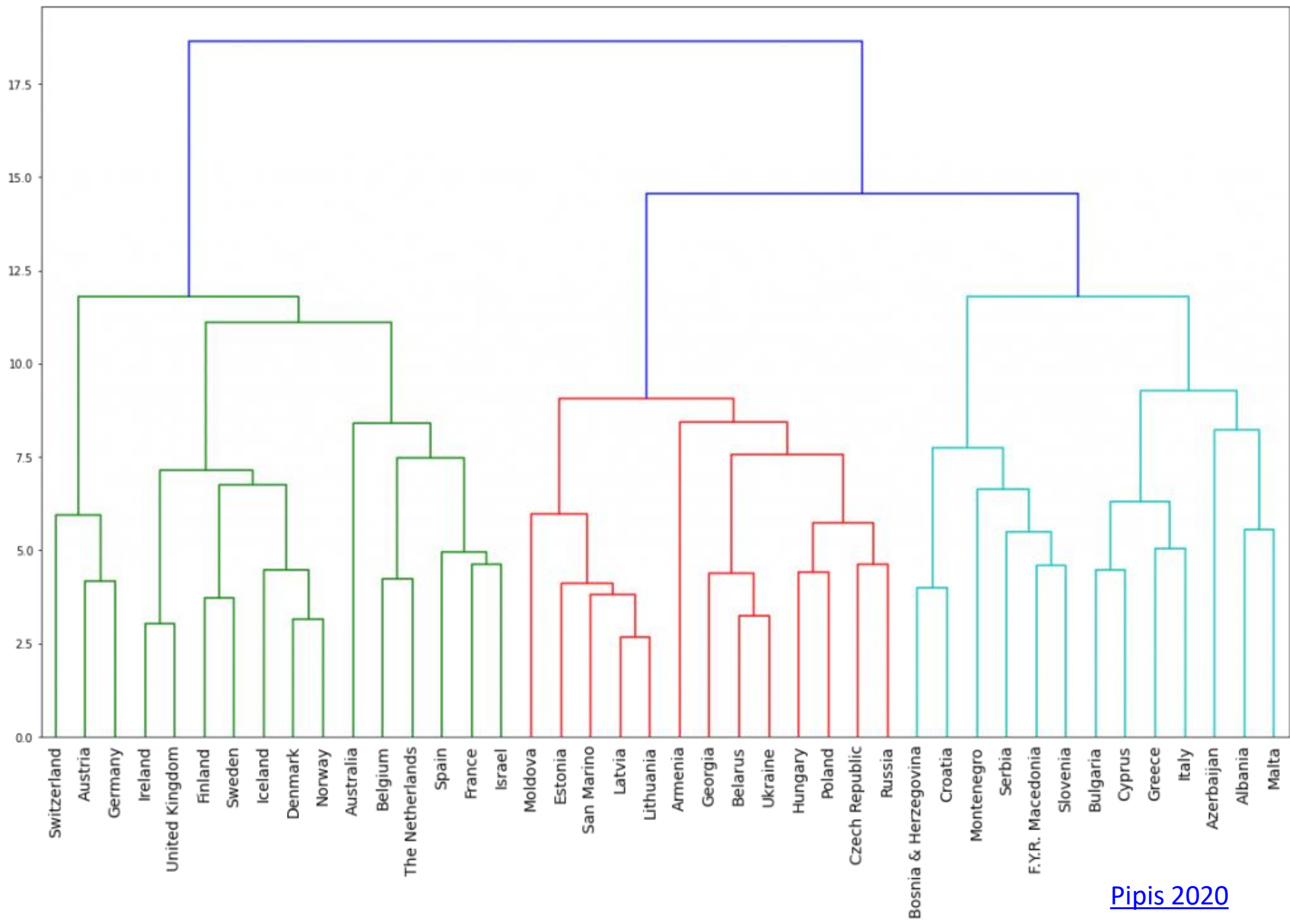
Rescaled dissimilarity coefficient

Articles



# Topic Dendrogram







# Cluster analysis: process

1. Sampling and data collection
2. Similarity measures
3. Clustering methods
4. Cluster solution interpretation
5. Cluster solution diagnostics

# 1. Sampling and data collection

- What is a **target population**?
- What is the **level of analysis**?
- What is the **unit of observation**?
- What **set of variables** are we interested in?
- Practical considerations

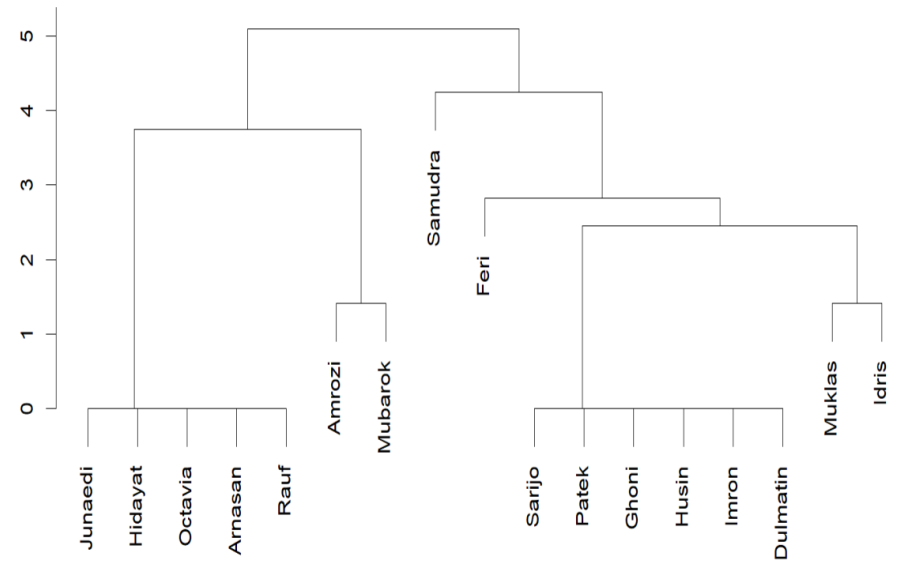
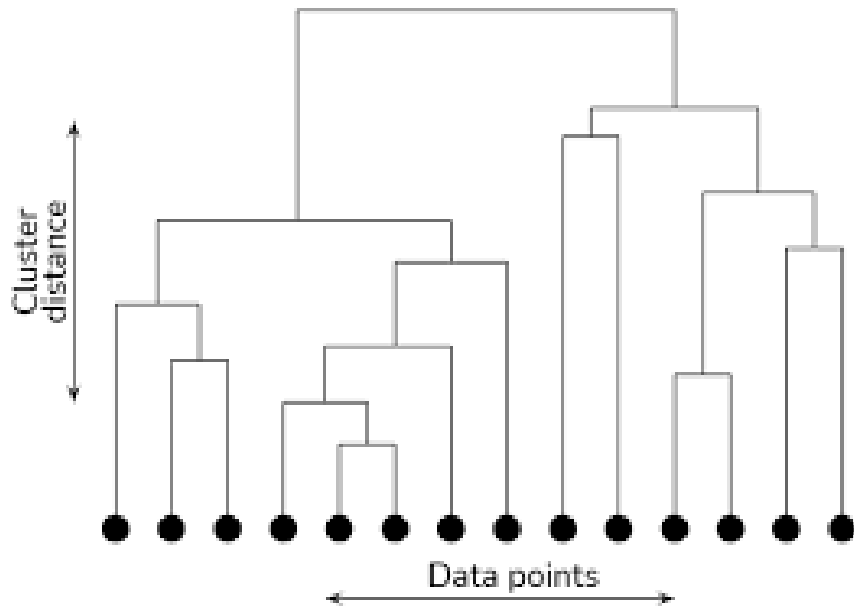
## 2. Similarity measures

- (Dis)similarity of objects is quantified by using **similarity measures**
- Choice of similarity measure needs to consider:
  1. **level of measurement**: categorical vs continuous
  2. **data dimensionality**: number of variables (vars)
  3. **scale sensitivity**: small vs large data, vars scales
- We distinguish between **association-based** and **distance-based** similarity measures (not exhaustive)

