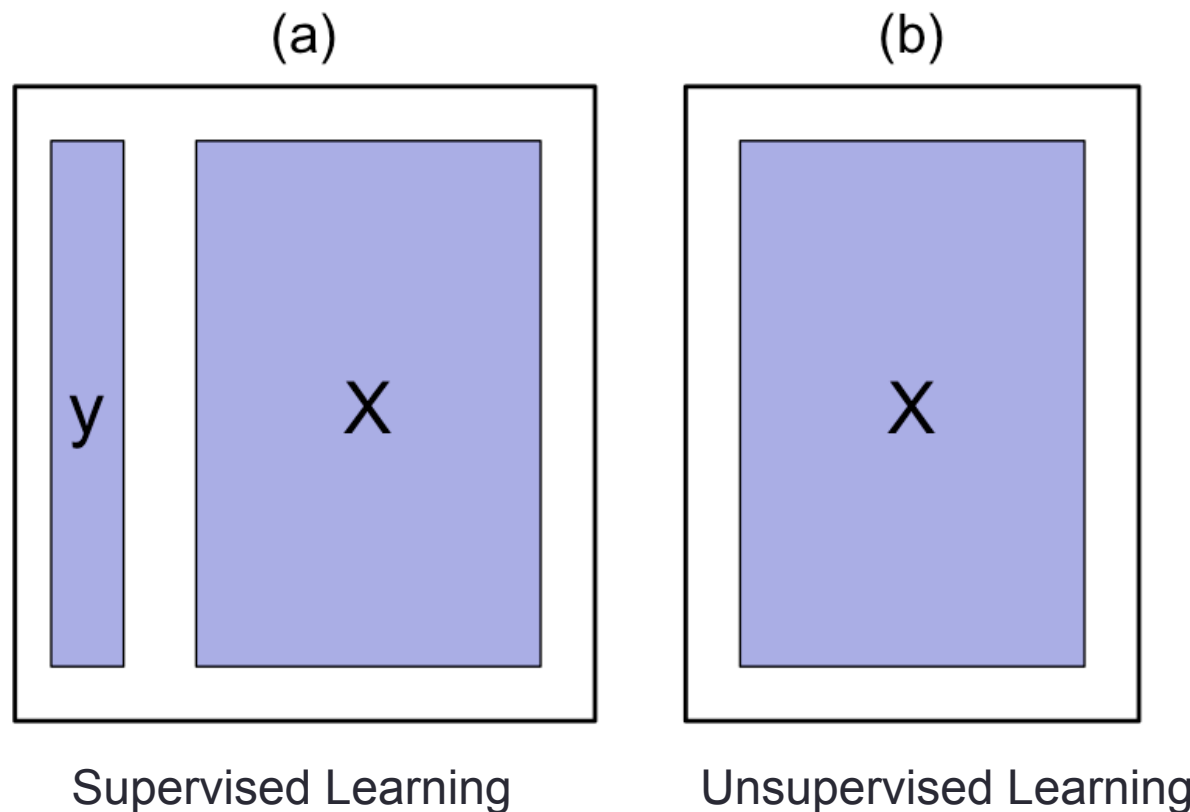# CLUSTERING

Chapter 10

# Outline

➢What is Clustering?

➢K-Means Clustering

➢Hierarchical Clustering

# WHAT IS CLUSTERING?

# Supervised vs. Unsupervised Learning

- Supervised Learning: both X and Y are known
- Unsupervised Learning: only X



Supervised Learning      Unsupervised Learning

# Clustering

- Clustering refers to a set of techniques for finding subgroups, or clusters, in a data set.

- A good clustering is one when the observations within a group are similar but between groups are very different

- For example, suppose we collect $p$ measurements on each of $n$ breast cancer patients. There may be different unknown types of cancer which we could discover by clustering the data
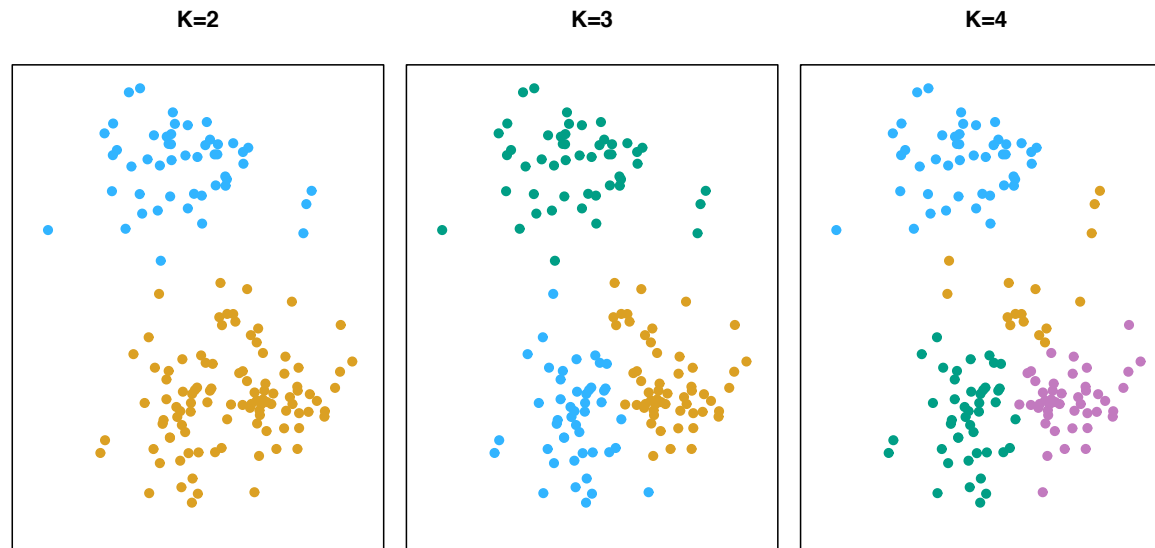
# Different Clustering Methods

