

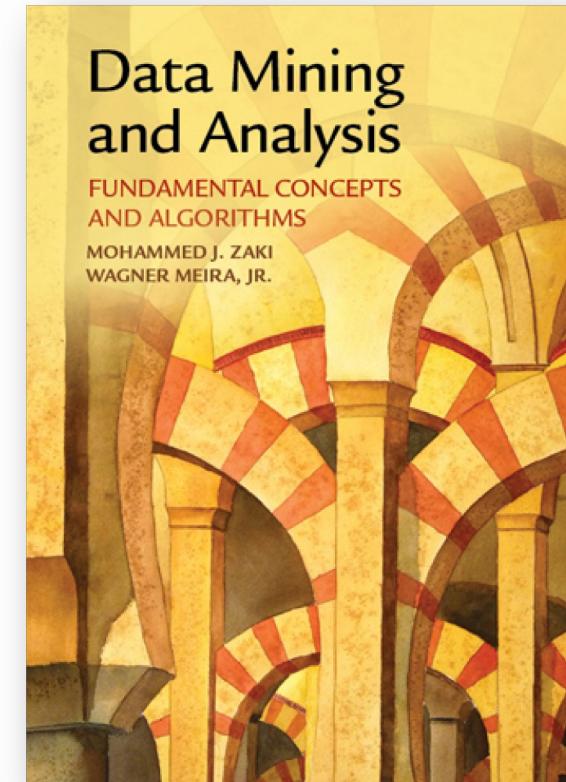


Hierarchical Clustering

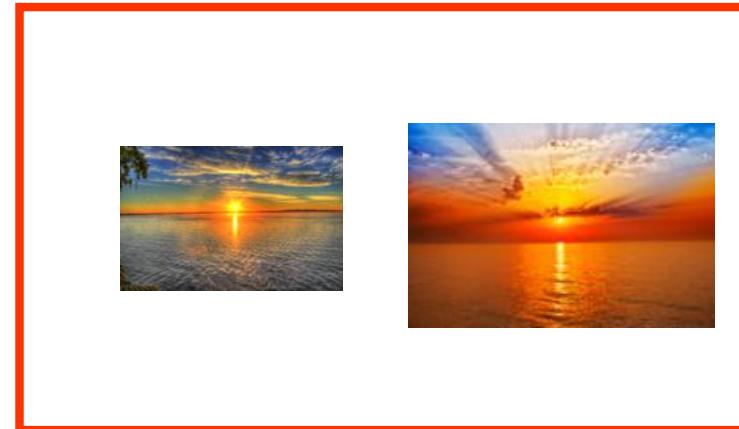
Data Science for Mobility

Readings

- “Data Mining and Analysis” by Zaki & Meira
 - Chapter 14
- <http://www.dataminingbook.info>







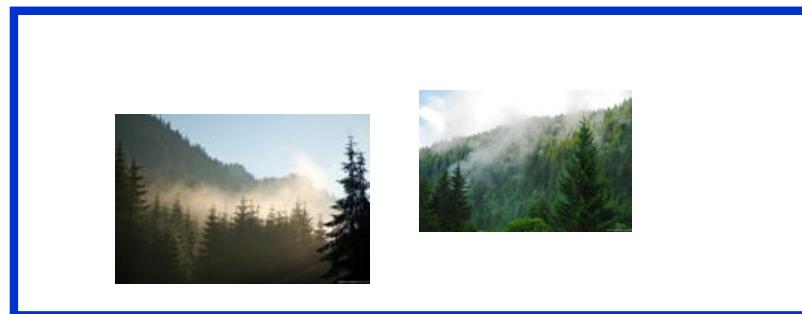
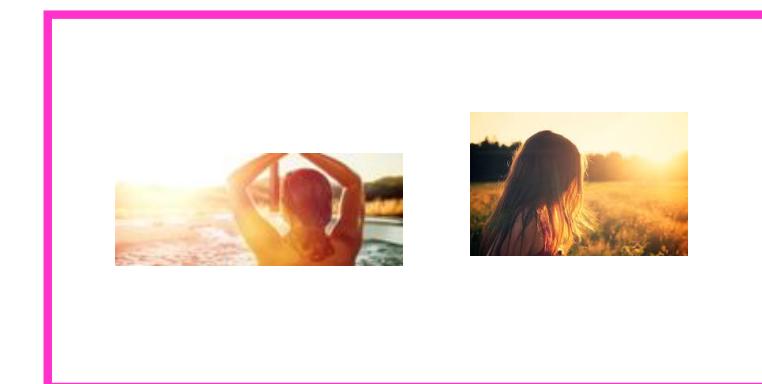
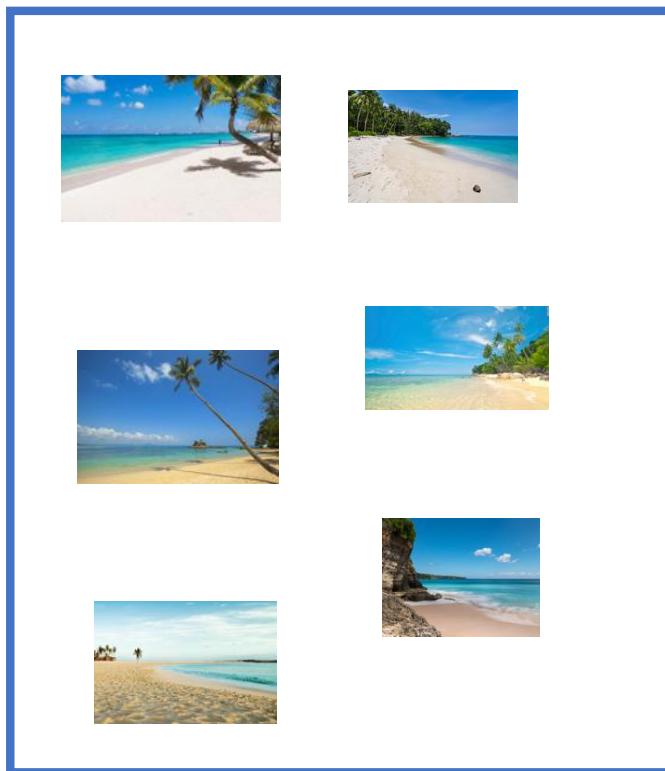
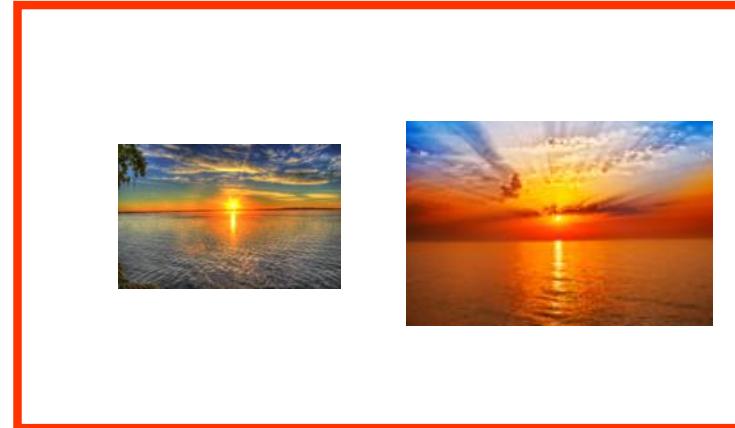
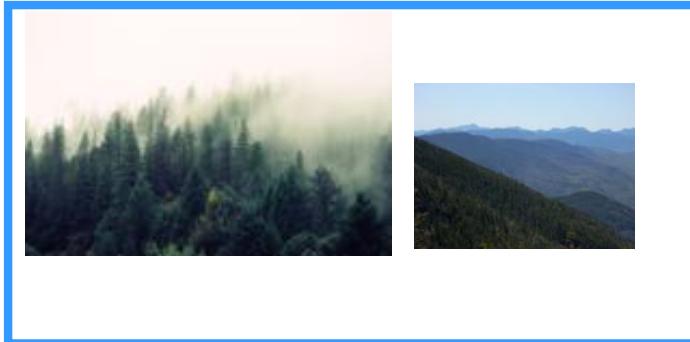








