

STA30005:

MULTIVARIATE ANALYSIS

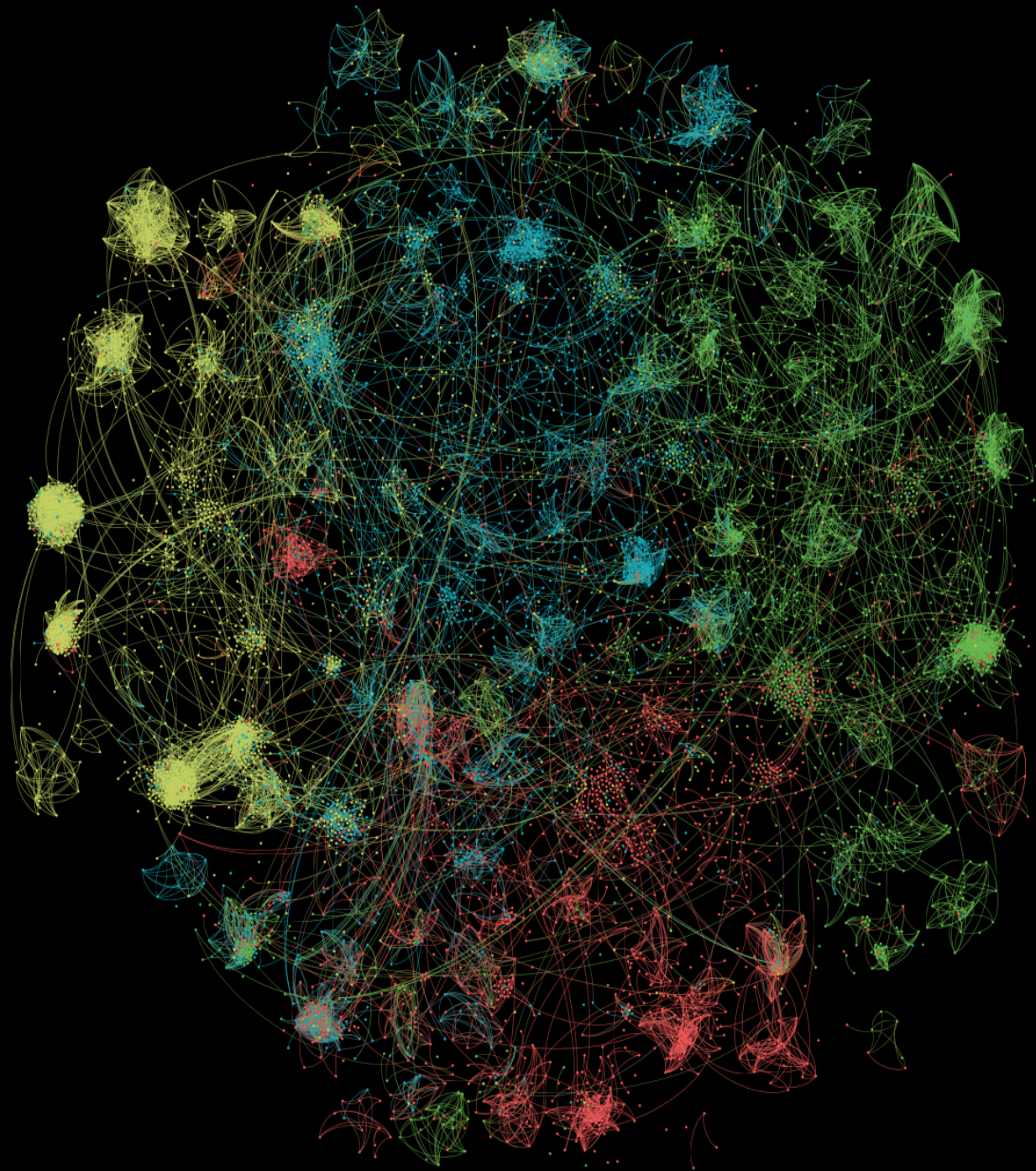


WEEK 8:

CLUSTER ANALYSIS

Tutorial Slides

STA30005: Multivariate Analysis



TASKS

Task A: K-means Clustering

Using the base R data file, USArrests, conduct a K-means cluster analysis and produce an appropriate visualisation to present your findings. Your visualisation should clearly display which US states belong to which cluster.

Task B: Hierarchical Clustering

Using the base R data file, USArrests, conduct a hierarchical (agglomerative) cluster analysis and produce a dendrogram to present your findings. You should compare different distance and linkage methods to find the optimal dendrogram for this scenario.

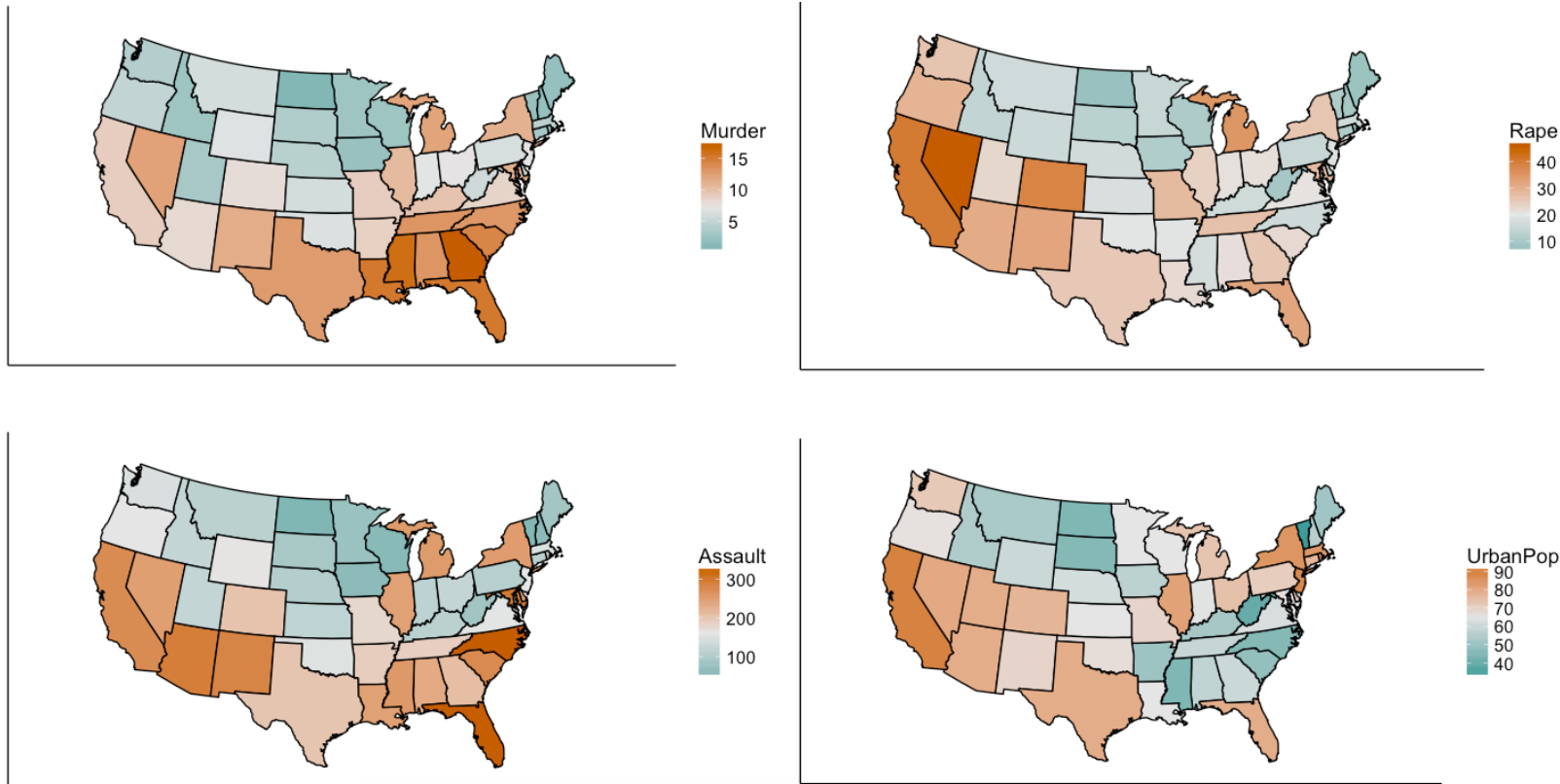
Task C: Enhancing your dendrograms

Recreate your dendrograms (from Task B) as a:

- Circular plot
- Phylogenetic tree plot
- Heatmap
- Interactive heatmap

THE DATA

The data file USArrests contains the number of arrests made (per 100,000) in US states for murder, assault and rape as well as the percent of urban population of those states in 1973. We would like to use cluster analysis to determine which states are similar to each other with regards to types of arrests.