- There are many different types of clustering methods

- We will concentrate on two of the most commonly used approaches
  - K-Means Clustering
  - Hierarchical Clustering

# K-MEANS CLUSTERING

# K-Means Clustering

- To perform K-means clustering, one must first specify the desired number of clusters K

- Then the K-means algorithm will assign each observation to exactly one of the K clusters

# How does K-Means work?

- We would like to partition that data set into K clusters
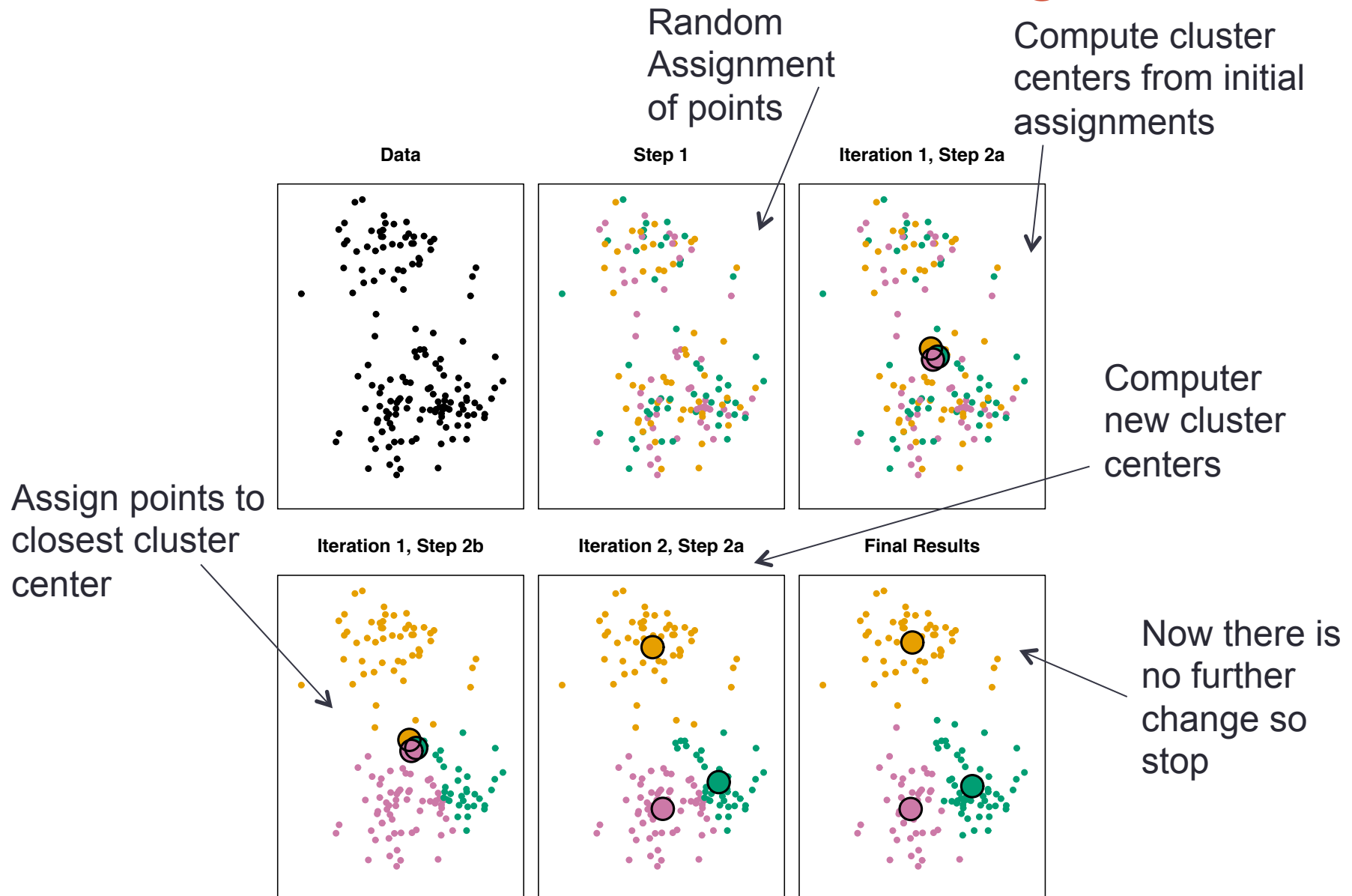
$$C_1, \ldots, C_K$$