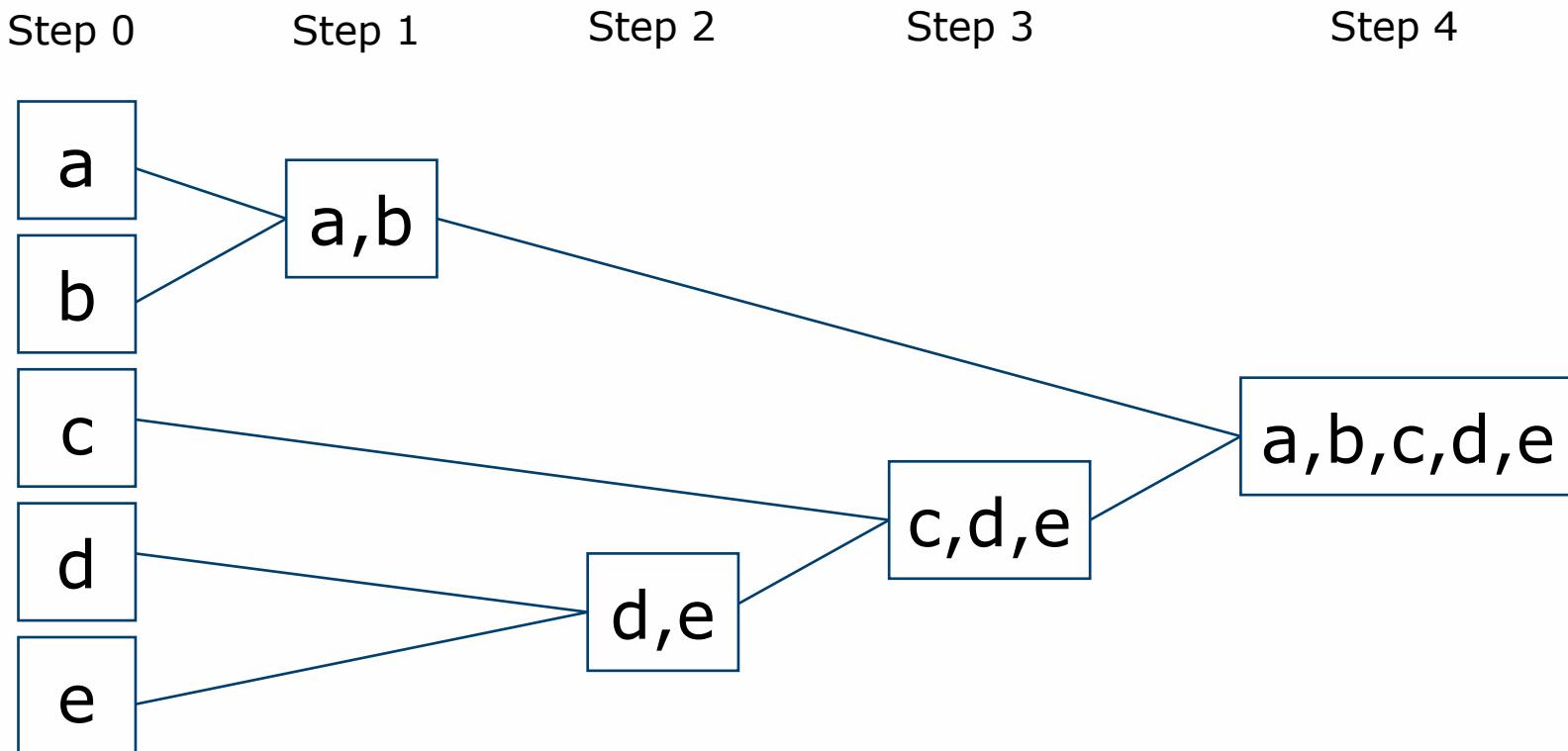




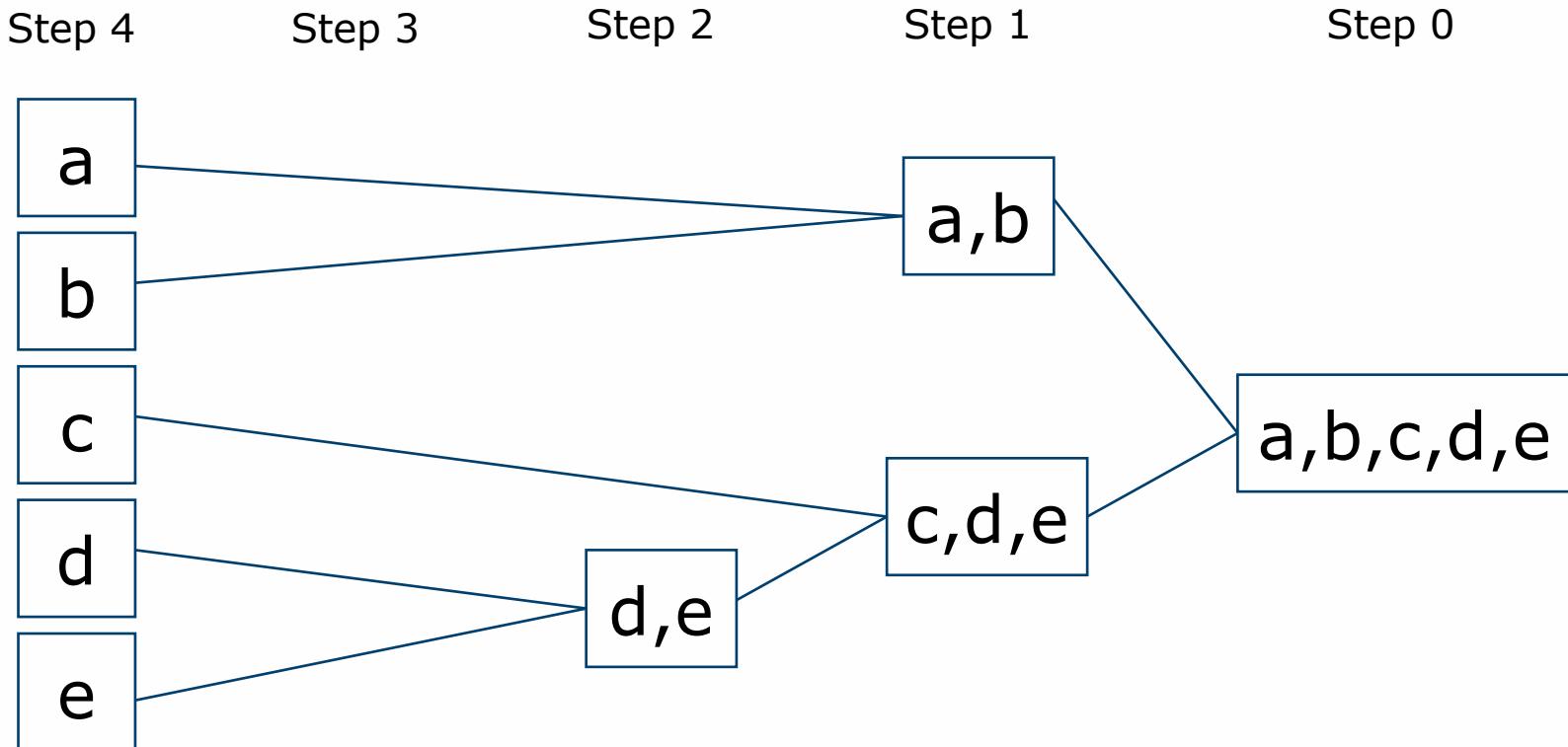
What is Hierarchical Clustering?

11

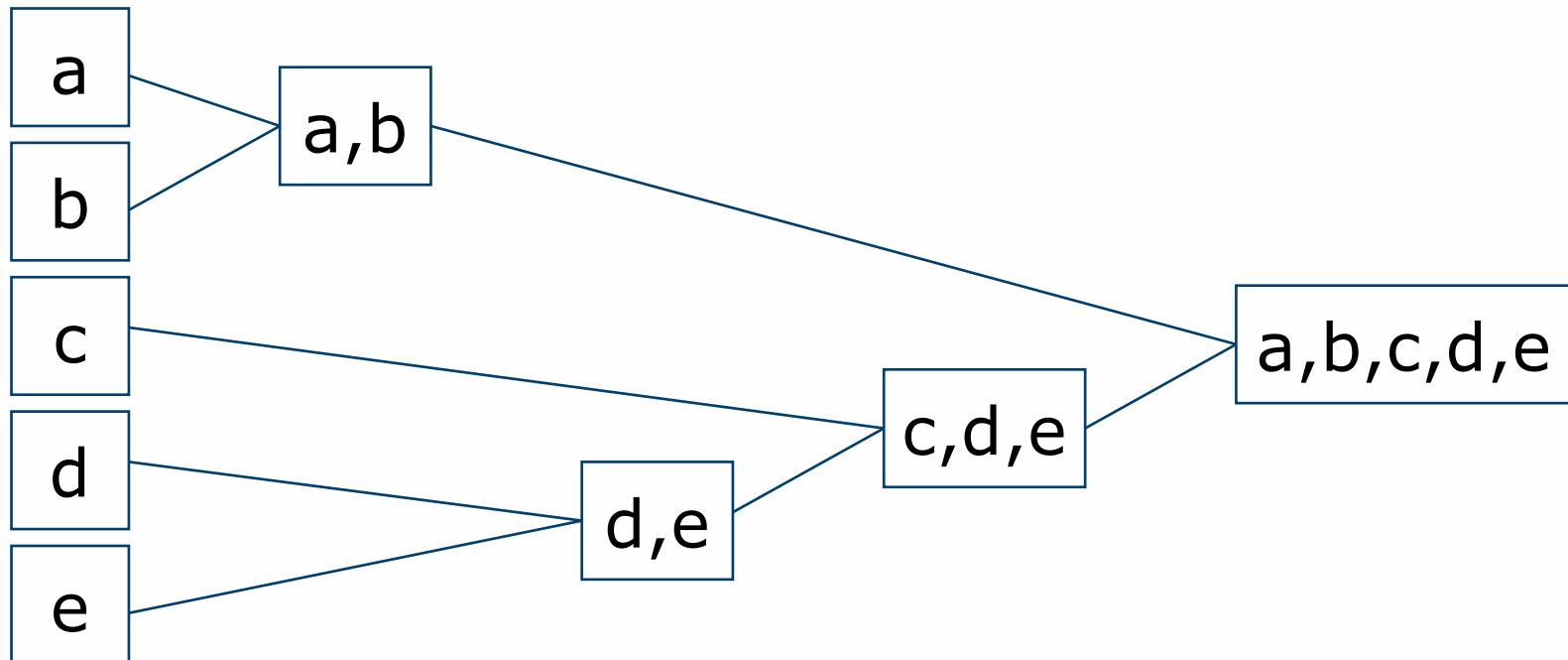
- Suppose we have five items, a, b, c, d, and e.
- Initially, we consider one cluster for each item
- Then, at each step we merge together the most similar clusters, until we generate one cluster



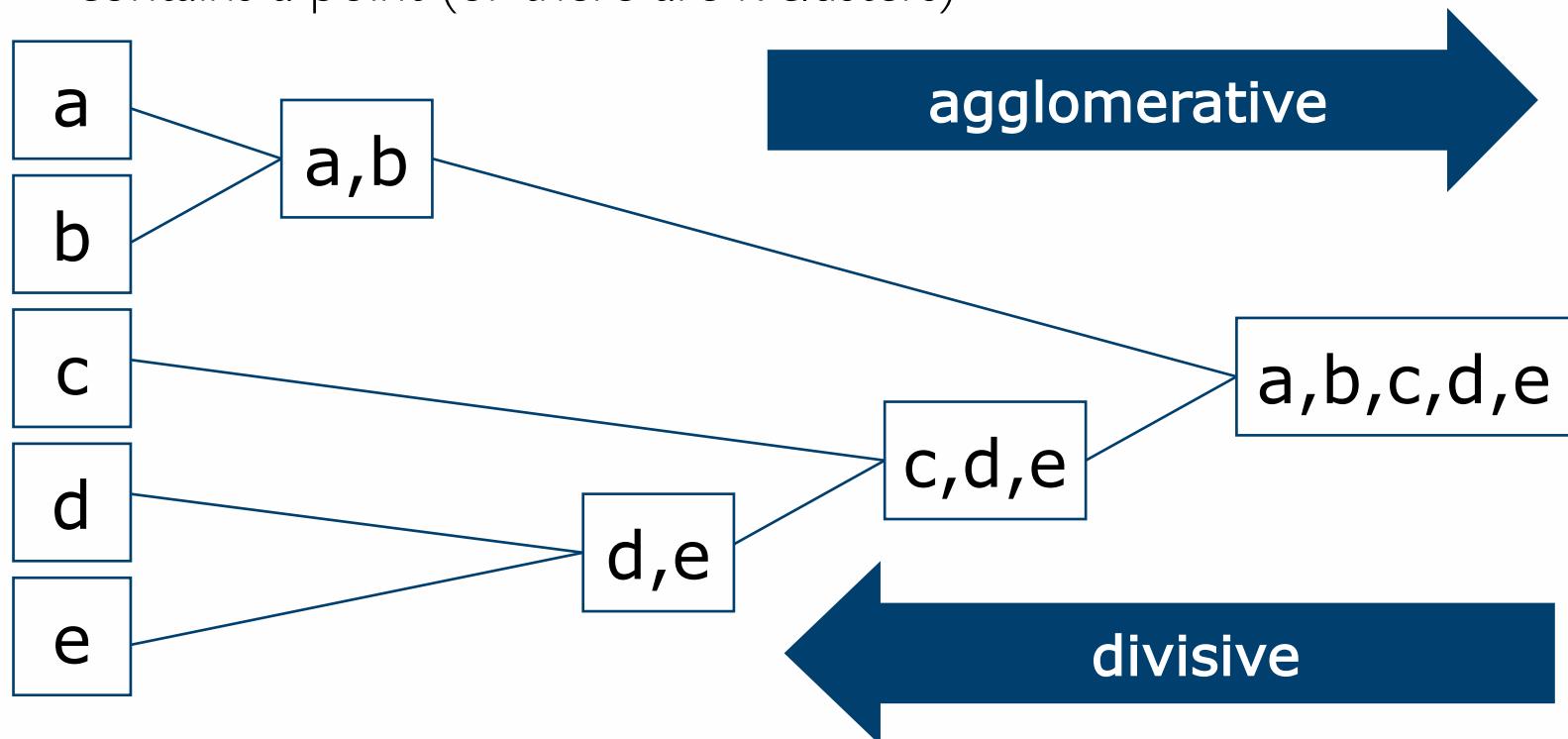
- Alternatively, we start from one cluster containing the five elements
- Then, at each step we split one cluster to improve introcluster similarity, until all the elements are contained in one cluster



