

# Statistical Computing with R: Masters in Data Sciences 503 (S27) Third Batch, SMS, TU, 2024

Shital Bhandary

Associate Professor

Statistics/Bio-statistics, Demography and Public Health Informatics

Patan Academy of Health Sciences, Lalitpur, Nepal

Faculty, Data Analysis and Decision Modeling, MBA, Pokhara University, Nepal

Faculty, FAIMER Fellowship in Health Professions Education, India/USA.

# Review Preview: Unsupervised models

- Clustering
  - K-means clustering
  - Clustering
    - Hierarchical clustering

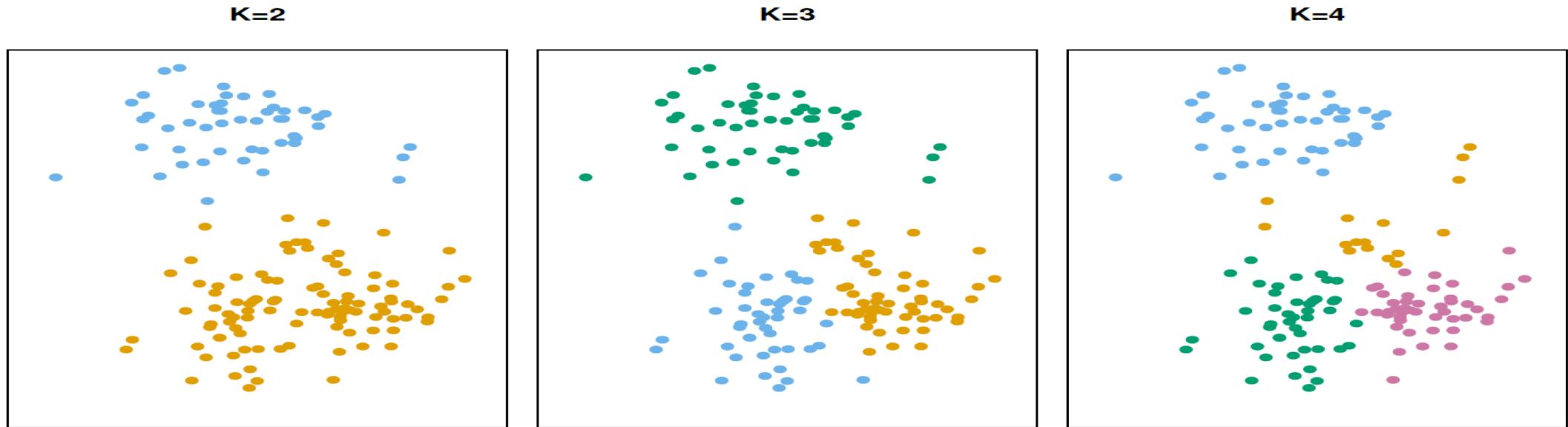
# Cluster analysis (Clustering): Chapter 12 “An Introduction to Statistical Learning” book!

- Clustering refers to a very broad set of techniques for finding subgroups, or clustering clusters, in a data set.
- When we cluster the observations of a data set, we seek to partition them into distinct groups so that the observations within each group are quite similar to each other, while observations in different groups are quite different from each other.
- Both clustering and PCA seek to simplify the data via a small number of summaries, but their mechanisms are different:
- PCA looks to find a low-dimensional representation of the observations that explain a good fraction of the **variance**;
- **Clustering looks to find homogeneous subgroups among the observations or cases.**

# Cluster analysis (Clustering): Chapter 12 “An Introduction to Statistical Learning” book!

- Since clustering is popular in many fields, there exist a great number of clustering methods.
- In this section we focus on perhaps the two best-known clustering approaches:
  - **K-means clustering** and
  - **hierarchical clustering**.
- In K-means clustering, we seek to partition the observations into a pre-specified number of clusters.
- On the other hand, in hierarchical clustering, we do not know in advance how many clusters we want and we use “**dendogram**” to find the number of clusters for the data

# k-means clustering: which “k” is the best?



**FIGURE 12.7.** A simulated data set with 150 observations in two-dimensional space. Panels show the results of applying  $K$ -means clustering with different values of  $K$ , the number of clusters. The color of each observation indicates the cluster to which it was assigned using the  $K$ -means clustering algorithm. Note that there is no ordering of the clusters, so the cluster coloring is arbitrary. These cluster labels were not used in clustering; instead, they are the outputs of the clustering procedure.

# k-means clustering with random data: ISLR

## #ISLR book

- set.seed (2)
- x <- matrix(rnorm (50 \* 2), ncol = 2)
- x[1:25, 1] <- x[1:25, 1] + 3
- x[1:25, 2] <- x[1:25, 2] - 4

#We are creating two group!

## #k-means clustering

- km.out <- kmeans(x, 2, nstart = 20)
- **km.out (check the variance explained!)**