  - Each observation belong to at least one of the K clusters
  - The clusters are non-overlapping, i.e. no observation belongs to more than one cluster

- The objective is to have a minimal "within-cluster-variation", i.e. the elements within a cluster should be as similar as possible

- One way of achieving this is to minimize the sum of all the pair-wise squared Euclidean distances between the observations in each cluster.

$$\underset{C_1,\ldots,C_K}{\text{minimize}} \left\{ \sum_{k=1}^{K} \frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^{p} (x_{ij} - x_{i'j})^2 \right\}$$
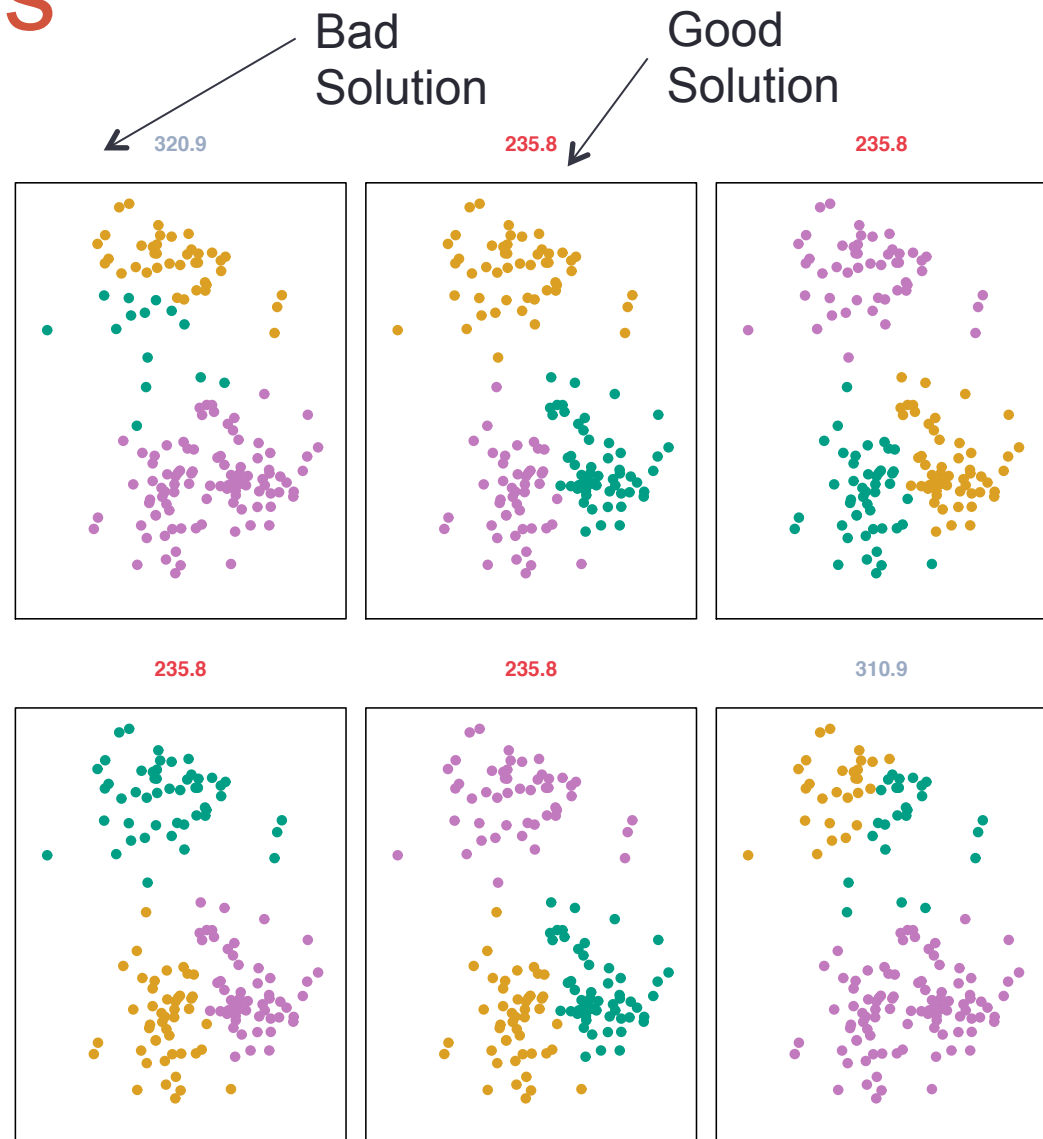
# K-Means Algorithm

- <u>Initial Step</u>: Randomly assign each observation to one of K clusters

- Iterate until the cluster assignments stop changing:
  - For each of the K clusters, compute the cluster centroid. The $k^{th}$ cluster centroid if the mean of the observations assigned to the $k^{th}$ cluster
  - Assign each observation to the cluster whose centroid is closest (where "closest" is defined using Euclidean distance.