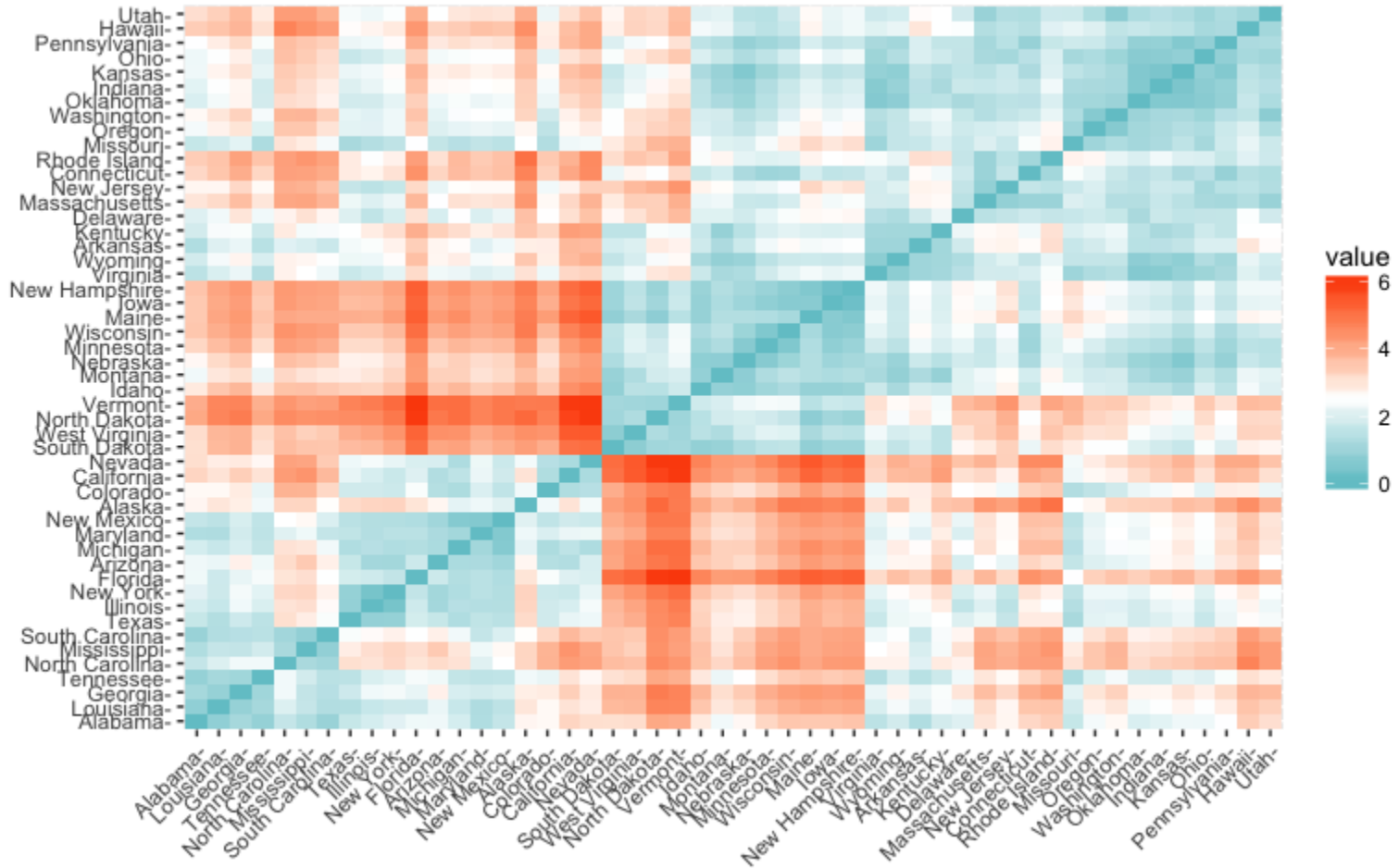
## 2. Similarity measures: representations

- The degree of (dis)similarity can be captured numerically or **graphically**.
- Dendrogram
- Heatmap
- Cluster profile

# Dendrogram

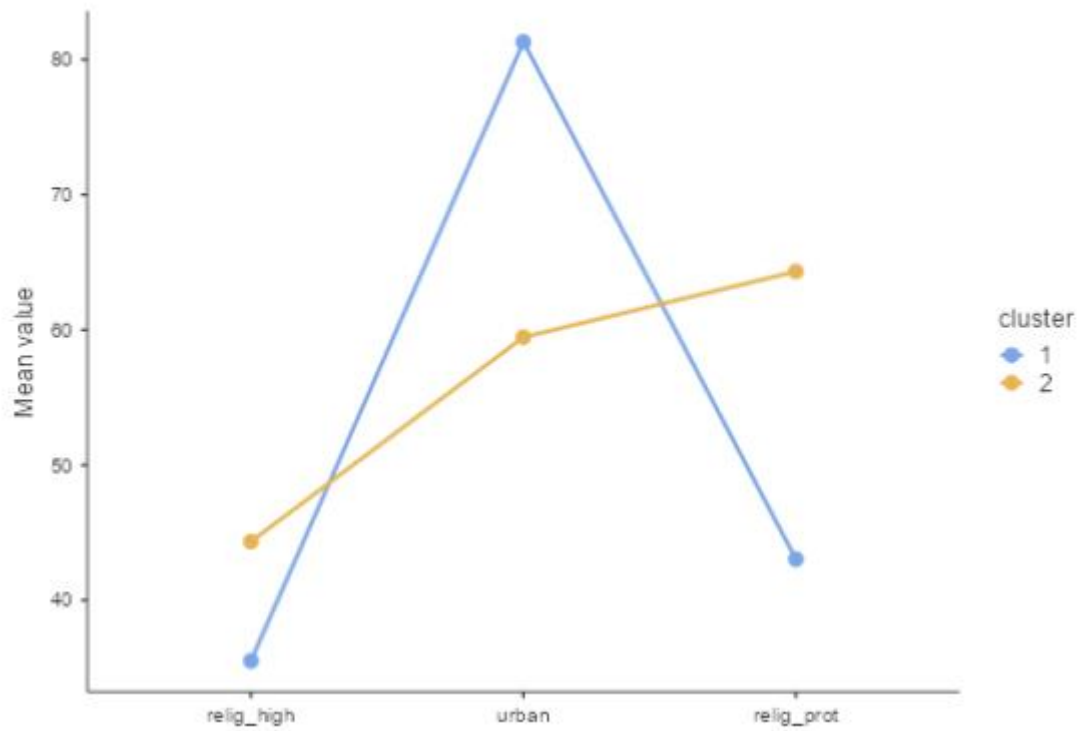


# Heatmap



# Cluster profile

Plot of means across clusters



## 2.1 Pearson's coefficient

- Pearson's  $r$  measures the existence, strength and direction of the **linear relationship** between two variables

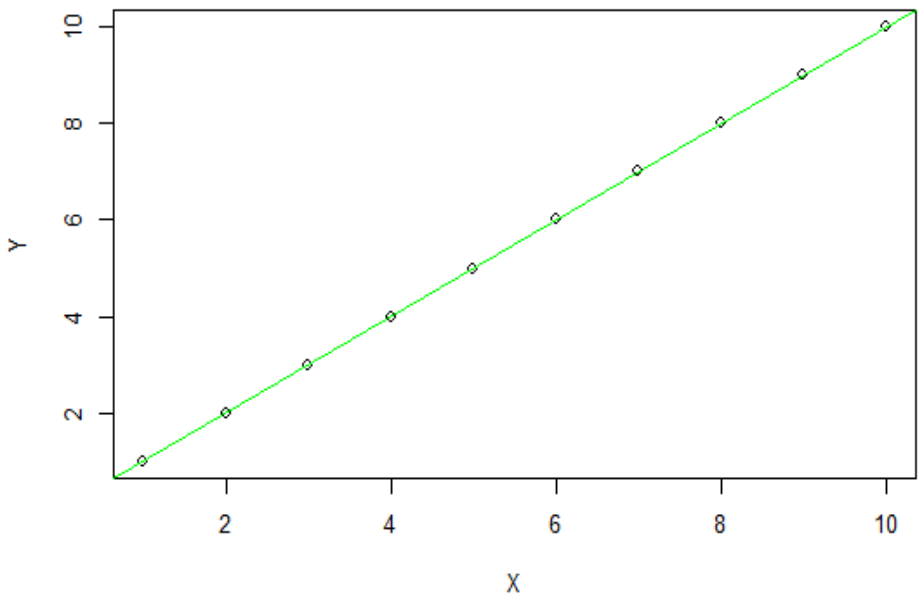
measurement level	number of values	range	coefficient
continuous-continuous	many-many	$\langle -1, 1 \rangle$	Pearson's $r$

Soukup et al. 2022

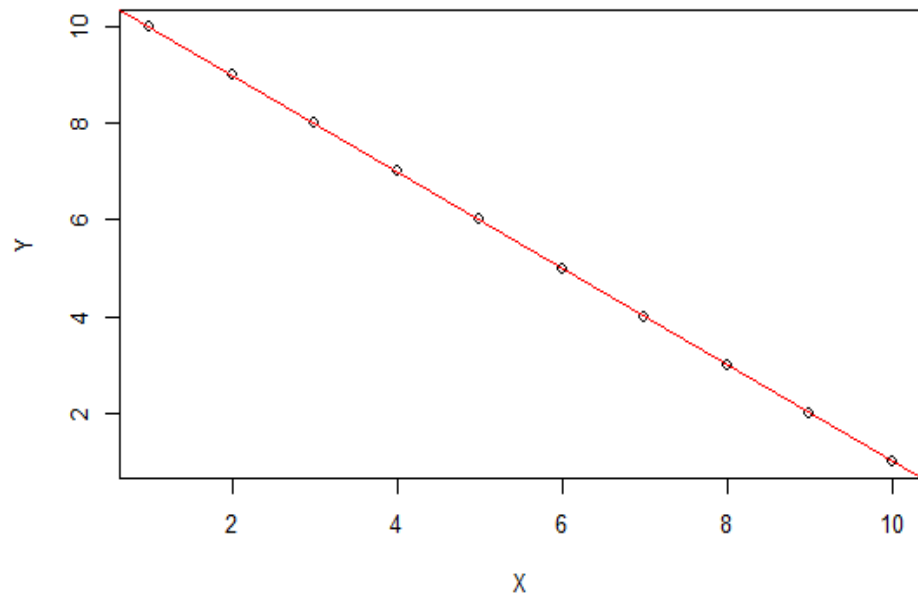
- Not suitable for heterogeneous data and/or nonlinear relationships between variables



