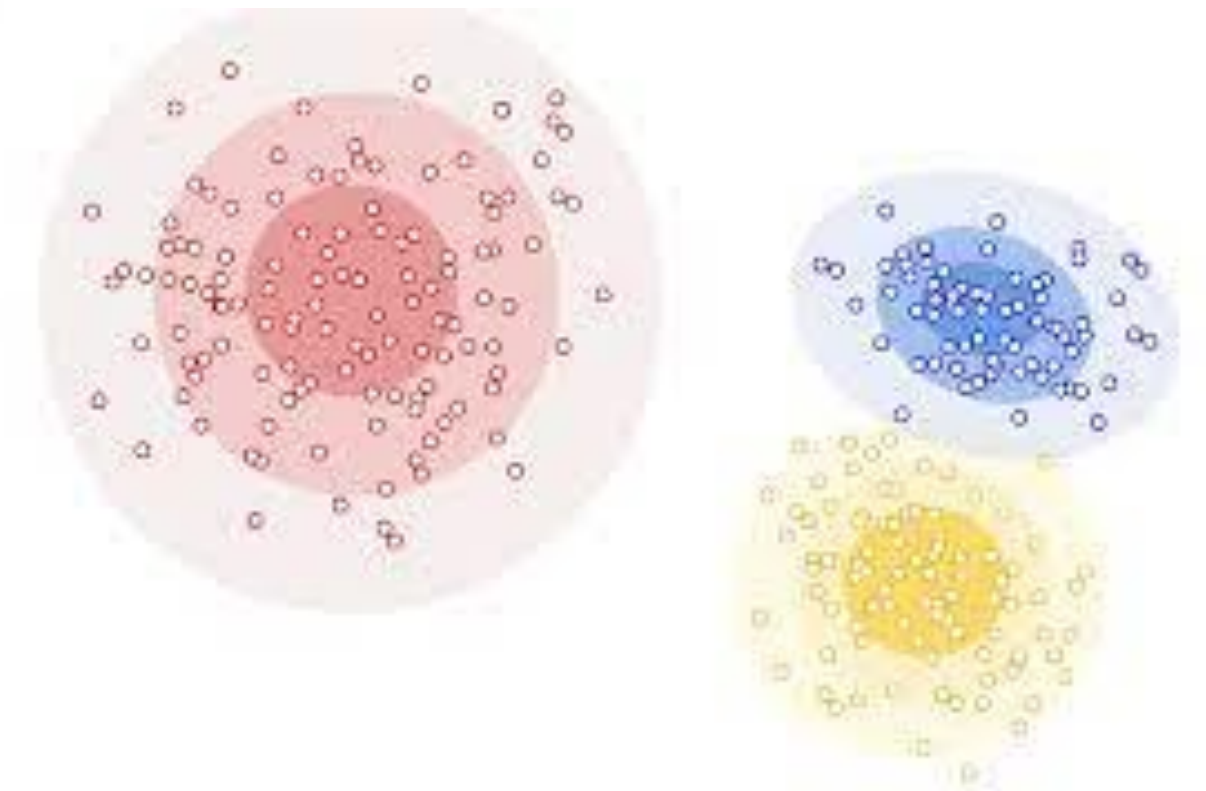
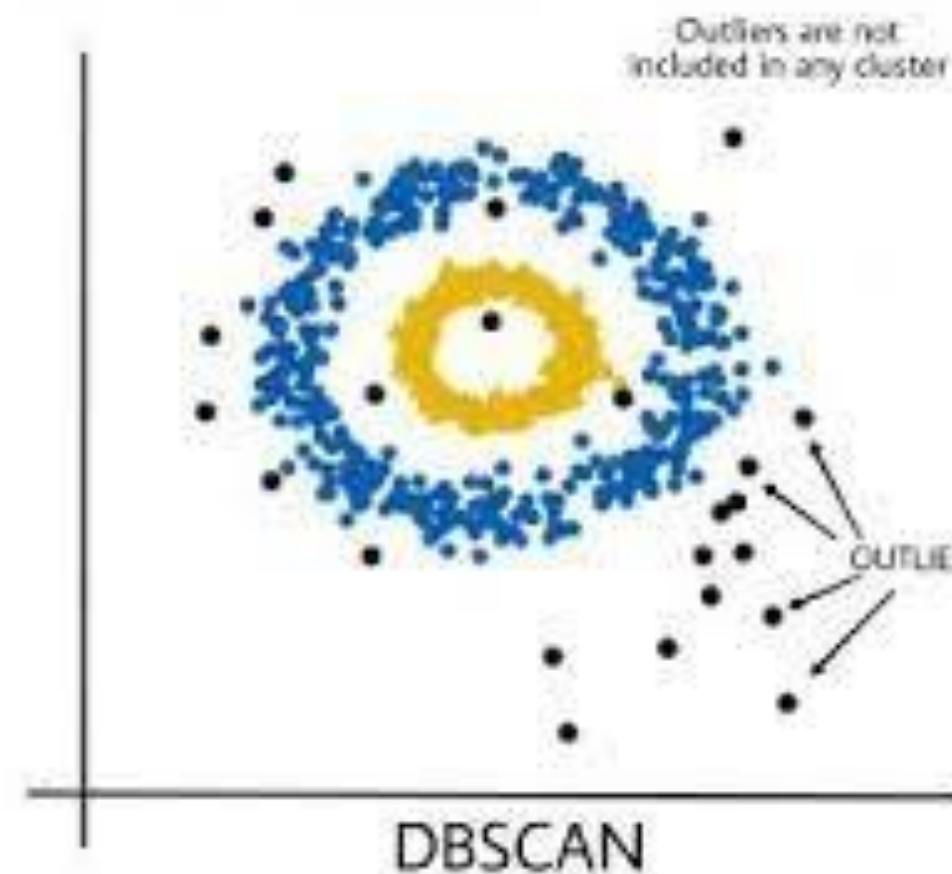
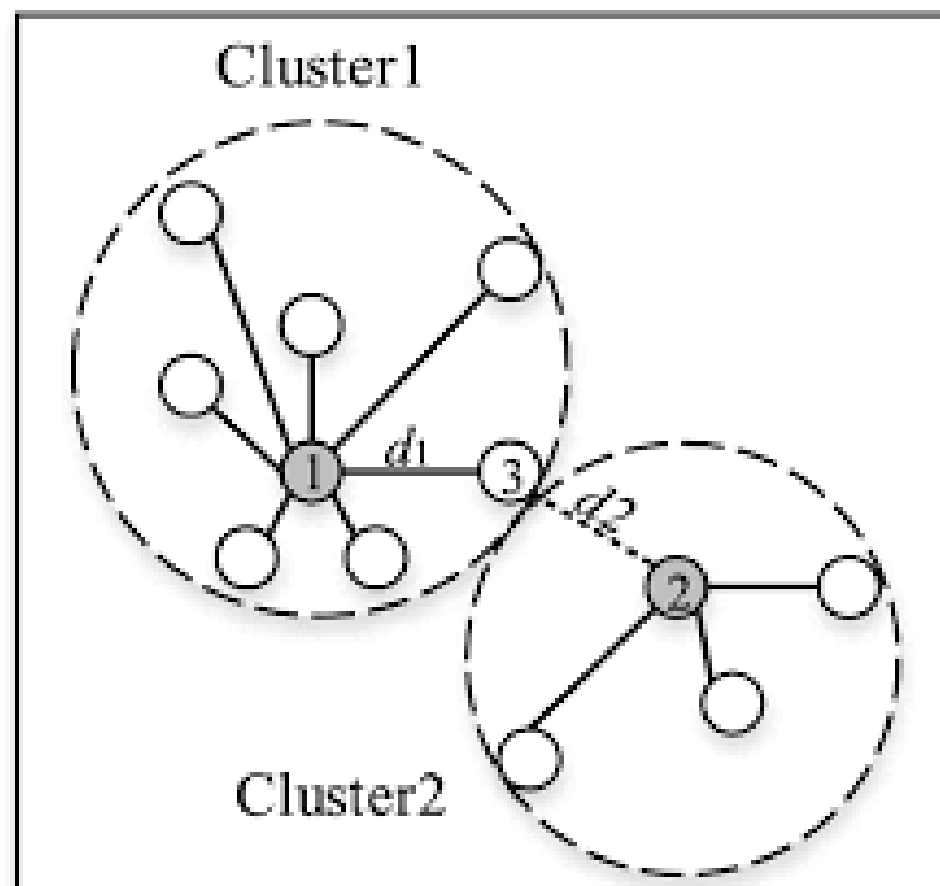
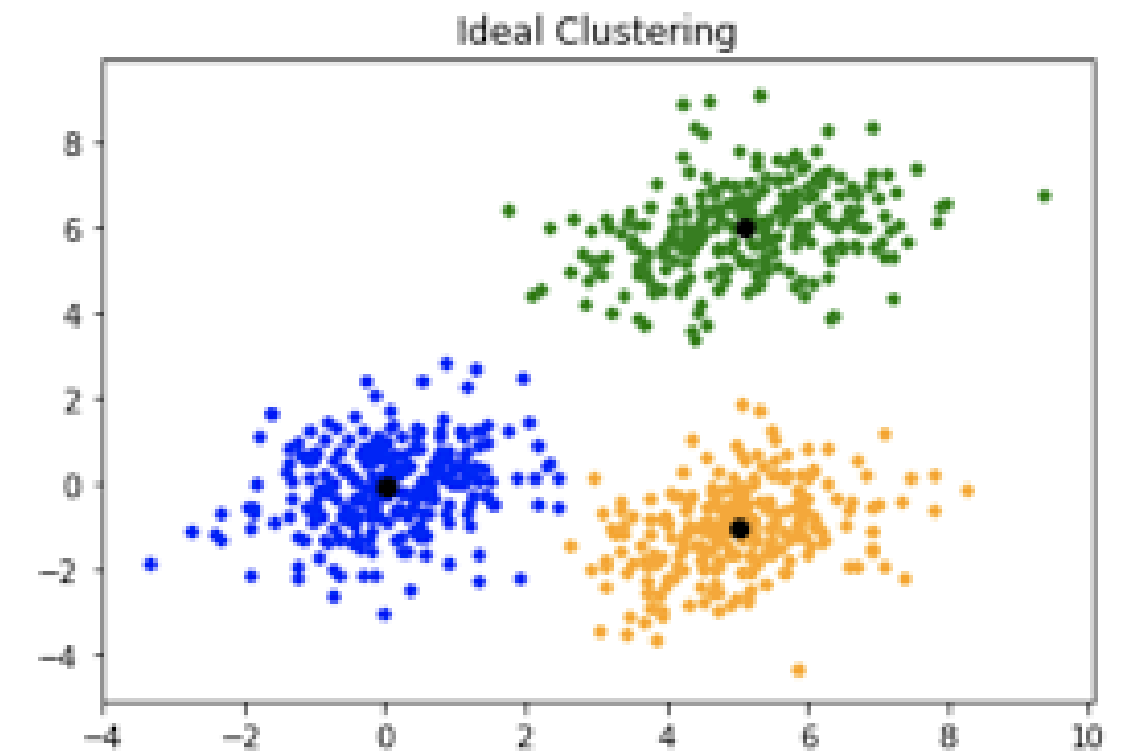