# An Illustration of the K-Means Algorithm



Random Assignment of points

Compute cluster centers from initial assignments

Data

Step 1

Iteration 1, Step 2a

Computer new cluster centers

Assign points to closest cluster center

Iteration 1, Step 2b

Iteration 2, Step 2a

Final Results

Now there is no further change so stop

# Local Optimums

- The K-means algorithm can get stuck in "local optimums" and not find the best solution

- Hence, it is important to run the algorithm multiple times with random starting points to find a good solution
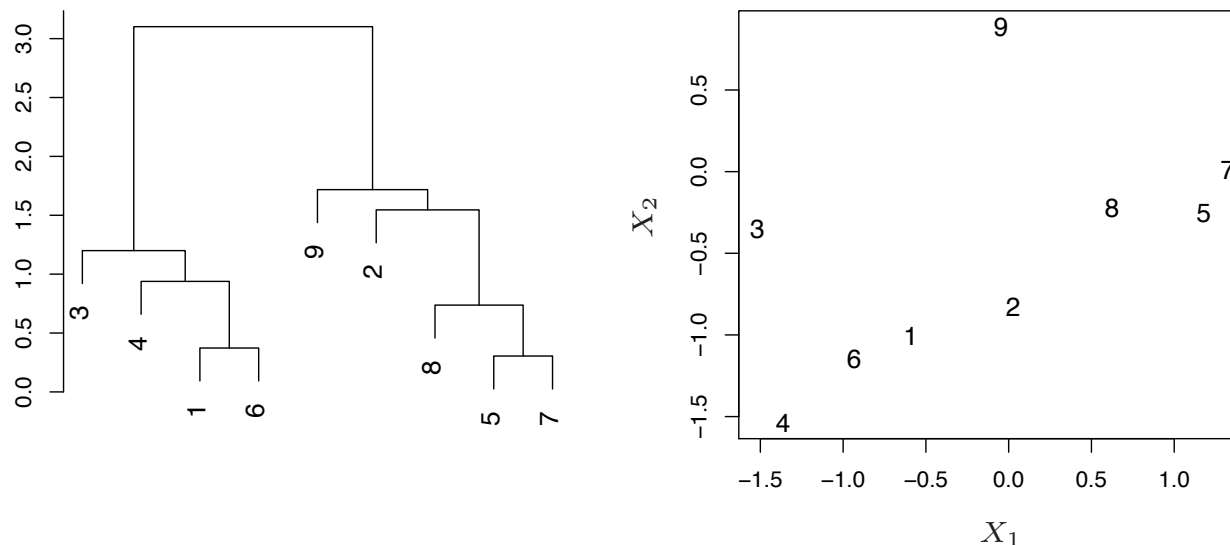
# HIERARCHICAL CLUSTERING

# Hierarchical Clustering

- K-Means clustering requires choosing the number of clusters.

- If we don't want to do that, an alternative is to use Hierarchical Clustering

- Hierarchical Clustering has an added advantage that it produces a tree based representation of the observations, called a Dendogram
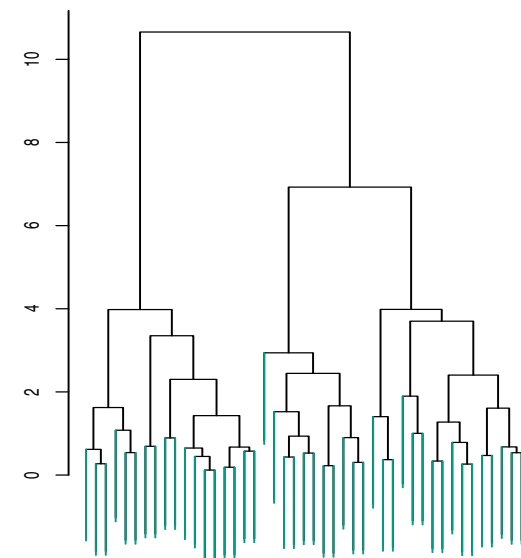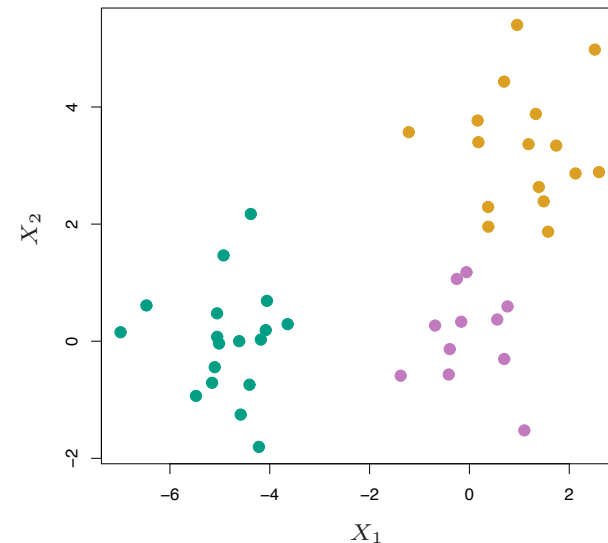
# Dendograms