TASKS

Task A: K-means Clustering

Using the base R data file, USArrests, conduct a K-means cluster analysis and produce an appropriate visualisation to present your findings. Your visualisation should clearly display which US states belong to which cluster.

Task B: Hierarchical Clustering

Using the base R data file, USArrests, conduct a hierarchical (agglomerative) cluster analysis and produce a dendrogram to present your findings. You should compare different distance and linkage methods to find the optimal dendrogram for this scenario.

Task C: Enhancing your dendrograms

Recreate your dendrograms (from Task B) as a:

- Circular plot
- Phylogenetic tree plot
- Heatmap
- Interactive heatmap

K-MEANS CLUSTERING

Data / Package Preparation

Determine K

Compute K-mean clusters

Visualise

Load the data (this is a base data set for R)

```
data("USArrests")
```

Show the data

```
View(USArrests)
```

*# Define (I like to use **df**) and Standardize the data*

```
df <- scale(USArrests)
```

Show the first 10 rows (round to 2 decimal places)

```
round(head(df, n = 10), 2)
```

Set a seed so that we can reproduce the same results

```
set.seed(123)
```

The screenshot shows the RStudio interface. The top pane displays the 'USArrests' dataset as a table with columns: Murder, Assault, UrbanPop, and Rape. The bottom pane shows the R console with the following code and output:

```
> data("USArrests")
> View(USArrests)
> df <- scale(USArrests)
> round(head(df, n = 10), 2)
```

	Murder	Assault	UrbanPop	Rape
Alabama	1.24	0.78	-0.52	0.00
Alaska	0.51	1.11	-1.21	2.48
Arizona	0.07	1.48	1.00	1.04
Arkansas	0.23	0.23	-1.07	-0.18
California	0.28	1.26	1.76	2.07
Colorado	0.03	0.40	0.86	1.86
Connecticut	-1.03	-0.73	0.79	-1.08

K-MEANS CLUSTERING

Data / Package Preparation

Determine K

Compute K-mean clusters

Visualise

Install / Load the factoextra package (this is used to generate cluster plots)

```
install.packages("factoextra")
```

```
library(factoextra)
```

Install / Load the cluster package (this is used to determine K)

```
install.packages("cluster")
```

```
library(cluster)
```


K-MEANS CLUSTERING

Data / Package Preparation

Determine K

Compute K-mean clusters

Visualise

There are several traditional methods for estimating K (such as Elbow and Silhouette plots). Modern methods have become more sophisticated due to advancements in computing. We will use the Gap Statistic method here

Compute the gap statistic

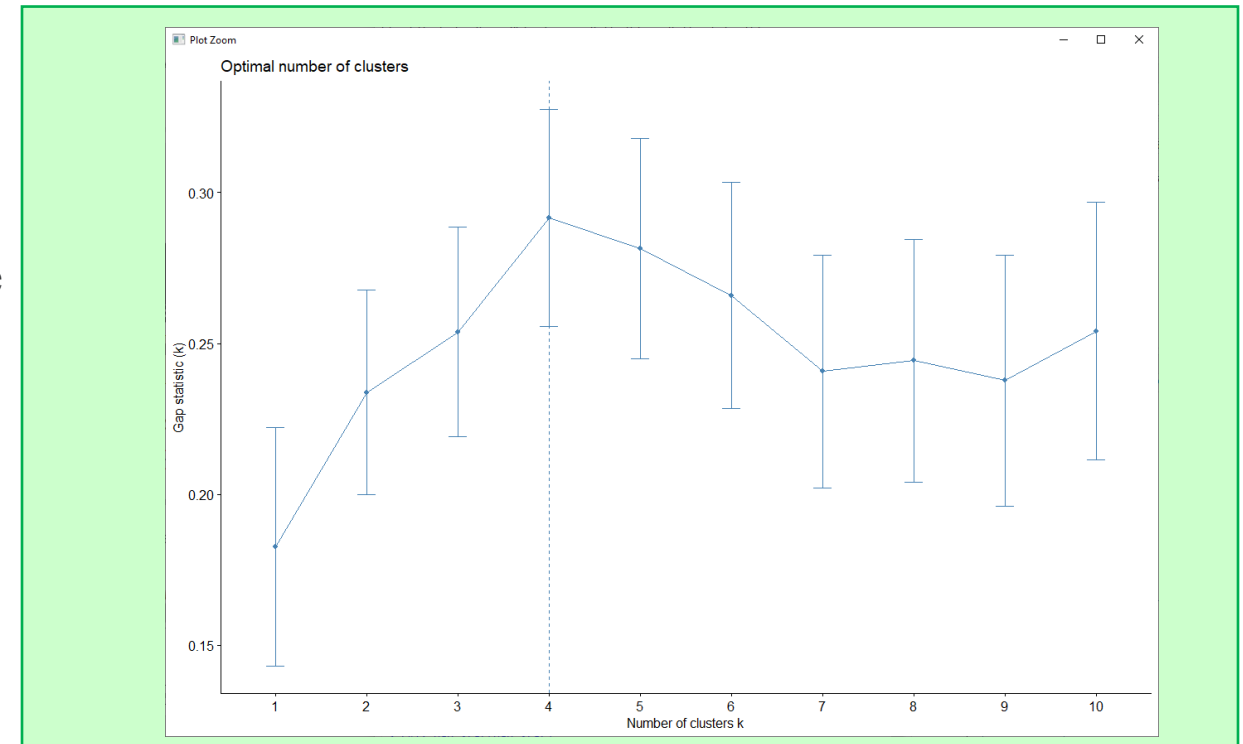
```
gap_stat <- clusGap(df, FUN = kmeans, nstart = 25, K.max =  
10, B = 500)
```

This function will conduct a k-means cluster analysis using the gap statistic method (with 25 random starting points) using 500 bootstrap samples

