



MONASH University

ETC3250

Business Analytics

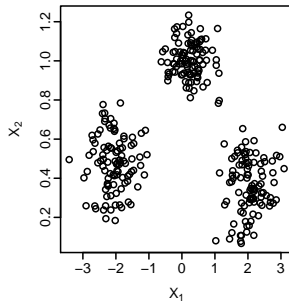
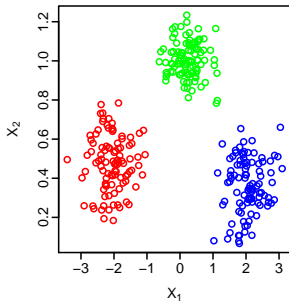
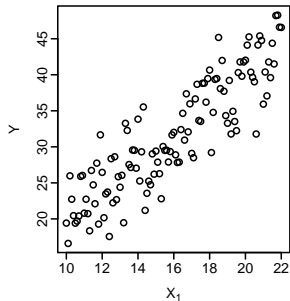
Week 5
Clustering

26 March 2018

Outline

Week	Topic	Chapter	Lecturer
1	Introduction	1	Souhaib
2	Statistical learning	2	Souhaib
3	Regression	3	Souhaib
4	Classification	4	Souhaib
5	Clustering	10	Souhaib
Semester break			
6	Model selection and resampling methods	5	Souhaib
7	Dimension reduction	6,10	Souhaib
8	Advanced regression	6	Souhaib
9	Advanced regression	6	Souhaib
10	Advanced classification	9	Souhaib
11	Tree-based methods	8	Souhaib
12	Project presentation		Souhaib

Statistical learning problems



Unsupervised learning

- **Unsupervised learning** is often performed as part of an **exploratory data analysis**.
- Unsupervised learning is often much **more challenging than supervised learning**. The exercise tends to be more **subjective**, and there is no simple goal for the analysis, such as prediction of a response
- **Hard to assess** the results obtained from unsupervised learning methods
- Techniques for unsupervised learning are of **growing importance** in a number of fields

Unsupervised learning methods

Both **PCA** and **clustering** seek to simplify the data via a small number of summaries, but their mechanisms are different:

- **PCA** (unsupervised dimension reduction method) looks to find a **low-dimensional representation** of the observations that explain a large fraction of the variance.
- **Clustering** looks to find **homogeneous subgroups** among the observations.

Clustering methods

■ K-means clustering

- We seek to partition the observations into K ($K \leq n$) clusters.

■ Hierarchical clustering

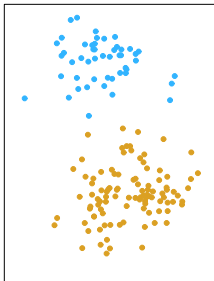
- We do not know in advance how many clusters we want. We consider all possible number of clusters, from 1 to n .

K-means clustering

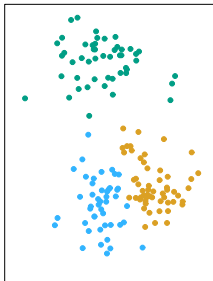
Find K clusters C_1, \dots, C_K where

- $C_1 \cup C_2 \cup \dots \cup C_K = \{1, \dots, n\}$
- $C_k \cap C_{k'} = \emptyset$ for all $k \neq k'$

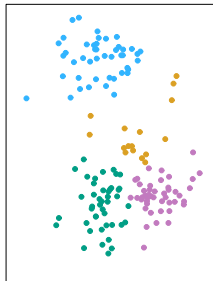
K=2



K=3



K=4



K-means clustering

For K-means, good clustering means small **total within-cluster variation**:

$$\underset{C_1, \dots, C_K}{\text{minimize}} \left\{ \sum_{k=1}^K W(C_k) \right\}$$

where $W(C_k)$ is the within-cluster variation for cluster C_k , i.e. the amount by which the *observations within a cluster* differ from each other.

K-means within-cluster variation

There are many possible ways to define the **within-cluster variation**, but by far the most common choice involves **squared Euclidean distance**:

$$W(C_k) = \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2$$

K-means optimization problem

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

The number of possible assignments of n data points into K clusters is

$$S(n, K) = \frac{1}{K!} \sum_{k=1}^K (-1)^{K-k} \binom{K}{k} k^n \approx K^n,$$

which is the *Stirling numbers of the second kind*.