- By far, it is the most common clustering technique
- Produces a hierarchy of nested clusters
- The hierarchy can be visualized as a dendrogram: a tree-like diagram that records the sequences of merges or splits



- Agglomerative
 - Start individual clusters, at each step, merge the closest pair of clusters until only one cluster (or k clusters) left
- Divisive
 - Start with one cluster, at each step, split a cluster until each cluster contains a point (or there are k clusters)



- No need to assume any particular number of clusters
- Any desired number of clusters can be obtained by 'cutting' the dendrogram at the proper level
- They may correspond to meaningful taxonomies
- Example in biological sciences include animal kingdom, phylogeny reconstruction, etc.
- Traditional hierarchical algorithms use a similarity or distance matrix to merge or split one cluster at a time

- More popular hierarchical clustering technique
- Compute the proximity matrix
- Let each data point be a cluster
- Repeat
 - Merge the two closest clusters
 - Update the proximity matrix
- Until only a single cluster remains
- Key operation is the computation of the proximity of two clusters
- Different approaches to defining the distance between clusters distinguish the different algorithms

Hierarchical Clustering: Time and Space Requirements

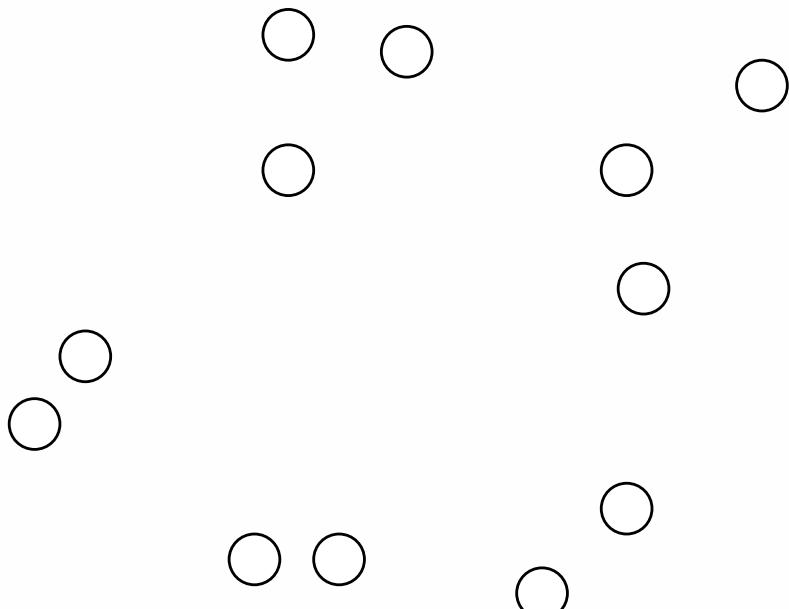
17

- $O(N^2)$ space since it uses the proximity matrix.
 - N is the number of points.
- $O(N^3)$ time in many cases
 - There are N steps and at each step the size, N^2 , proximity matrix must be updated and searched
 - Complexity can be reduced to $O(N^2 \log(N))$ time for some approaches

- Compute the distance between all pairs of points [$O(N^2)$]
- Insert the pairs and their distances into a priority queue to find the min in one step [$O(N^2)$]
- When two clusters are merged, we remove all entries in the priority queue involving one of these two clusters [$O(N \log N)$]
- Compute all the distances between the new cluster and the remaining clusters [$O(N \log N)$]
- Since the last two steps are executed at most N time, the complexity of the whole algorithms is $O(N^2 \log N)$

Distance Between Clusters

- Start with clusters of individual points and the distance matrix

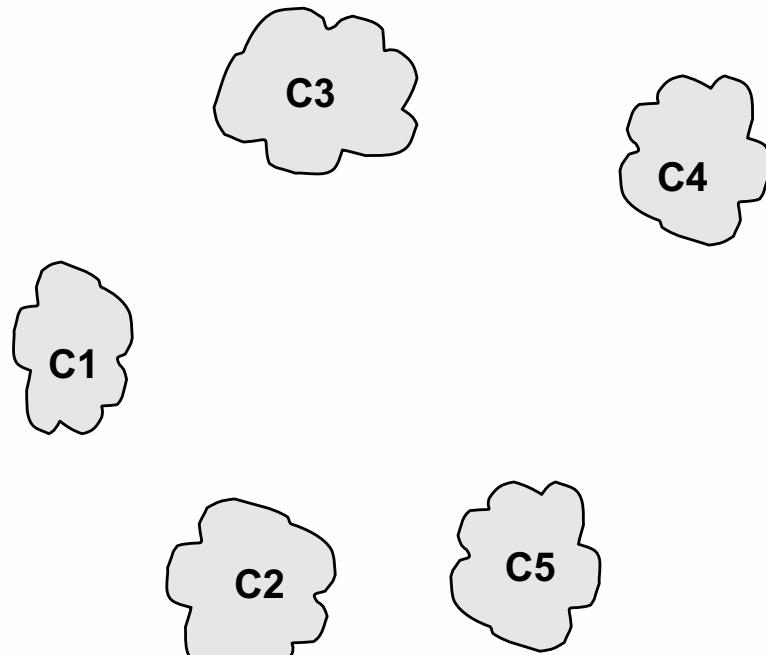


	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

Distance Matrix

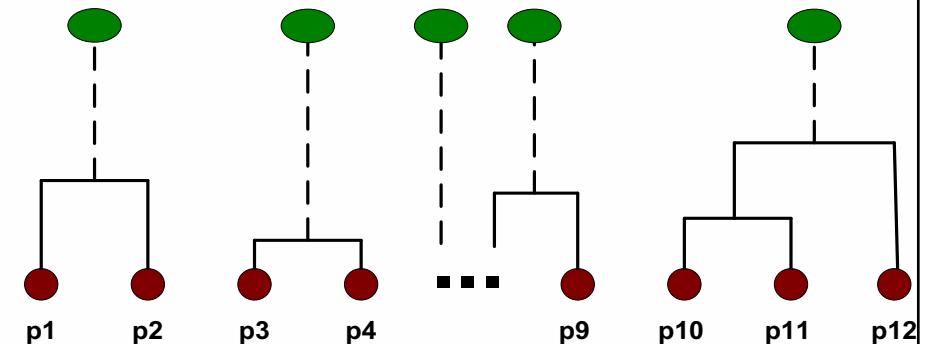
$p_1 \quad p_2 \quad p_3 \quad p_4 \quad \dots \quad p_9 \quad p_{10} \quad p_{11} \quad p_{12}$

- After some merging steps, we have some clusters

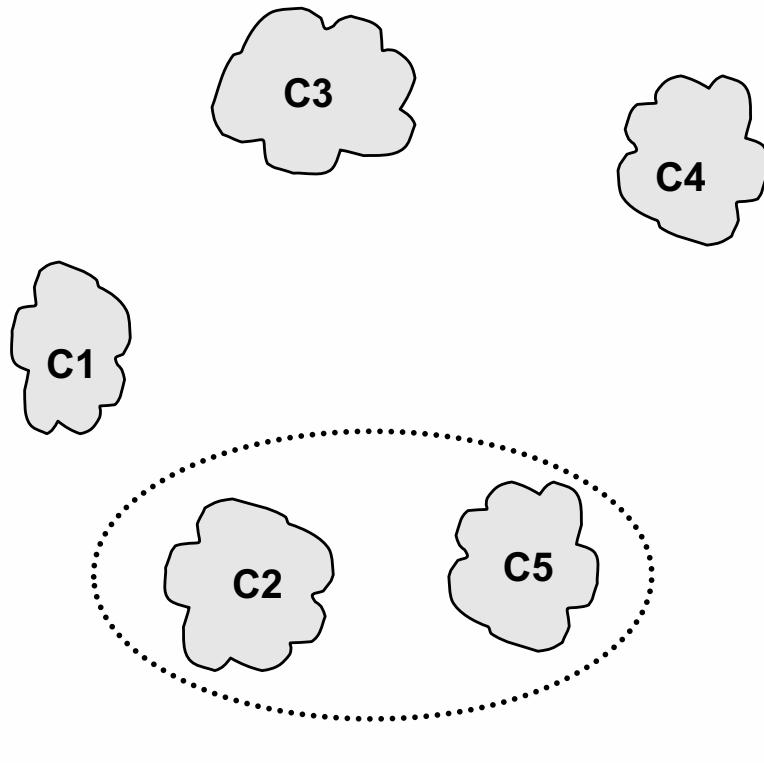


	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Distance Matrix

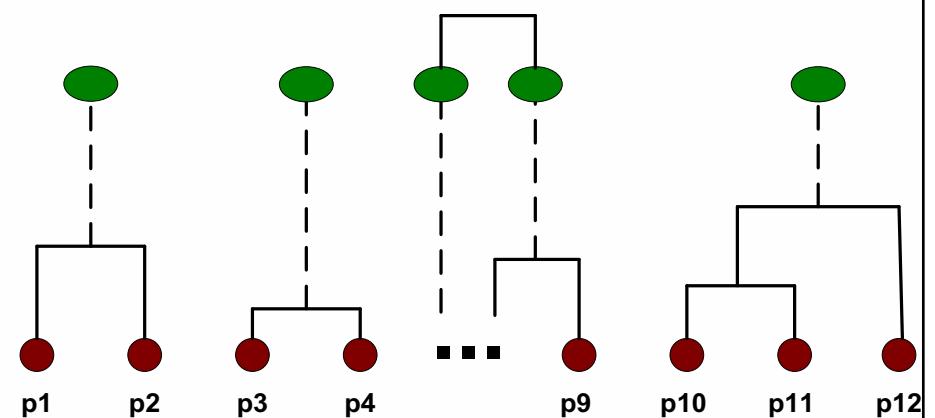


- We want to merge the two closest clusters (C_2 and C_5) and update the proximity matrix.