Visualise the gap statistic plot

```
fviz_gap_stat(gap_stat)
```

Based upon the gap statistic method, $K = 4$ is optimal



K-MEANS CLUSTERING

Data / Package Preparation

Determine K

Compute K-mean clusters

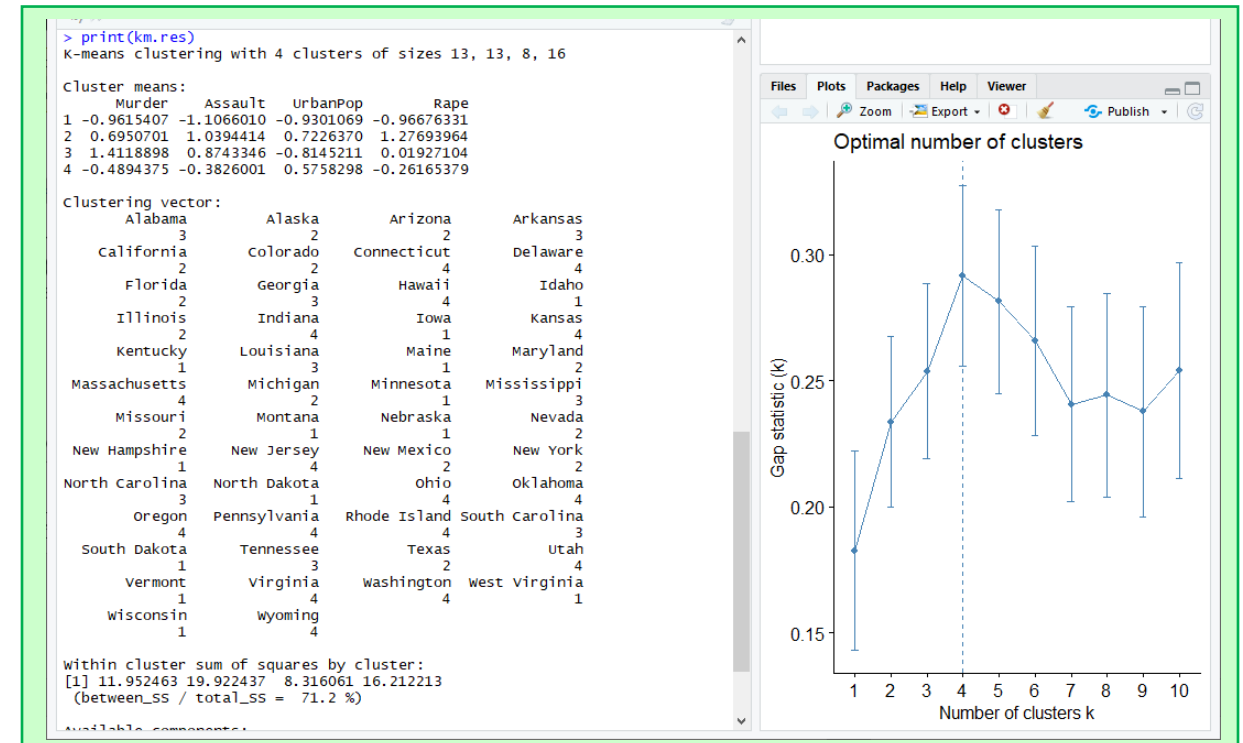
Visualise

As the final result of k-means clustering result is sensitive to the random starting assignments, we specify `nstart = 25`. This means that R will try 25 different random starting assignments and then select the best results corresponding to the one with the lowest within cluster variation. The default value of `nstart` in R is one. But, it's strongly recommended to compute k-means clustering with a large value of `nstart` such as 25 or 50, in order to have a more stable result.

```
# Compute k-means with K = 4
km.res <- kmeans(df, 4, nstart = 25)
```

```
# Print the k-means results
print(km.res)
```

The amount of variation within each of the four clusters is provided with these results



K-MEANS CLUSTERING

Data / Package Preparation

Determine K

Compute K-mean clusters

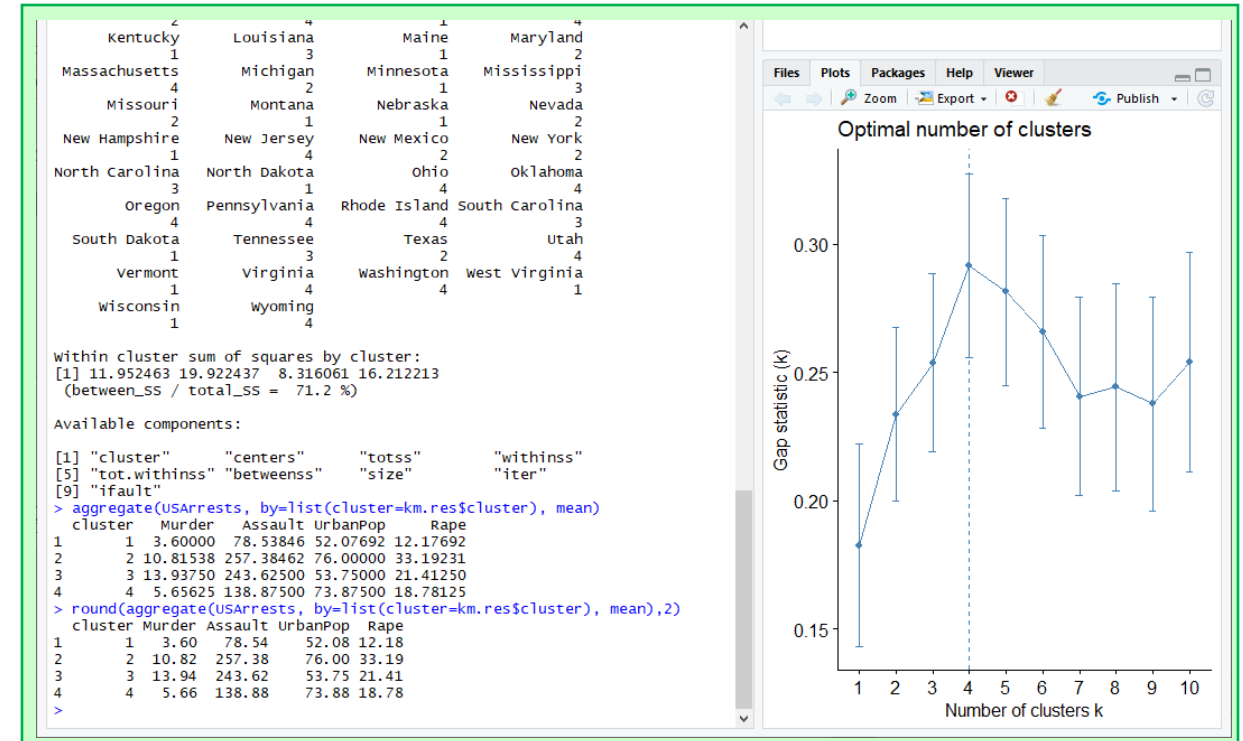
Visualise

View the cluster means by the original data

```
round(aggregate(USArrests, by=list(cluster=km.res$cluster),  
mean),2)
```

Based upon these results, we can see that:

- US states in cluster 1 are low across all variables
- US states in cluster 2 are high across all variables
- US states in cluster 3 are high across all variables
- US states in cluster 4 are low across all variables



K-MEANS CLUSTERING

Data / Package Preparation

Determine K

Compute K-mean clusters

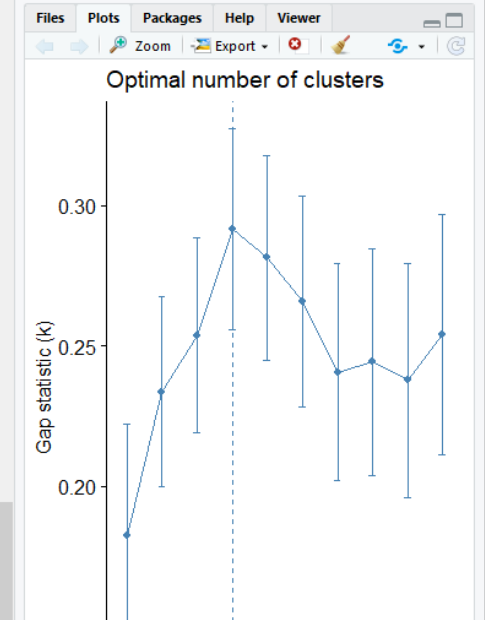
Visualise

View which US state was assigned to which cluster

```
dd <- cbind(USArrests, cluster = km.res$cluster)
show(dd)
```

```
> dd <- cbind(USArrests, cluster = km.res$cluster)
> show(dd)
```

	Murder	Assault	UrbanPop	Rape	cluster
Alabama	13.2	236	58	21.2	3
Alaska	10.0	263	48	44.5	2
Arizona	8.1	294	80	31.0	2
Arkansas	8.8	190	50	19.5	3
California	9.0	276	91	40.6	2
Colorado	7.9	204	78	38.7	2
Connecticut	3.3	110	77	11.1	4
Delaware	5.9	238	72	15.8	4
Florida	15.4	335	80	31.9	2
Georgia	17.4	211	60	25.8	3
Hawaii	5.3	46	83	20.2	4
Idaho	2.6	120	54	14.2	1
Illinois	10.4	249	83	24.0	2
Indiana	7.2	113	65	21.0	4
Iowa	2.2	56	57	11.3	1
Kansas	6.0	115	66	18.0	4
Kentucky	9.7	109	52	16.3	1
Louisiana	15.4	249	66	22.2	3
Maine	2.1	83	51	7.8	1
Maryland	11.3	300	67	27.8	2
Massachusetts	4.4	149	85	16.3	4
Michigan	12.1	255	74	35.1	2
Minnesota	2.7	72	66	14.9	1
Mississippi	16.1	259	44	17.1	3
Missouri	9.0	178	70	28.2	2
Montana	6.0	109	53	16.4	1
Nebraska	4.3	102	62	16.5	1
Nevada	12.2	252	81	46.0	2
New Hampshire	2.1	57	56	9.5	1
New Jersey	7.4	159	89	18.8	4
New Mexico	11.4	285	70	32.1	2
New York	11.1	254	86	26.1	2
North Carolina	13.0	337	45	16.1	3