- $S(n, K)$ is a huge number unless K and n are tiny.
- Fortunately, a very simple algorithm can be shown to provide a **local optimum**—a pretty good solution—to the K-means optimization problem.

K-means optimization problem

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

The number of possible assignments of n data points into K clusters is

$$S(n, K) = \frac{1}{K!} \sum_{k=1}^K (-1)^{K-k} \binom{K}{k} k^n \approx K^n,$$

which is the *Stirling numbers of the second kind*.

- $S(n, K)$ is a huge number unless K and n are tiny.
- Fortunately, a very simple algorithm can be shown to provide a **local optimum**—a pretty good solution—to the K-means optimization problem.

Rewriting the objective

$$\frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = 2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2$$

where

$$\bar{x}_{kj} = \frac{1}{|C_k|} \sum_{i \in C_k} x_{ij}$$

$$\Rightarrow \underset{C_1, \dots, C_K}{\text{minimize}} \left\{ \sum_{k=1}^K \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2 \right\}$$

$$\equiv \underset{C_1, \dots, C_K}{\text{minimize}} \left\{ \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \bar{x}_k\|_2^2 \right\}$$

K-means optimization problem

For any $x_1, \dots, x_m \in \mathbb{R}^p$, the quantity $\sum_{i=1}^m \|x_i - c\|_2^2$ is minimized by $c = \bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$.

So our problem

$$\underset{C_1, \dots, C_K}{\text{minimize}} \left\{ \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \bar{x}_k\|_2^2 \right\}$$

is the same as the following enlarged criterion:

$$\underset{C_1, \dots, C_K, \mathbf{s}_1, \dots, \mathbf{s}_K}{\text{minimize}} \left\{ \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mathbf{s}_k\|_2^2 \right\}.$$

The K-means clustering algorithm **approximately** minimizes the previous objective by **alternately minimizing** over C_1, \dots, C_K and $\mathbf{s}_1, \dots, \mathbf{s}_K$.

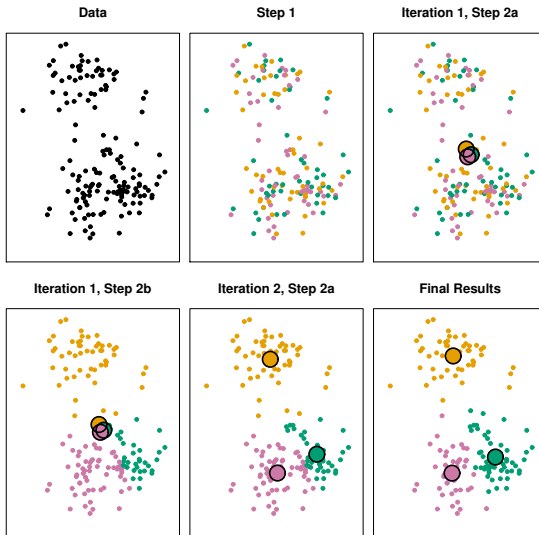
K-means optimization problem

Algorithm 10.1 *K-Means Clustering*

1. Randomly assign a number, from 1 to K , to each of the observations. These serve as initial cluster assignments for the observations.
 2. Iterate until the cluster assignments stop changing:
 - (a) For each of the K clusters, compute the cluster *centroid*. The k th cluster centroid is the vector of the p feature means for the observations in the k th cluster.
 - (b) Assign each observation to the cluster whose centroid is closest (where *closest* is defined using Euclidean distance).
-

The K-means algorithm is **guaranteed to decrease the value of the objective for each iteration.**

Example



Random initial clusters



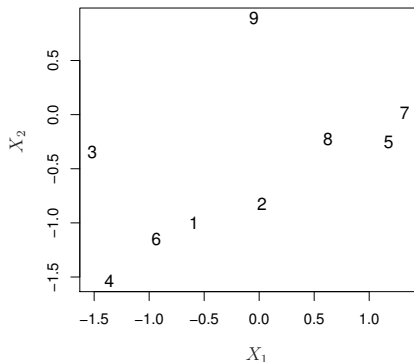
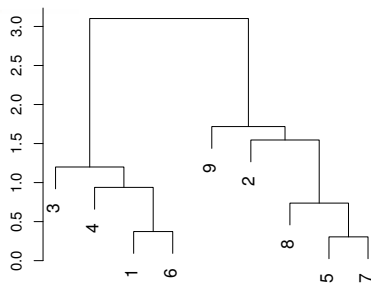
Practical issues with K-means

- Should we standardise the data?
- How many clusters should we use?
 - Objective function is minimal with $K = n$.
- How far from optimal solution?
 - Arbitrarily far with random initial clusters.
 - The clustering quality depends heavily on initial clusters.
 - Better approach: K-means++.