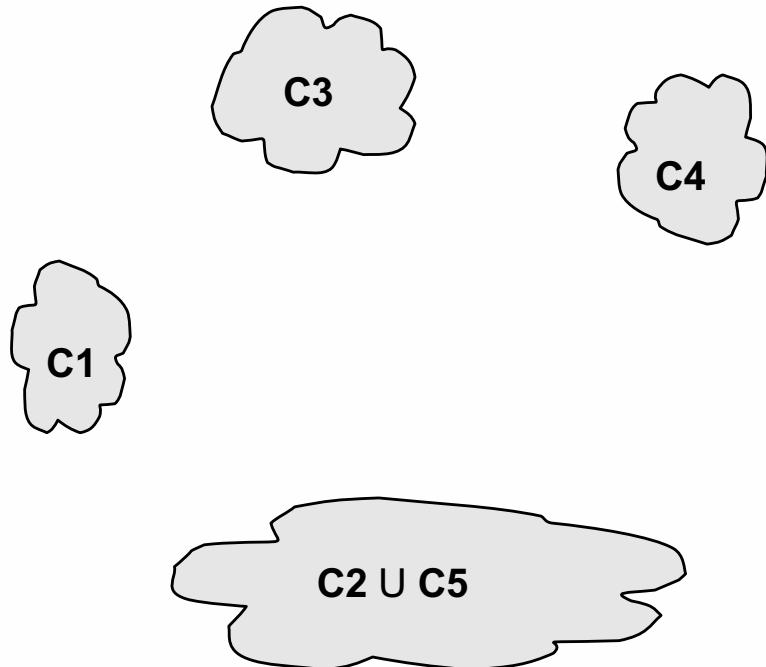


	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Distance Matrix

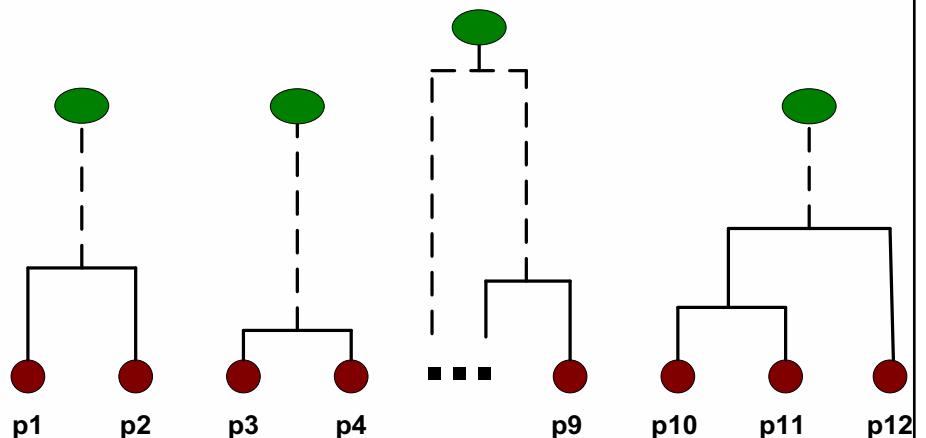


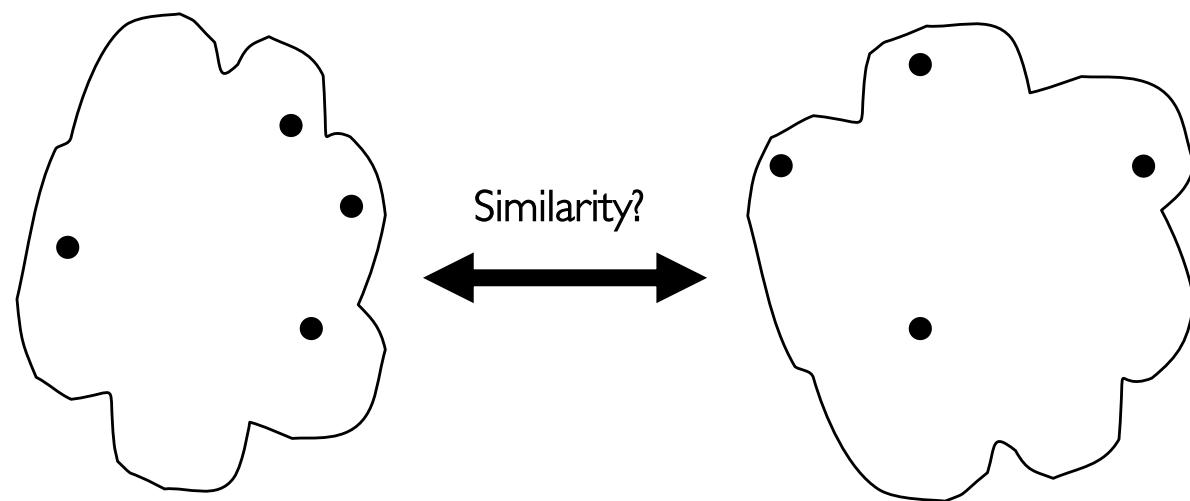
- The question is “How do we update the proximity matrix?”

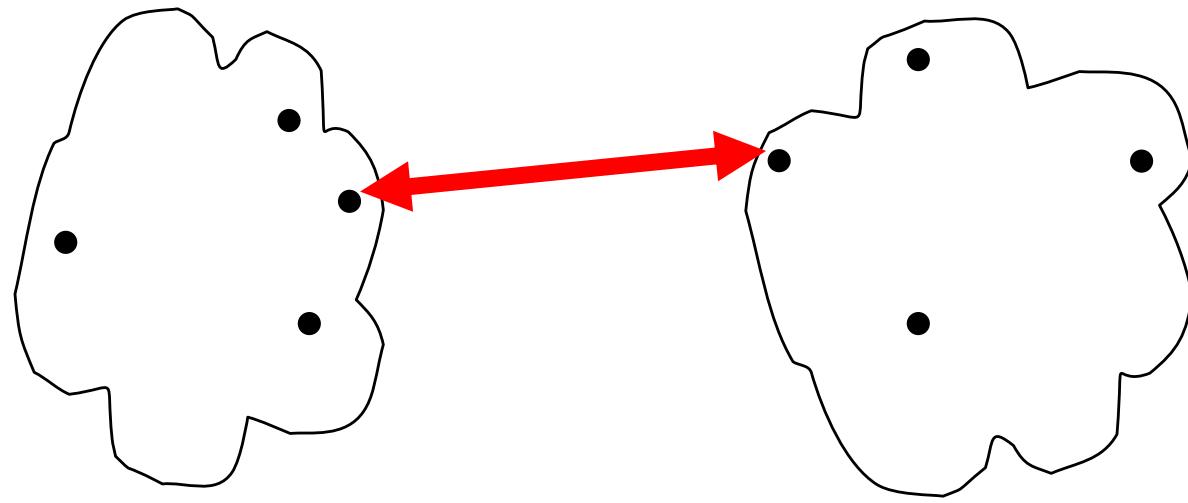


	C1	C2 U C5	C3	C4
C1		?		
C2 U C5	?		?	?
C3		?		
C4		?		