**r=1**

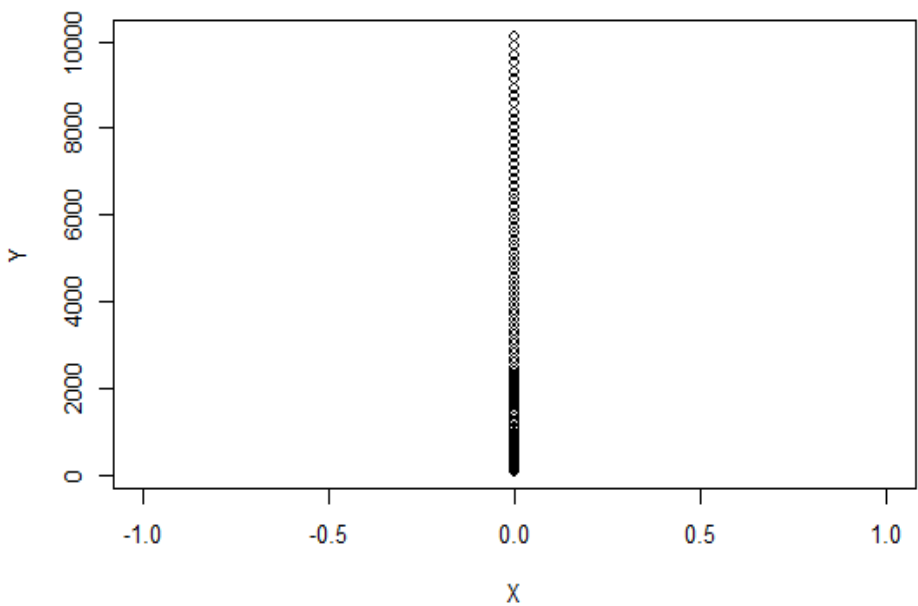


**r=-1**

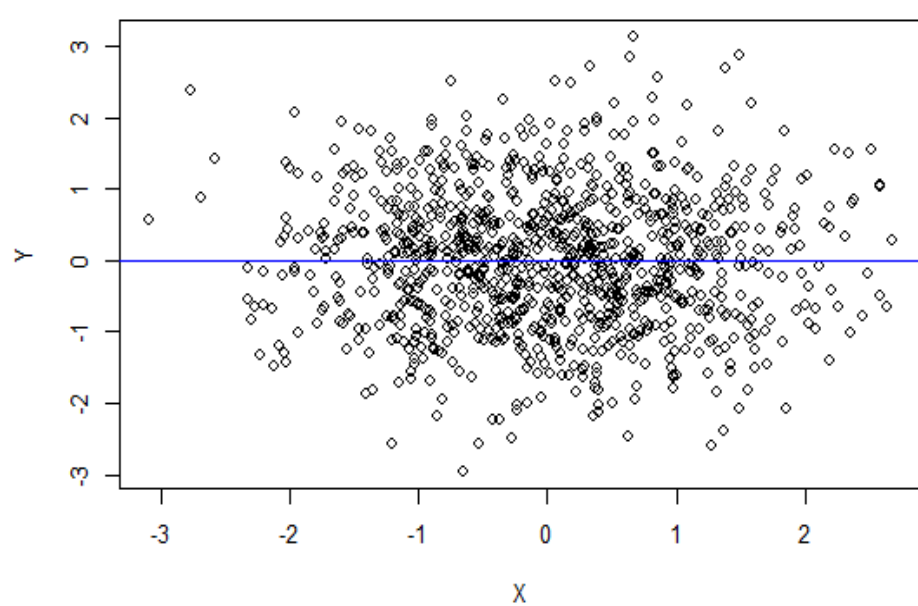


<http://guessthecorrelation.com/>

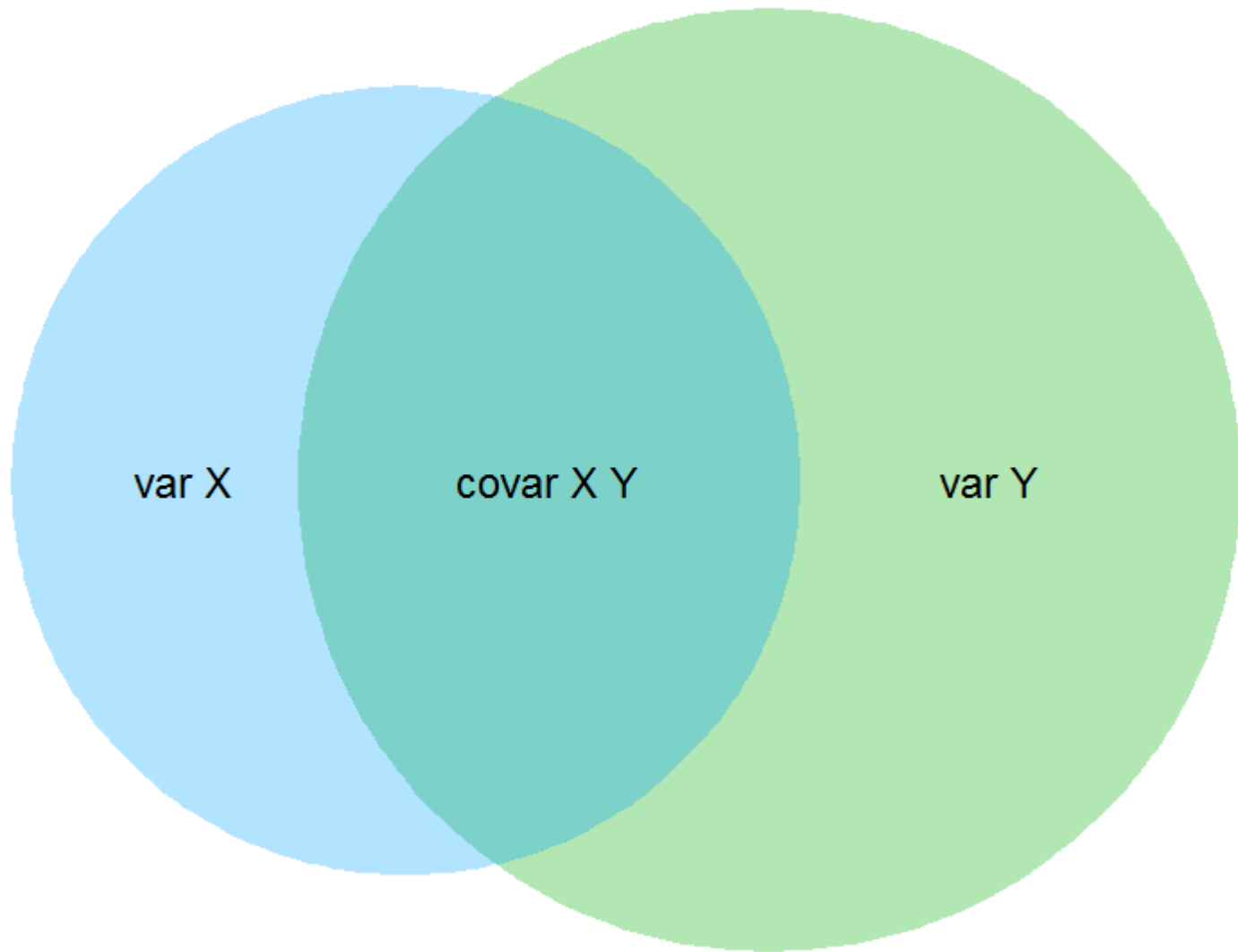
**r=0**

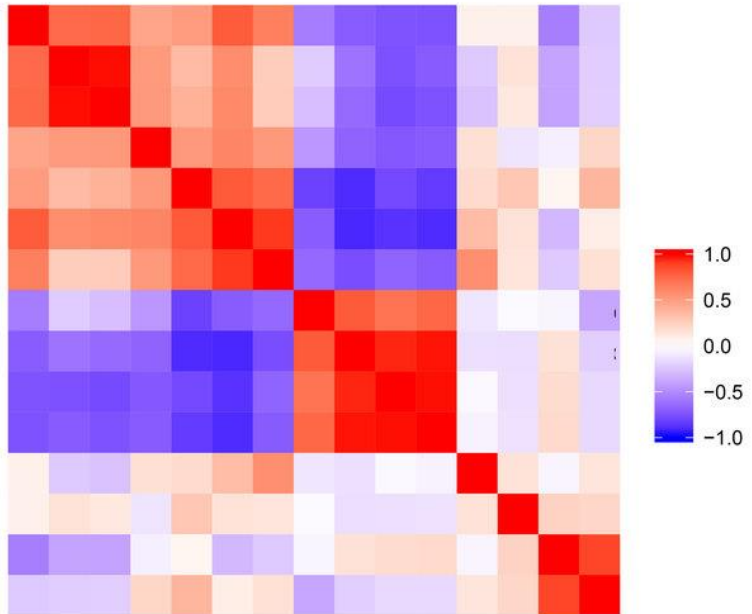
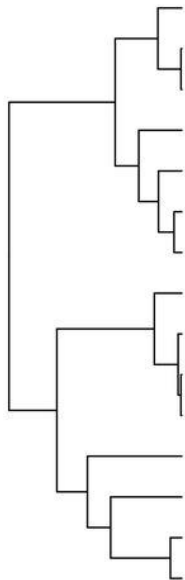
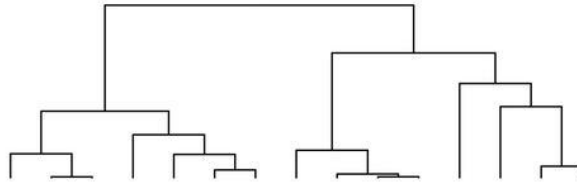