## #Checking the clusters

- km.out\$cluster

#We have used K=2 as we have created a random data with 2 groups

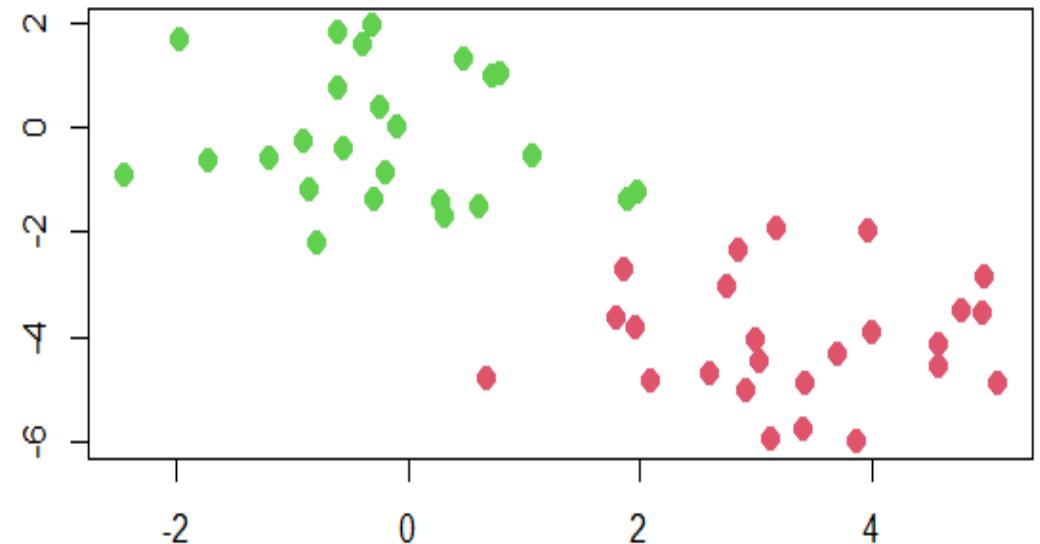
ISLR book: We strongly recommend always running K-means clustering with a large value of nstart, such as 20 or 50, since otherwise an undesirable local optimum may be obtained.

# Plot the clusters:

```
#Plot
```

```
plot(x, col = (km.out$cluster + 1),  
main = "K-Means Clustering  
Results with K = 2",  
xlab = "", ylab = "", pch = 20, cex =  
2)
```

K-Means Clustering Results with K = 2



Let us use k=3 in the data and see what happens:

**#Clustering with 3 clusters in the random data**

- set.seed (4)
- km.out <- kmeans(x, 3, nstart = 20)
- km.out

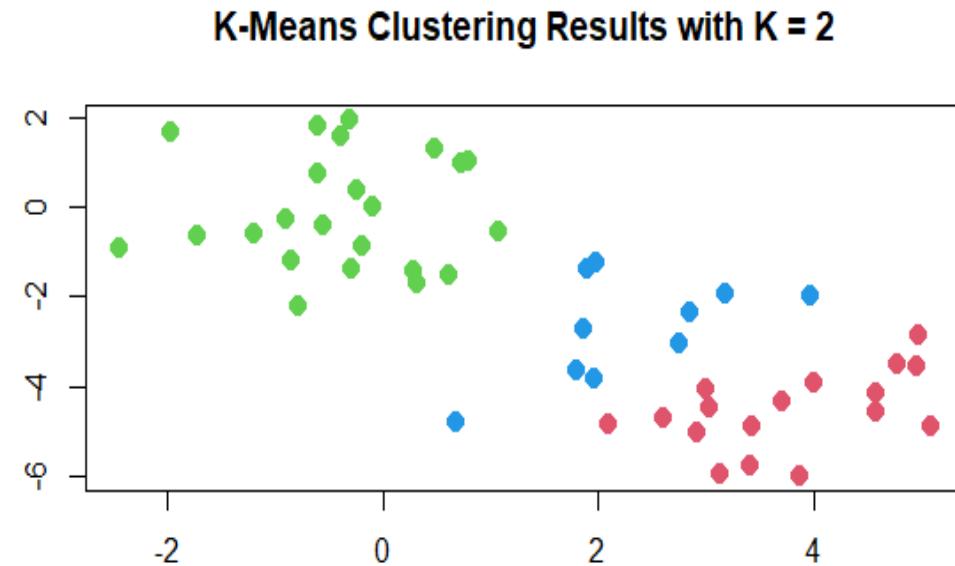
**K-means clustering with 3 clusters of sizes 17, 23, 10**

- Cluster means:
  - [,1] [,2]
  - 1 3.7789567 -4.56200798
  - 2 -0.3820397 -0.08740753
  - 3 2.3001545 -2.69622023
- Within cluster sum of squares by cluster:
  - [1] 25.74089 52.67700 19.56137
  - (between\_SS / total\_SS = 79.3 %)

Plot:

#Plot

```
plot(x, col = (km.out$cluster + 1),  
main = "K-Means Clustering  
Results with K = 3",  
xlab = "", ylab = "", pch = 20, cex =  
2)
```



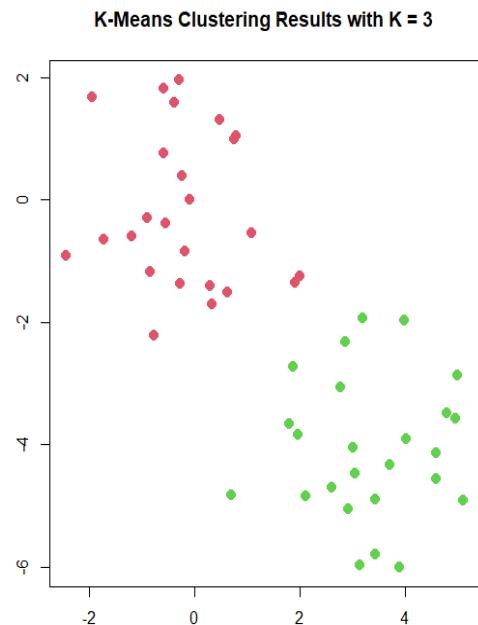
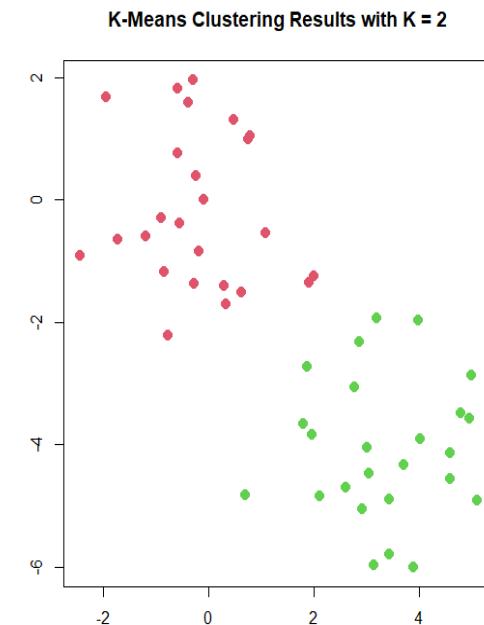
# Comparing two plots: 2-cluster and 3-cluster

## #Plot the clusters

- `par(mfrow = c(1, 2))`
- Code of plot with 2 clusters
- Code of plot with 3 clusters

How to decide: which k is best?