K-MEANS CLUSTERING

Data / Package Preparation

Determine K

Compute K-mean clusters

Visualise

Visualise the clustering
`fviz_cluster(km.res, data = df)`



Source: Kassambara (2017). Practical Guide To Cluster Analysis in R

K-MEANS CLUSTERING

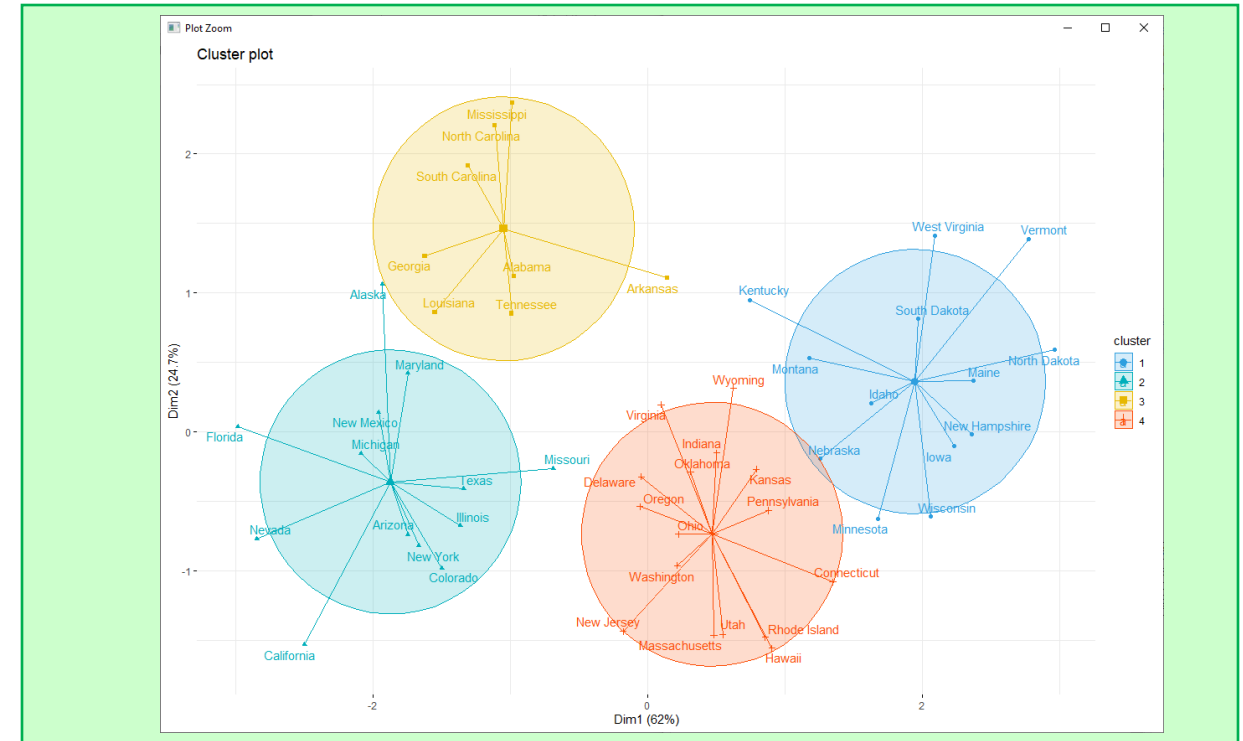
Data / Package Preparation

Determine K

Compute K-mean clusters

Visualise

```
# Visualise the clustering (enhanced)
fviz_cluster(km.res, data = df, palette = c("#2E9FDF",
"#00AFBB", "#E7B800", "#FC4E07"), ellipse.type = "euclid",
star.plot = TRUE, repel = TRUE, ggtheme =
theme_minimal())
```



TASKS

Task A: K-means Clustering

Using the base R data file, USArrests, conduct a K-means cluster analysis and produce an appropriate visualisation to present your findings. Your visualisation should clearly display which US states belong to which cluster.

Task B: Hierarchical Clustering

Using the base R data file, USArrests, conduct a hierarchical (agglomerative) cluster analysis and produce a dendrogram to present your findings. You should compare different distance and linkage methods to find the optimal dendrogram for this scenario.

Task C: Enhancing your dendrograms

Recreate your dendrograms (from Task B) as a:

- Circular plot
- Phylogenetic tree plot
- Heatmap
- Interactive heatmap

AGGLOMERATIVE CLUSTERING

Data Preparation

Similarity Measures

Linkage

Dendrogram

Verification

Cutting

Load the data (this is a base data set for R)

```
data("USArrests")
```

Show the data

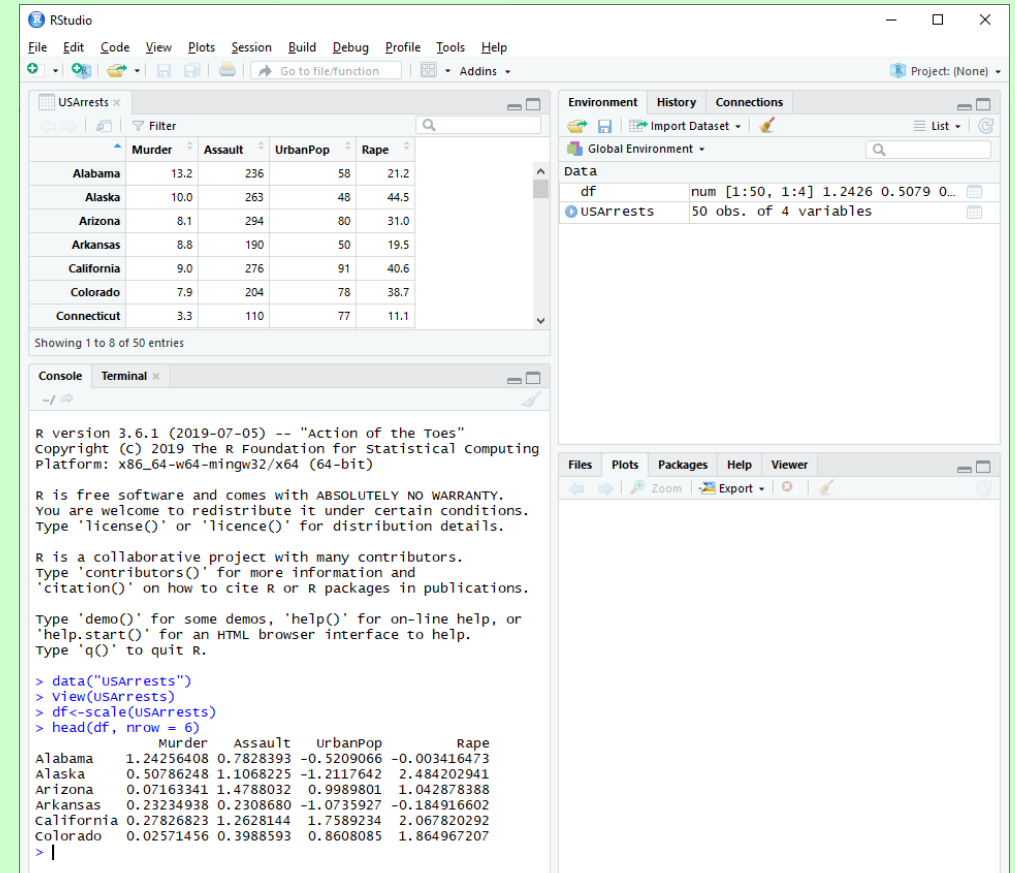
```
View(USArrests)
```

*# Define (I like to use **df**) and Standardize the data*

```
df <- scale(USArrests)
```

Show the first 6 rows (round to 2 decimal places)

```
round(head(df, nrow = 6), 2)
```



The screenshot shows the RStudio interface. The 'Environment' pane on the right lists 'df' as a numeric matrix with dimensions [1:50, 1:4] and 'USArrests' as a data frame with 50 observations and 4 variables. The 'Console' pane at the bottom shows the execution of the following R code:

```
> data("USArrests")
> View(USArrests)
> df <- scale(USArrests)
> head(df, nrow = 6)
```

The output of the `head` function is displayed in the console, showing the first 6 rows of the scaled data (df) with columns Murder, Assault, UrbanPop, and Rape. The values are rounded to 2 decimal places.