**r=0**



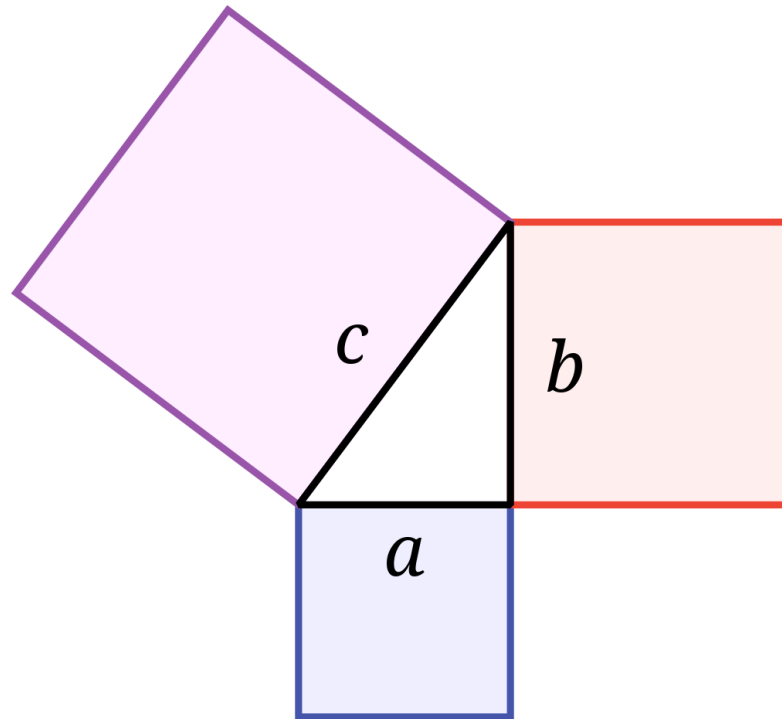
- $r = \text{covariance} / \text{combined total variance}$





## 2.2 Euclidean distance

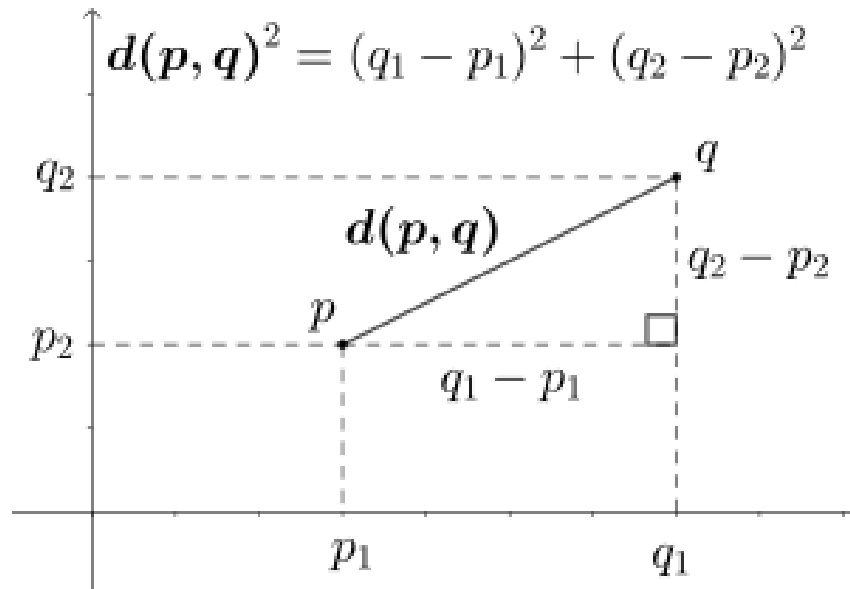
- Do you recall this?



$$a^2 + b^2 = c^2$$

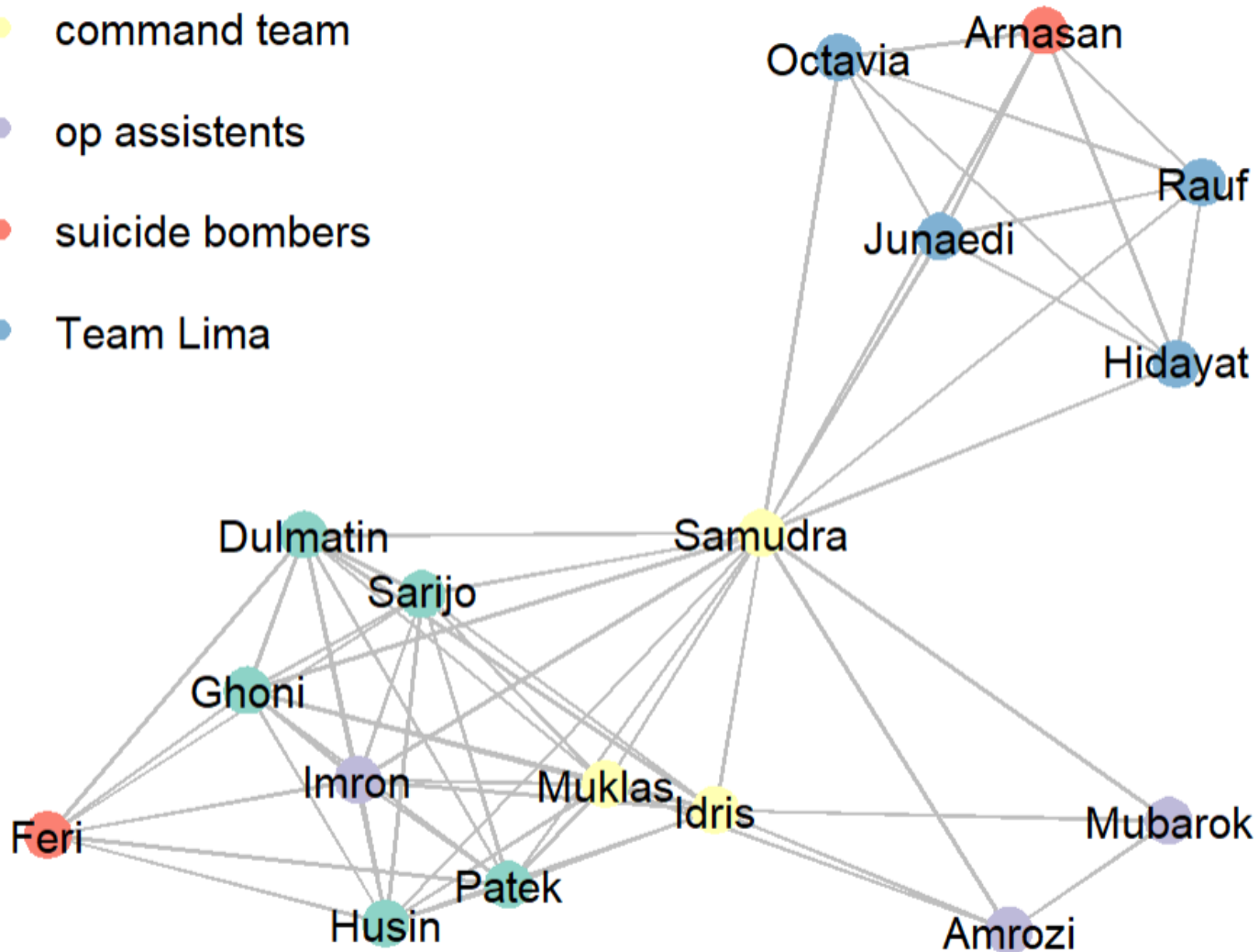
## 2.2 Euclidean distance

- **Euclidean distance (d)** is a generalization of the Pythagorean theorem
- $d(p, q)$ , of two objects  $(p, q)$  equals **the length of the straight line between them**



