

Cluster Analysis

- **CLUSTERING ANALYSIS** IS ABOUT **DISCOVERING GROUPINGS IN THE DATA**.
- **CLUSTERING METHODS** ARE **UNSUPERVISED METHODS**.
 - * GIVE AN **UNLABELED** DATASET
 - * GROUP DATA INTO **CLUSTERS**
- PART OF THE ANALYSIS IS \rightarrow **CHOOSING K** WITH K BEING **# OF CLUSTERS**.
- **VISUALIZING** DATA SOMETIMES CAN BE **MORE EFFECTIVE**.

Clustering Methods

- WE SEEK TO PARTITION DATA INTO A **# OF CLUSTERS** $\rightarrow C_1, \dots, C_k$.
- **ASSUMPTION**: EACH **DATAPoint** CAN **ONLY** BELONG TO A **SINGLE CLUSTER**.
- DATAPoints IN THE **SAME CLUSTER** ARE **SIMILAR** & DATA POINTS IN \neq CLUSTERS ARE **DISSIMILAR!**
- METHODS TO COVER:
 - * **K-MEANS CLUSTERING**
 - * **HIERARCHICAL CLUSTERING**

K-MEAN CLUSTERING

- **First** CHOOSE A **WT (K)**
- **PARTITION DATA** INTO **K CLUSTERS** C_1, \dots, C_k
- All clusters ARE **NON-EMPTY** & AN DATA POINTS BELONGS TO **EXACTLY ONE CLUSTER!**
- * SEEK A **PARTITIONING** THAT **MINIMIZES** **TOTAL WITHIN-CLUSTER VARIATION**:

$$\text{minimize}_{C_1, \dots, C_k} \sum_{k=1}^K W(C_k)$$

- * USUALLY USE **"SIMILARITY"** TO CORRELATE WITH **DISTANCE** SO POPULAR $W(C_k)$ IS:

$$W(C_k) = \frac{1}{|C_k|} \sum_{i: x_i \in C_k} \|x_i - x_i^*\|^2$$

What about Centroids?

- WE DO SOMETHING A BIT DIFFERENT:

- ↳ WE MEASURE THE **WITHIN-CLUSTER VARIATION** AS THE **SUMMED DISTANCE BETWEEN DATA POINTS & CLUSTER CENTROID r_k** :

$$W_{C_k}(r_k) = \sum_{i: x_i \in C_k} \|x_i - r_k\|^2$$

- ↳ IF WE KNOW WHICH OBS BELONG TO A C_k , THEN THE CENTROID IS CHOSEN AS **CLUSTER MEAN**:

$$r_k = \frac{1}{|C_k|} \sum_{i: x_i \in C_k} x_i$$



The K-means algorithm

Initialize the K cluster centres r_1, \dots, r_k .
(e.g. by choosing K random points in the data)

1. Assign all data points to their nearest cluster center.

$$x_i \in C_k \Leftrightarrow \forall j \neq k : \|x_i - r_k\|^2 < \|x_i - r_j\|^2$$

2. Move the cluster center to the mean of the cluster.

$$r_k = \frac{1}{|C_k|} \sum_{i: x_i \in C_k} x_i$$

This is iterated until centres stop moving.

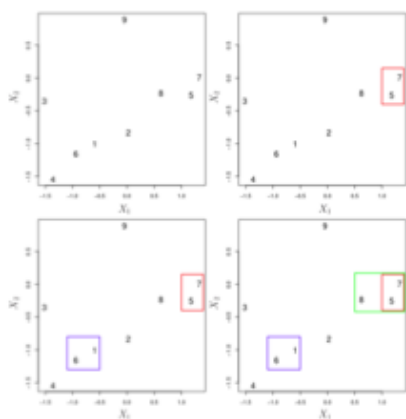
- The algo stops in a finite # of iterations
- We tend to stop when nothing changes.
- The algo can converge to local opt.
- Choice of K highly impacts which data points are clustered together
- Anything that impacts distance, impacts clustering:
 - * High dimensionality.
 - * Outliers.
 - * Scale of variables

Hierarchical Methods

- Have the following property:
 - * Clustering will look similar if you vary the # of clusters by a little.
- All k clusters will be subsets of clusters in a clustering with fewer clusters
- Two obvious strategies:
 - Having k clusters, merge two most similar ones to get $k-1$ clusters.
 - Having k clusters, divide one to get $k+1$ clusters.
- High dissimilarity is best for the clusters.

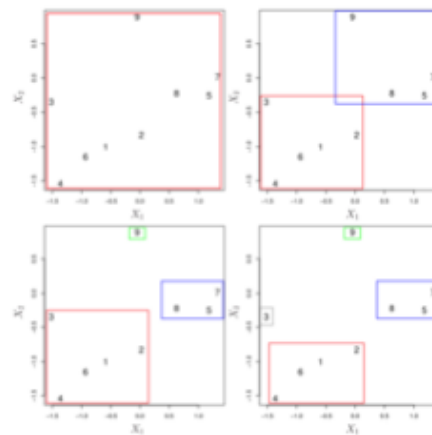
Agglomerative/bottom-up/merging

Agglomerative strategies start with each observation as a separate cluster and repeatedly merge the two clusters with the smallest dissimilarity.