# Lecture 16: Hierarchical Clustering

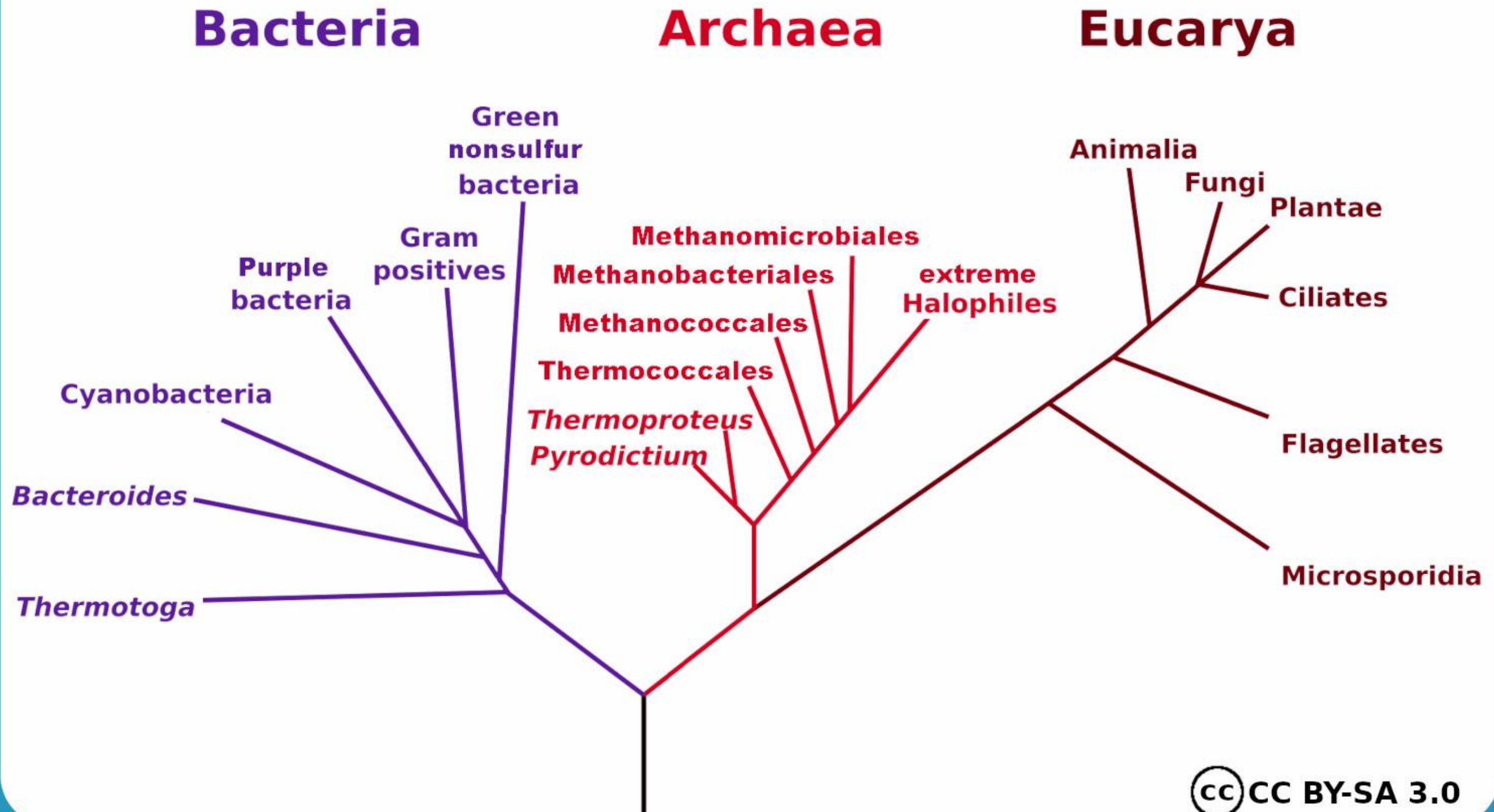
# Cluster Analysis - Recap

- Centroid based clustering
- Distribution based clustering
- Density based clustering
- Connectivity based clustering

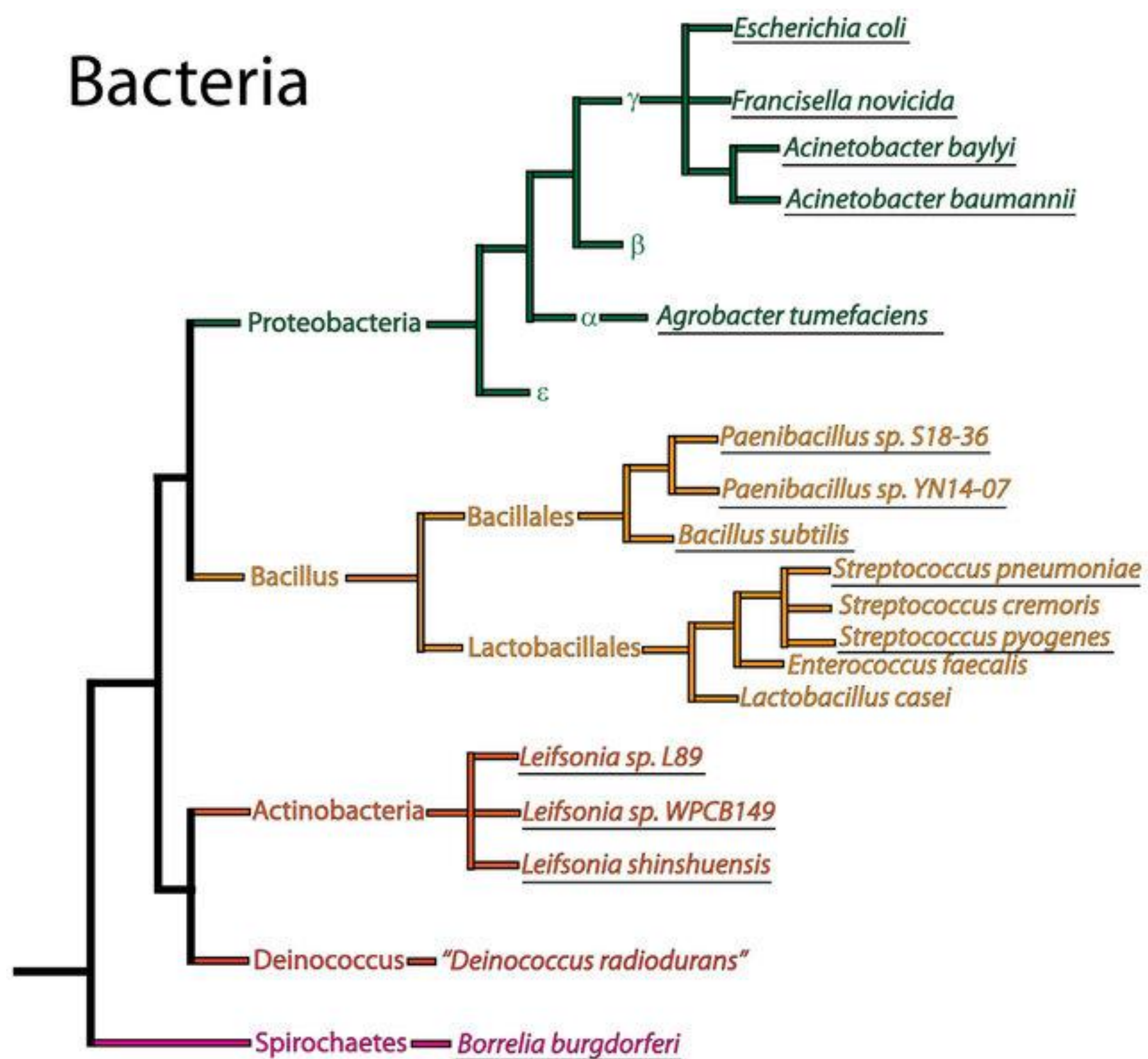




# Phylogenetic Tree of Life



# Bacteria







# Intro to hierarchical clustering

- How many clusters in the data is
  - Technical question
  - Centroid, Distribution & Density answers this question
- Business questions
  - How are clusters related?
  - How to say 2 clusters are different but are also similar?
- We need a hierarchical organization of clusters
  - Based on high level effects / granular sub groupings
- Very natural for some problems
  - Taxonomy