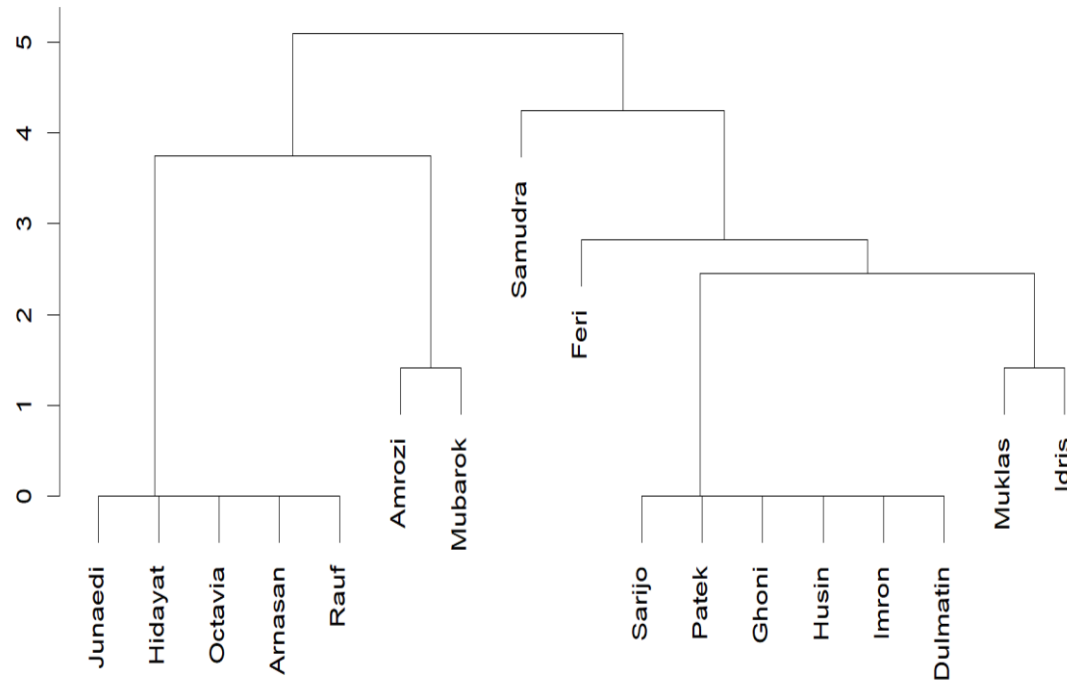
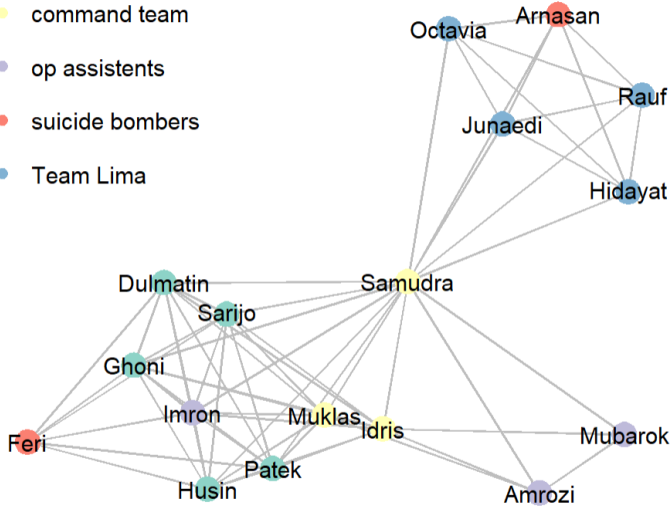
$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2}.$$

- bomb makers
- command team
- op assistants
- suicide bombers
- Team Lima



# 2.2 Euclidean distance

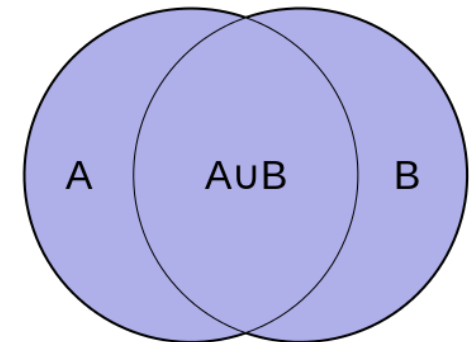
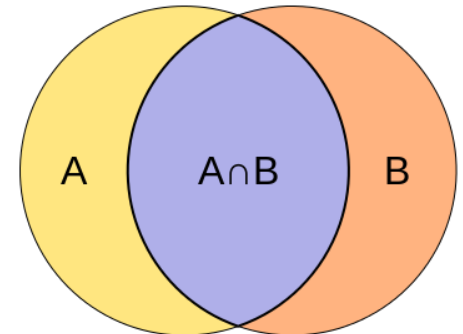
- bomb makers
- command team
- op assistants
- suicide bombers
- Team Lima



## 2.3 Jaccard's coefficient

measurement level	number of values	range	coefficient
categorical-categorical	binary	$\langle 0, 1 \rangle$	Jaccard's

		sample B	
		present	absent
sample A	present	$A \cap B$	$A - B$
	absent	$B - A$	$\notin A \cup B$

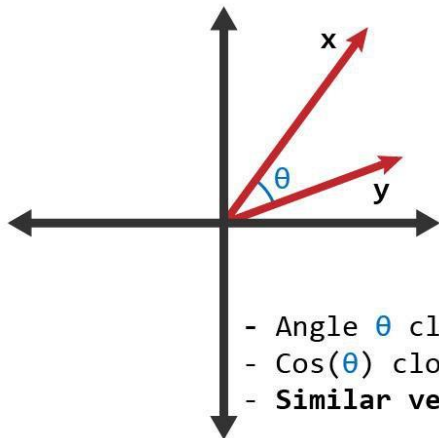


- $J$  = the **size of the intersection** ( $A \cap B$ ) by the **size of the union** ( $A + B = A \cup B$ ) of the samples
- $J = A \cap B / (A \cup B)$
- Does not account for observations missing in both samples ( $\notin A \cup B$ )

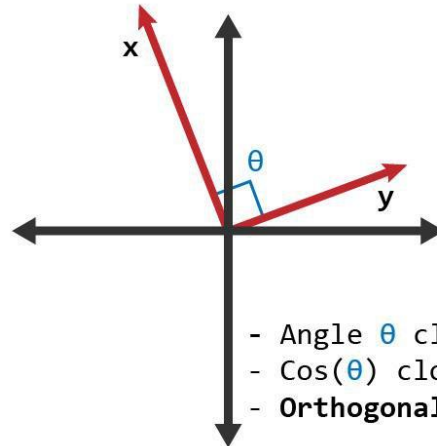
wikimedia commons



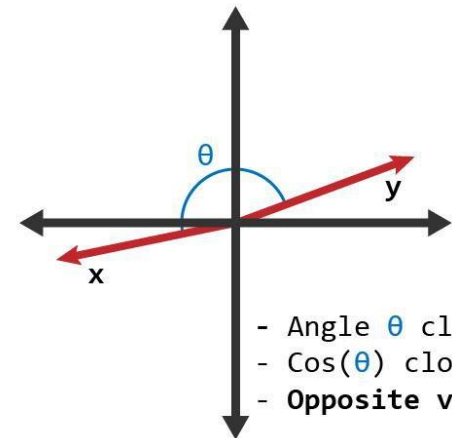
## 2.4 Cosine distance



- Angle  $\theta$  close to  $0$
- $\text{Cos}(\theta)$  close to 1
- **Similar vectors**



- Angle  $\theta$  close to  $90$
- $\text{Cos}(\theta)$  close to  $0$
- **Orthogonal vectors**



- Angle  $\theta$  close to  $180$
- $\text{Cos}(\theta)$  close to  $-1$
- **Opposite vectors**

## 2.4 Cosine distance (CD)

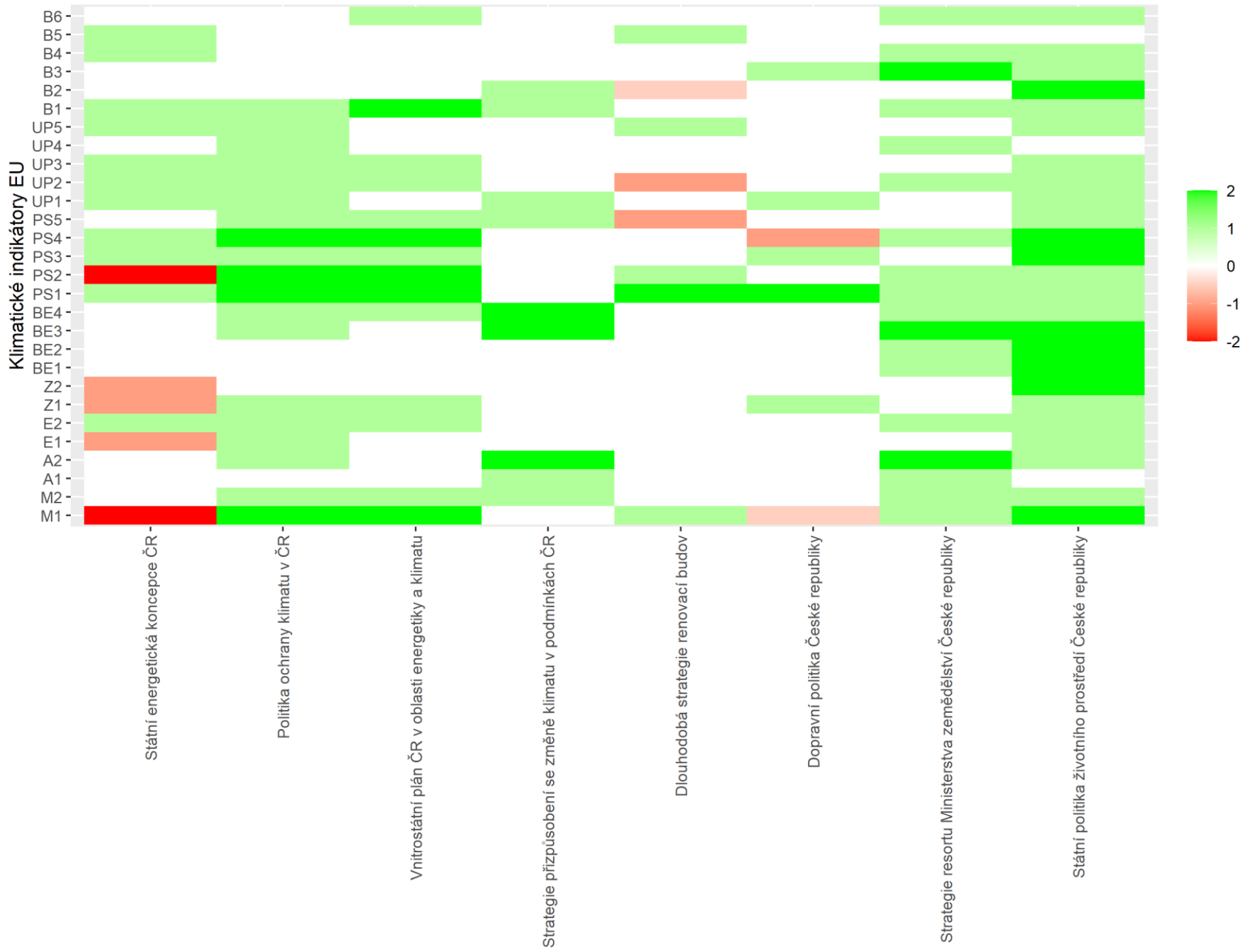
	doc_A	doc_B
geopolitics	4	0
climate change	0	1
Brexit	0	12
Euro	3	9
sovereignty	5	2