Hierarchical clustering

- One potential disadvantage of **K-means clustering** is that it requires us to **pre-specify the number of clusters** K .
- **Hierarchical clustering** is an alternative approach which **does not require** that we commit to a particular choice of K
- Hierarchical clustering has an added advantage over K-means clustering in that it results in an attractive **tree-based representation** of the observations, called a **dendrogram**.

Dendrogram



- Each **leaf** of the dendrogram represents one observation
- Leaves **fuse** into branches and branches **fuse**, either with leaves or other branches.
- Fusions **lower in the tree** means the groups of observations are more similar to each other.

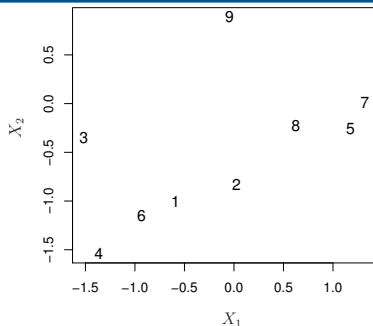
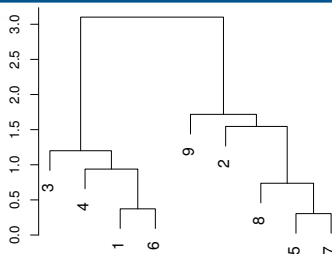
Dendrogram interpretation

- Observations **5 and 7**, and **1 and 6** are quite similar to each other
- Observations **9 and 2** are located near each other on the dendrogram but incorrect to conclude from the figure that they are quite similar
- Actually, observation 9 is **no more similar** to observation 2 than it is to observations 8, 5, and 7.
- There are 2^{n-1} possible reorderings of the dendrogram, where n is the number of leaves. This is because at each of the $n - 1$ points where fusions occur, **the positions of the two fused branches could be swapped without affecting the meaning of the dendrogram.**
- Therefore, we **cannot** draw conclusions about the similarity of two observations based on their proximity along the **horizontal axis**.

Dendrogram interpretation

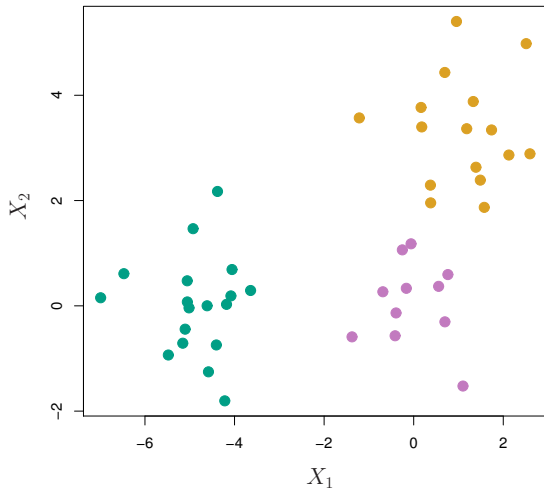
- Rather, we draw conclusions about the similarity of two observations based on the location on the **vertical axis** where branches containing those two observations first are fused.
- For any two observations, we can look for the point in the tree where branches containing those two observations are first fused. **The height of this fusion**, as measured on the **vertical axis**, indicates **how different** the two observations are.
- Thus, observations that fuse at the very **bottom** of the tree are **quite similar** to each other, whereas observations that fuse **close to the top** of the tree will tend to be **quite different**