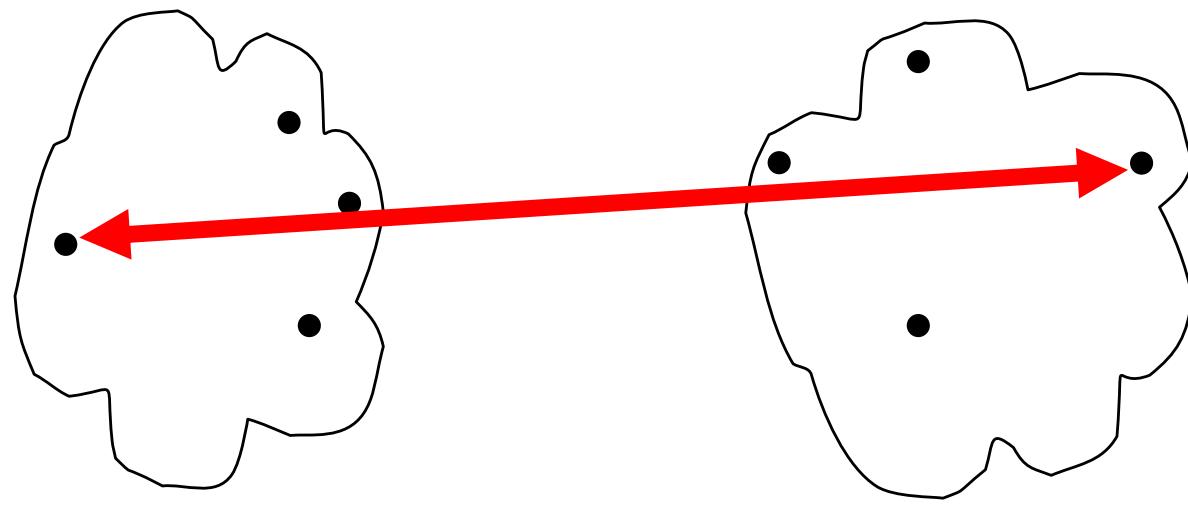
Distance Matrix



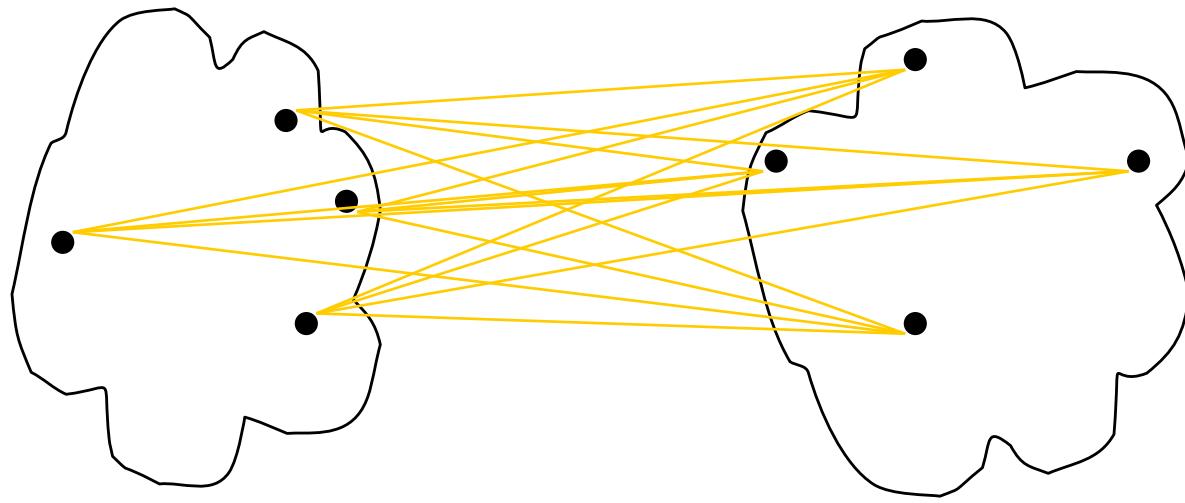




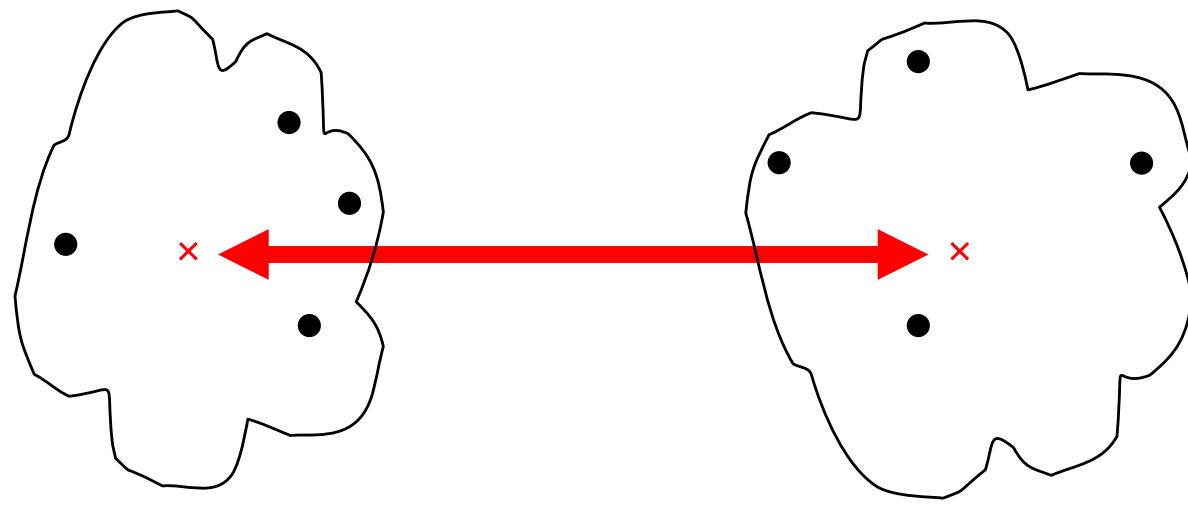
Single Linkage or MIN



Complete Linkage or MAX



Average or Group Average



Distance between Centroids

- Single link (or MIN)
 - smallest distance between an element in one cluster and an element in the other, i.e., $d(C_i, C_j) = \min(t_{i,p}, t_{j,q})$
- Complete link (or MAX)
 - largest distance between an element in one cluster and an element in the other, i.e., $d(C_i, C_j) = \max(t_{i,p}, t_{j,q})$
- Average (or group average)
 - average distance between an element in one cluster and an element in the other, i.e., $d(C_i, C_j) = \text{avg}(d(t_{i,p}, t_{j,q}))$
- Centroid
 - distance between the centroids of two clusters, i.e.,
 $d(C_i, C_j) = d(\mu_i, \mu_j)$ where μ_i and μ_j are the centroids
- ...

Example

- Suppose we have five items, a, b, c, d, and e.
- We want to perform hierarchical clustering on five instances following an agglomerative approach
- First: we compute the distance or similarity matrix
- D_{ij} is the distance between instance “i” and “j”

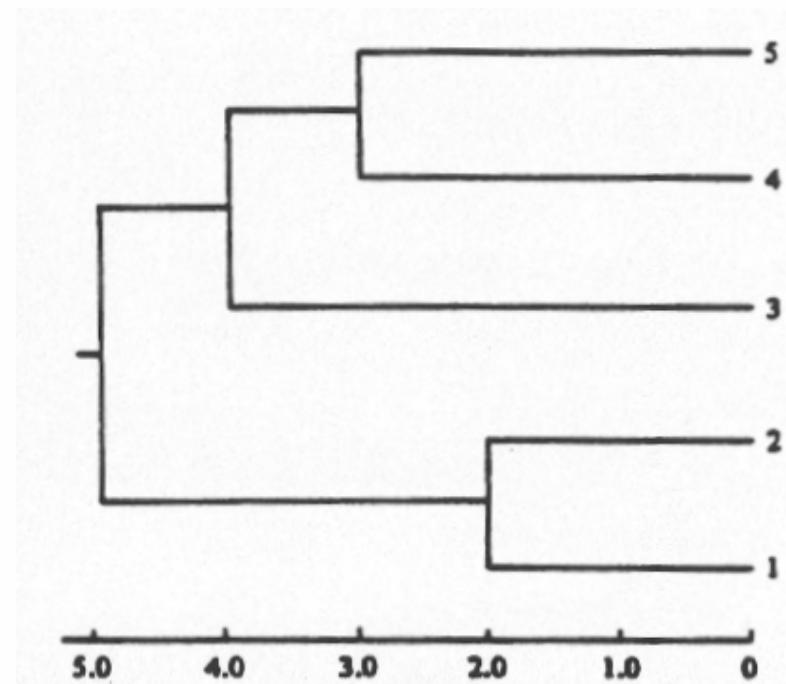
$$D = \begin{pmatrix} 0.0 & & & & \\ 2.0 & 0.0 & & & \\ 6.0 & 5.0 & 0.0 & & \\ 10.0 & 9.0 & 4.0 & 0.0 & \\ 9.0 & 8.0 & 5.0 & 3.0 & 0.0 \end{pmatrix}$$

- Group the two instances that are closer
- In this case, a and b are the closest items ($D_{2,1}=2$)
- Compute again the distance matrix, and start again.
- Suppose we apply single-linkage (MIN), we need to compute the distance between the new cluster $\{1,2\}$ and the others
 - $d(12)3 = \min[d_{13}, d_{23}] = d_{23} = 5.0$
 - $d(12)4 = \min[d_{14}, d_{24}] = d_{24} = 9.0$
 - $d(12)5 = \min[d_{15}, d_{25}] = d_{25} = 8.0$

Example

- The new distance matrix is,
- At the end, we obtain the following dendrogram

$$D = \begin{pmatrix} 0.0 \\ 5.0 & 0.0 \\ 9.0 & 4.0 & 0.0 \\ 8.0 & 5.0 & 3.0 & 0.0 \end{pmatrix}$$



Determining the Number of Clusters

hierarchical clustering generates
a set of N possible partitions

which one should I choose?

From the previous lecture we know ideally
a good cluster should partition points so that ...

Data points in the same cluster should have
a small distance from one another

Data points in different clusters should be at
a large distance from one another.

- Within-cluster sum of squares

$$\text{WSS}(C) = \sum_{i=1}^k \sum_{x_j \in C_i} d(x_j, \mu_i)^2$$

where μ_i is the centroid of cluster C_i (in case of Euclidean spaces)

- Between-cluster sum of squares

$$\text{BSS}(C) = \sum_{i=1}^k |C_i| d(\mu, \mu_i)^2$$

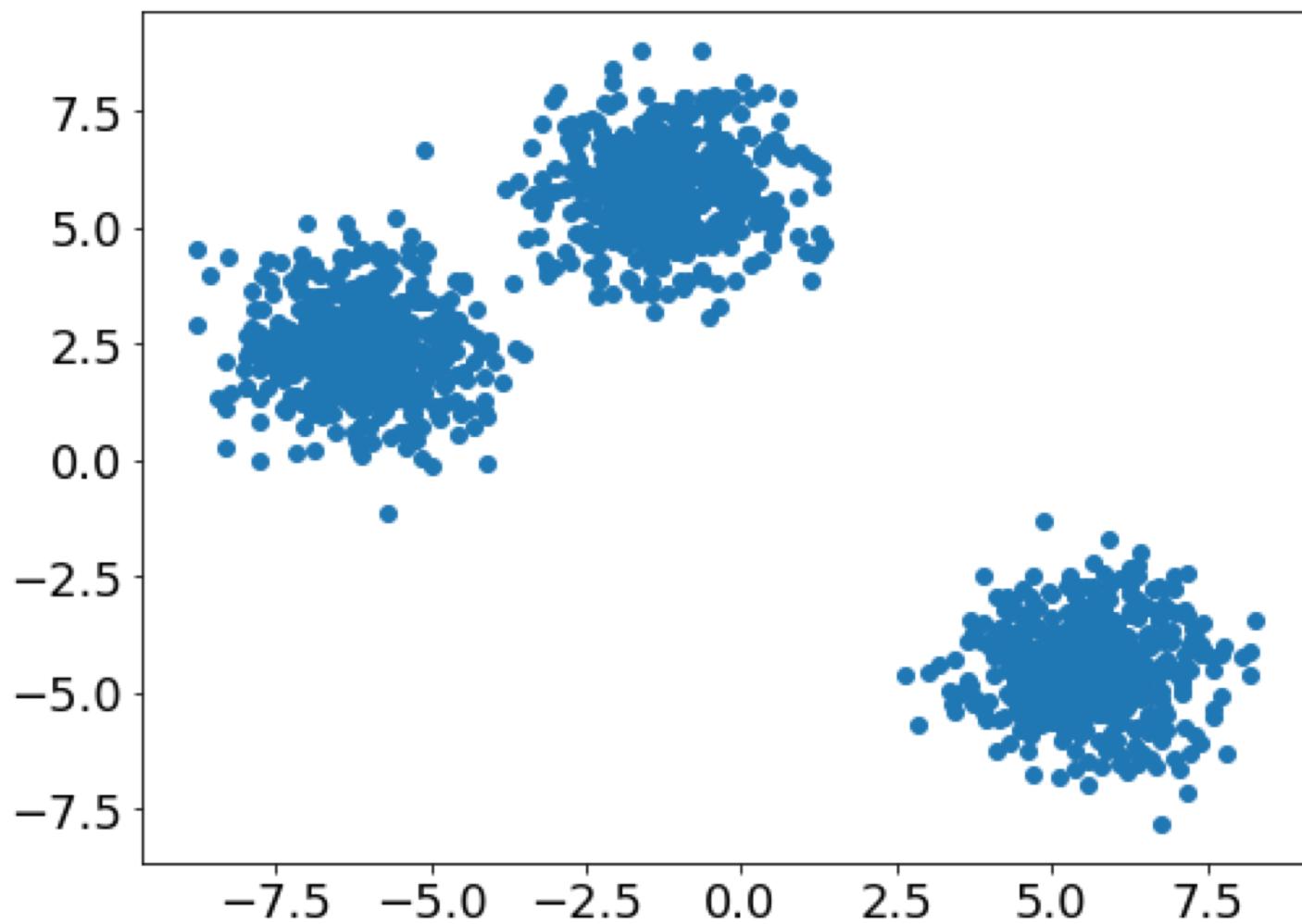
where μ is the centroid of the whole dataset