- First join closest points (5 and 7)
- Height of fusing/merging  (on vertical axis) indicates how similar the points are
- After the points are fused they are treated as a single observation and the algorithm continues
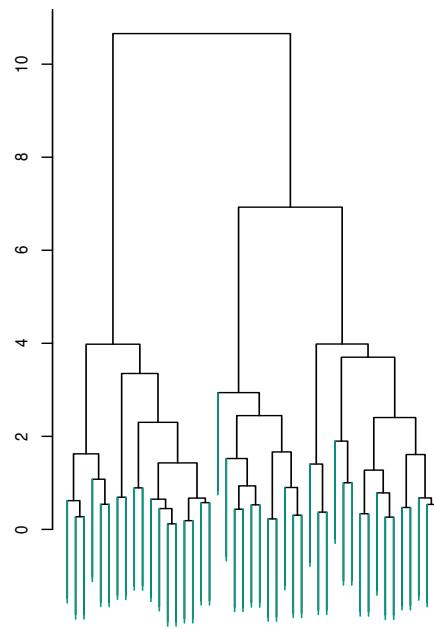
# Interpretation

- Each "leaf" of the dendogram represents one of the 45 observations

- At the bottom of the dendogram, each observation is a distinct leaf. However, as we move up the tree, some leaves begin to fuse. These correspond to observations that are similar to each other.

- As we move higher up the tree, an increasing number of observations have fused. The earlier (lower in the tree) two observations fuse, the more similar they are to each other.

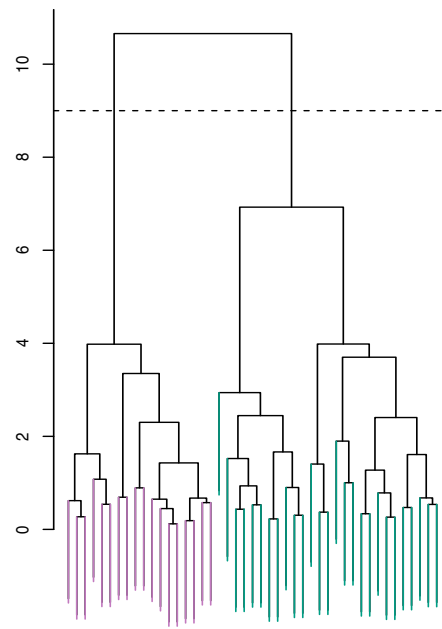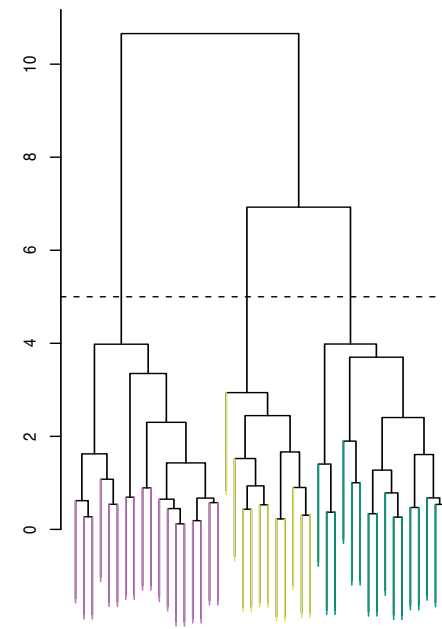- Observations that fuse later are quite different