	Murder	Assault	UrbanPop	Rape
Alabama	1.24256408	0.7828393	-0.5209066	-0.003416473
Alaska	0.50786248	1.1068225	-1.2117642	2.484202941
Arizona	0.07163341	1.4788032	0.9989801	1.042878388
Arkansas	0.23234938	0.2308680	-1.0735927	-0.184916602
California	0.27826823	1.2628144	1.7589234	2.067820292
Colorado	0.02571456	0.3988593	0.8608085	1.864967207

AGGLOMERATIVE CLUSTERING

Data Preparation

Similarity Measures

Linkage

Dendrogram

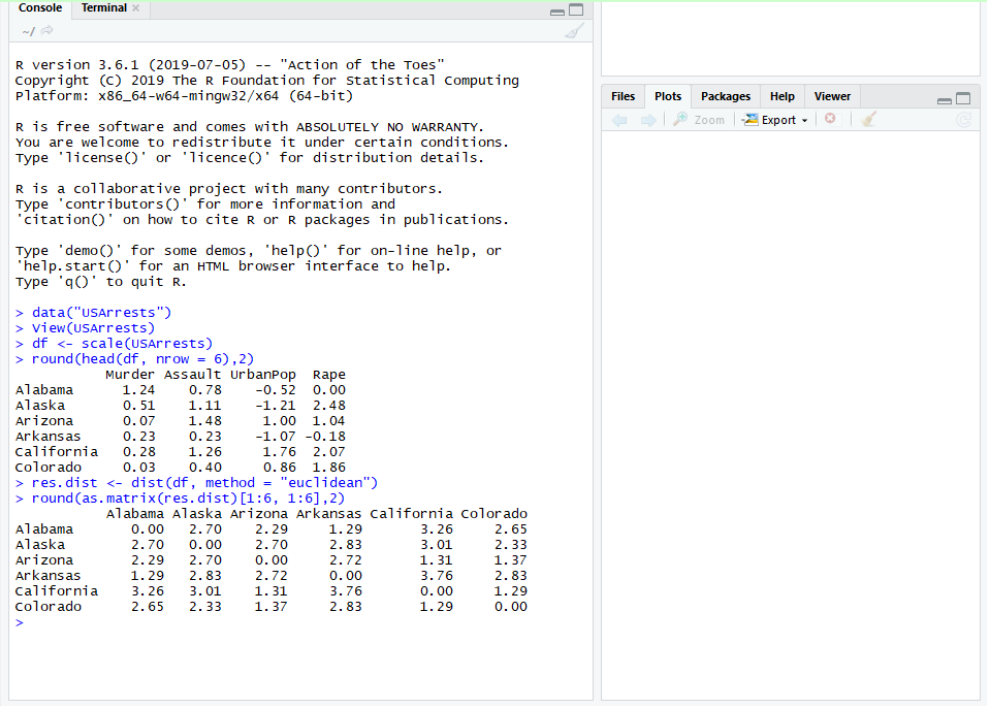
Verification

Cutting

There are many methods to calculate the (dis)similarity information (e.g. Euclidean, Manhattan, etc.). In R, we can use the function `dist()` to compute distances between every pair of objects in a data set

```
# Define and compute the dissimilarity matrix  
res.dist <- dist(df, method = "euclidean")
```

```
# View distance information in matrix form (first 6)  
round(as.matrix(res.dist)[1:6, 1:6], 2)
```



```
R version 3.6.1 (2019-07-05) -- "Action of the Toes"  
Copyright (C) 2019 The R Foundation for Statistical Computing  
Platform: x86_64-w64-mingw32/x64 (64-bit)  
  
R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.  
  
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.  
  
Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.  
  
> data("USArrests")  
> view(USArrests)  
> df <- scale(USArrests)  
> round(head(df, nrow = 6), 2)  
      Murder Assault UrbanPop Rape  
Alabama   1.24    0.78   -0.52  0.00  
Alaska    0.51    1.11  -1.21  2.48  
Arizona   0.07    1.48    1.00  1.04  
Arkansas  0.23    0.23  -1.07 -0.18  
California 0.28    1.26    1.76  2.07  
Colorado  0.03    0.40    0.86  1.86  
> res.dist <- dist(df, method = "euclidean")  
> round(as.matrix(res.dist)[1:6, 1:6], 2)  
      Alabama Alaska Arizona Arkansas California Colorado  
Alabama  0.00  2.70  2.29  1.29  3.26  2.65  
Alaska   2.70  0.00  2.70  2.83  3.01  2.33  
Arizona  2.29  2.70  0.00  2.72  1.31  1.37  
Arkansas 1.29  2.83  2.72  0.00  3.76  2.83  
California 3.26  3.01  1.31  3.76  0.00  1.29  
Colorado 2.65  2.33  1.37  2.83  1.29  0.00  
>
```

AGGLOMERATIVE CLUSTERING

Data Preparation

Similarity Measures

Linkage

Dendrogram

Verification

Cutting

The linkage function takes the distance information, returned by the function `dist()`, and groups pairs of objects into clusters based on their similarity. Next, these newly formed clusters are linked to each other to create bigger clusters. We can use `hclust()` to create the hierarchical tree

Define and apply a hierarchical clustering linkage (in this example we will use Ward's method)
`res.hc <- hclust(d = res.dist, method = "ward.D2")`

Common linkage methods are described below:

- **Maximum or complete linkage:** The distance between two clusters is defined as the maximum value of all pairwise distances between the elements in cluster 1 and the elements in cluster 2. It tends to produce more compact clusters.
- **Minimum or single linkage:** The distance between two clusters is defined as the minimum value of all pairwise distances between the elements in cluster 1 and the elements in cluster 2. It tends to produce long, “loose” clusters.
- **Mean or average linkage:** The distance between two clusters is defined as the average distance between the elements in cluster 1 and the elements in cluster 2.
- **Centroid linkage:** The distance between two clusters is defined as the distance between the centroid for cluster 1 (a mean vector of length p variables) and the centroid for cluster 2.
- **Ward's minimum variance method:** It minimizes the total within-cluster variance. At each step the pair of clusters with minimum between-cluster distance are merged.

AGGLOMERATIVE CLUSTERING

Data Preparation

Similarity Measures

Linkage

Dendrogram

Verification

Cutting

Dendrograms correspond to the graphical representation of the hierarchical tree generated by the function `hclust()`. Here, we'll use the function `fviz_dend()` [in factoextra R package] to produce our dendrograms.