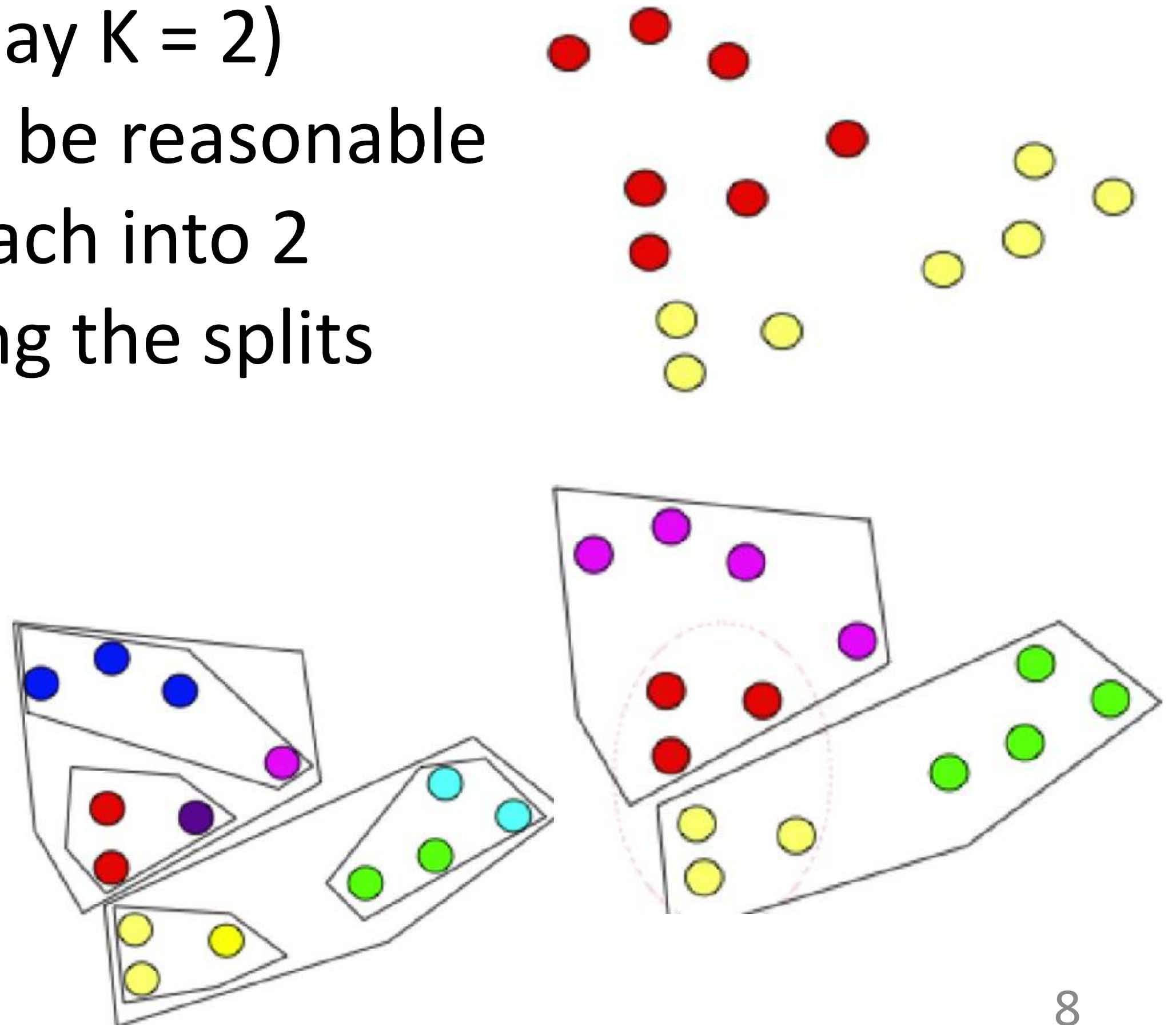
# Hierarchical clustering

- Top-down
  - Start with one cluster, split recursively
  - Divisive clustering
- Bottom-up
  - Each point is a cluster, merge clusters
  - Agglomerative clustering
- Tree in both cases
  - All data in one cluster at top
  - Each data in its own cluster at bottom



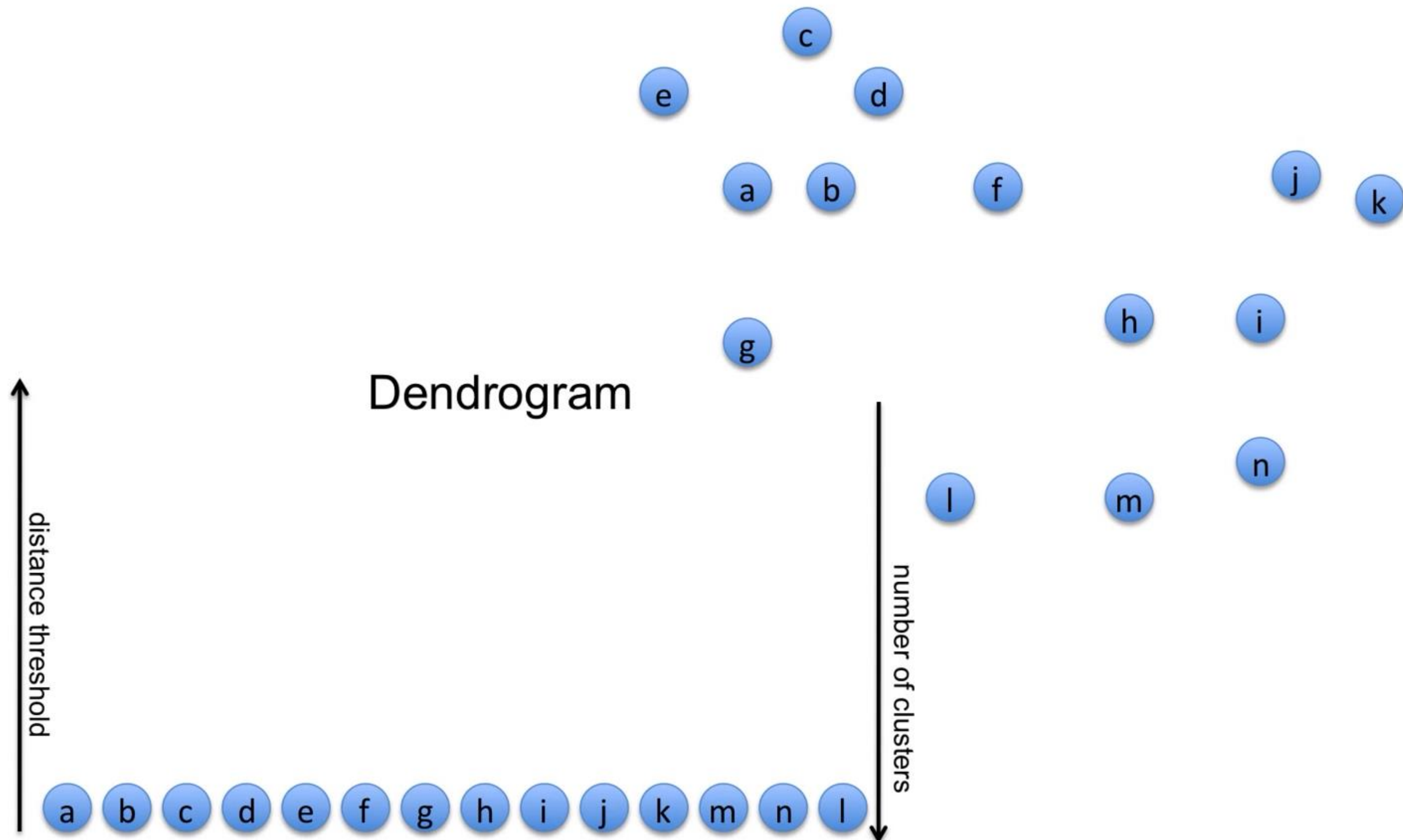
# Divisive clustering

- Recursively run K means (say  $K = 2$ )
- Partition into two may not be reasonable
- Split into 2. Further split each into 2
- No going back and changing the splits
  - Greedy
  - Sometimes problematic
- Computationally efficient