# Choosing Clusters

- To choose clusters we draw lines across the dendogram
- We can form any number of clusters depending on where we draw the break point.
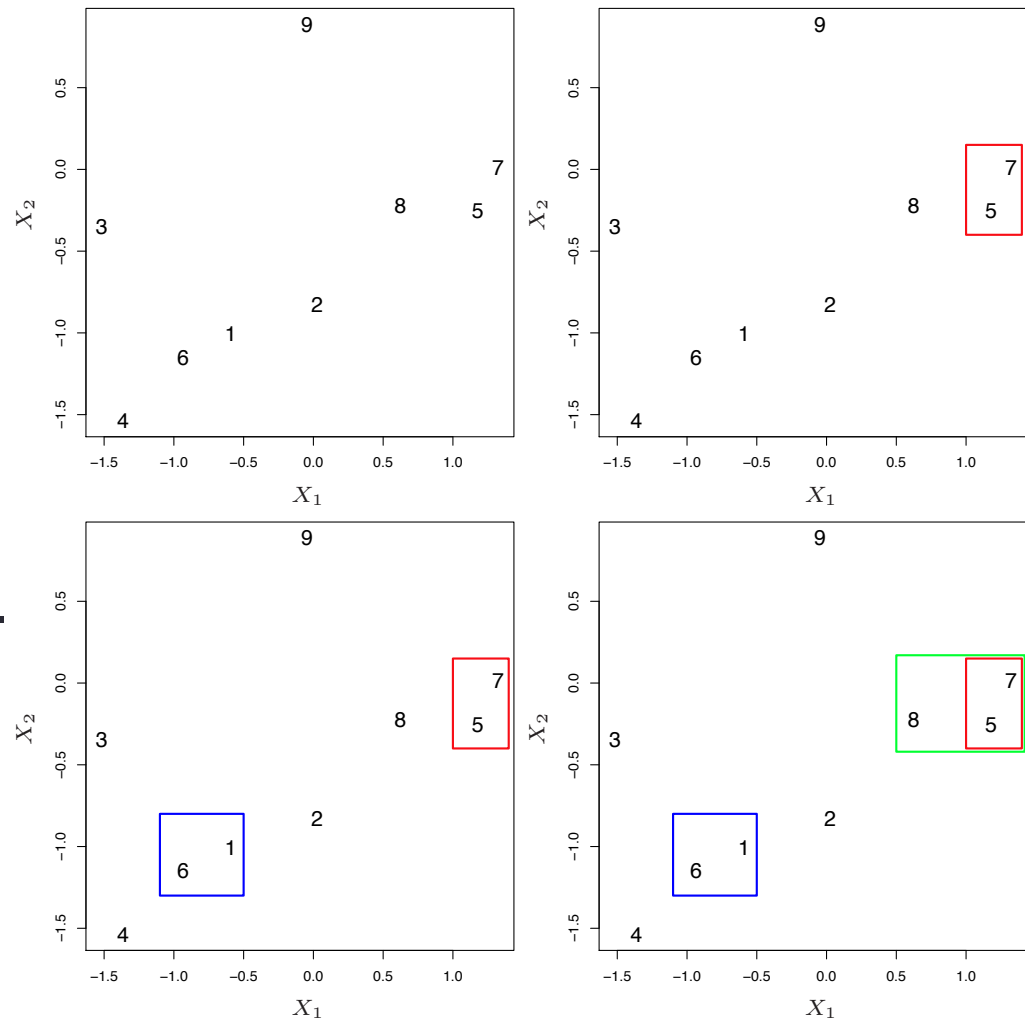


One Cluster          Two Clusters          Three Clusters

# Algorithm (Agglomerative Approach)

- The dendogram is produced as follows:
  - Start with each point as a separate cluster (n clusters)

  - Calculate a measure of dissimilarity between all points/ clusters

  - Fuse two clusters that are most similar so that there are now n-1 clusters

  - Fuse next two most similar clusters so there are now n-2 clusters

  - Continue until there is only 1 cluster

# An Example

- Start with 9 clusters
- Fuse 5 and 7
- Fuse 6 and 1
- Fuse the (5,7) cluster with 8.
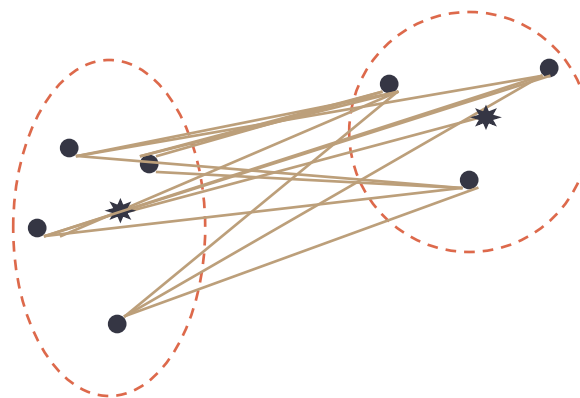- Continue until all observations are fused.

# How do we define dissimilarity?

- Implementing hierarchical clustering involves one obvious issue

- How do we define the dissimilarity, or linkage, between the fused (5,7) cluster and 8?