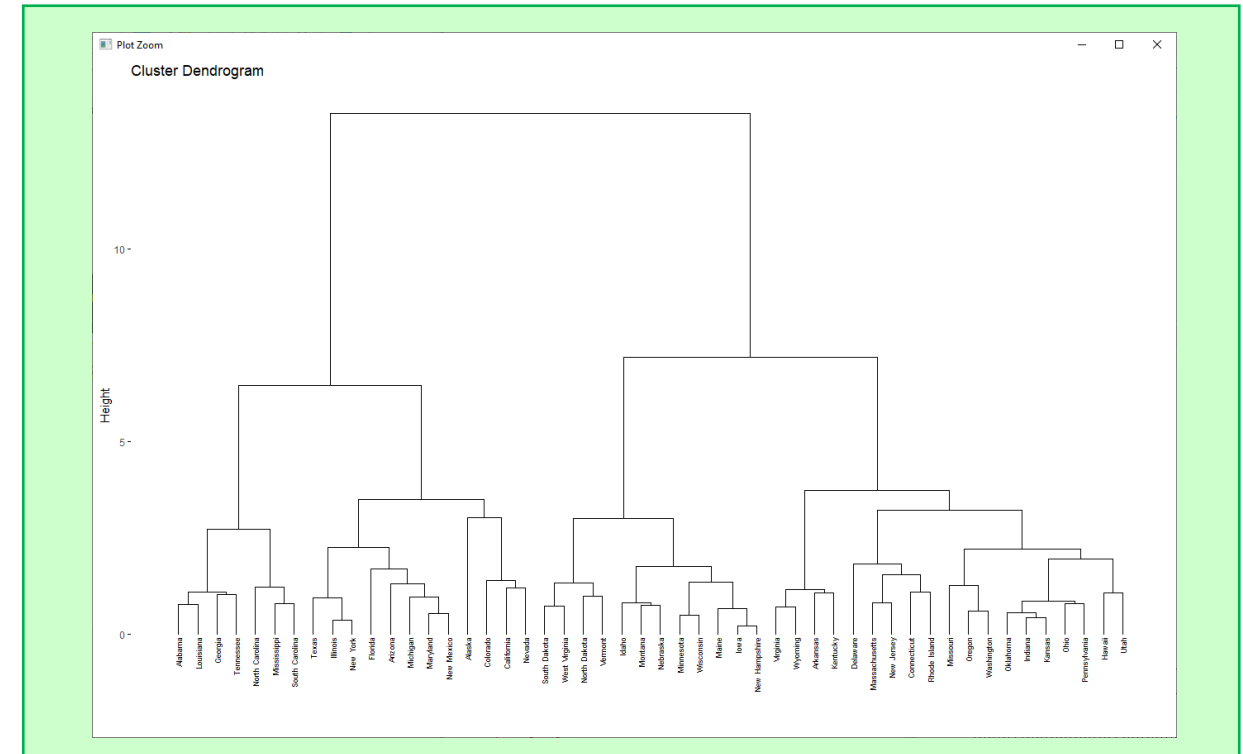
Install and Load the factoextra package

```
install.packages(factoextra)
```

```
library(factoextra)
```

Plot the dendrogram

```
fviz_dend(res.hc, cex = 0.5)
```



AGGLOMERATIVE CLUSTERING

Data Preparation

Similarity Measures

Linkage

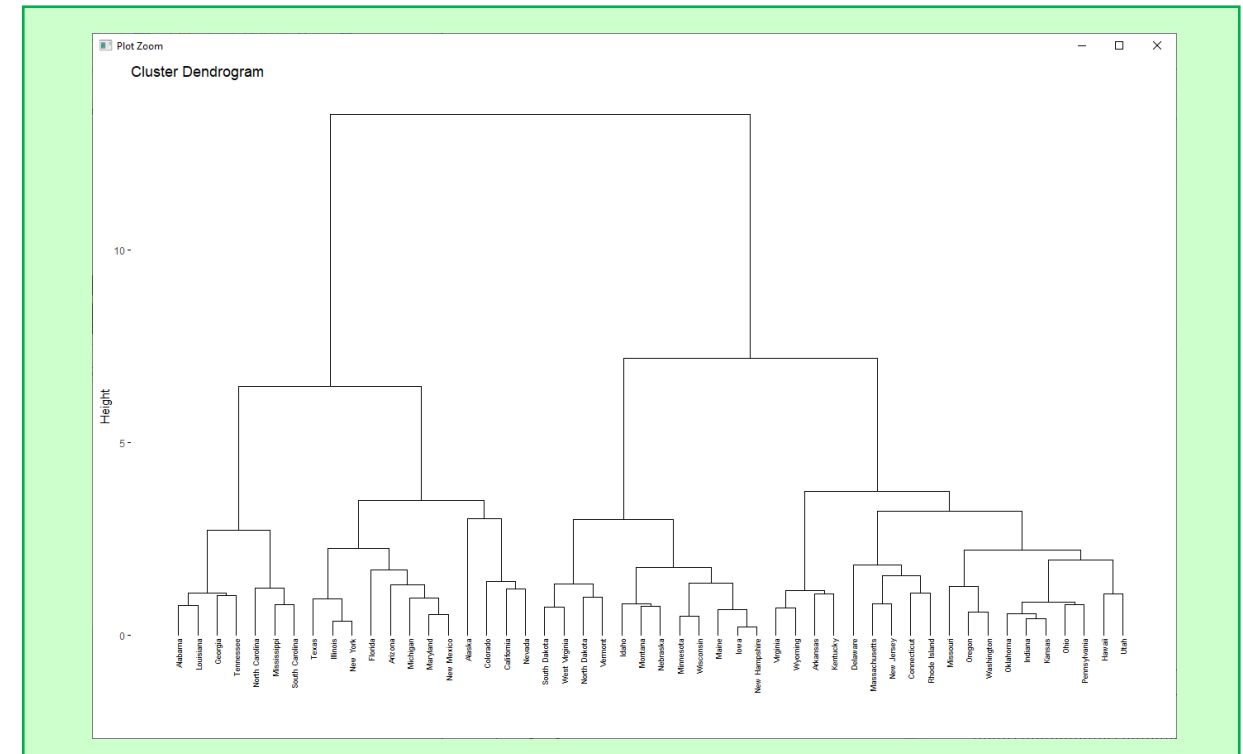
Dendrogram

Verification

Cutting

Notes about this dendrogram:

- In this dendrogram, each leaf corresponds to one object.
- As we move up the tree, objects that are similar to each other are combined into branches, which are themselves fused at a higher height.
- The height of the fusion, provided on the vertical axis, indicates the (dis)similarity/distance between two objects/clusters.
- The higher the height of the fusion, the less similar the objects are. This height is known as the *cophenetic distance* between the two objects.



AGGLOMERATIVE CLUSTERING

Data Preparation

Similarity Measures

Linkage

Dendrogram

Verification

Cutting

We should now assess the dendrogram heights against the original distances.

Compute the cophenetic distance

```
res.coph <- cophenetic(res.hc)
```

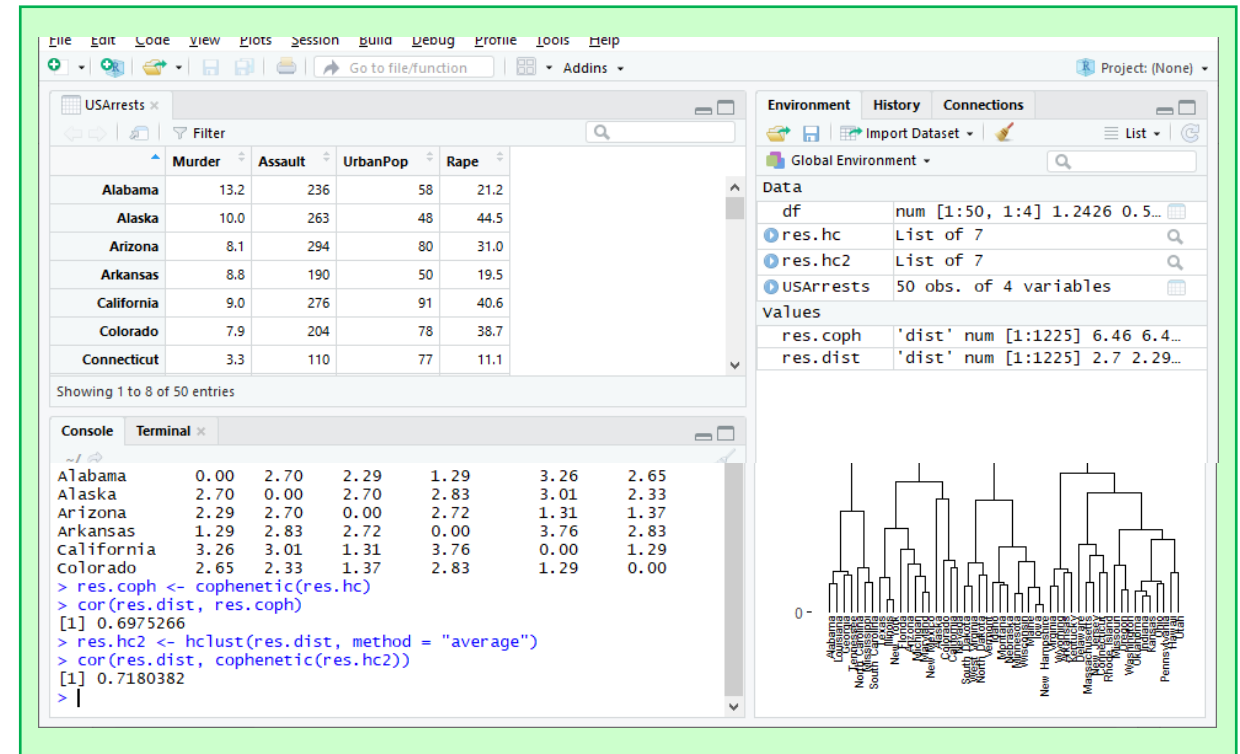
Correlate cophenetic distance and original distance

```
cor(res.dist, res.coph)
```