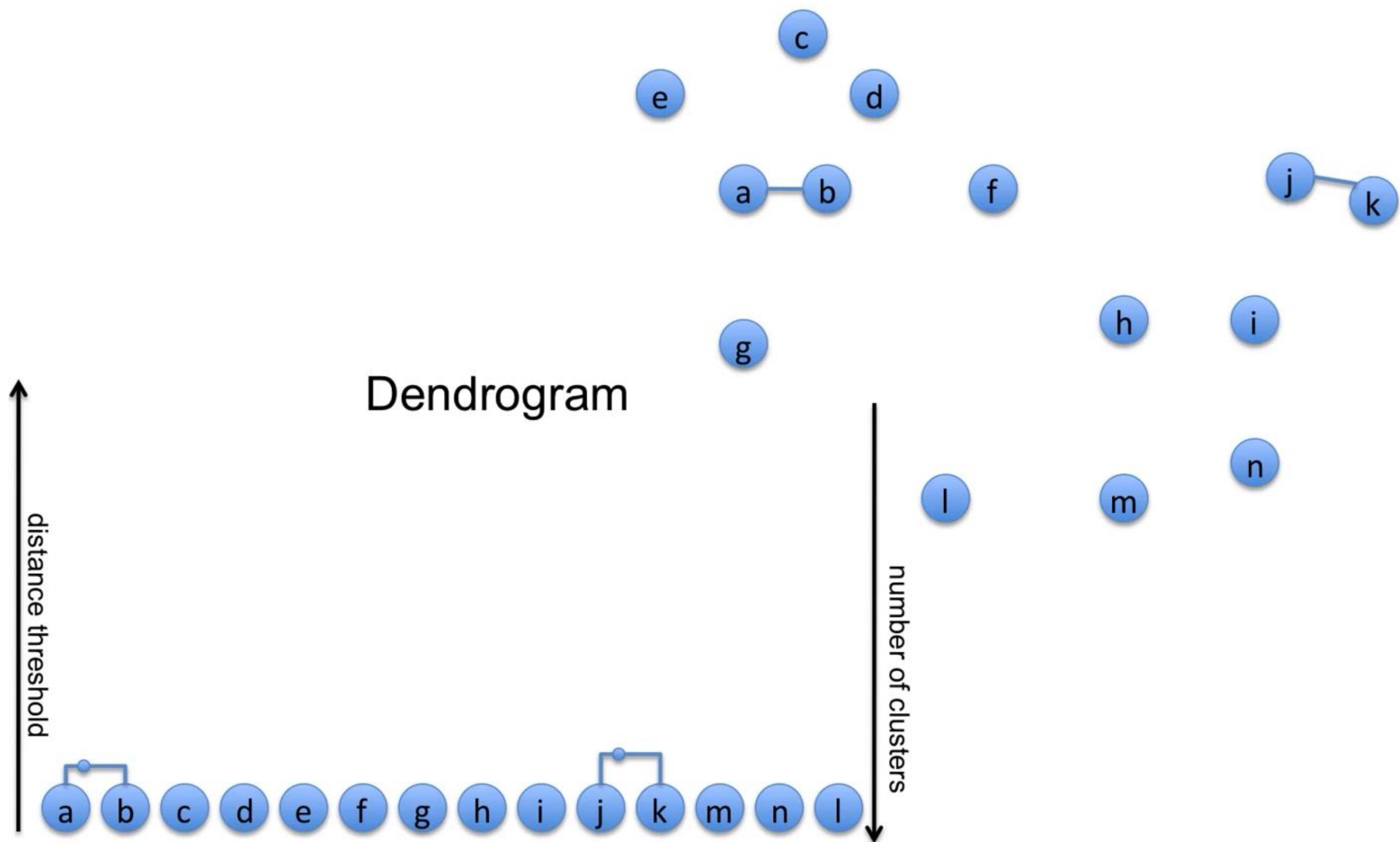


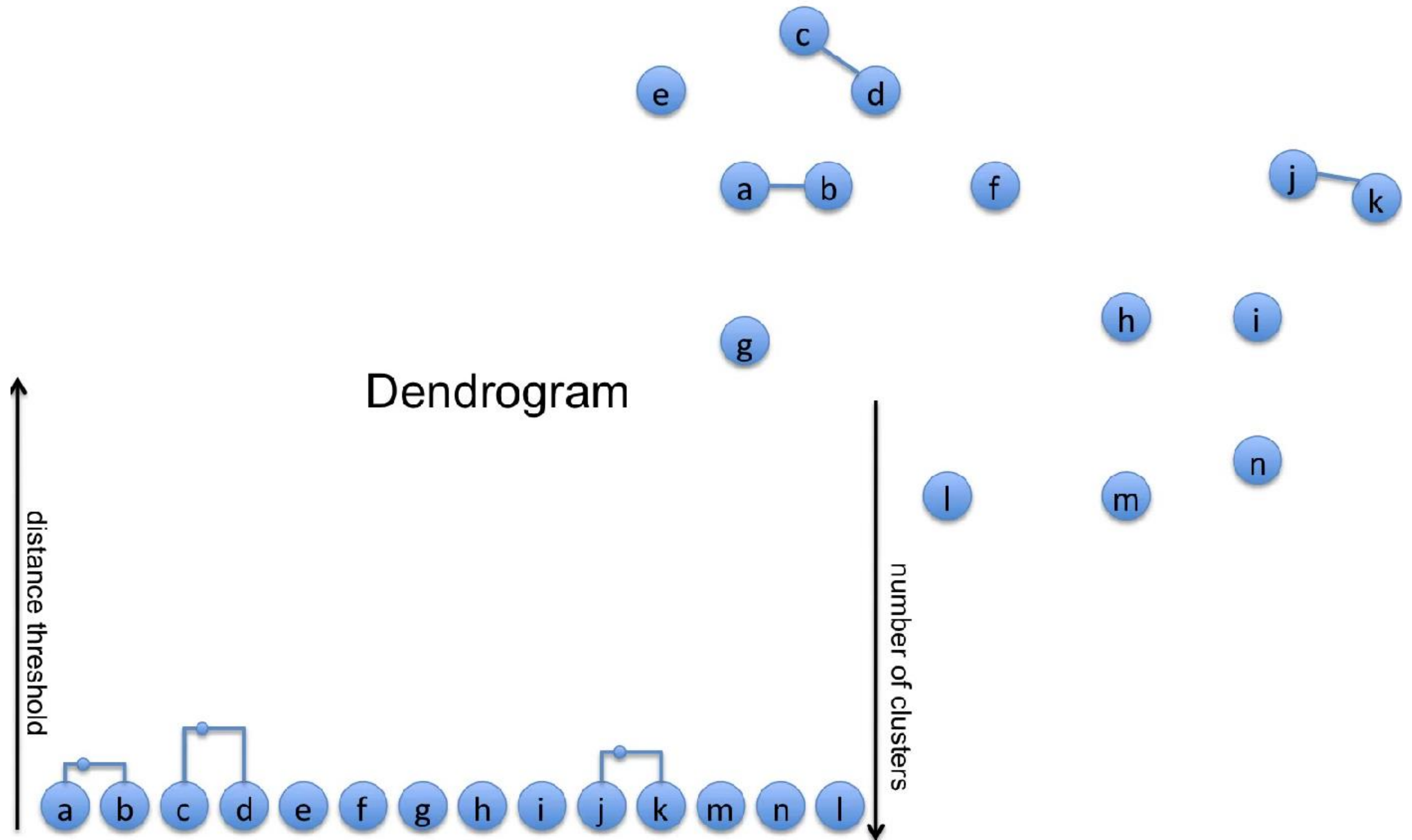














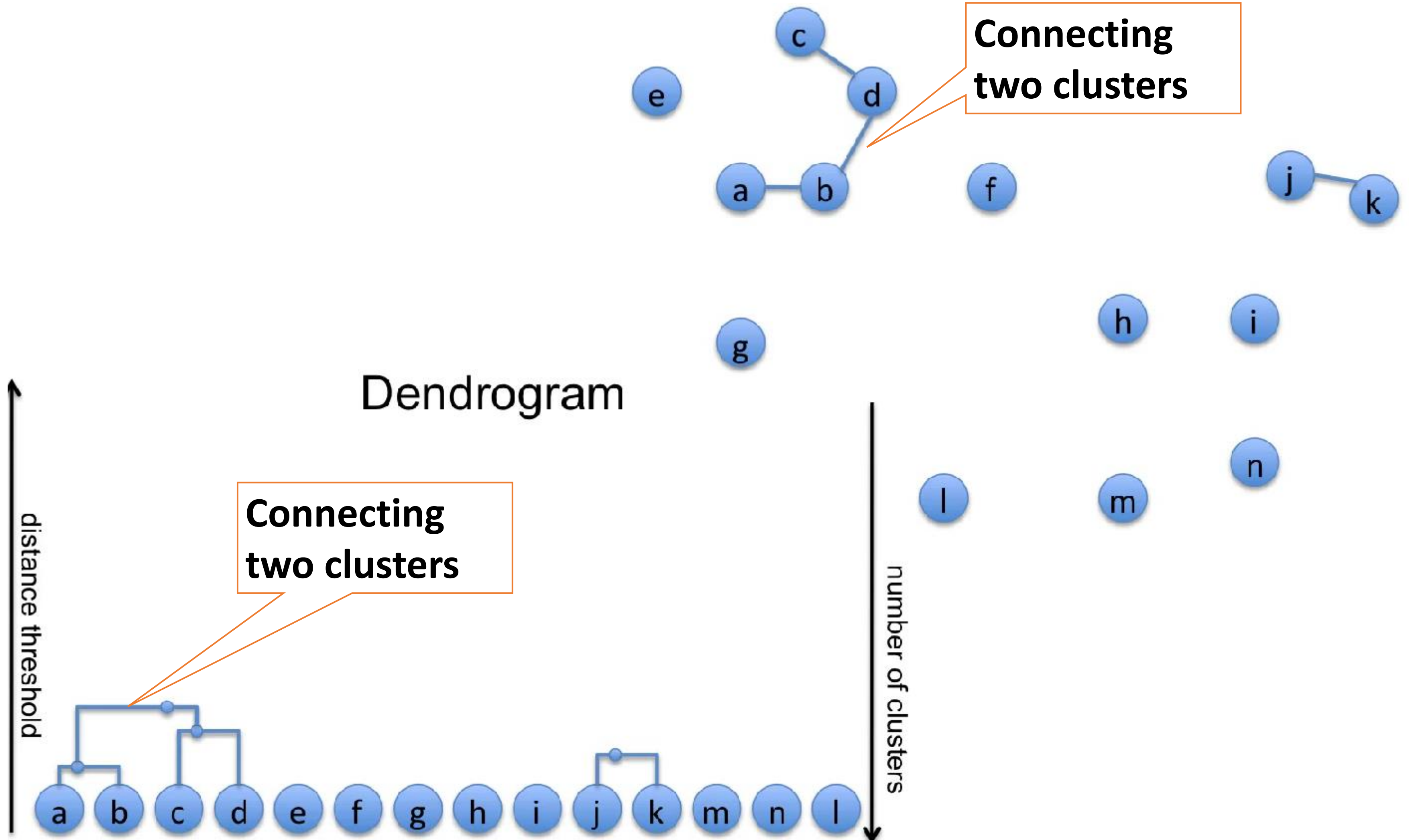