We need to use hierarchical clustering!



# Comparing different nstarts: Plot them and see the changes!

**#nstart = 1**

- set.seed (4)
- km.out <- kmeans(x, 3, nstart = 1)
- km.out\$tot.withinss

#Km within ss for

- 97.97927
- (between\_SS / total\_SS = 79.3 %)

**#nstart =20**

- km.out <- kmeans(x, 3, nstart = 20)
- km.out\$tot.withinss

#Km within ss

- 104.3319
- (between\_SS / total\_SS = 78.0 %)

# Question:

- What is the difference when nstart = 1 and nstart = 20 used?
- Which one should we use?
- Why nstart=1 has more variance than nstart=20?
- <https://datascience.stackexchange.com/questions/11485/k-means-in-r-usage-of-nstart-parameter>

# Let's do k-means clustering with "iris" data:

```
#Load two packages for special  
plot
```

- library(ClusterR)
- **library(cluster)**

```
#Get, check and make data
```

- data(iris)
- str(iris)
- iris\_1 <- iris[,-5]

```
# Fitting K-Means clustering  
Model to training dataset
```

- set.seed(240)
- kmeans.res <- kmeans(iris\_1,  
**centers = 3, nstart = 20**)
- kmeans.res

We have used k=3 as we know that there are 3 types of flowers!

# k-means fit:

- K-means clustering with **3 clusters** of sizes 50, 62, 38
- Cluster means:

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
• 1	5.006000	3.428000	1.462000	0.246000
• 2	5.901613	2.748387	4.393548	1.433871
• 3	6.850000	3.073684	5.742105	2.071053

- Within cluster sum of squares by cluster:
- [1] 15.15100 39.82097 23.87947
- ( $\text{between\_SS} / \text{total\_SS} = \mathbf{88.4\%}$ )

# Confusion matrix: Possible here as we also have dependent variable to compare!

## # Confusion Matrix (**not usual!**)

- cm <- table(iris\$Species, kmeans.res\$cluster)
  - cm
    - setosa
    - versicolor
    - virginica
  - #Accuracy
  - (accuracy <- sum(diag(cm))/sum(cm))
  - (mce <- 1 - accuracy)
- [1] 0.8933333
  - [1] 0.1066667

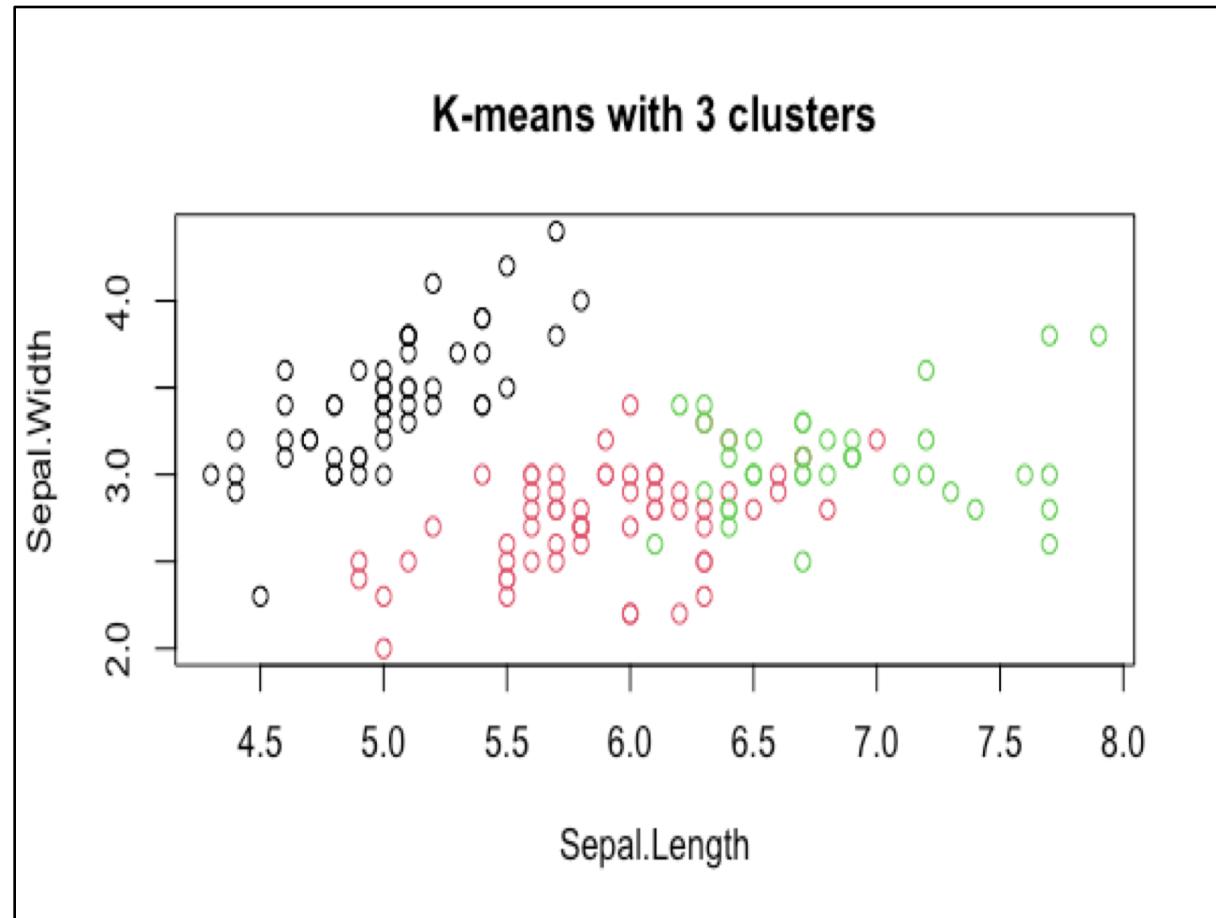
	1	2	3
setosa	50	0	0
versicolor	0	48	2
virginica	0	14	36

Do the same for the Decision Tree based models fitted with CTG data in the previous class and compare!