Divisive/top-down/splitting

Divisive strategies work in the opposite direction: they start with all observations in one cluster and recursively split one cluster.



How to measure dissimilarity between clusters?

- Most of the time dissimilarity is based on distance between two data points

Single-Link Clustering

$$D(C_i, C_j) = \min \{d(x, y) \mid x \in C_i, y \in C_j\}$$

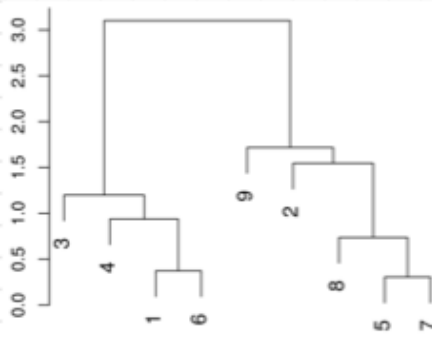
Group-Average Clustering

$$D(C_i, C_j) = \text{mean} \{d(x, y) \mid x \in C_i, y \in C_j\}$$

Complete-Link Clustering

$$D(C_i, C_j) = \max \{d(x, y) \mid x \in C_i, y \in C_j\}$$

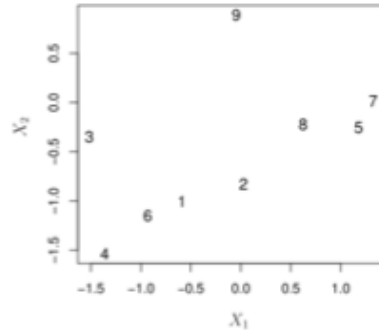
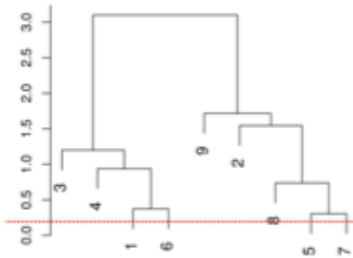
- Hierarchical clustering gives us DENDROGRAMS



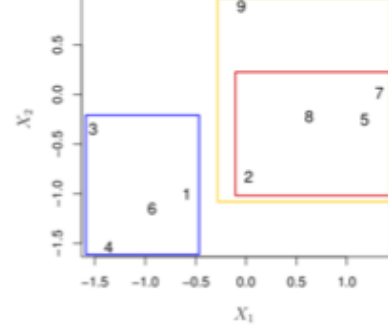
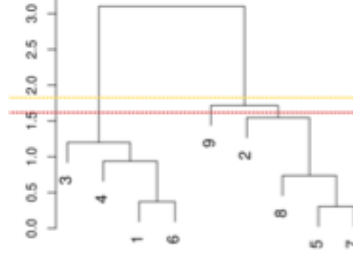
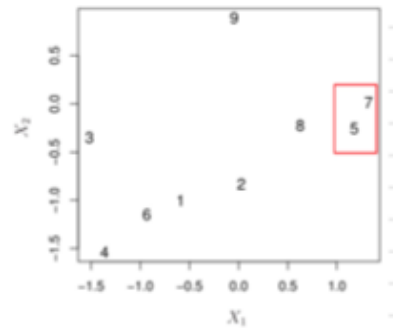
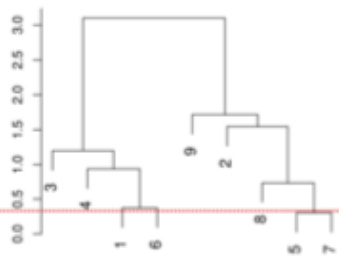
- A DENDROGRAM visualizes PARTITIONS FOR CHOICES OF:
 - $k=1$ (top)
 - $k=n$ (bottom)

- Moving FROM THE BOTTOM TOWARDS THE TOP, PAIRS OF CLUSTERS FUSE.

- A HORIZONTAL CUT IN A DENDROGRAM CORRESPONDS TO A PARTITIONING



IF A CUT INTERSECTS A LINK, ALL NODES "DOWNSTREAM" FROM THAT INTERSECTION ARE GROUPED.



Divisive Strategy!

- Among all clusters & all possible splits of the cluster:
 - Choose comb of cluster & split that gives the GREATEST DISSIMILARITY BETWEEN THE TWO.
- More complex than agglomerative methods cause:
 - DECIDE WHICH TO DIVIDE
 - DECIDE HOW TO DIVIDE

There are:

$$\frac{2^N - 2}{2} = 2^{N-1} - 1$$

ways of splitting a cluster with N obs in two non-empty a.