- There are four options:
  - Complete Linkage
  - Single Linkage
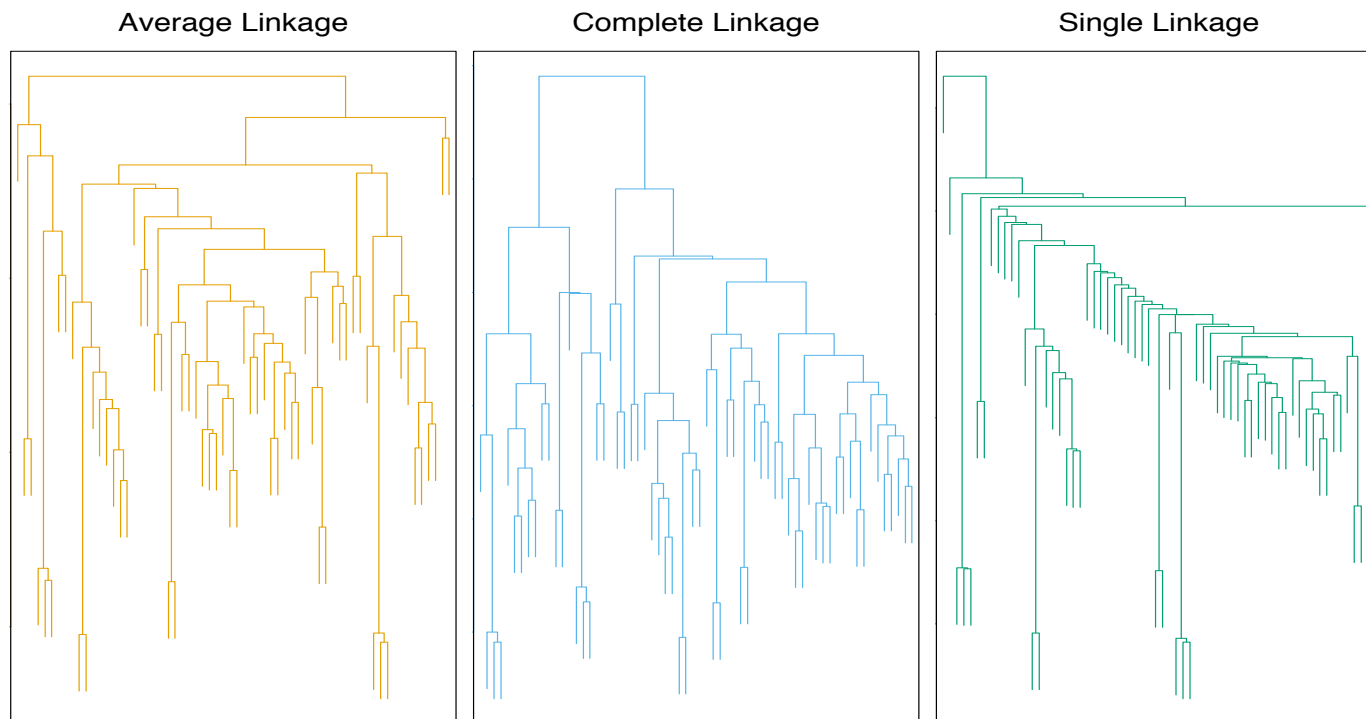  - Average Linkage
  - Centriod Linkage

# Linkage Methods: Distance Between Clusters

- **<u>Complete Linkage:</u>** Largest distance between observations

- **<u>Single Linkage:</u>** Smallest distance between observations

- **<u>Average Linkage</u>:** Average distance between observations

- **<u>Centroid:</u>** distance between centroids of the observations

# Linkage Can be Important

- Here we have three clustering results for the same data
- The only difference is the linkage method but the results are very different
- Complete and average linkage tend to yield evenly sized clusters whereas single linkage tends to yield extended clusters to which single leaves are fused one by one.



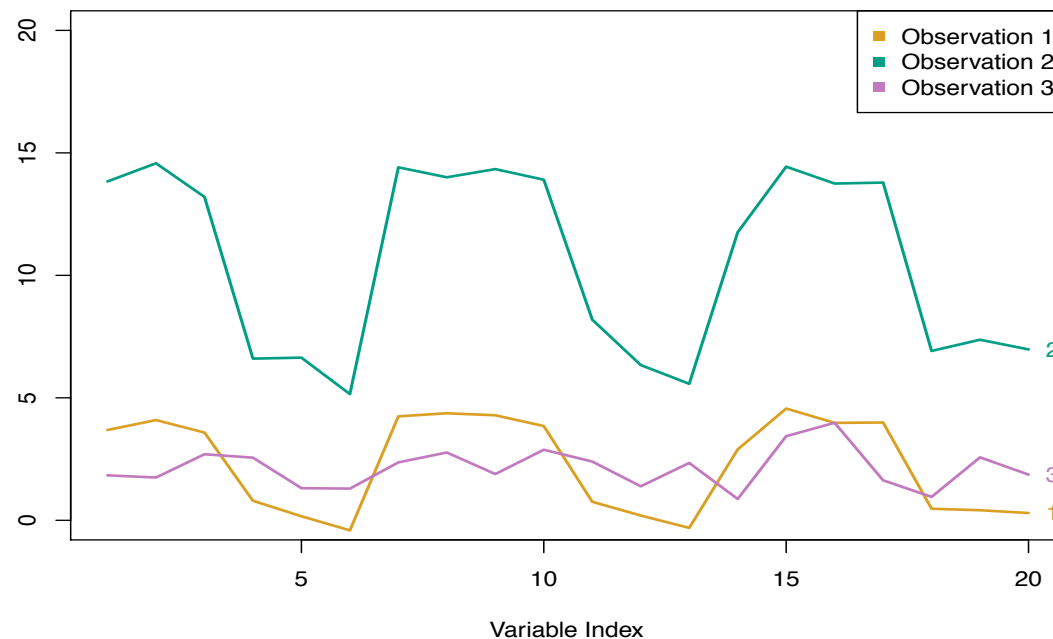Average Linkage          Complete Linkage          Single Linkage

