

# 3 Clustering

Clustering refers to a broad set of techniques for finding *subgroups* in a data set.

We seek to partition observations into distinct groups so that

- observations within a group are similar
  - observations in different groups are dissimilar
- need to define depends on the domain!

For instance, suppose we have a set of  $n$  observations, each with  $p$  features. The  $n$  observations could correspond to tissue samples for patients with breast cancer and the  $p$  features could correspond to *measurements collected for each tissue sample.*

- clinical measurements, e.g. tumor stage or grade.
- gene expression measurements.

We may have reason to believe there is heterogeneity among the  $n$  observations.

e.g. *different unknown subtype of cancer.*

This is *unsupervised* because

We are trying to discover structure (distinct clusters).

This is different from a supervised problem (no goal of prediction).

Both clustering and PCA seek to simplify the data via a small number of summaries.

- PCA - find low dimensional representation of observations that explain a good fraction of variability
- Clustering = find homogeneous subgroups among observations.

Since clustering is popular in many fields, there are many ways to cluster.

we focus on the 2 best-known clustering approaches.

- *K*-means clustering  
we seek to partition the observations into a pre-specified # of clusters.
- Hierarchical clustering  
we do NOT know in advance how many clusters we want.  
we obtain clusterings for  $1, \dots, n$  clusters and view these in a dendrogram.

In general, we can cluster observations on the basis of features or we can cluster features on the basis of observations.

① identify subgroups among observations.

② discover subgroups among features.

we will focus on ①, but can perform ② just by transposing our data matrix.

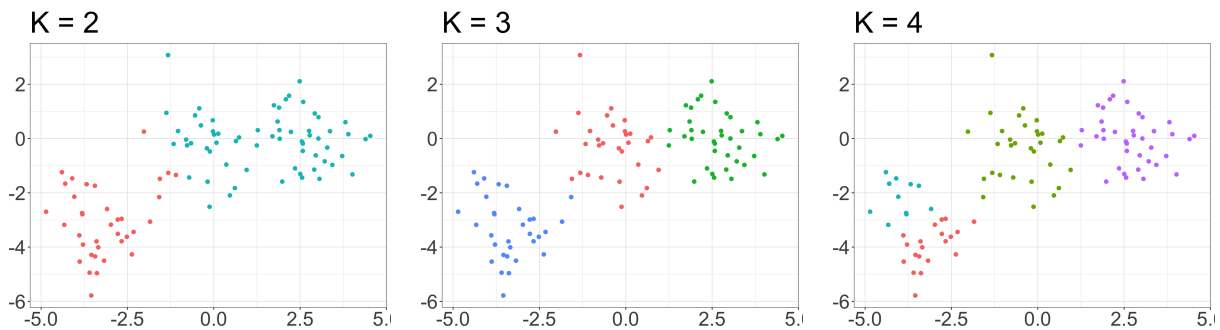
## 3.1 K-Means Clustering

Simple and elegant approach to partition a data set into  $K$  distinct, non-overlapping clusters.

We must first specify how many clusters  $K$ .

then  $K$ -means assigns each observation to one of  $K$  clusters.

e.g. clustering  $n=100$  observations into  $K$  clusters using  $p=2$  features.



The  $K$ -means clustering procedure results from a simple and intuitive mathematical problem. Let  $C_1, \dots, C_K$  denote sets containing the indices of observations in each cluster. These satisfy two properties:

1.

2.

**Idea:**



The *within-cluster variation* for cluster  $C_k$  is a measure of the amount by which the observations within a cluster differ from each other.

To solve this, we need to define within-cluster variation.

This results in the following optimization problem that defines  $K$ -means clustering:

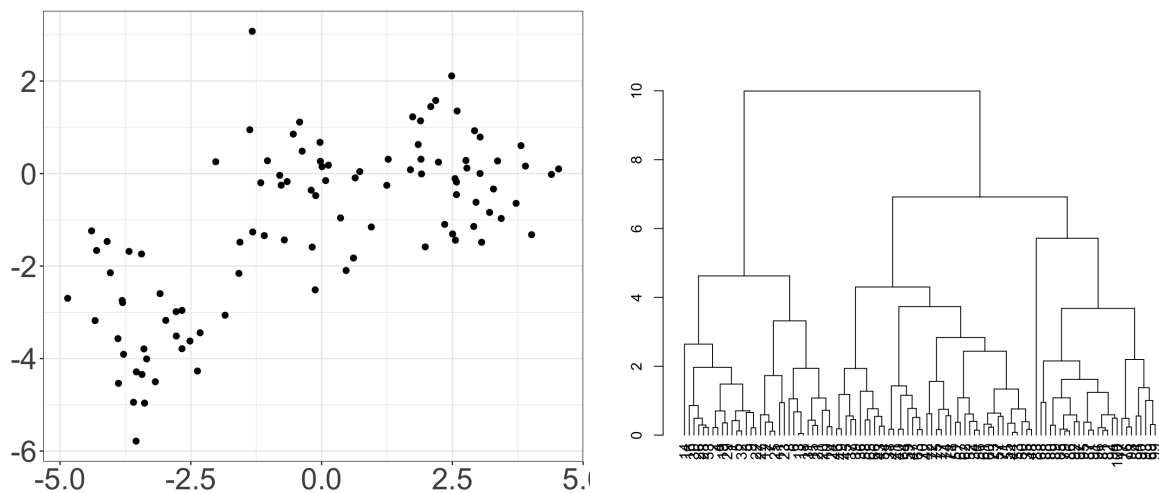
A very simple algorithm has been shown to find a local optimum to this problem:

## 3.2 Hierarchical Clustering

One potential disadvantage of  $K$ -means clustering is that it requires us to specify the number of clusters  $K$ . *Hierarchical clustering* is an alternative that does not require we commit to a particular  $K$ .