$$\text{cosine similarity} = |A||B|\cos\theta = \frac{A \cdot B}{|A||B|} = \frac{\sum_i^n A_i B_i}{\sqrt{\sum_i^n A_i^2} \sqrt{\sum_i^n B_i^2}}$$

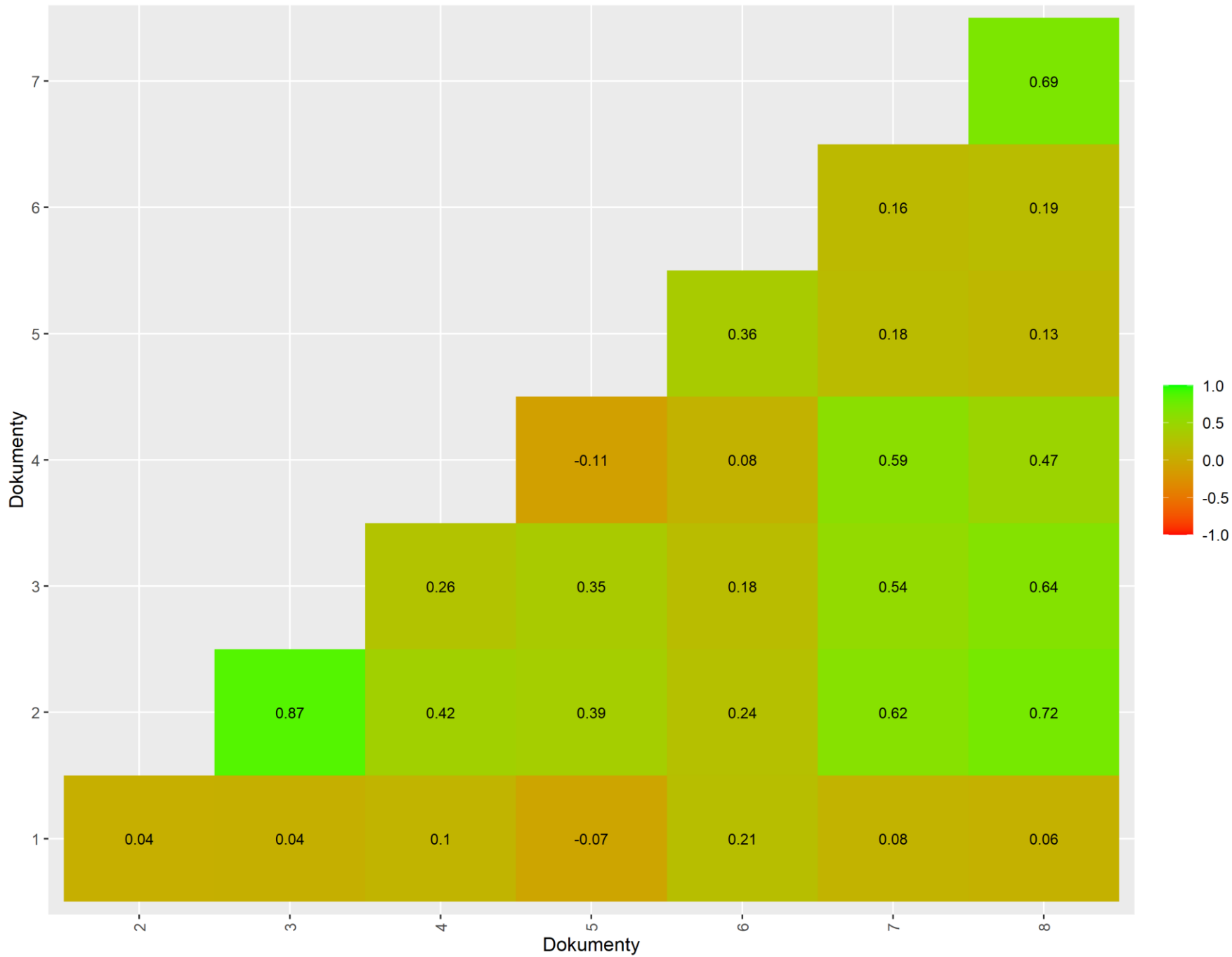
- $CS = \frac{(4*0 + 0*1 + 0*12 + 3*9 + 5*2)}{\sqrt{(16 + 0 + 0 + 9 + 2)} * \sqrt{(0 + 1 + 144 + 81 + 4)}} = 0.345$

- $CD = 1 - 0.345 = 0.655$

# Koherence dokumentů s indikátory EU



Kosinová vzdálenost koncepčních dokumentů



## 2. Similarity measures: summary

- Cosine and Euclidean distances appropriate for higher-dimensional data
- Cosine distance used for text-based data
- Euclidean distances and Jaccard coefficient used for network data
- Pearson coefficient appropriate for continuous data and linear relationships
- There are many more measures of similarity

# Cluster analysis: process

1. Sampling
2. Similarity measure
- 3. Clustering method**
4. Cluster solution interpretation
5. Cluster solution diagnosis

# 3. Clustering methods

- After the (dis)similarities between objects are calculated, we need to select a particular **clustering method** that **partitions** (clusters) **the data** according specific rules
- There are several clustering approaches, **k-means clustering** and **hierarchical clustering** belong to the most common

# 3.1 k-means clustering

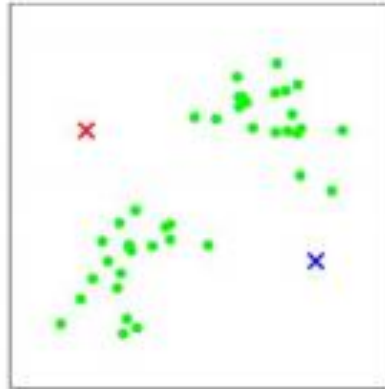
- **k-means clustering** partitions the data into  $k$  clusters based on the mean distances in each cluster
  1. The number of clusters  $k$  needs to be **pre-selected**
  2. The algorithm starts randomly assign cluster centers (**centroids**)
  3. Each object is assigned to the nearest centroid based on a particular similarity measure
  4. Within the clusters, **the centroids are updated** based on the mean similarity of the objects classified to the respective cluster
  5. The steps 3-4 repeat until the centroids no longer change → solution



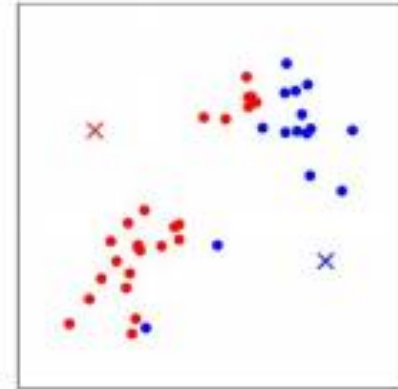
# 2-cluster solution



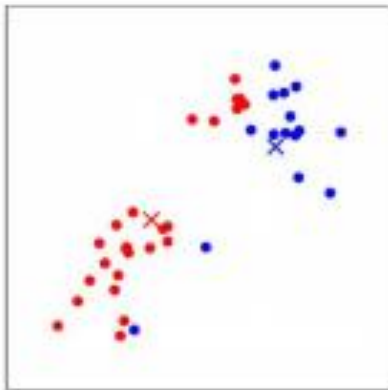
(a)



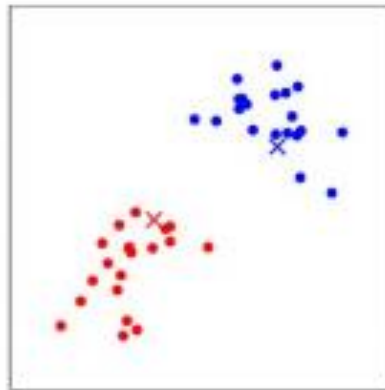
(b)