Clustering with a dendrogram

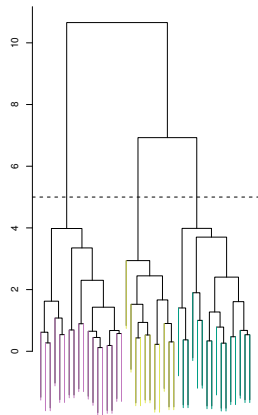
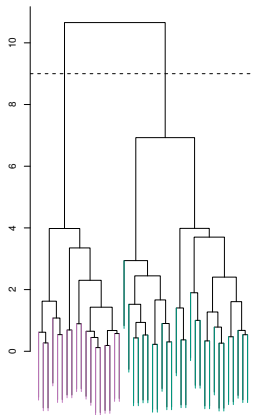
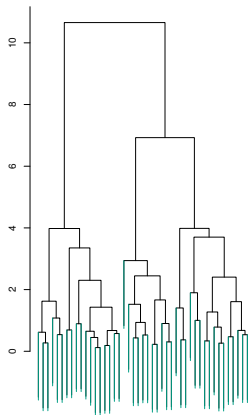


- The distinct sets of observations **beneath a cut** can be interpreted as clusters
- How many clusters if we cut at a height of two or one?
- Between 1 (no cut) and n (a cut at height 0) clusters.
- The **height of the cut** controls the number of clusters obtained as K in K-means clustering.

Example



Example



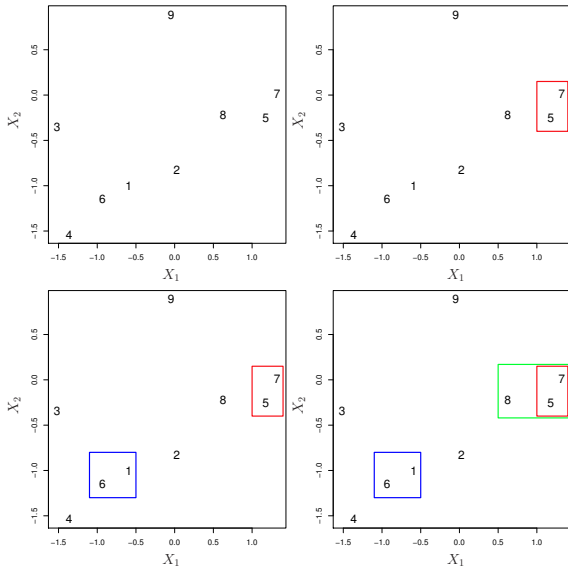
Challenges

- Often the choice of **where to cut** the dendrogram **is not so clear**. In practice, people select **by eye** a sensible number of clusters, based on the heights of the fusion and the number of clusters desired.
- Hierarchical clustering can provide bad clustering when the true clusters are **not nested**
- **Example:** a group of people with a 50–50 split of males and females, evenly split among Americans, Japanese, and French.
- Best division into **two groups** might split these people by gender, and the best division into **three groups** might split them by nationality.
 - The best division into three groups does not result from taking the best division into two groups and splitting up one of those groups (**not nested**).

Hierarchical clustering algorithm

- 1 We begin by defining some sort of **dissimilarity measure** between each pair of observations. Most often, Euclidean distance is used
- 2 Each of the n observations is treated as **its own cluster**
- 3 The two clusters that are most similar to each other are then fused so that there now are $n - 1$ clusters
- 4 Next the two clusters that are most similar to each other are fused again, so that there now are $n - 2$ clusters
- 5 The algorithm proceeds in this fashion **until all of the observations belong to one single cluster**, and the dendrogram is complete

Hierarchical clustering



Linkage

We have a concept of the **dissimilarity between pairs of observations**, but how do we define the **dissimilarity between two clusters** if one or both of the clusters contains multiple observations?

→ The concept of dissimilarity between a pair of observations needs to be extended to a pair of groups of observations. We introduce the concept of *linkage*.

Linkage

We have a concept of the **dissimilarity between pairs of observations**, but how do we define the **dissimilarity between two clusters** if one or both of the clusters contains multiple observations?