We will discuss *bottom-up* or *agglomerative* clustering.

### 3.2.1 Dendrograms

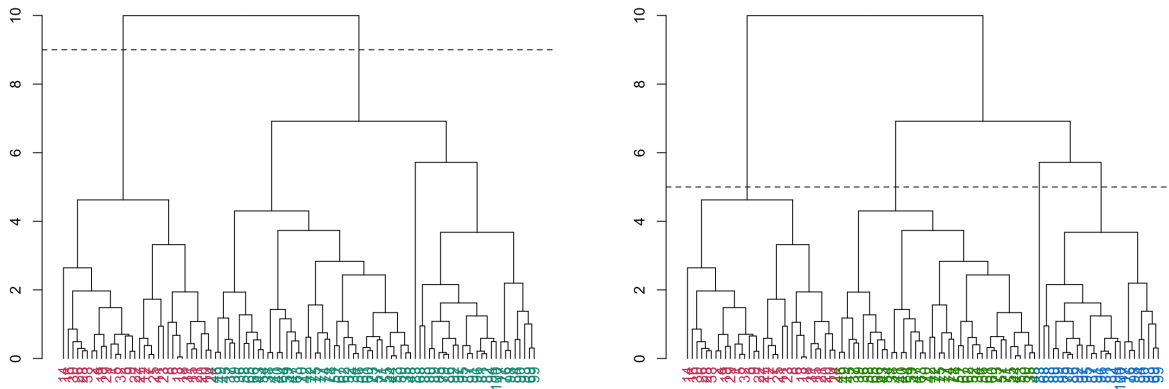


Each *leaf* of the dendrogram represents one of the 100 simulated data points.

As we move up the tree, leaves begin to fuse into branches, which correspond to observations that are similar to each other.

For any two observations, we can look for the point in the tree where branches containing those two observations are first fused.

How do we get clusters from the dendrogram?

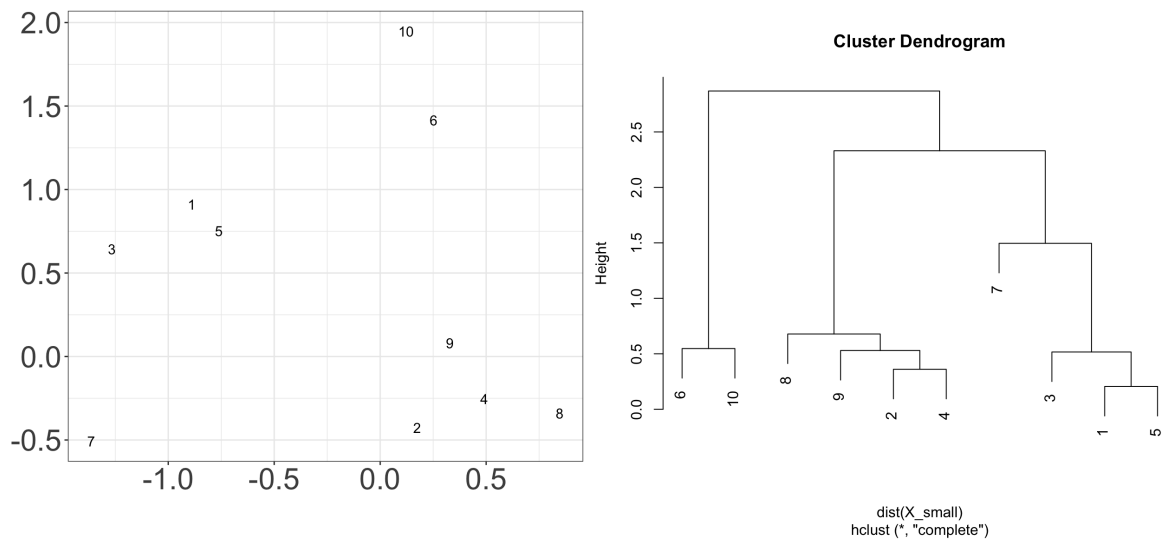


The term *hierarchical* refers to the fact that clusters obtained by cutting the dendrogram at a given height are necessarily nested within the clusters obtained by cutting the dendrogram at a greater height.

### 3.2.2 Algorithm

First, we need to define some sort of *dissimilarity* metric between pairs of observations.

Then the algorithm proceeds iteratively.





More formally,

One issue has not yet been addressed.

How do we determine the dissimilarity between two clusters if one or both of them contains multiple observations?

1.

2.

3.

4.



### 3.2.3 Choice of Dissimilarity Metric

## 3.3 Practical Considerations in Clustering

In order to perform clustering, some decisions should be made.

- 
- 
- 

Each of these decisions can have a strong impact on the results obtained. What to do?