Evaluation of Hierarchical Clustering using Knee/Elbow Analysis

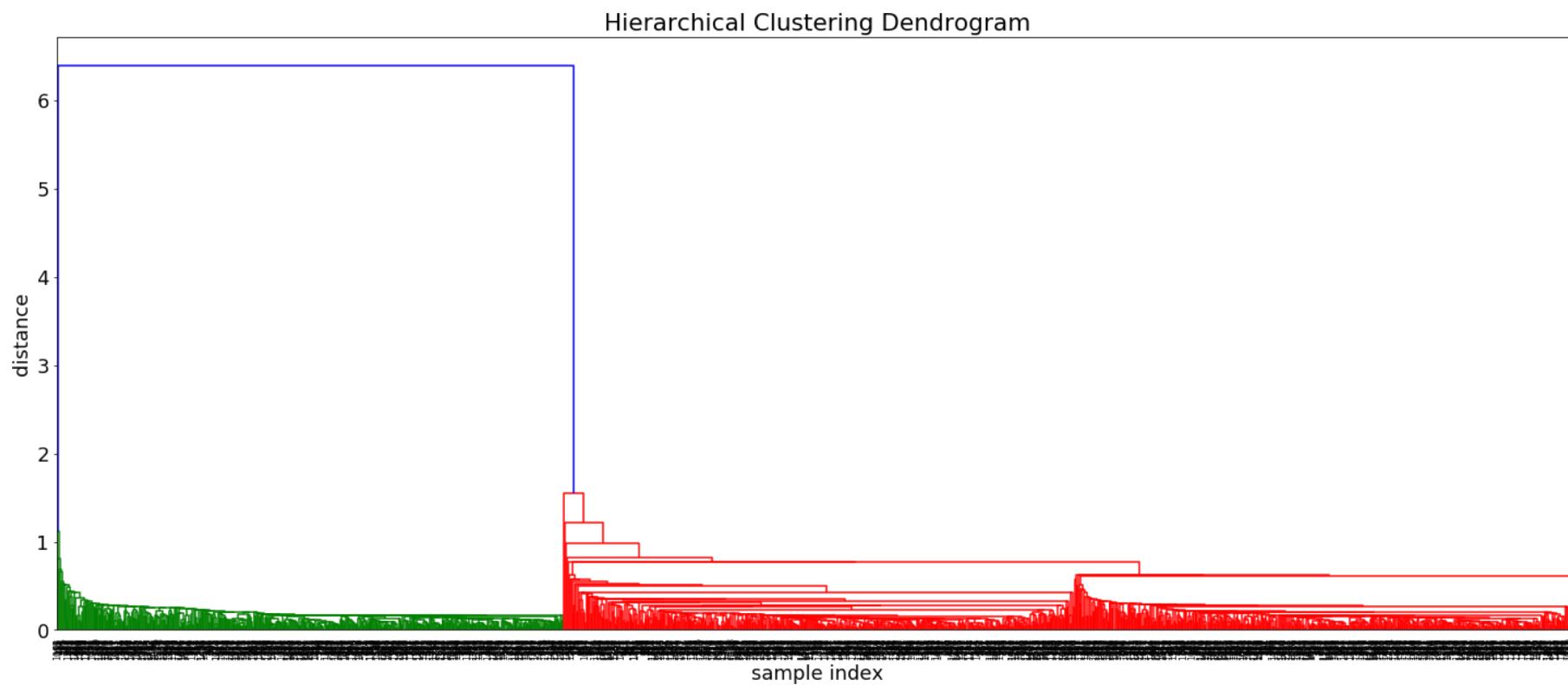
plot the WSS and BSS for every clustering and look
for a knee in the plot that show a significant
modification in the evaluation metrics

Run the Python notebook
for hierarchical clustering

Data Points

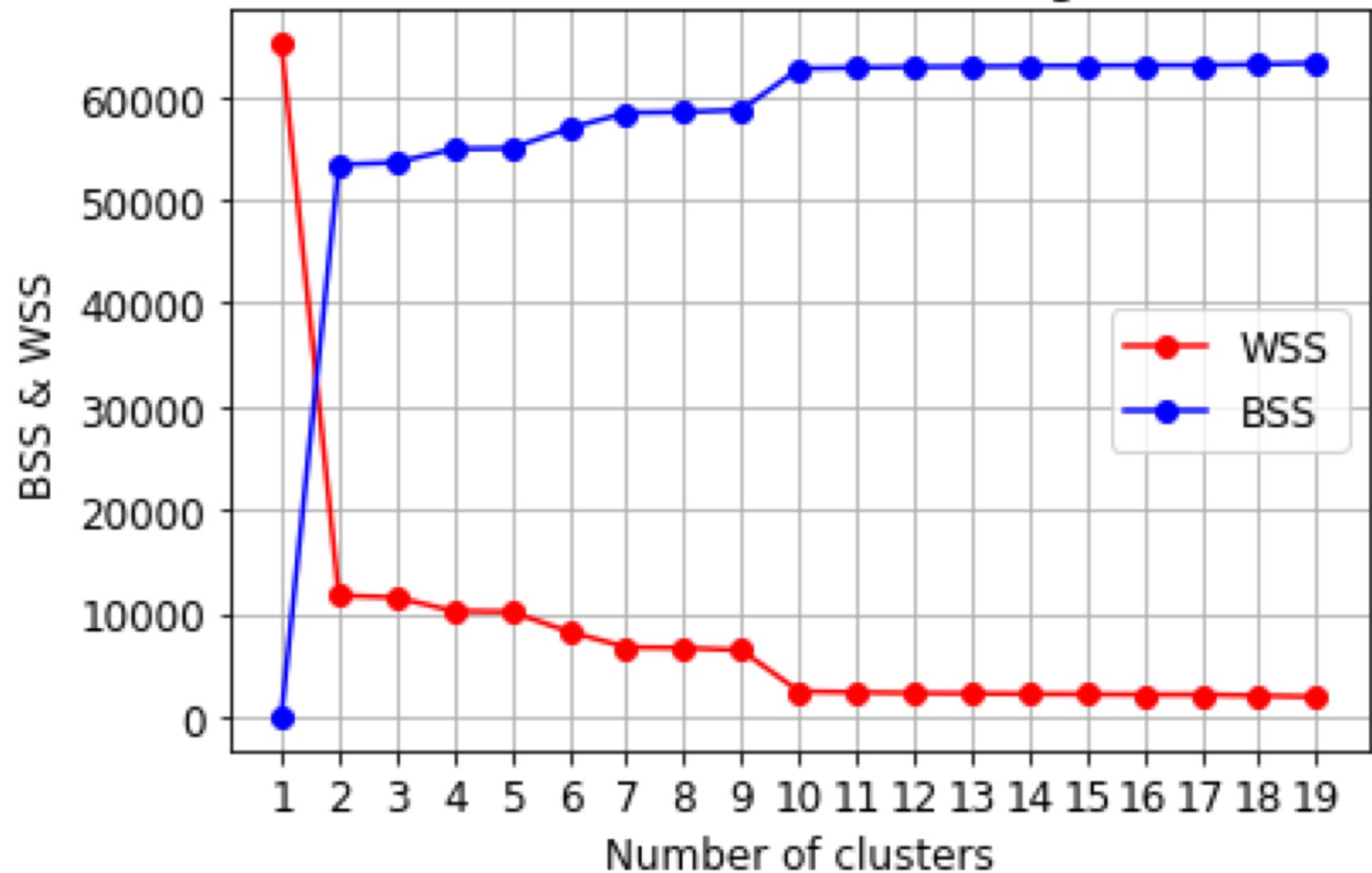


Example data generated using the `make_blob` function of Scikit-Learn



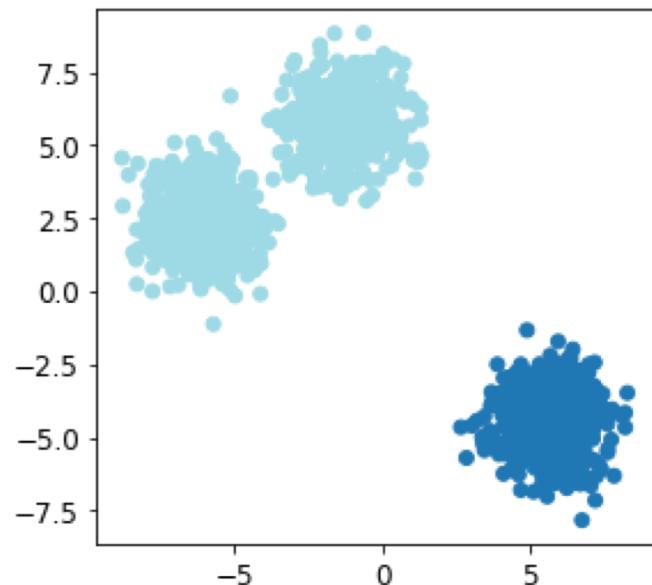
Dendrogram computed using single linkage.

Hierarchical Clustering

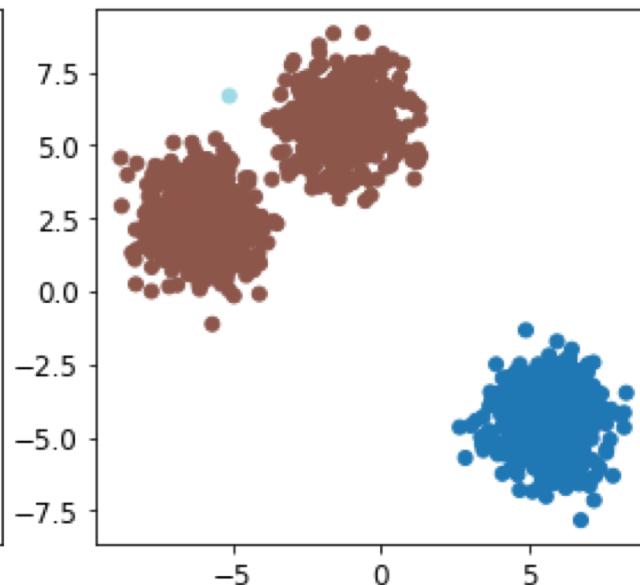


BSS and WSS for values of k from 1 until 19.