# Model Evaluation and Visualization:

```
# Model Evaluation and visualization
```

- `plot(iris_1[c("Sepal.Length", "Sepal.Width")])`
- `plot(iris_1[c("Sepal.Length", "Sepal.Width")], col = kmeans.res$cluster)`
- `plot(iris_1[c("Sepal.Length", "Sepal.Width")], col = kmeans.res$cluster, main = "K-means with 3 clusters")`



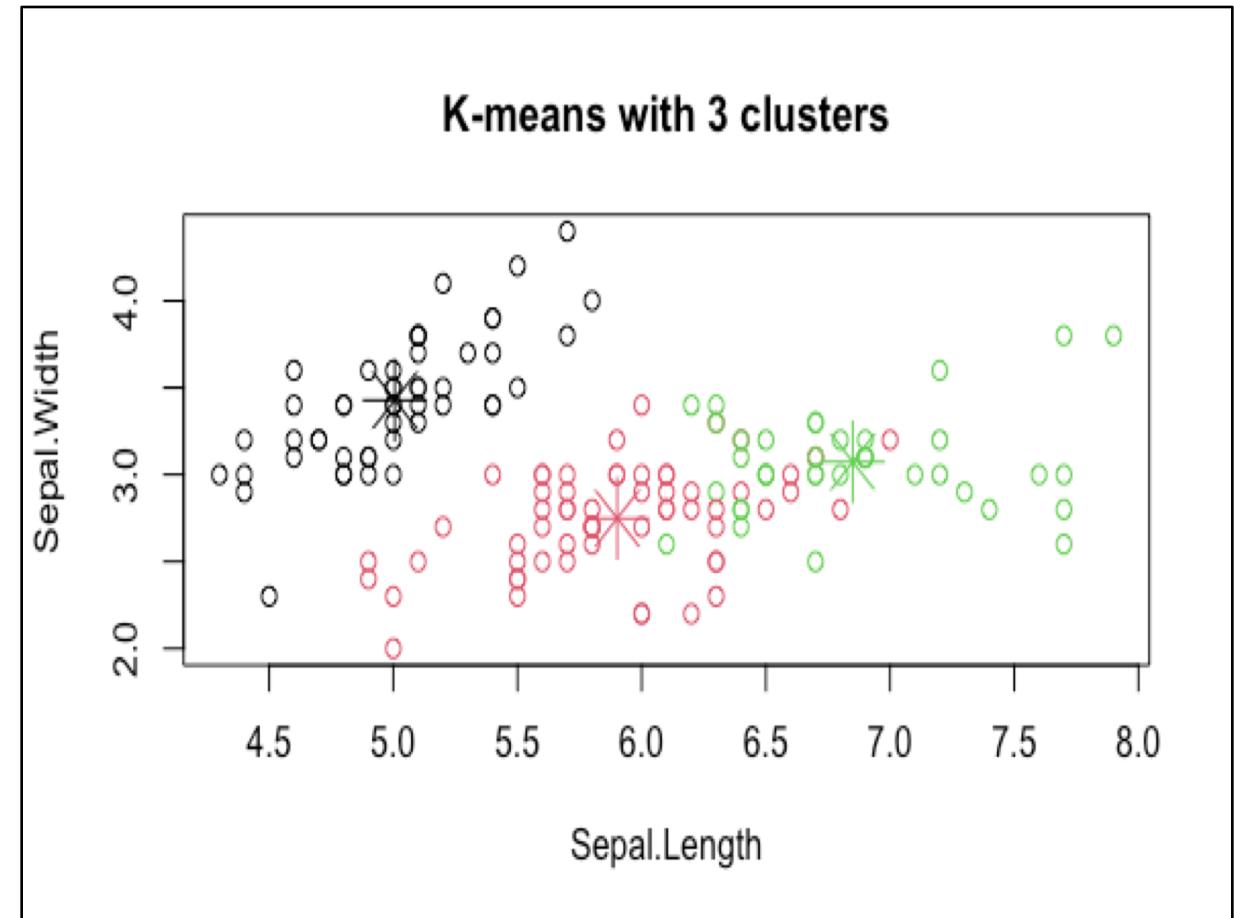
# Adding cluster centers ( use points after plot):

## # Getting cluster centers

- kmeans.res\$centers
- kmeans.res\$centers[,  
c("Sepal.Length",  
"Sepal.Width")]