# Choice of Dissimilarity Measure

- So far, we have considered using Euclidean distance as the dissimilarity measure

- However, an alternative measure that could make sense in some cases is the correlation based distance

# Comparing Dissimilarity Measures

- In this example, we have 3 observations and p = 20 variables

- In terms of Euclidean distance obs. 1 and 3 are similar

- However, obs. 1 and 2 are highly correlated so would be considered similar in terms of correlation measure
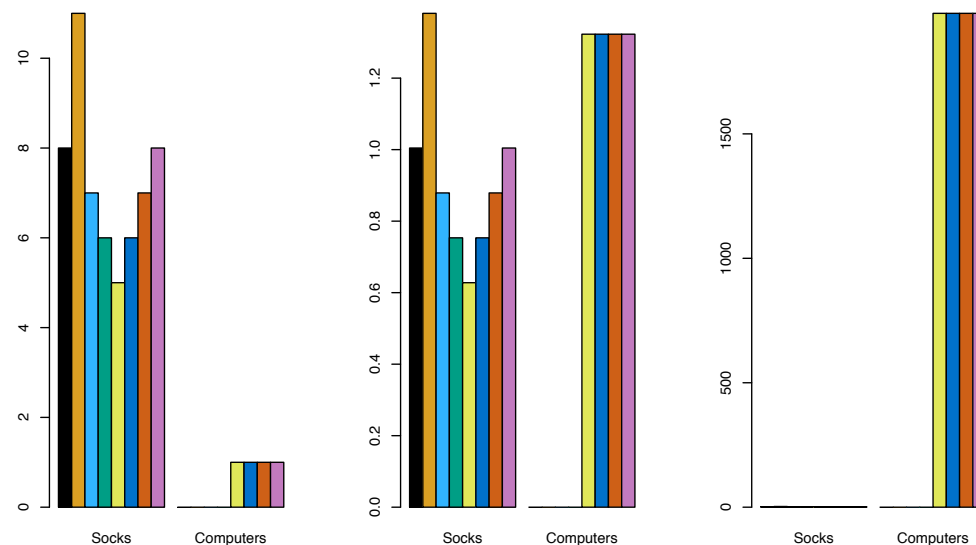
# Online Shopping Example

- Suppose we record the number of purchases of each item (columns) for each customer (rows)

- Using Euclidean distance, customers who have purchases very little will be clustered together

- Using correlation measure, customers who tend to purchase the same types of products will be clustered together even if the magnitude of their purchase may be quite different

# Standardizing the Variables

- Consider an online shop that sells two items: socks and computers
  - Left: In terms of quantity, socks have higher weight
  - Center: After standardizing, socks and computers have equal weight
  - Right: In terms of dollar sales, computers have higher weight

# FINAL THOUGHTS

# Practical Issues in Clustering

- In order to perform clustering, some decisions must be made:
  - Should the features first be standardized? i.e. Have the variables centered to have a mean of zero and standard deviation of one.
  - In case of hierarchical clustering:
    - What dissimilarity measure should be used?
    - What type of linkage should be used?
    - Where should we cut the dendogram in order to obtain clusters?
  - In case of K-means clustering:
    - How many clusters should we look for the data?
- In practice, we try several different choices, and look for the one with the most useful or interpretable solution. There is no single right answer!

# Final Thoughts

- Most importantly, one must be careful about how the results of a clustering analysis are reported

- These results should not be taken as the absolute truth about a data set

- Rather, they should constitute a starting point for the developments of a scientific hypothesis and further study, preferably on independent data

# Exercise

- Suppose that we have 5 observations, for which we compute a similarity (distance) matrix as follows:

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| A | 0 |   |   |   |   |
| B | 9 | 0 |   |   |   |
| C | 3 | 7 | 0 |   |   |
| D | 6 | 5 | 9 | 0 |   |
| E | 11 | 10 | 2 | 8 | 0 |

- On the basis of the similarity matrix, sketch the dendogram that results from hierarchically clustering these 5 observations using complete linkage.