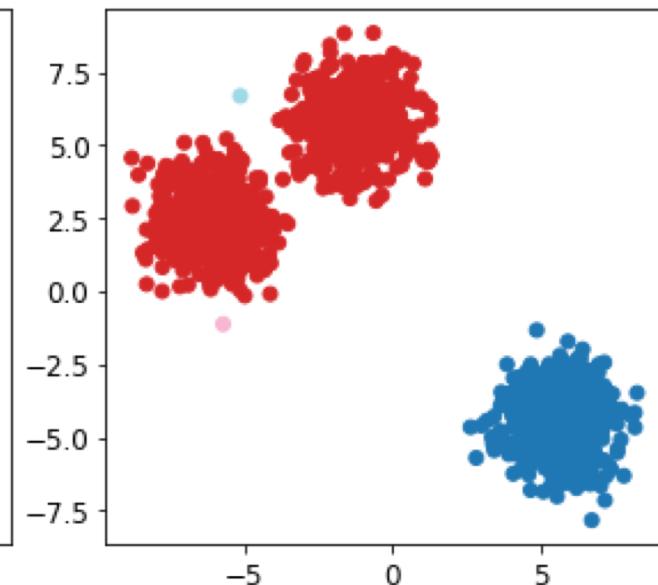
$k=2$



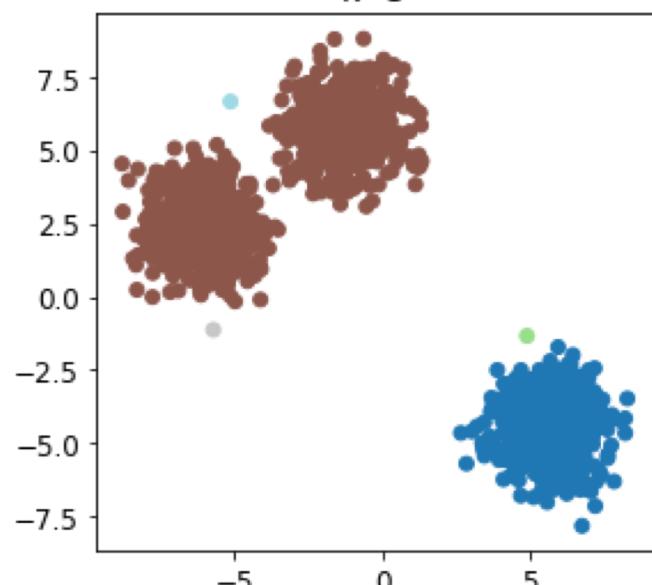
$k=3$



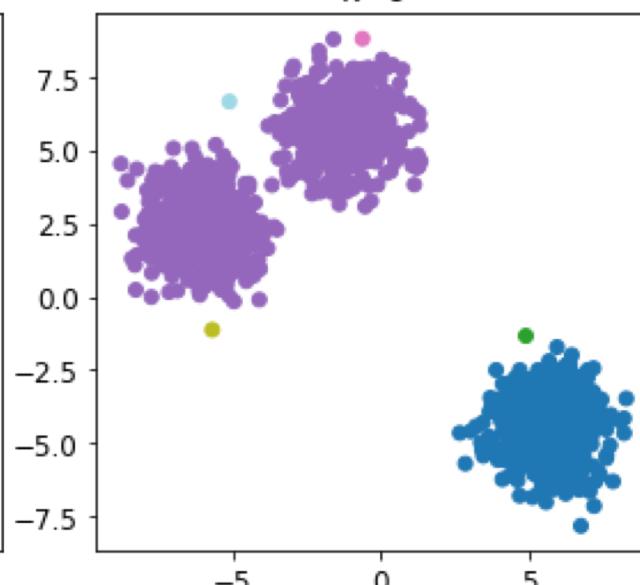
$k=4$



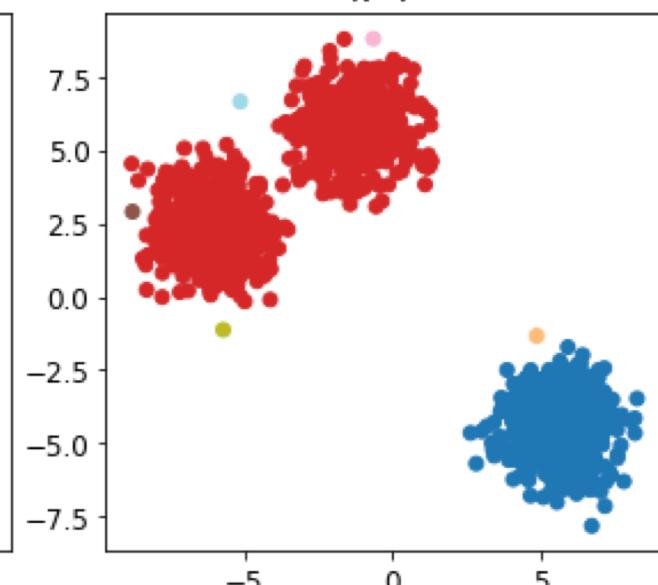
$k=5$



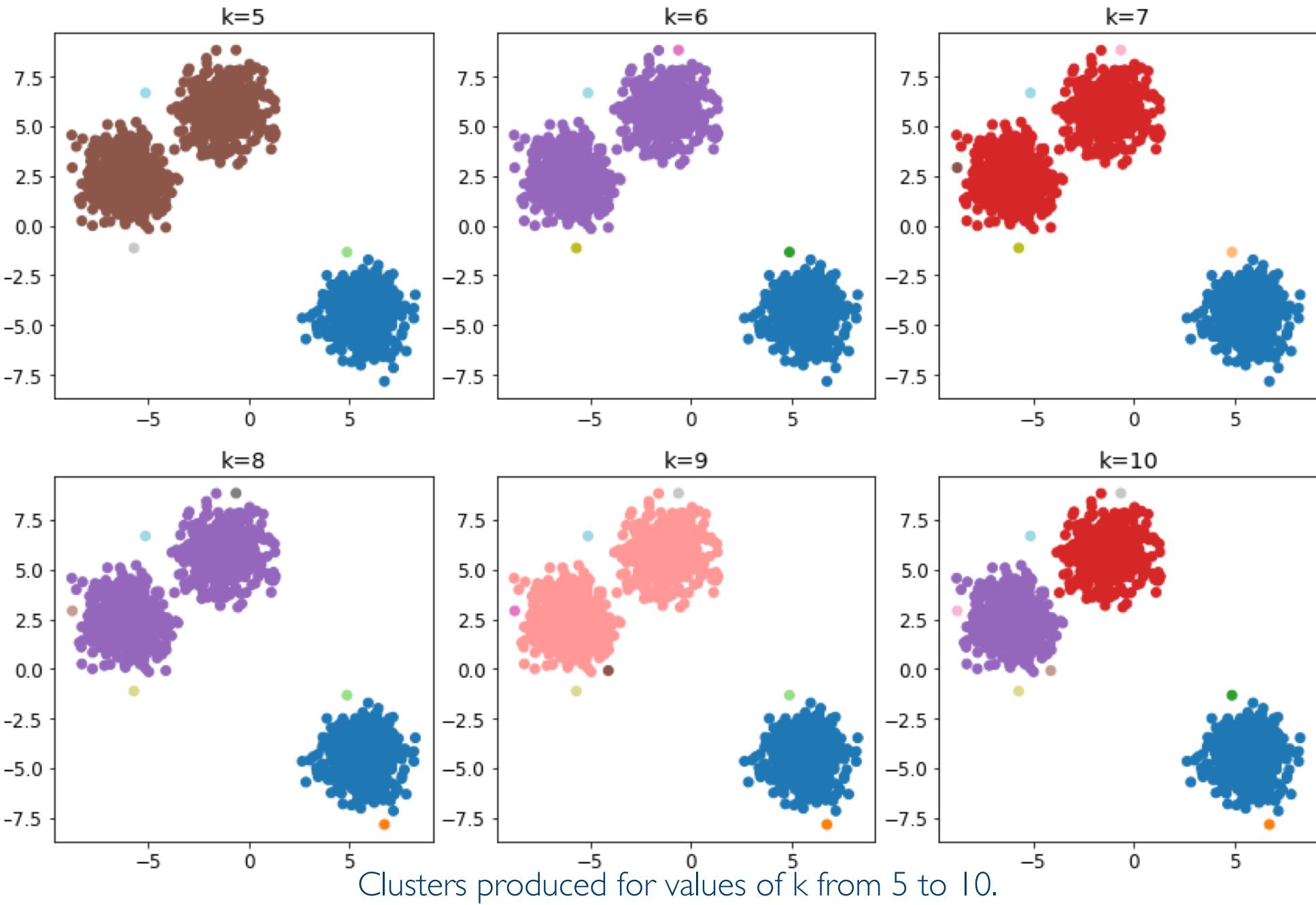
$k=6$



$k=7$



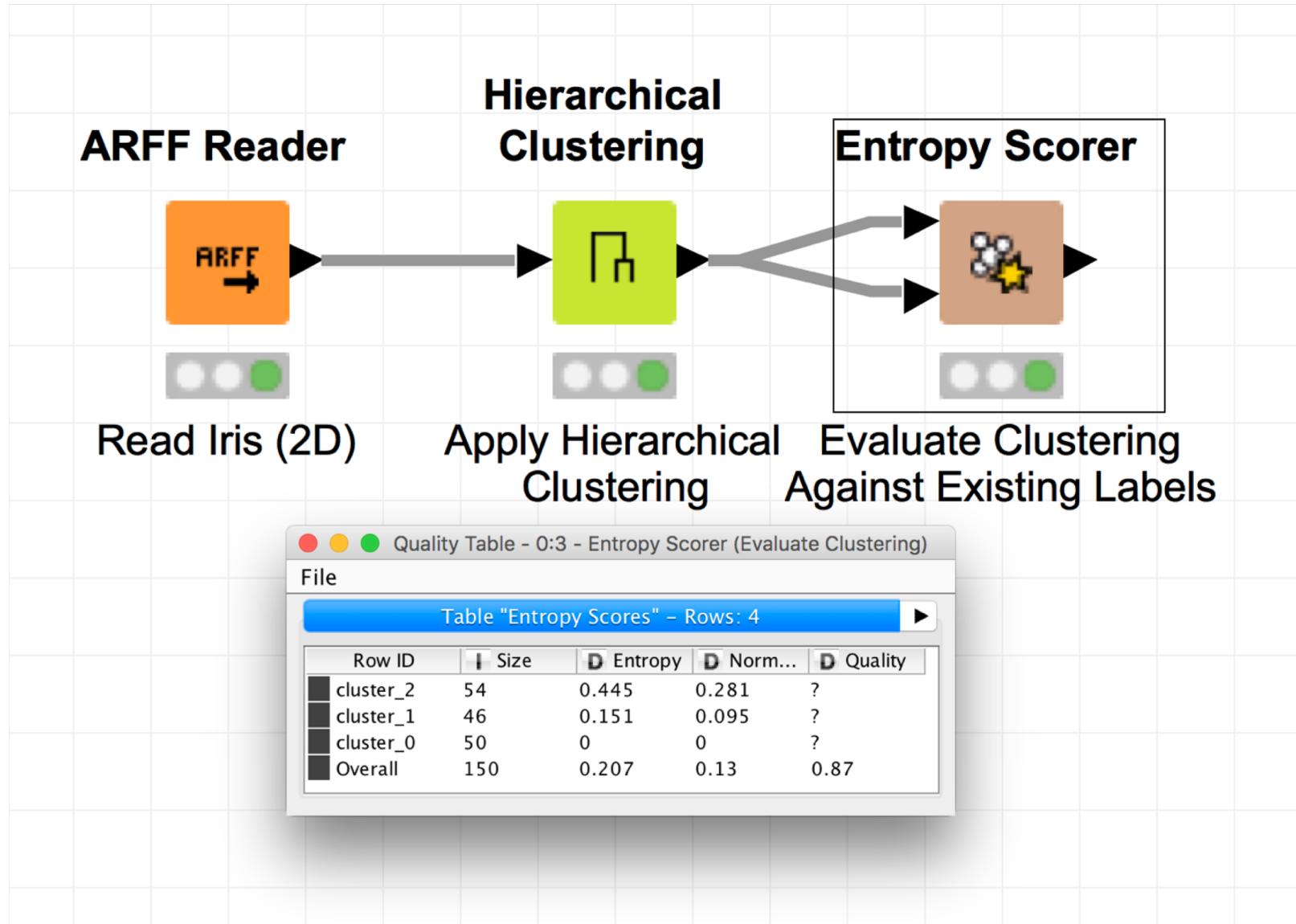
Clusters produced for values of k from 2 to 7.



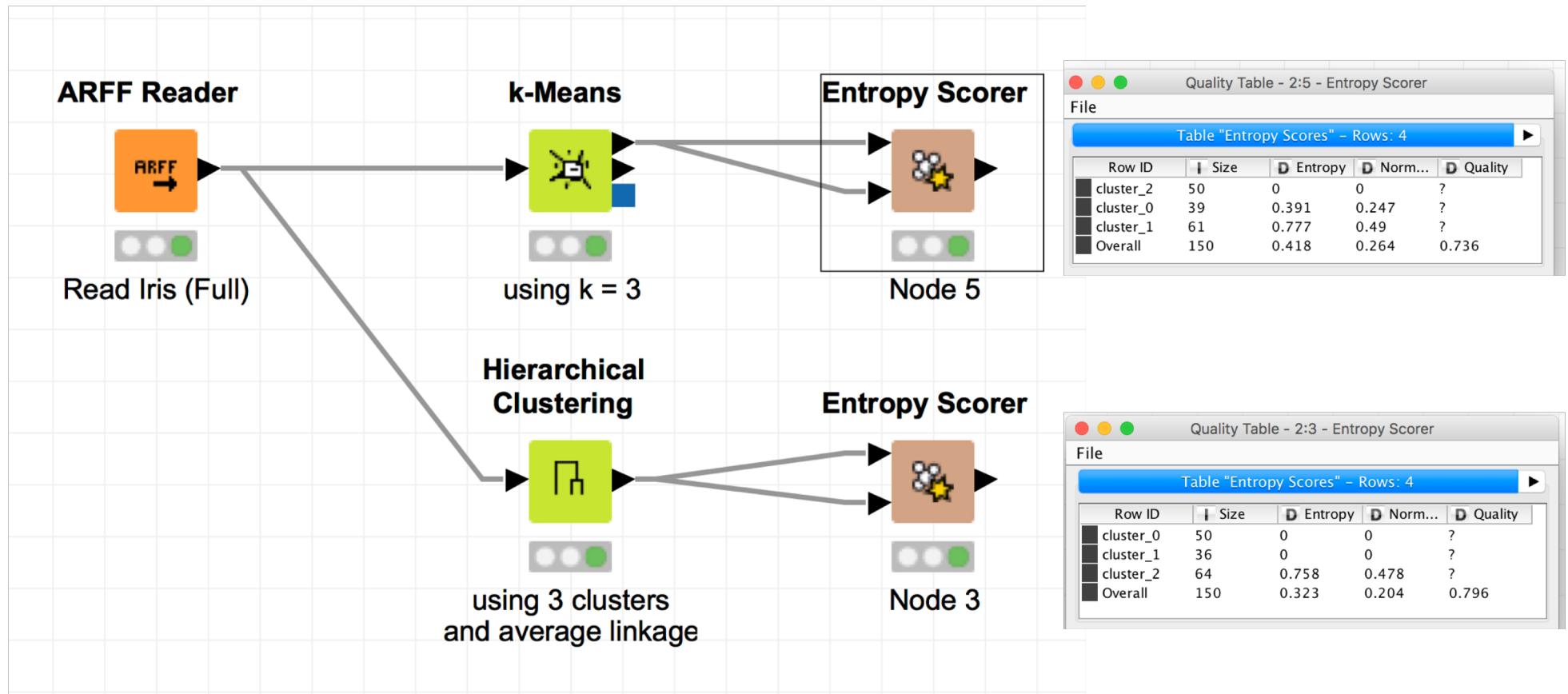
How can we represent clusters?

- Euclidean Spaces
 - We can identify a cluster using for instance its centroid (e.g. computed as the average among all its data points)
 - Alternatively, we can use its convex hull
- Non-Euclidean Spaces
 - We can define a distance (jaccard, cosine, edit)
 - We cannot compute a centroid and we can introduce the concept of clustroid
- Clustroid
 - An existing data point that we take as a cluster representative
 - It can be the point that minimizes the sum of the distances to the other points in the cluster
 - Or, the one minimizing the maximum distance to another point
 - Or, the sum of the squares of the distances to the other points in the cluster