## # Plotting cluster centers

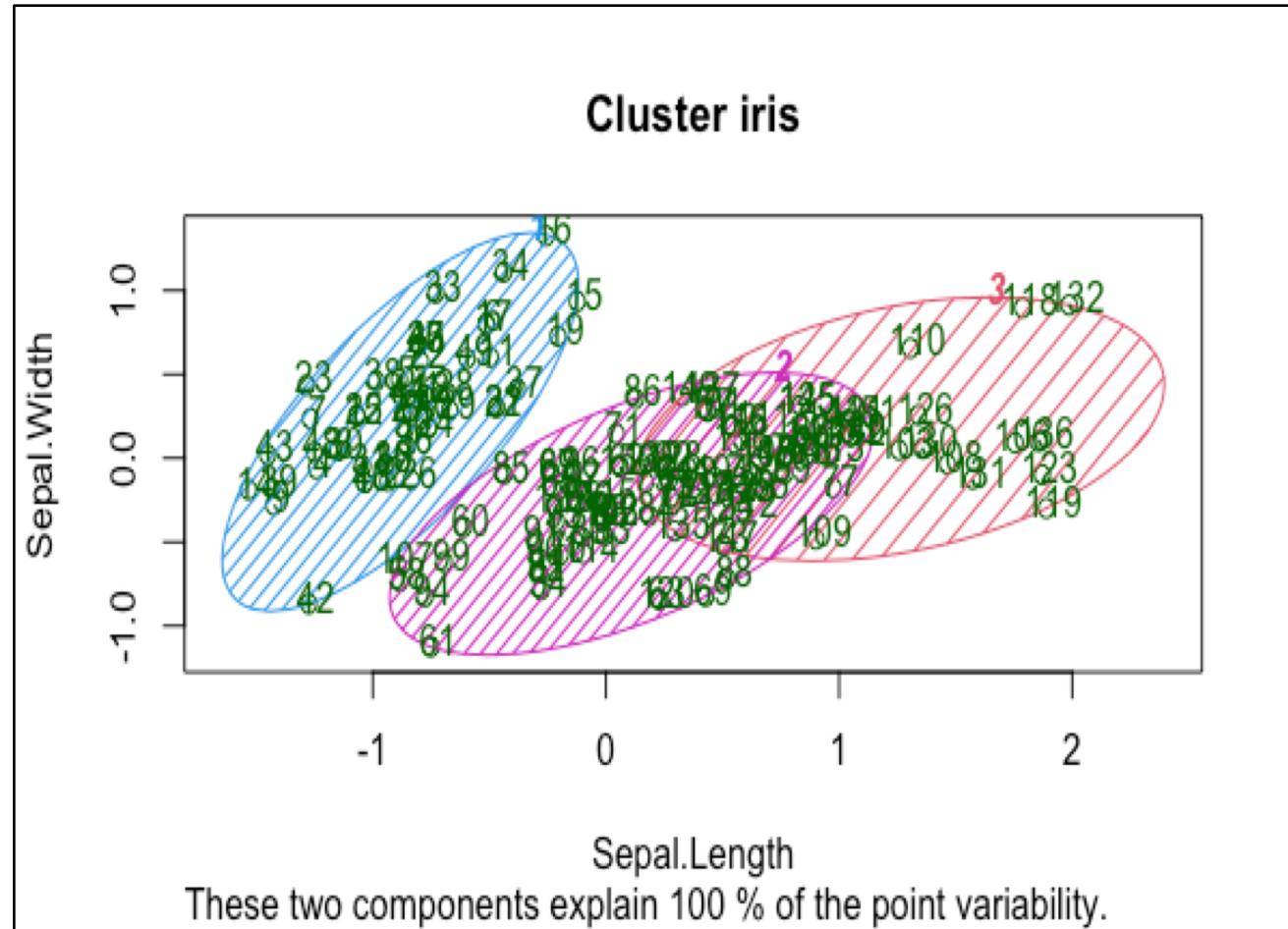
- points(kmeans.res\$centers[,  
c("Sepal.Length",  
"Sepal.Width")],  
col = 1:3, pch = 8, cex = 3)



# Visualizing clusters:

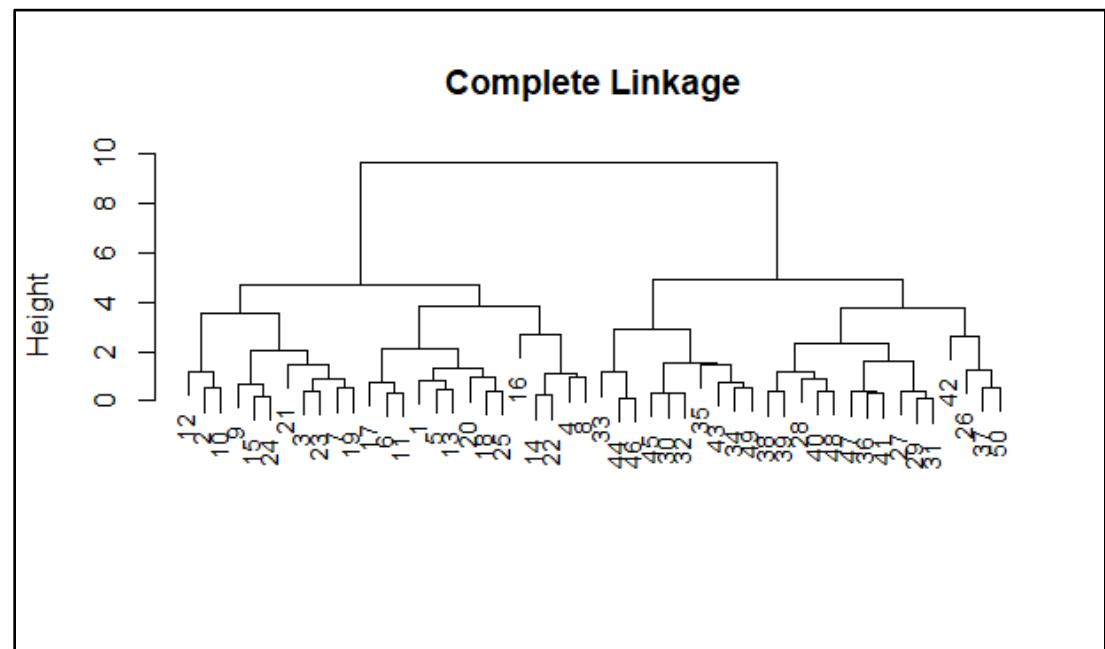
## # Visualizing clusters

- y\_kmeans <- kmeans.res\$cluster
- library(cluster)
- clusplot(iris\_1[, c("Sepal.Length", "Sepal.Width")],  
y\_kmeans,  
lines = 0,  
shade = TRUE, color = TRUE,  
labels = 2,  
plotchar = FALSE, span = TRUE,  
main = paste("Cluster iris"),  
xlab = 'Sepal.Length',  
ylab = 'Sepal.Width')



# Hierarchical cluster analysis (HCA):

- 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**.