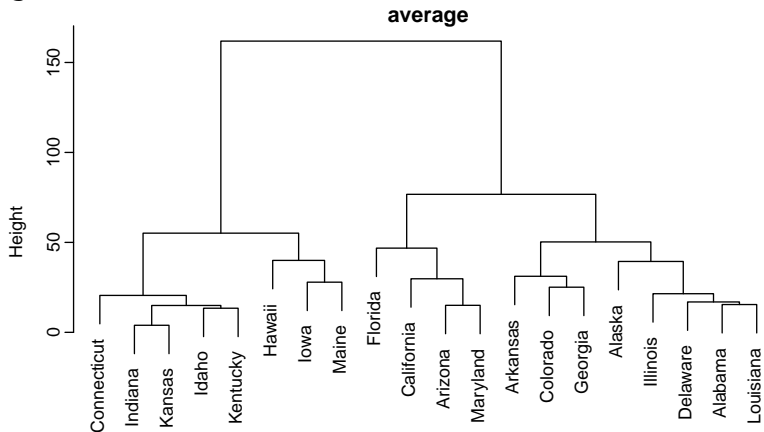
→ The concept of dissimilarity between a pair of observations needs to be extended to a pair of groups of observations. We introduce the concept of *linkage*.

Linkage

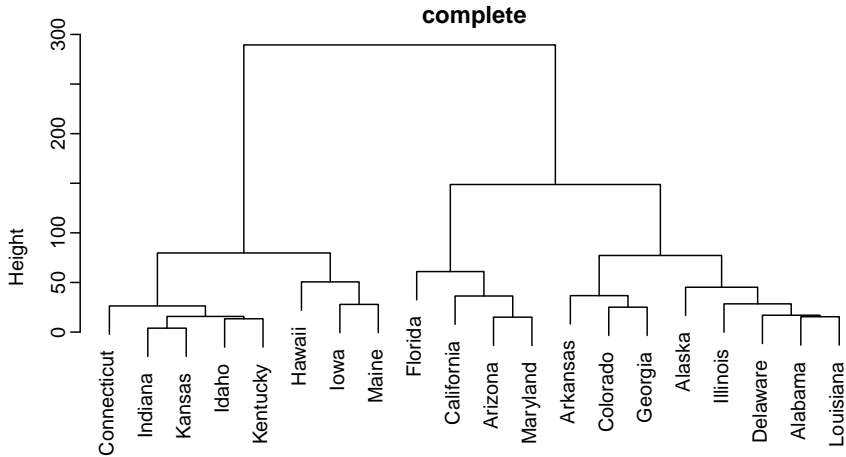
<i>Linkage</i>	<i>Description</i>
Complete	Maximal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>largest</i> of these dissimilarities.
Single	Minimal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>smallest</i> of these dissimilarities. Single linkage can result in extended, trailing clusters in which single observations are fused one-at-a-time.
Average	Mean intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>average</i> of these dissimilarities.
Centroid	Dissimilarity between the centroid for cluster A (a mean vector of length p) and the centroid for cluster B. Centroid linkage can result in undesirable <i>inversions</i> .

Example

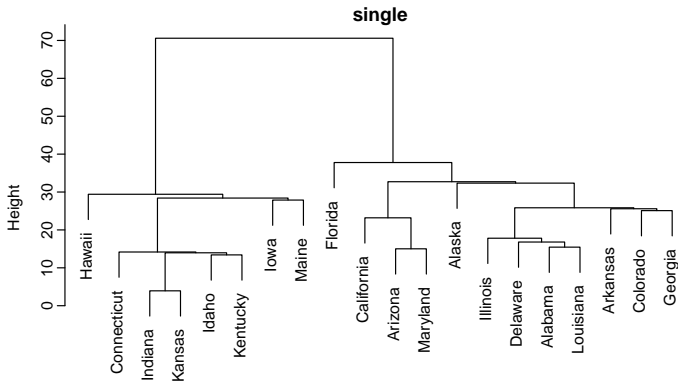
The data set contains statistics, in arrests per 100,000 residents for assault, murder, and rape in each of the 50 US states in 1973. Also given is the percent of the population living in urban areas.



Example

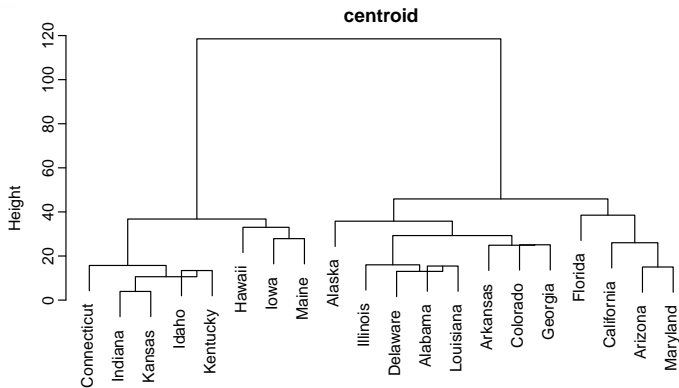


Example



Average and complete linkage are generally preferred over **single** linkage, as they tend to yield more balanced dendrograms.