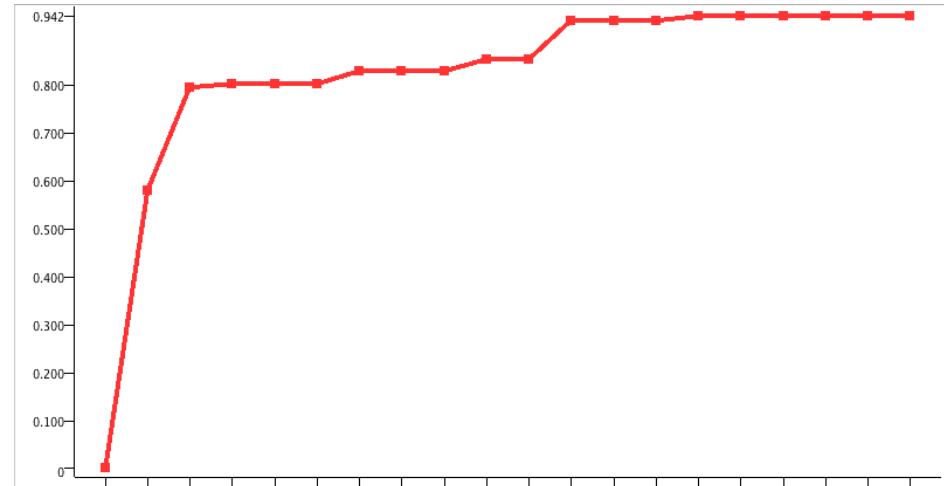
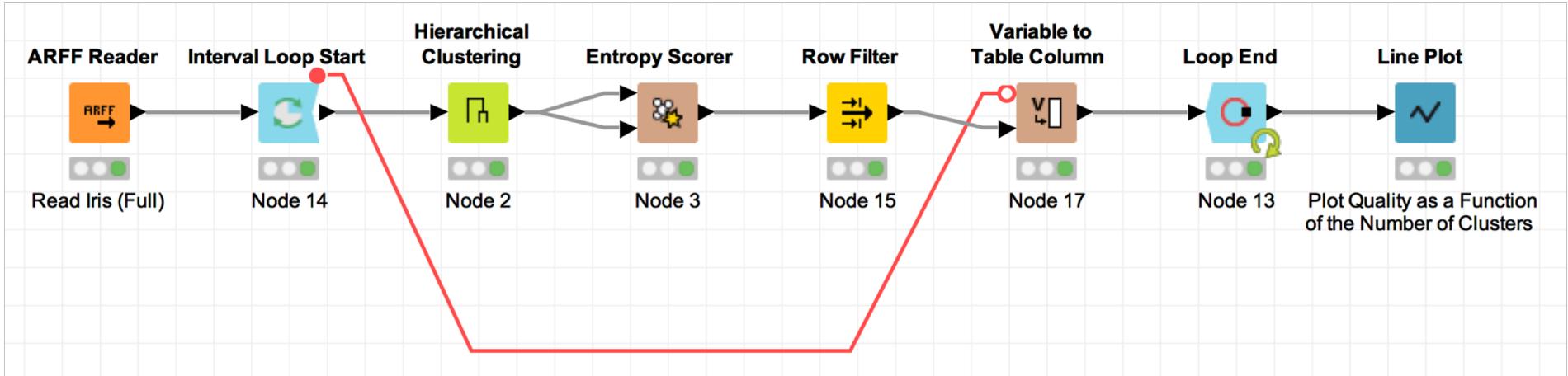
Examples using KNIME



Evaluation of the result from hierarchical clustering with 3 clusters and average linkage against existing labels



Comparison of hierarchical clustering with 3 clusters
and average linkage against k-Means with $k=3$



Computing cluster quality from one to 20 clusters
using the entropy scorer

Summary

Hierarchical Clustering: Problems and Limitations

- Once a decision is made to combine two clusters, it cannot be undone
- No objective function is directly minimized
- Different schemes have problems with one or more of the following:
 - Sensitivity to noise and outliers
 - Difficulty handling different sized clusters and convex shapes
 - Breaking large clusters
- Major weakness of agglomerative clustering methods
 - They do not scale well: time complexity of at least $O(n^2)$, where n is the number of total objects
 - They can never undo what was done previously

Run the python notebooks and the
KNIME workflows for this lecture