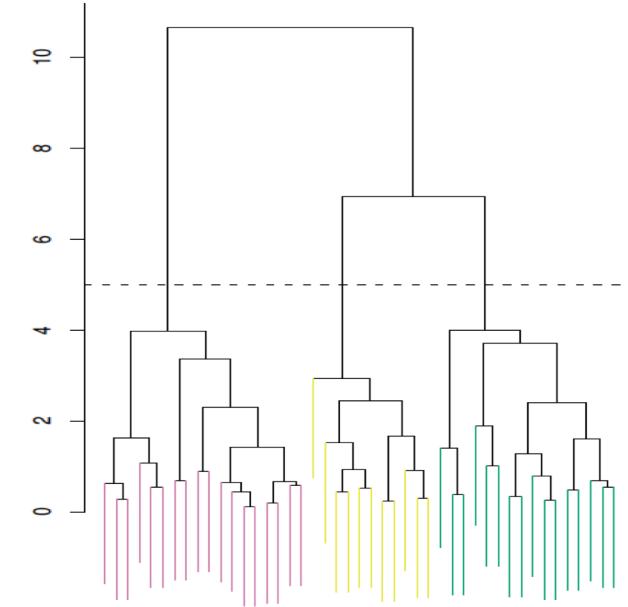
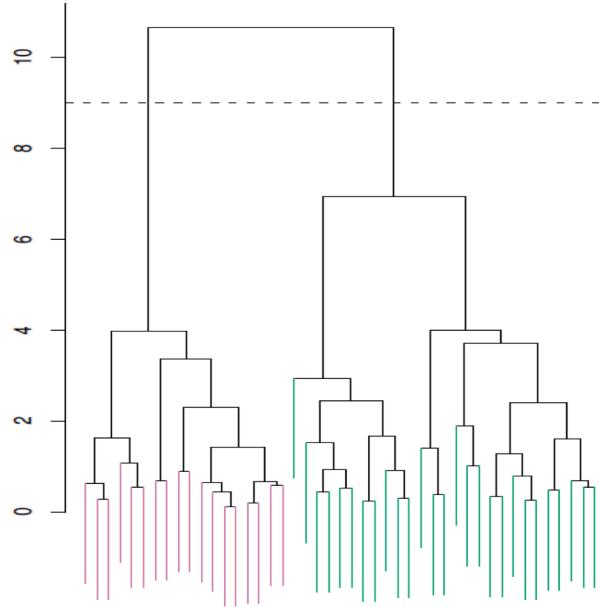
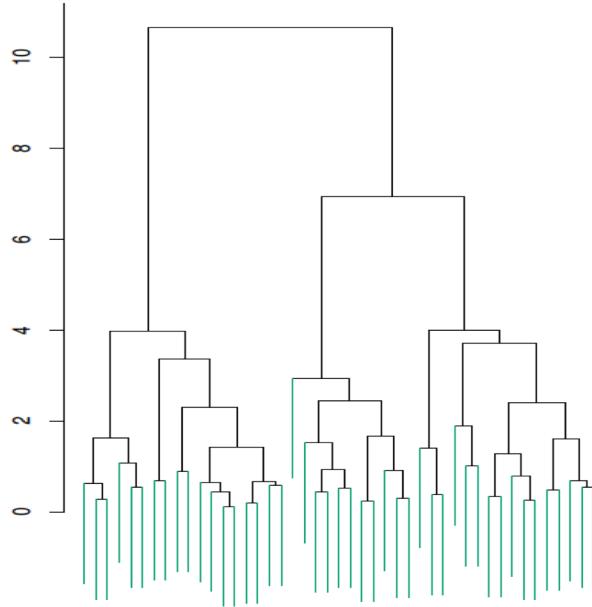
# HCA algorithm:

- The hierarchical clustering dendrogram is obtained via an extremely simple algorithm.
- We begin by defining some sort of dissimilarity measure between each pair of observations. Most often, Euclidean distance is used.
- The algorithm proceeds iteratively.
- Starting out at the bottom of the dendrogram, each of the  $n$  observations is treated as its own cluster.

# HCA algorithm:

- The two clusters that are most similar to each other are then fused so that there now are  $n-1$  clusters.
- Next the two clusters that are most similar to each other are fused again, so that there now are  $n - 2$  clusters.
- The algorithm proceeds in this fashion until all of the observations belong to one single cluster, and the dendrogram is complete.

# Hierarchical clustering: How many clusters? Need to find k using a vertical “cut” line/s!



**FIGURE 12.11.** Left: dendrogram obtained from hierarchically clustering the data from Figure 12.10 with complete linkage and Euclidean distance. Center: the dendrogram from the left-hand panel, cut at a height of nine (indicated by the dashed line). This cut results in two distinct clusters, shown in different colors. Right: the dendrogram from the left-hand panel, now cut at a height of five. This cut results in three distinct clusters, shown in different colors. Note that the colors were not used in clustering, but are simply used for display purposes in this figure.

# HCA: Linkage methods for choosing “dissimilarity” measure

<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> .