**Connecting  
two clusters**

## Dendrogram

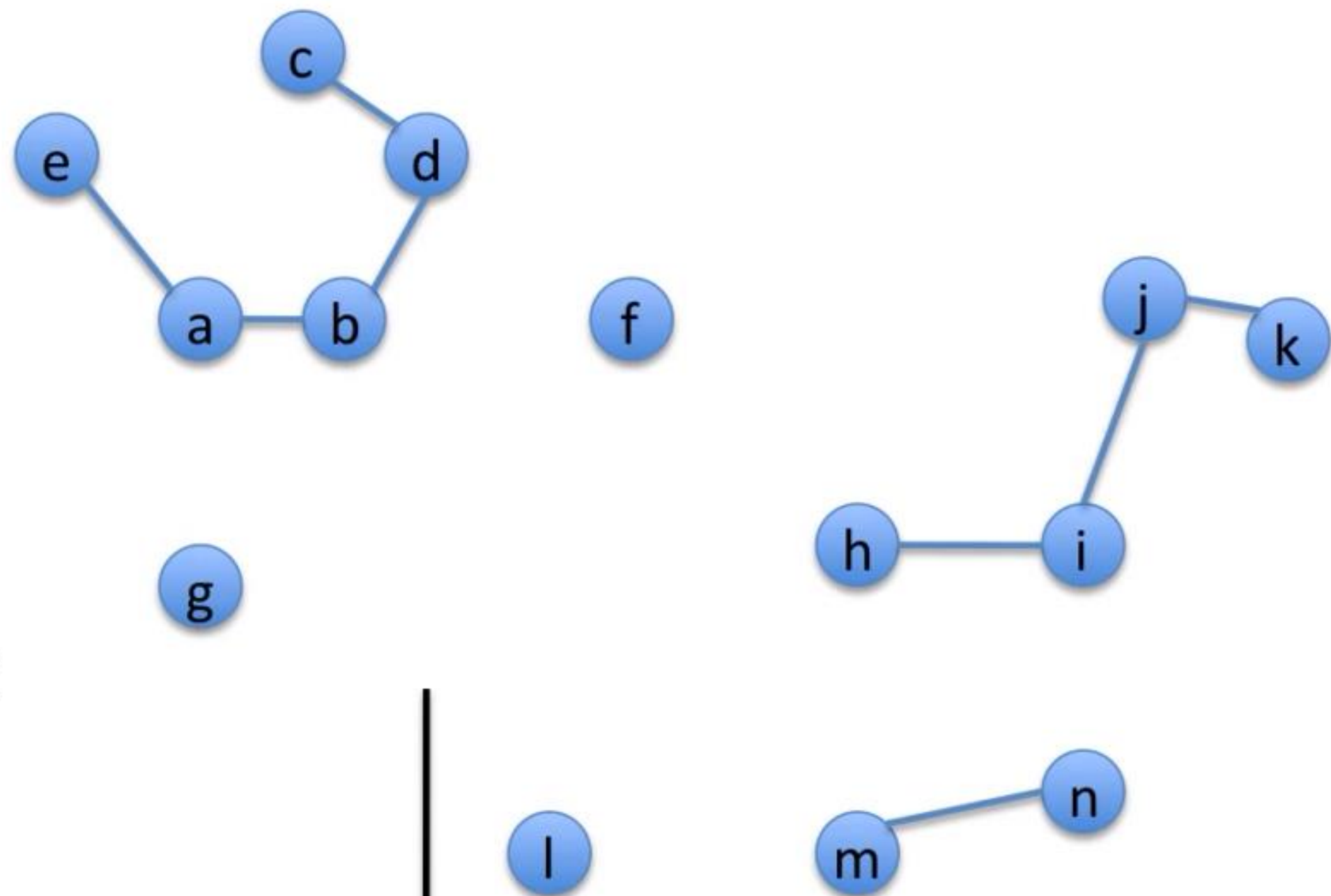
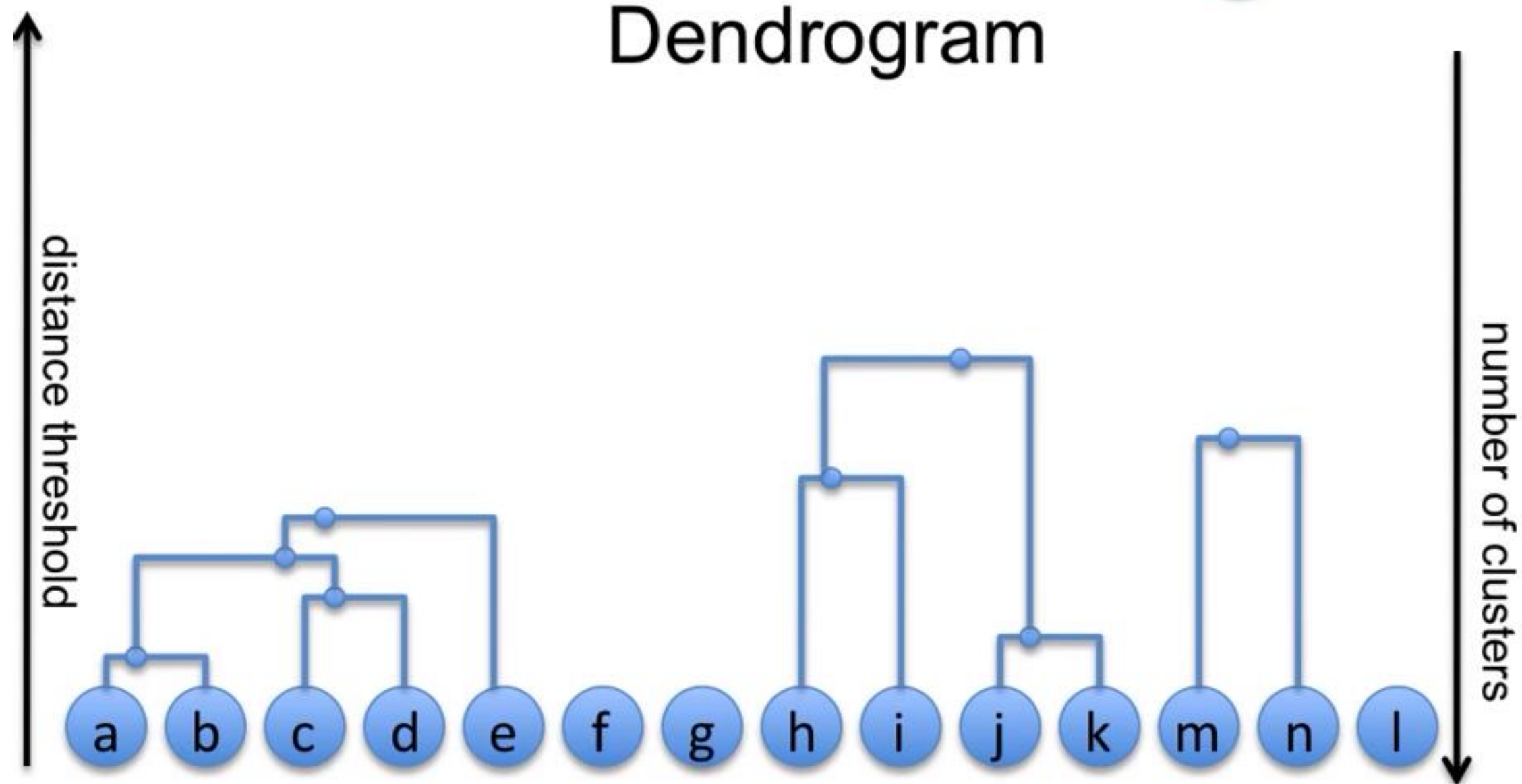
**Connecting  
two clusters**