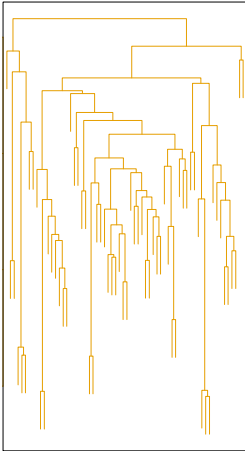
Example



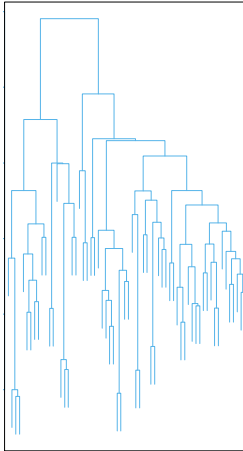
Centroid linkage is often used in genomics, but suffers from *inversions*: two clusters are fused at a height below either of the individual clusters in the dendrogram. This can lead to difficulties in visualization as well as in interpretation of the dendrogram.

Linkage

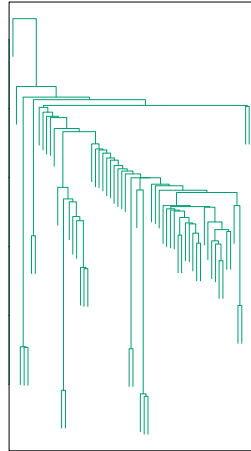
Average Linkage



Complete Linkage



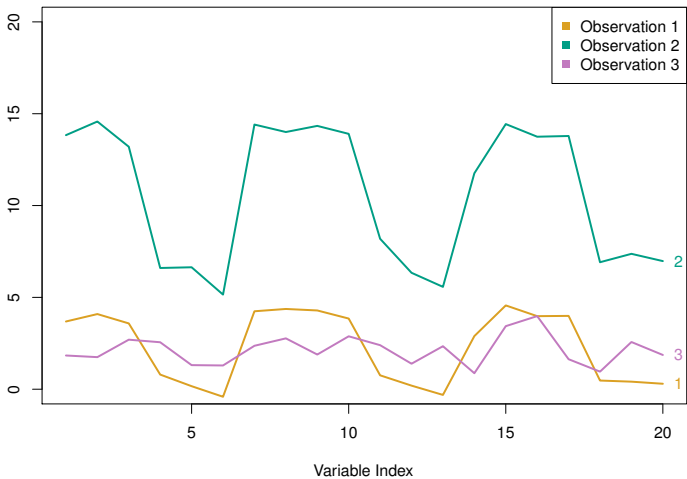
Single Linkage



Dissimilarity measure

- The choice of dissimilarity measure is very important, as it has a strong effect on the resulting dendrogram.
- In general, careful attention should be paid to the **type of data** being clustered and the **scientific question** at hand.

Dissimilarity measure



Hierarchical clustering algorithm

Algorithm 10.2 *Hierarchical Clustering*

1. Begin with n observations and a measure (such as Euclidean distance) of all the $\binom{n}{2} = n(n-1)/2$ pairwise dissimilarities. Treat each observation as its own cluster.
 2. For $i = n, n-1, \dots, 2$:
 - (a) Examine all pairwise inter-cluster dissimilarities among the i clusters and identify the pair of clusters that are least dissimilar (that is, most similar). Fuse these two clusters. The dissimilarity between these two clusters indicates the height in the dendrogram at which the fusion should be placed.
 - (b) Compute the new pairwise inter-cluster dissimilarities among the $i-1$ remaining clusters.
-

Practical issues in clustering

- Should we standardise the data?
- How many clusters should we use?
- For hierarchical clustering:
 - Which dissimilarity measure?
 - What type of linkage?
 - Where should we cut the dendrogram?
- Other considerations
 - Soft clustering (e.g. using mixture models)
 - Clustering high-dimensional data