**TABLE 12.3.** A summary of the four most commonly-used types of linkage in hierarchical clustering.

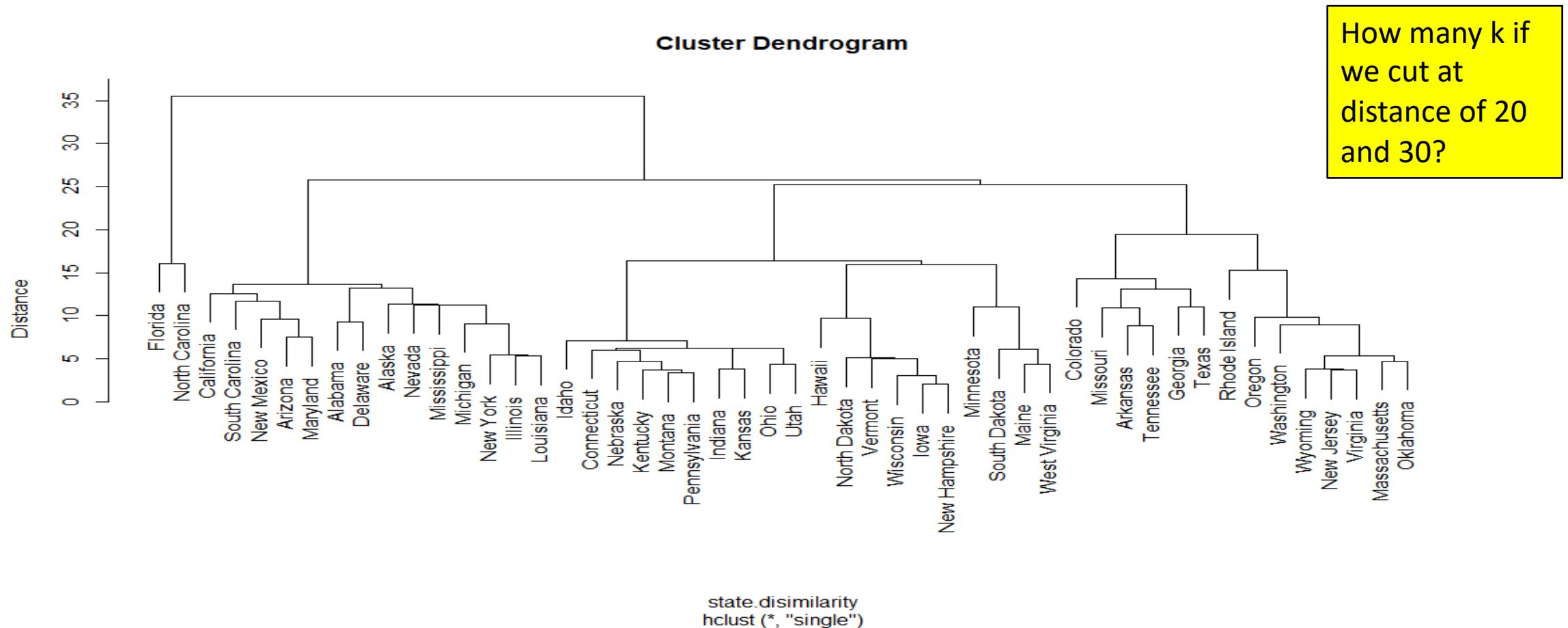
# HCA with “single” linkage: USArests.1 data:

## #Hierarchical clustering with single linkage

- #US Arrests data
- USArests.1 <- USArests[,-3]
- state.disimilarity <- dist(USArests.1)
- hirar.1 <- hclust(state.disimilarity, method='single')
- **plot(hirar.1, labels=rownames(USArests.1), ylab="Distance")**

- hirar.1
- Call:  
• hclust(d = state.disimilarity, method = "single")
- Cluster method : single
- Distance : euclidean
- Number of objects: 50

# HCA with “single” linkage in USArests.1 data



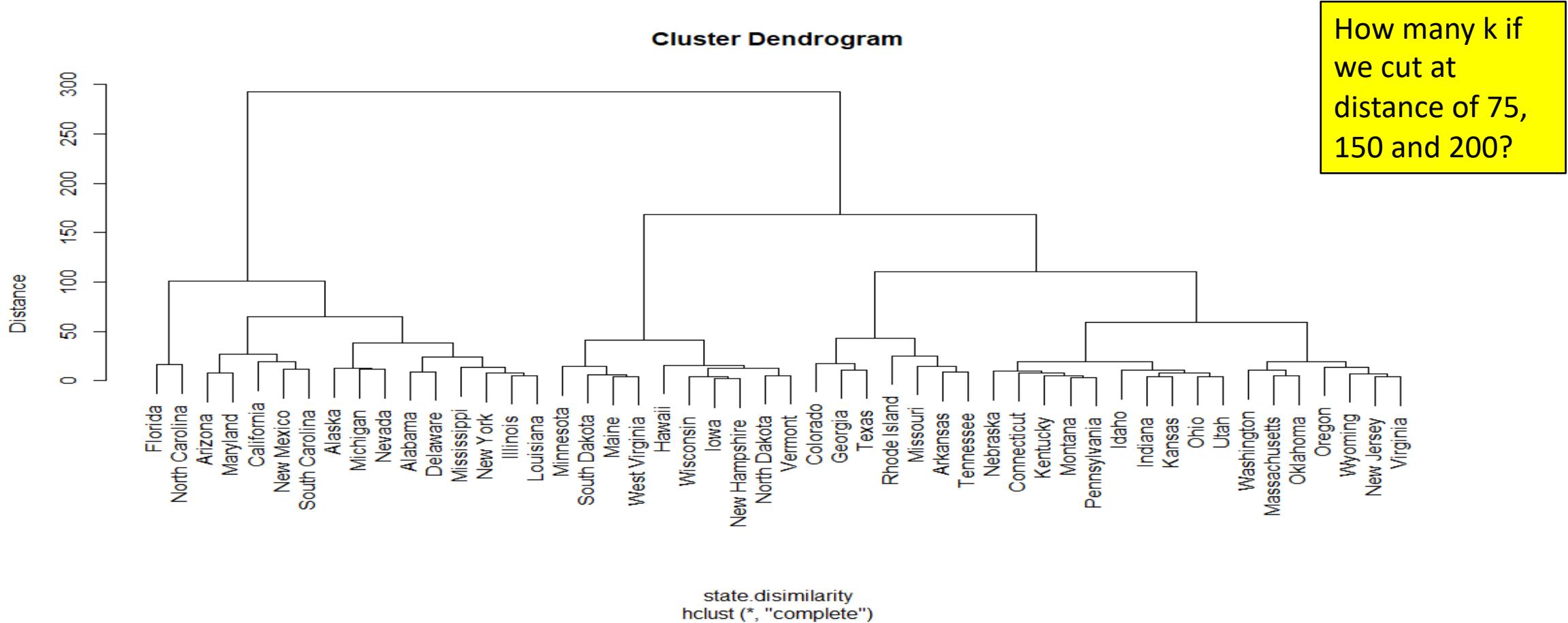
# HCA with “single” linkage: USArrests.1 data:

## #Hierarchical clustering with complete linkage

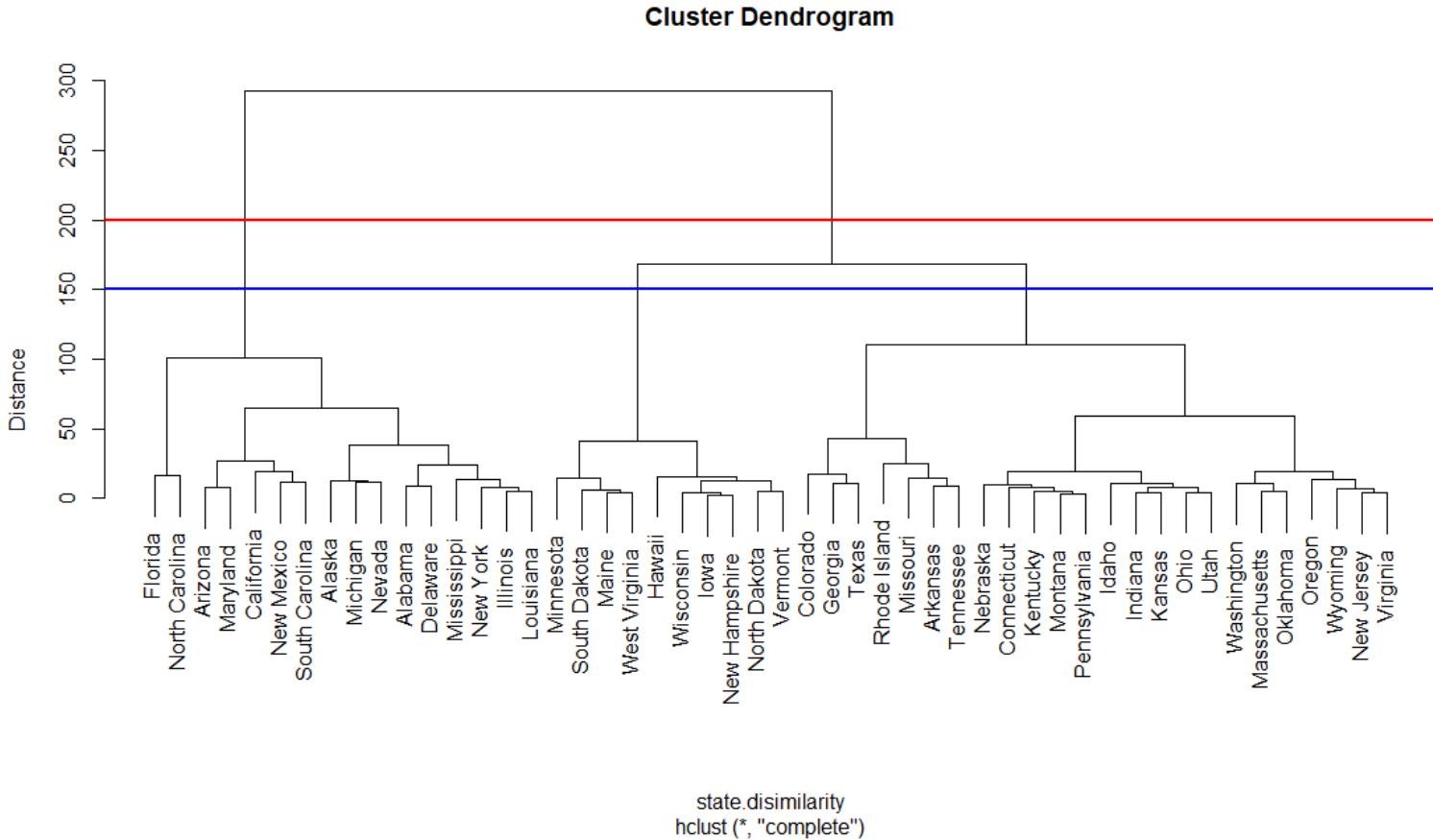
- #US Arrests data
- hirar.2 <-  
hclust(state.disimilarity,  
method='complete')
- **plot(hirar.2,**  
**labels=rownames(USArrests.1),**  
**ylab="Distance")**

- Call:  
hclust(d = state.disimilarity,  
method = "complete")
- Cluster method : complete
- Distance : euclidean
- Number of objects: 50

# HCA with “complete” linkage in USArests.1 data



# HCA with “complete” linkage in USArests.1 data with cut at distance of 200 and 150!



So, it is always better to use the HCA to determine the K and then use it to fit the k-means clustering.

In Data Science, we need to use all the four methods and find best k and the fit k-mean for each of the best K's. Then select the best clustering model based on the highest R-square value.

# Question/queries?

- Next two classes:
  1. Association rules
  2. Monte Carlo Simulations
- Final class: **Projects in R**
  - Install “git” on Windows so that we can use “github” in R Studio while creating online projects
  - We will need to use the offline projects in R Studio too

# Thank you!

@shitalbhandary