number of clusters

distance threshold

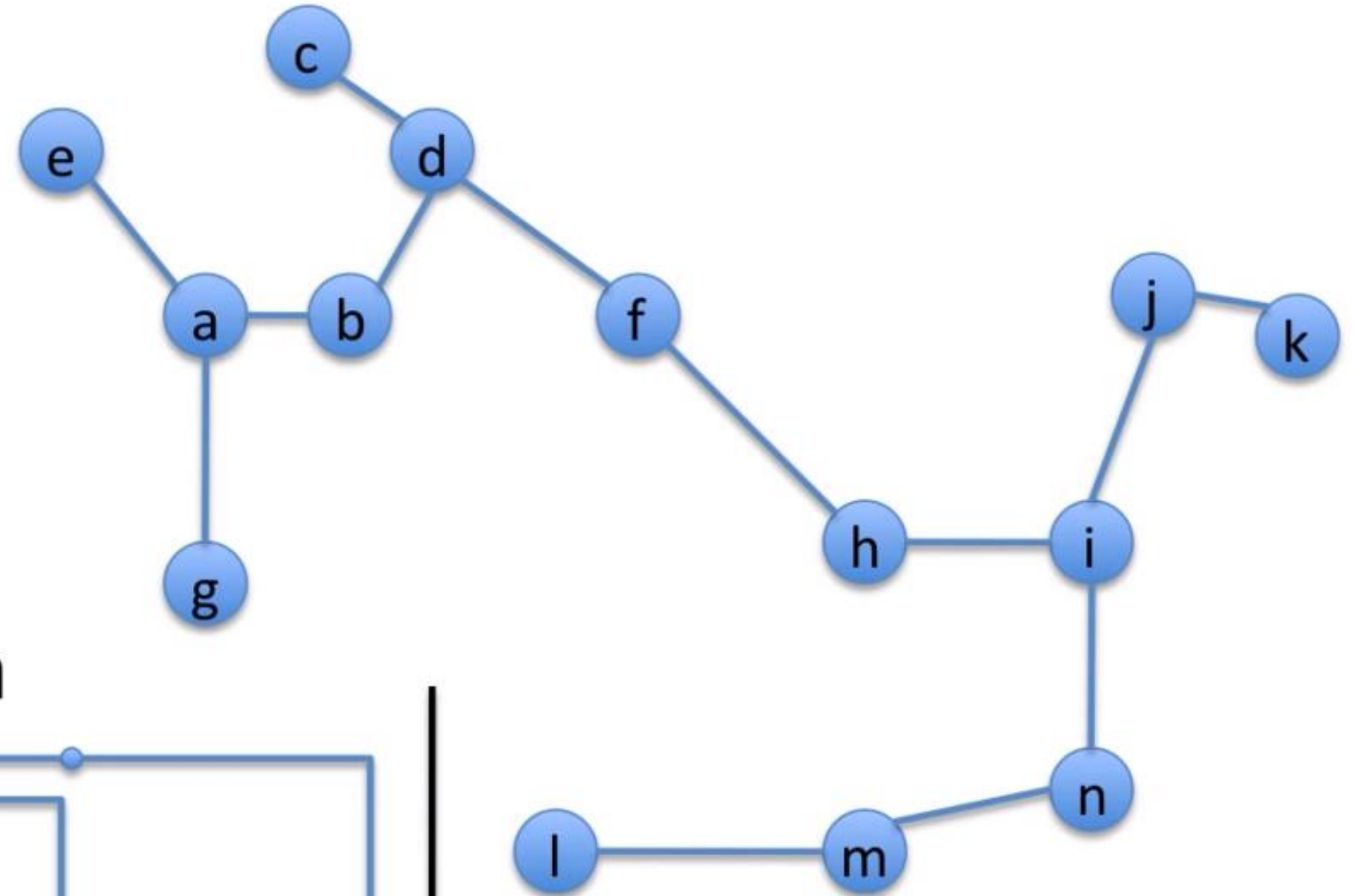


# Dendrogram

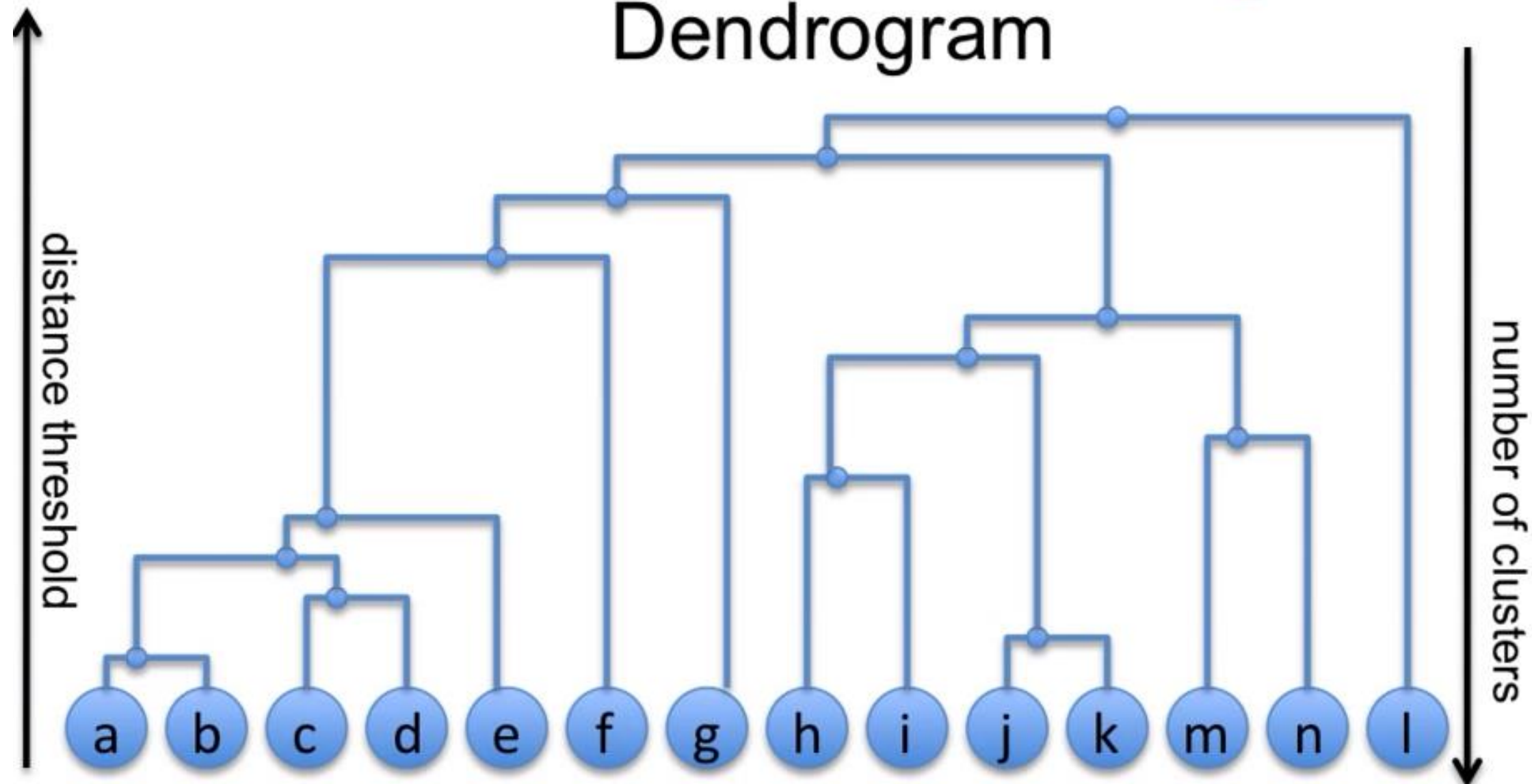


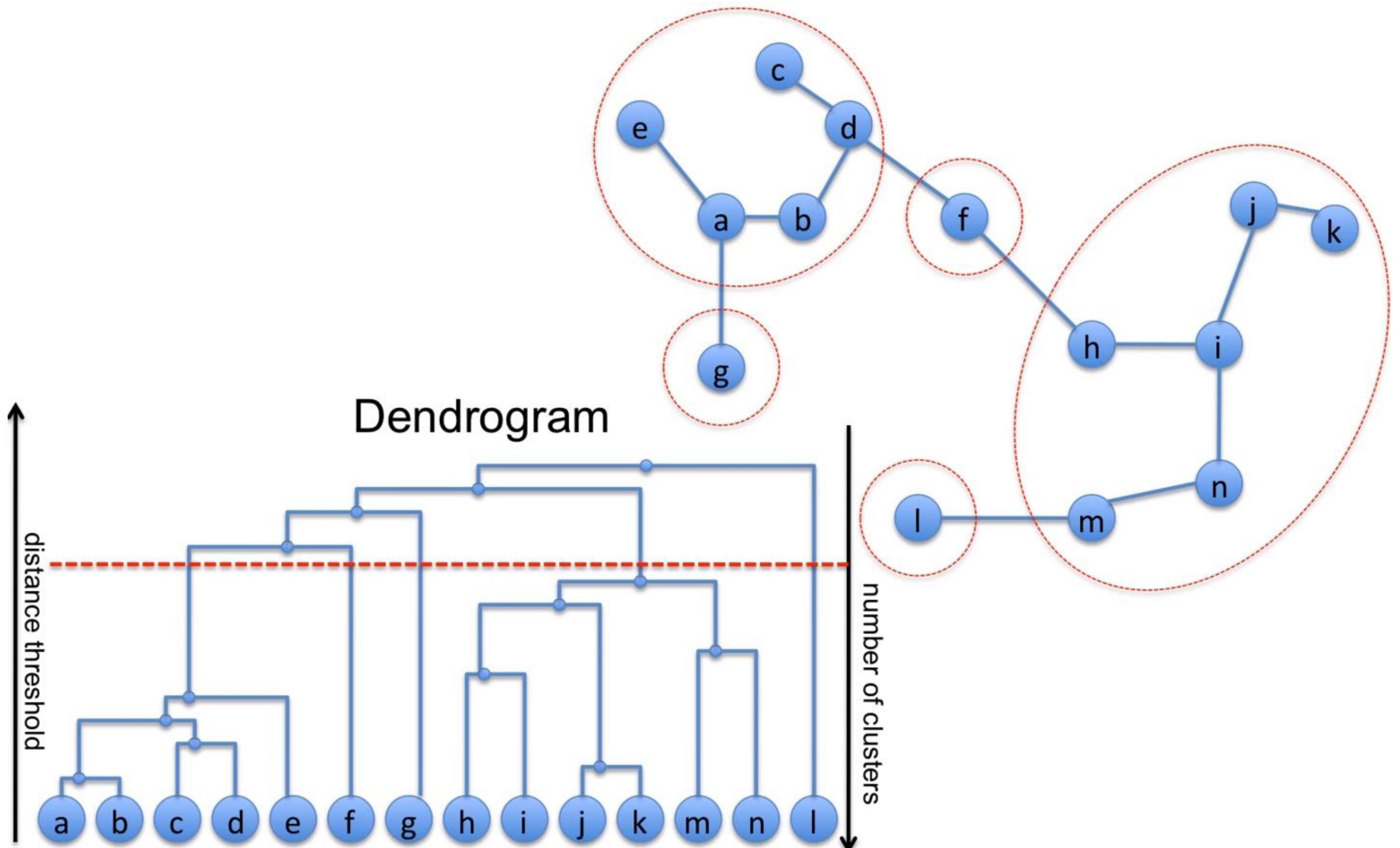


How many  
clusters is  
good?



Dendrogram





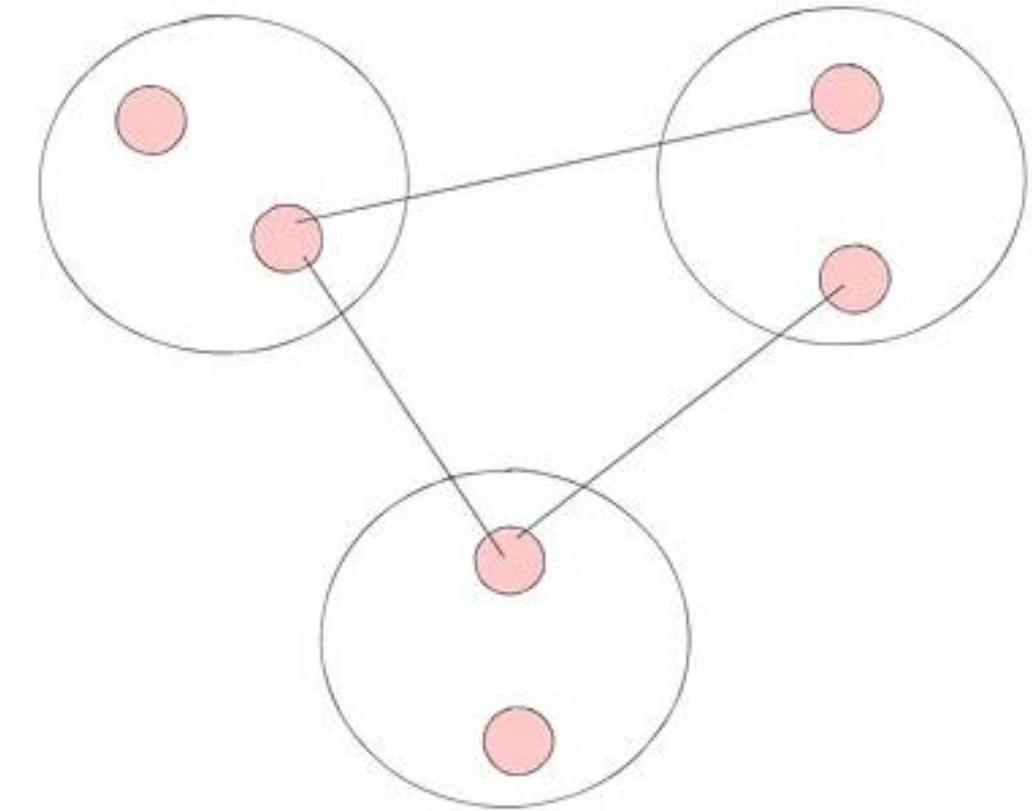


# Agglomerative clustering

- Idea: Nearby points end up in same cluster
- Algorithm
  - Start with collection  $C$  of  $N$  singleton clusters
  - $c_j = \{x_j\}$
  - Find another cluster that is closest  $\min_{ij} (c_i, c_j)$
  - Merge  $c_i$  and  $c_j$  into new cluster  $c_{i+j}$
  - Remove  $c_i$  and  $c_j$  from  $C$ . Add  $c_{i+j}$
  - Repeat
- Slower algorithm
- Understand with dendrogram (connectivity tree)

# Cluster distance metric for hierarchical clustering

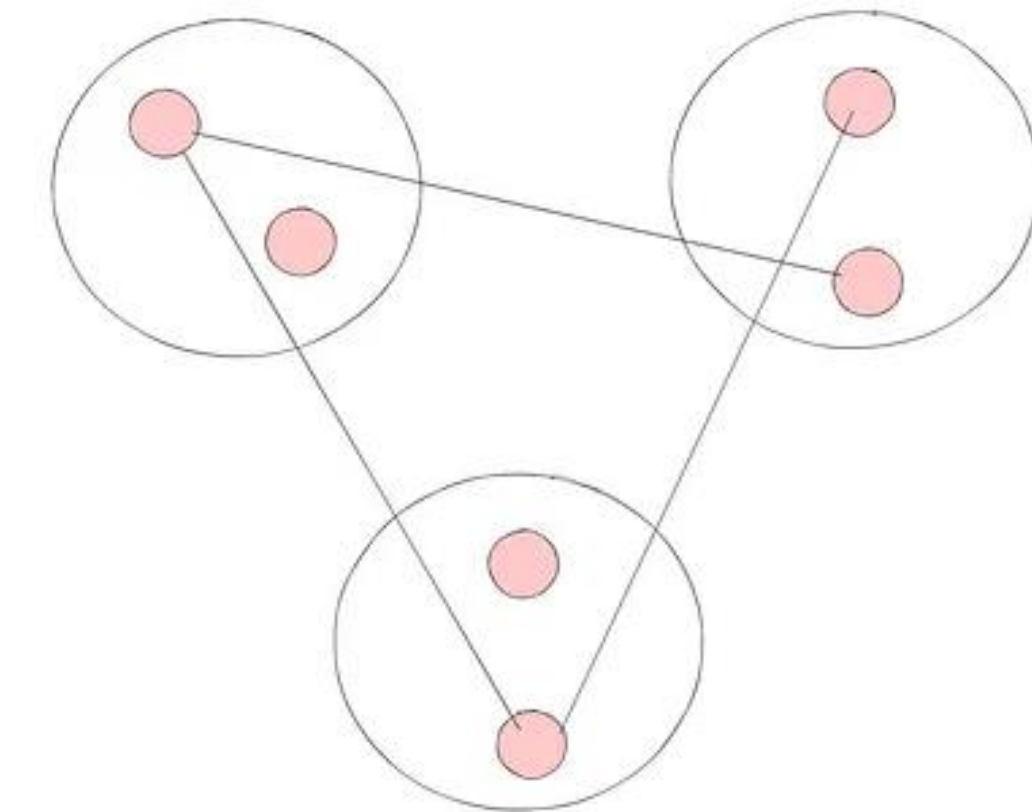
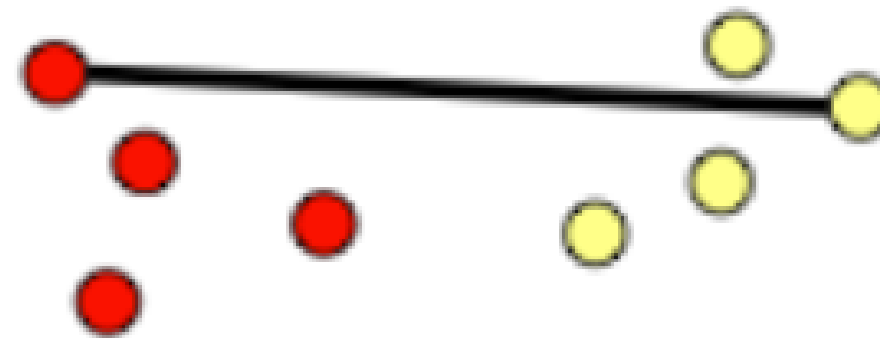
- Distance between clusters
- Not individual points
- Single Link