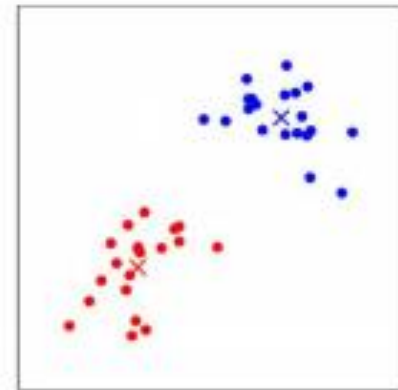
(c)



(d)

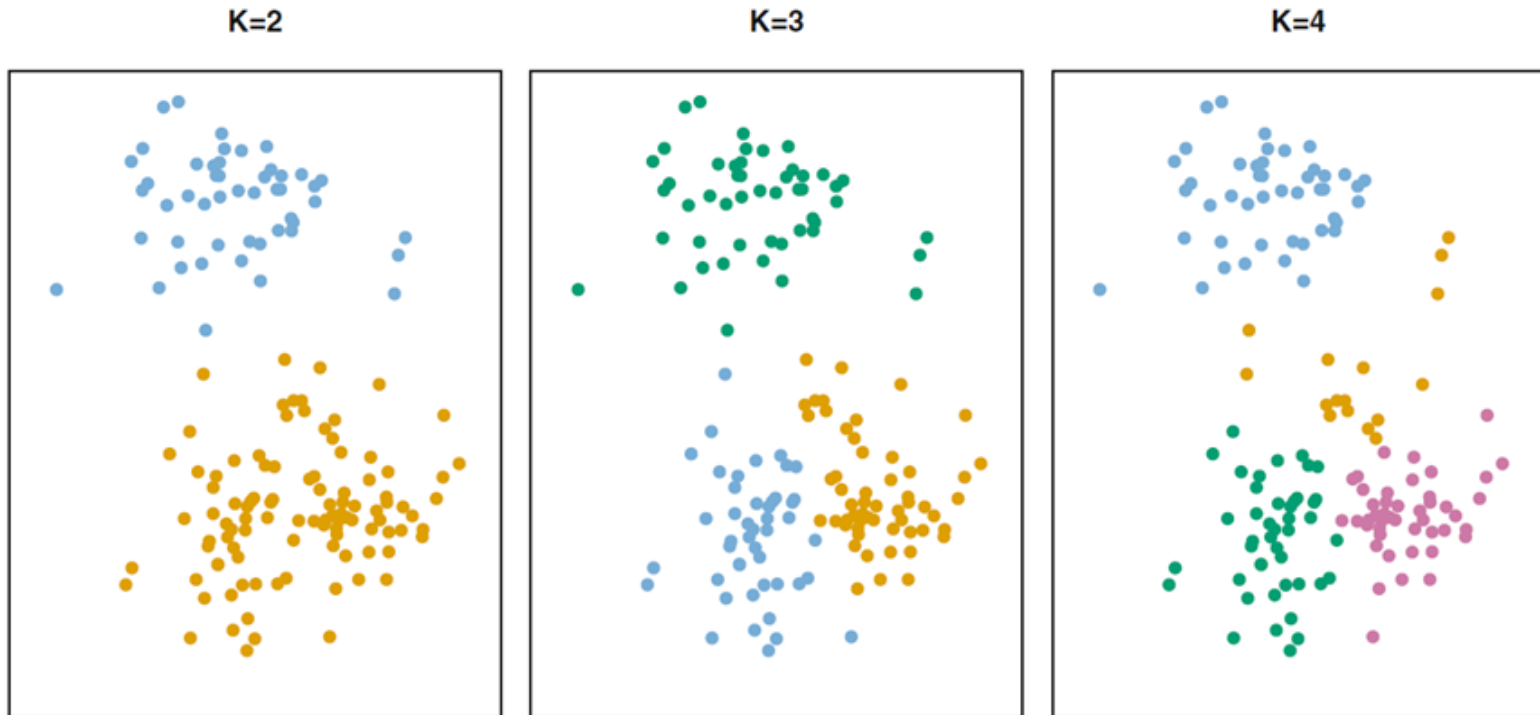


(e)



(f)

# Determining $k$



*A simulated data set with 150 observations in two-dimensional space. Each figure show the results of applying  $K$ -means clustering with different values of  $K$ , the number of clusters. The color of each observation indicates the cluster to which it was assigned using the  $K$ -means clustering algorithm. Note that there is no ordering of the clusters, so the cluster coloring is arbitrary. These cluster labels were not used in clustering; instead, they are the outputs of the clustering procedure.*

# 3.1 Hierarchical clustering

- **Hierarchical clustering** can be performed in **agglomerative** and **divisive** sequence
  - The number of clusters is **not pre-selected**
  - The **linkage method** (how are clusters merged/split) needs to be defined
1. Treat each objects as a separate cluster (**agglomerative** sequence)
  2. The average distances between the clusters are calculated (**average linkage**)
  3. The clusters with the lowest average distance are merged
  4. The steps 2-3 are repeated until there is only a single cluster
  5. The process is represented by **dendrogram**
  6. Considering substantive insights, the  $k$ -cluster solution is identified

- **Single Linkage**

$$D(c_1, c_2) = \min D(x_1, x_2)$$

Minimum distance or distance between closest elements in clusters



- **Complete Linkage**

$$D(c_1, c_2) = \max D(x_1, x_2)$$

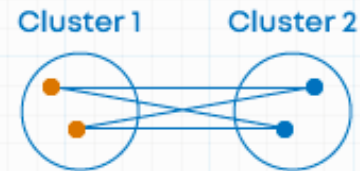
Maximum distance between elements in clusters



- **Average Linkage**

$$D(c_1, c_2) = \frac{1}{|c_1|} \frac{1}{|c_2|} \sum \sum D(x_1, x_2)$$

Average of the distances of all pairs



- **Centroid Method**

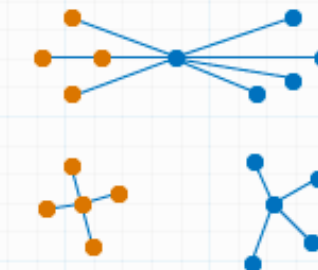
Combining clusters with minimum distance between the centroids of the two clusters

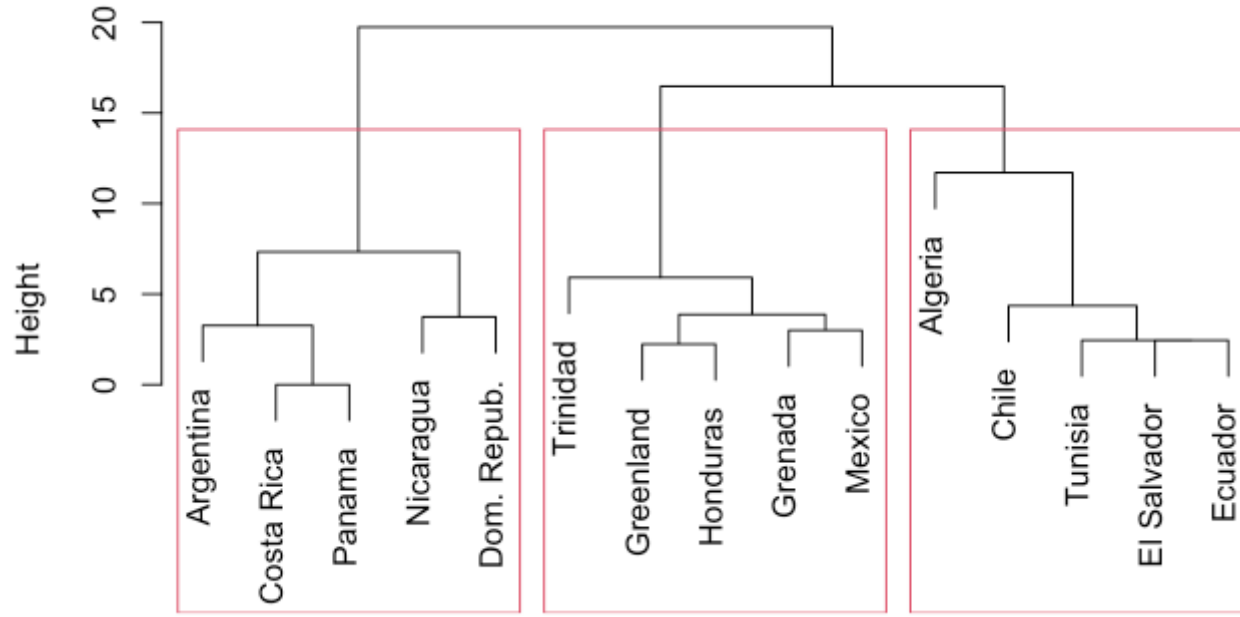


- **Ward's Method**

- Combining clusters where increase in within cluster variance is to the smallest degree.