- Distance b/w closest elements of clusters

$$\mathcal{D}(c_1, c_2) = \min_{x_1 \in c_1, x_2 \in c_2} \mathcal{D}(x_1, x_2)$$

- Produces long chains
- Complete Link



- Distance b/w farthest elements of clusters

$$\mathcal{D}(c_1, c_2) = \max_{x_1 \in c_1, x_2 \in c_2} \mathcal{D}(x_1, x_2)$$

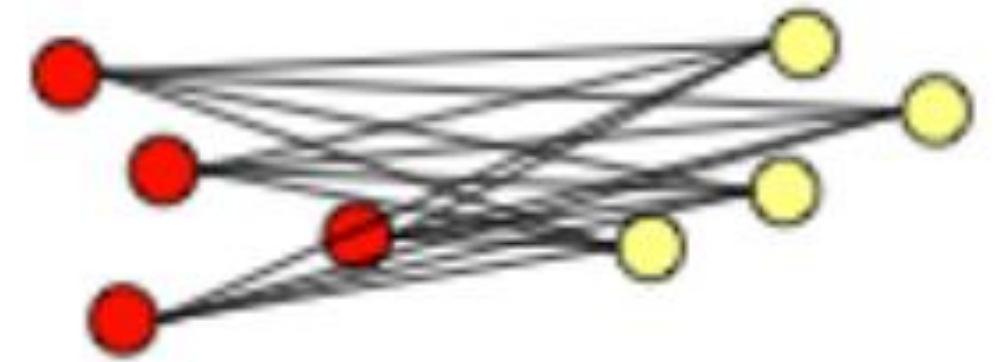
- Produces spherical clusters



# Cluster distance metric for hierarchical clustering

- Average Link

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

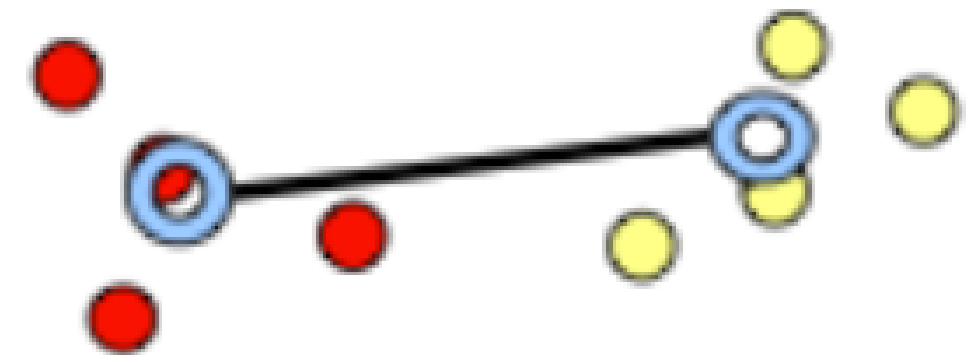


- Avg of pairwise distance b/w cluster elements

- Less affected by outliers

- Centroid Link

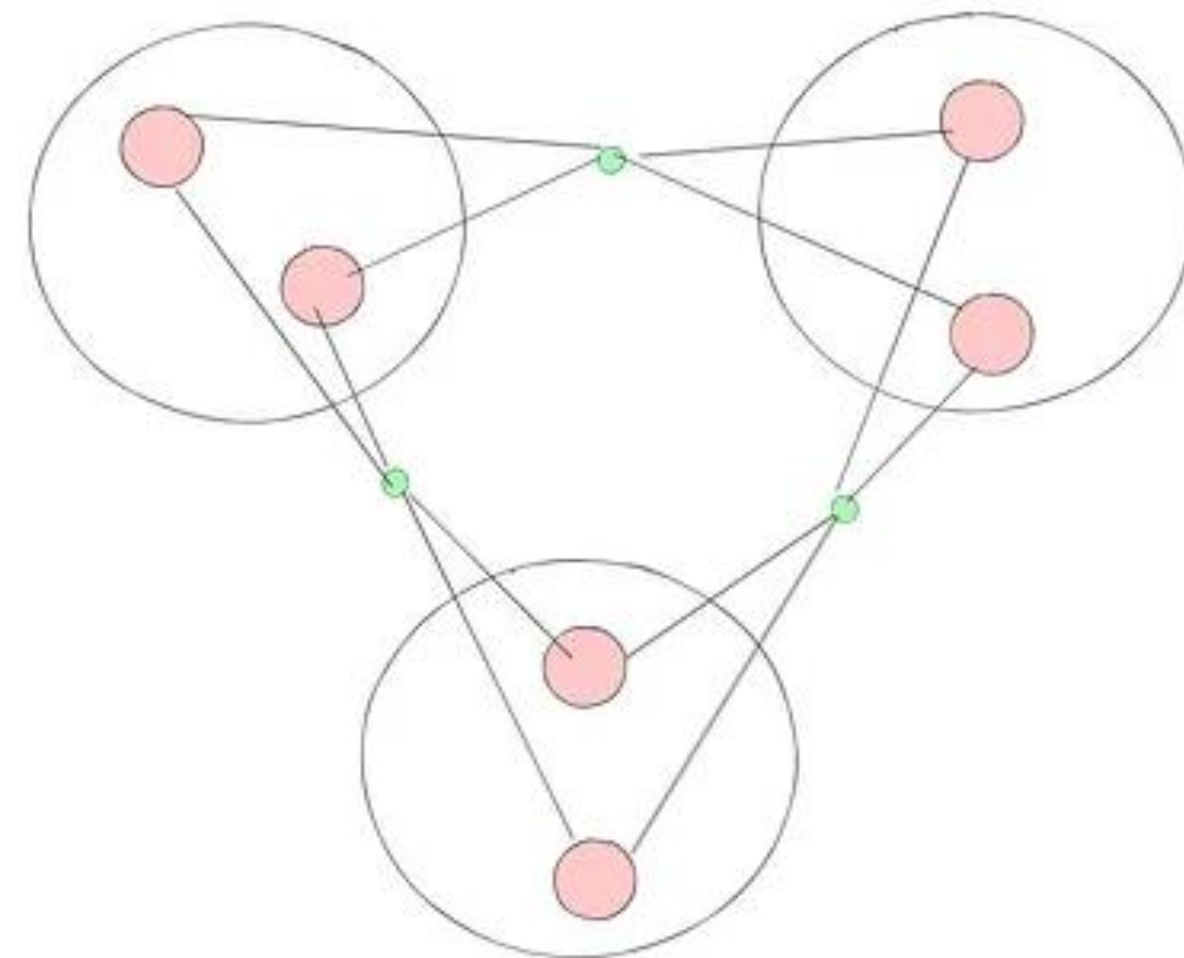
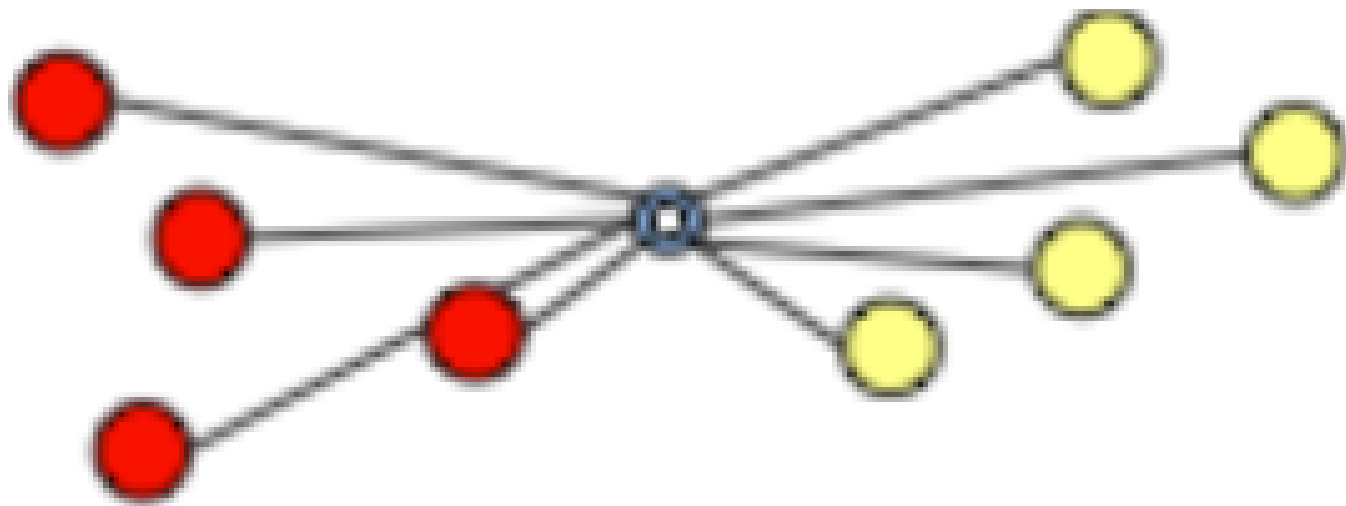
- Distance b/w cluster centroids



$$\mathcal{D}(c_1, c_2) = \mathcal{D}\left(\left(\frac{1}{|c_1|} \sum_{x_1 \in c_1}\right), \left(\frac{1}{|c_2|} \sum_{x_2 \in c_2}\right)\right)$$

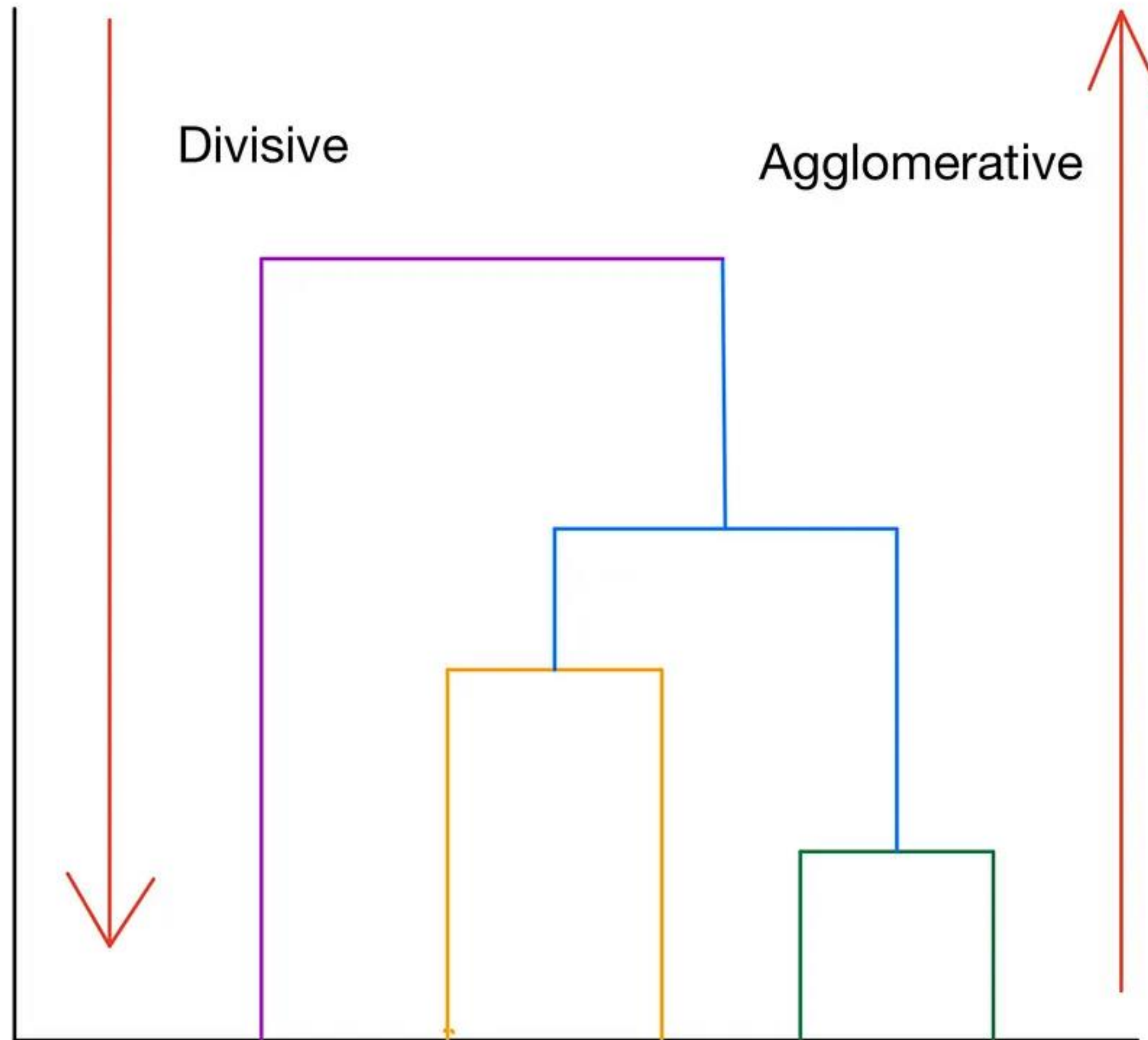
# Cluster distance metric for hierarchical clustering

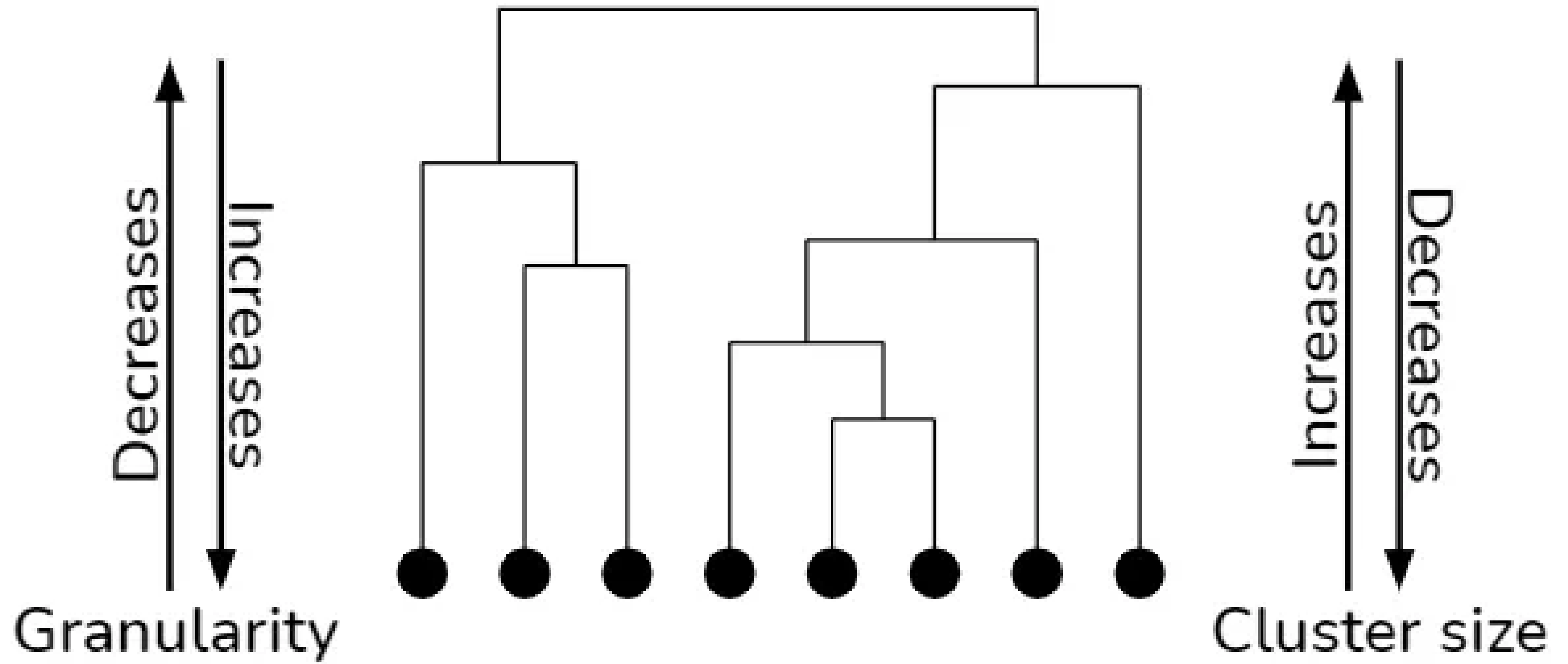
- Ward Link
  - If the two clusters were to join
  - What will be their mean and standard deviation?
  - The two clusters with minimum SD will be merged

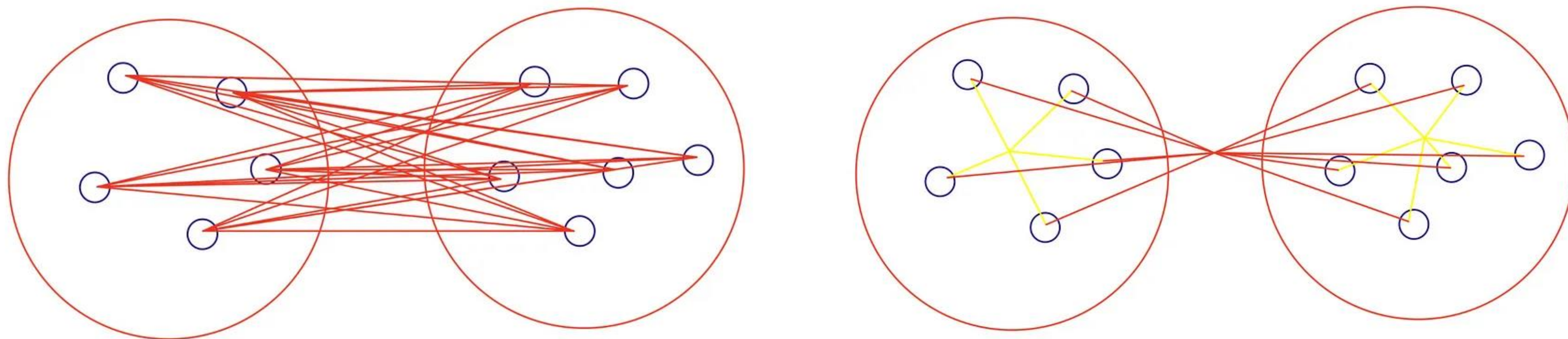
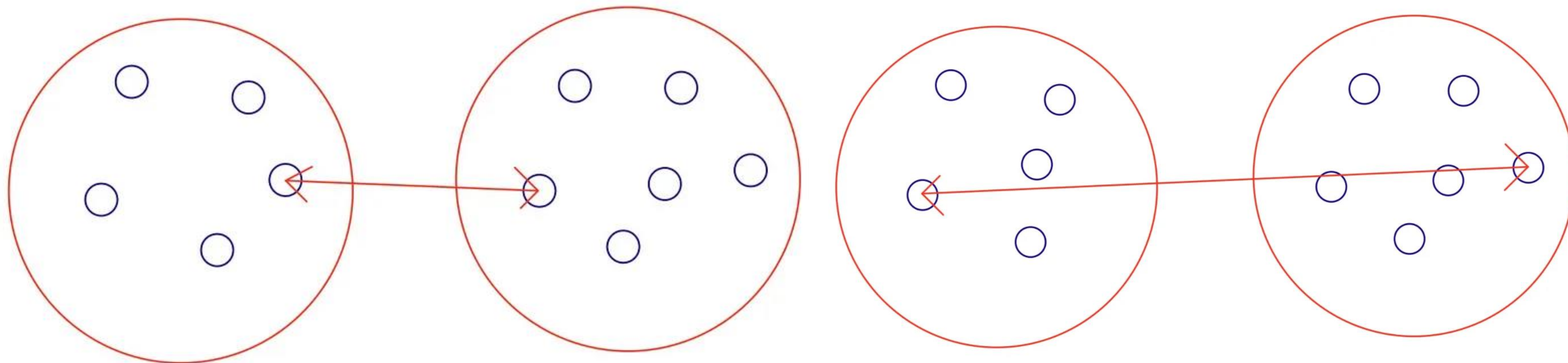














# Connectivity v/s other clustering algorithms

- Compactness (Cluster cohesion)
- Separation

$$Index = \frac{(\alpha \times Separation)}{(\beta \times Compactness)}$$

- Density – looks for a point within a epsilon distance
- Connectivity – No upper limit on nearest point distance
  - No good evaluation metric
  - Empirical or distance metric should capture domain knowledge



QUESTIONS