- Objective is to minimize the total within cluster variance





# Cluster analysis: process

1. Sampling
2. Similarity measure
3. Clustering method
4. **Cluster solution interpretation**
5. **Cluster solution diagnosis**

# 4. Cluster solution interpretation

- Do the clusters have **prima facie** validity (**eyeballing** test)?
- **Substantive** and **theoretical insights** are vital – to what extent the solution aligns with our expectations?
- **Size of the clusters** – do clusters markedly differ in their size?
- **Outliers** – are there any?
- Do clusters **reduce our data well?** → cluster **solution diagnostics**

# 5. Cluster solution diagnostics

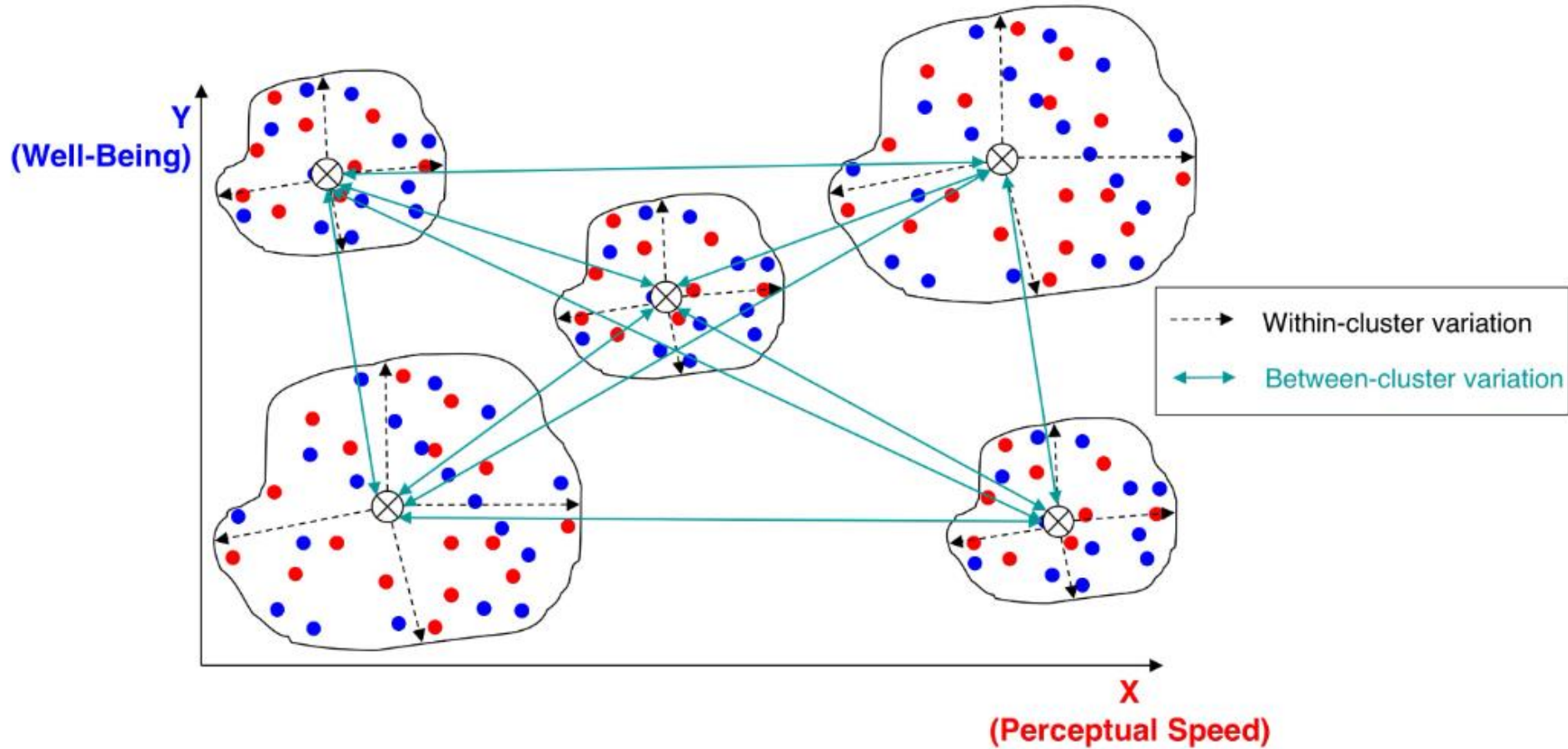
- Important to assess the **quality of our solution**
- Are the clusters internally **cohesive**?
- Are the clusters well **separated**?
- What is the optimal **number of clusters**?



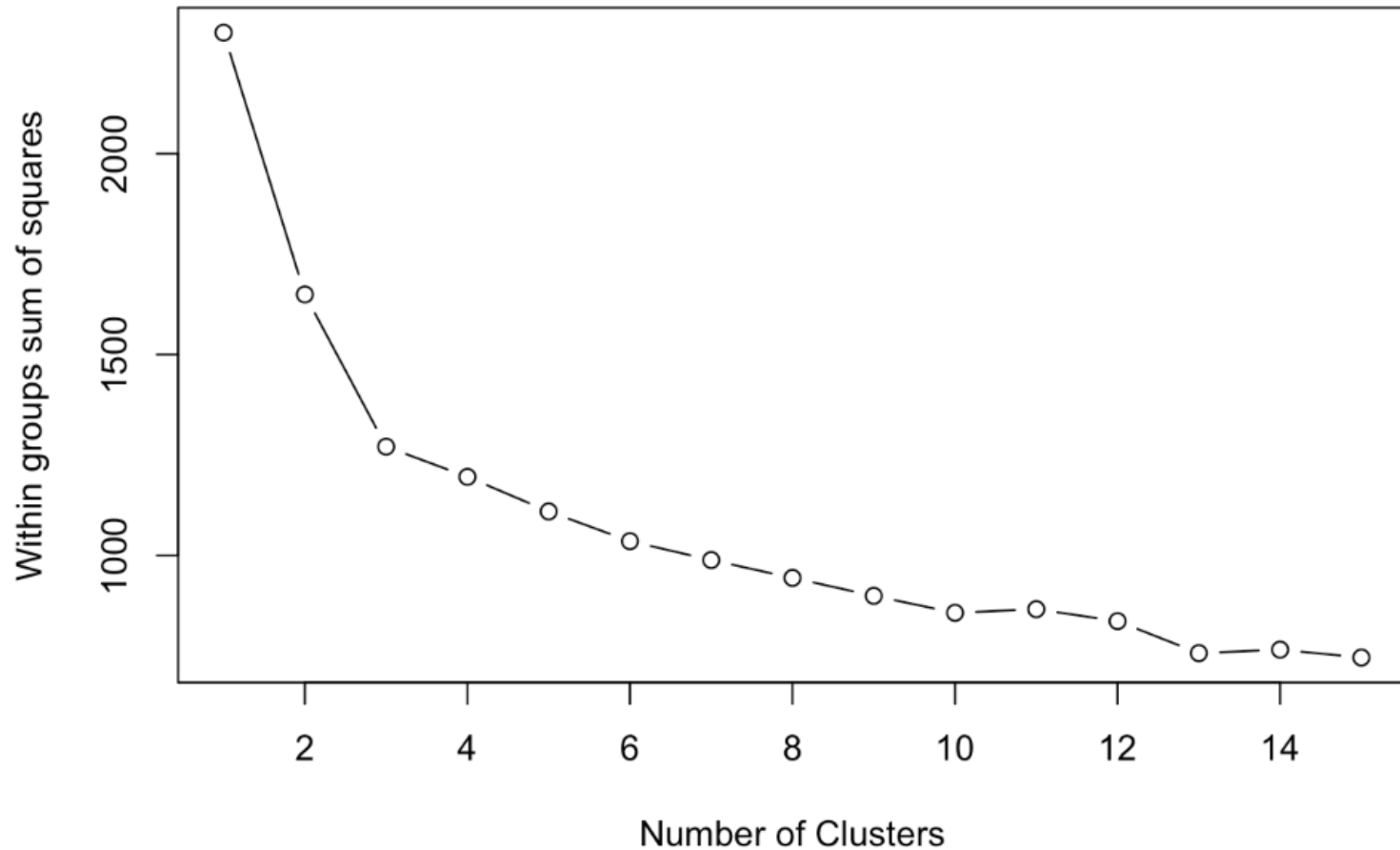
# 5.1 Within/between sum of squares

- The **within-cluster sum of squares (WCSS)** calculate the sum of squared distances between each object and its cluster centroid (k-means clustering approach)
- The **between-cluster sum of squares (BCSS)** measure the sum of squared distance between centroid of each cluster and the overall centroid of all objects
- The **WCSS** capture the **cohesiveness** of clusters, while **BCSS** measure the **separation between clusters**

# 5.1 WCSS and BCSS



# Elbow graph



## 5.2 Silhouette score

- The **Silhouette score** ranges  $\langle -1, 1 \rangle$ ; where high positive values indicate a good fit of the object within its own cluster, zero indicates a borderline position, and negative values indicate that the object is misclassified
- $s = \frac{b - a}{\max(b - a)}$ ;  $a$  = average within-cluster d,  $b$  = average between-cluster d
- The Silhouette is calculated for each object and then average is taken to evaluate the cluster solution
- The Silhouette score can be applied both to the **k-means** and **hierarchical clustering**
- Rule of thumb:  $s > 0.5$  good solution;  $s < 0.25$  bad solution

# Exercise

# Exercise

- Open Jamovi and install **snowCluster** extension
- Load into Jamovi file “state.csv”
- Check variables relig\_prot, urban, and relig\_high in the codebook
- Describe the variables
- Apply **(1) k-means** clustering and **(2) hierarchical** clustering
- Interpret the results