Note: the closer the value of the correlation coefficient is to 1, the more accurate the clustering solution reflects your data (with values above 0.75 being very good). You may want to try different linkage methods and evaluate the correlation again:

```
res.hc2 <- hclust(res.dist, method = "average")  
cor(res.dist, cophenetic(res.hc2))
```

Source: Kassambara (2017). Practical Guide To Cluster Analysis in R



AGGLOMERATIVE CLUSTERING

Data Preparation

Similarity Measures

Linkage

Dendrogram

Verification

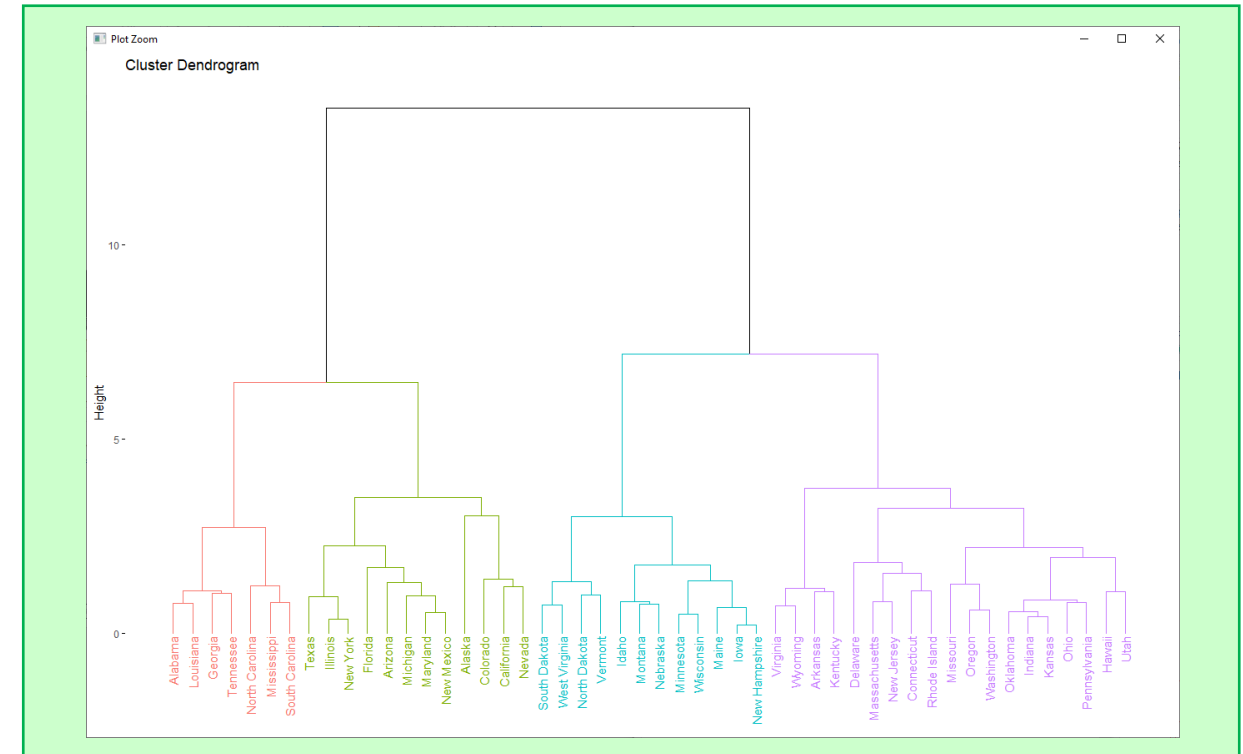
Cutting

The R base function `cutree()` can be used to cut a tree, generated by the `hclust()` function, into several groups either by specifying the desired number of groups or the cut height.

```
# Cut the dendrogram into 4 groups  
grp <- cutree(res.hc, k = 4)
```

```
# Plot this new dendrogram  
fviz_dend(res.hc, k = 4)
```

Try recreating the dendrogram with different values of k or cut points



TASKS

Task A: K-means Clustering

Using the base R data file, USArrests, conduct a K-means cluster analysis and produce an appropriate visualisation to present your findings. Your visualisation should clearly display which US states belong to which cluster.

Task B: Hierarchical Clustering

Using the base R data file, USArrests, conduct a hierarchical (agglomerative) cluster analysis and produce a dendrogram to present your findings. You should compare different distance and linkage methods to find the optimal dendrogram for this scenario.

Task C: Enhancing your dendrograms

Recreate your dendrograms (from Task B) as a:

- Circular plot
- Phylogenetic tree plot
- Heatmap
- Interactive heatmap

ENHANCING YOUR DENDROGRAMS

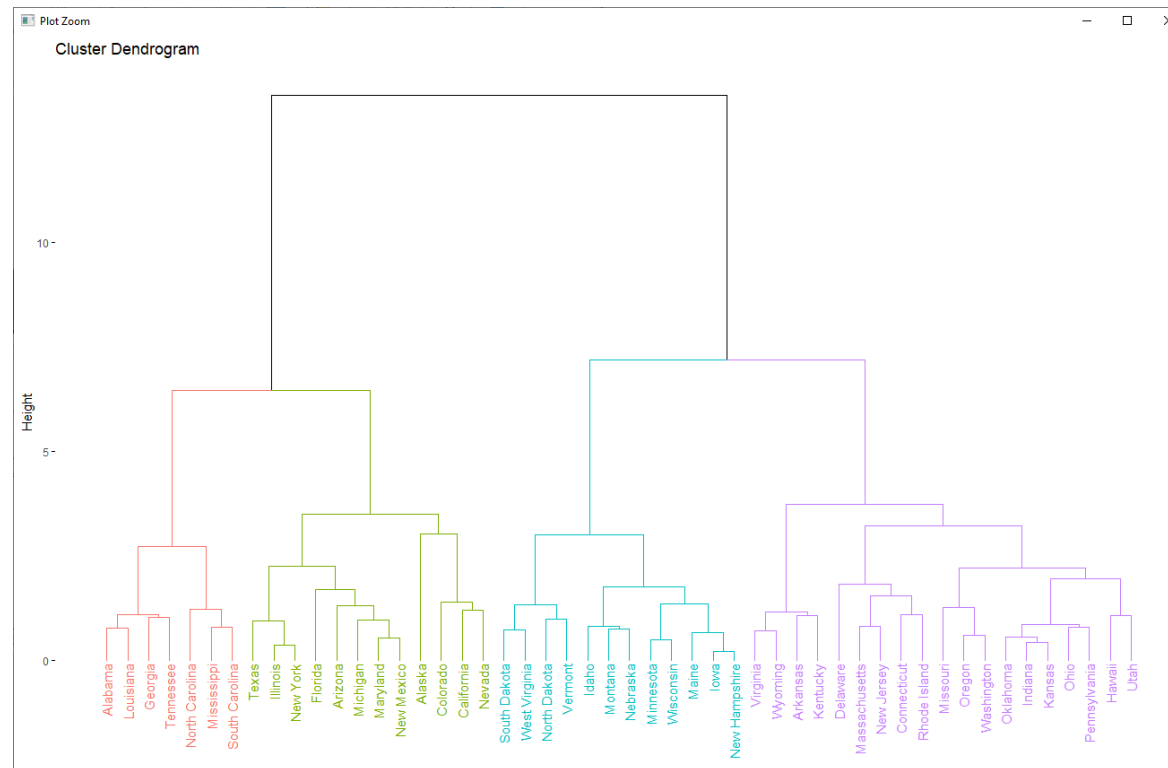
Circular plot

Phylogenetic tree plot

Heatmap

Interactive heatmap

In this task we will enhance the dendrogram we produced earlier with 4 new plots:



Source: Kassambara (2017). Practical Guide To Cluster Analysis in R

ENHANCING YOUR DENDROGRAMS

Circular plot

Phylogenetic tree plot

Heatmap

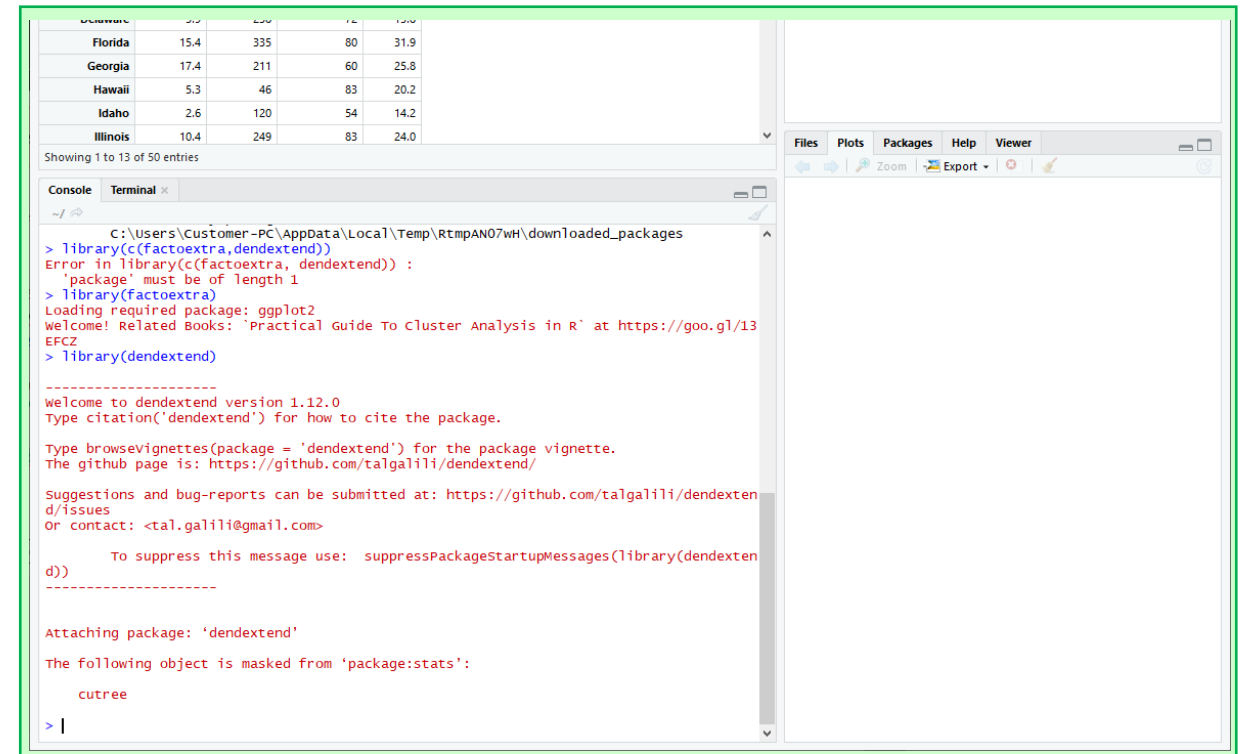
Interactive heatmap

Suppose we are starting from scratch and have not produced the outputs from the previous task

```
# Load the data (this is a base data set for R)  
data("USArrests")
```

```
# Install / Load the required packages  
install.packages(c("factoextra", "dendextend"))  
library(factoextra)  
library(dendextend)
```

```
# Compute distances and hierarchical clustering  
dd <- dist(scale(USArrests), method = "euclidean")  
hc <- hclust(dd, method = "ward.D2")
```



```
C:\Users\Customer-PC\AppData\Local\Temp\RtmpAN07wH\downloaded_packages  
> library(c(factoextra, dendextend))  
Error in library(c(factoextra, dendextend)) :  
  'package' must be of length 1  
> library(factoextra)  
Loading required package: ggplot2  
Welcome! Related Books: 'Practical Guide To Cluster Analysis in R' at https://goo.gl/13EFCZ  
> library(dendextend)  
-----  
welcome to dendextend version 1.12.0  
Type citation('dendextend') for how to cite the package.  
  
Type browsevignettes(package = 'dendextend') for the package vignette.  
The github page is: https://github.com/talgalili/dendextend/  
  
Suggestions and bug-reports can be submitted at: https://github.com/talgalili/dendextend/issues  
Or contact: <tal.galili@gmail.com>  
  
To suppress this message use: suppressPackageStartupMessages(library(dendextend))  
-----  
  
Attaching package: 'dendextend'  
  
The following object is masked from 'package:stats':  
  
    cutree  
  
> |
```

ENHANCING YOUR DENDROGRAMS

Circular plot

Phylogenetic tree plot

Heatmap

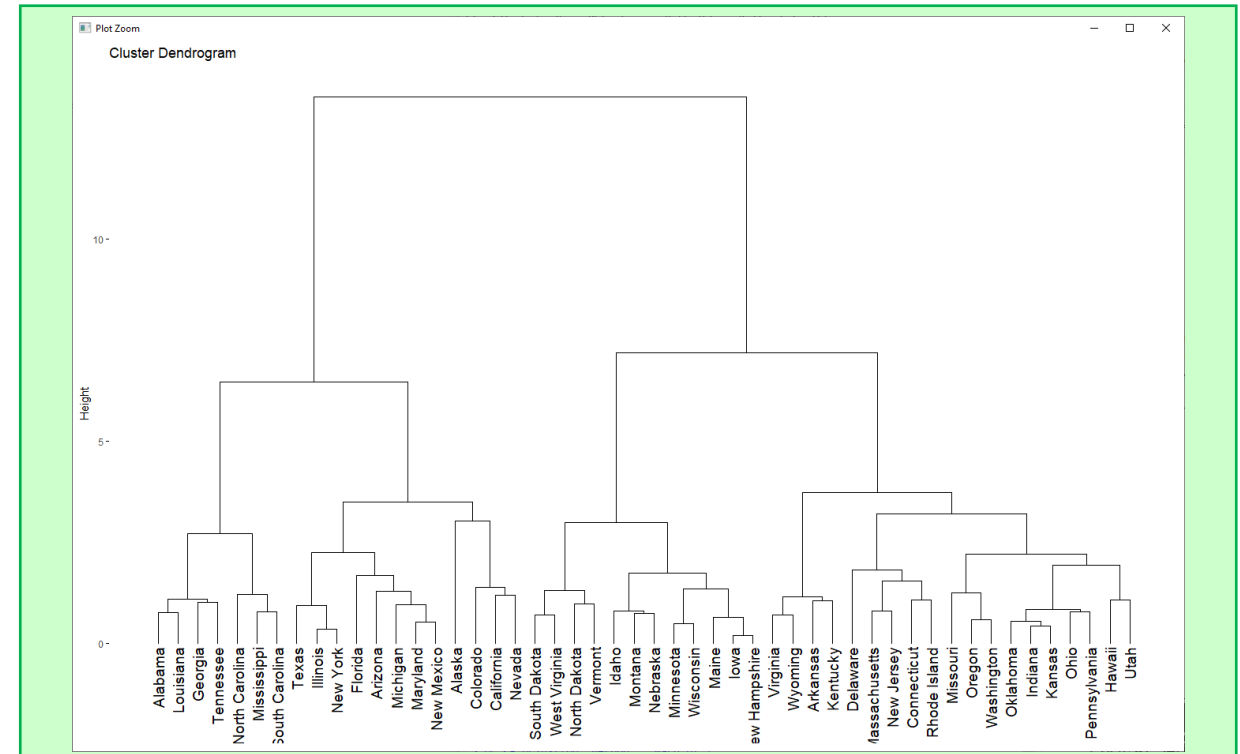
Interactive heatmap

We will begin by recreating the dendrogram with the *fviz_dend* function

Create the default dendrogram

`fviz_dend(hc, cex = 1)`

This changes the size of the labels



ENHANCING YOUR DENDROGRAMS

Circular plot

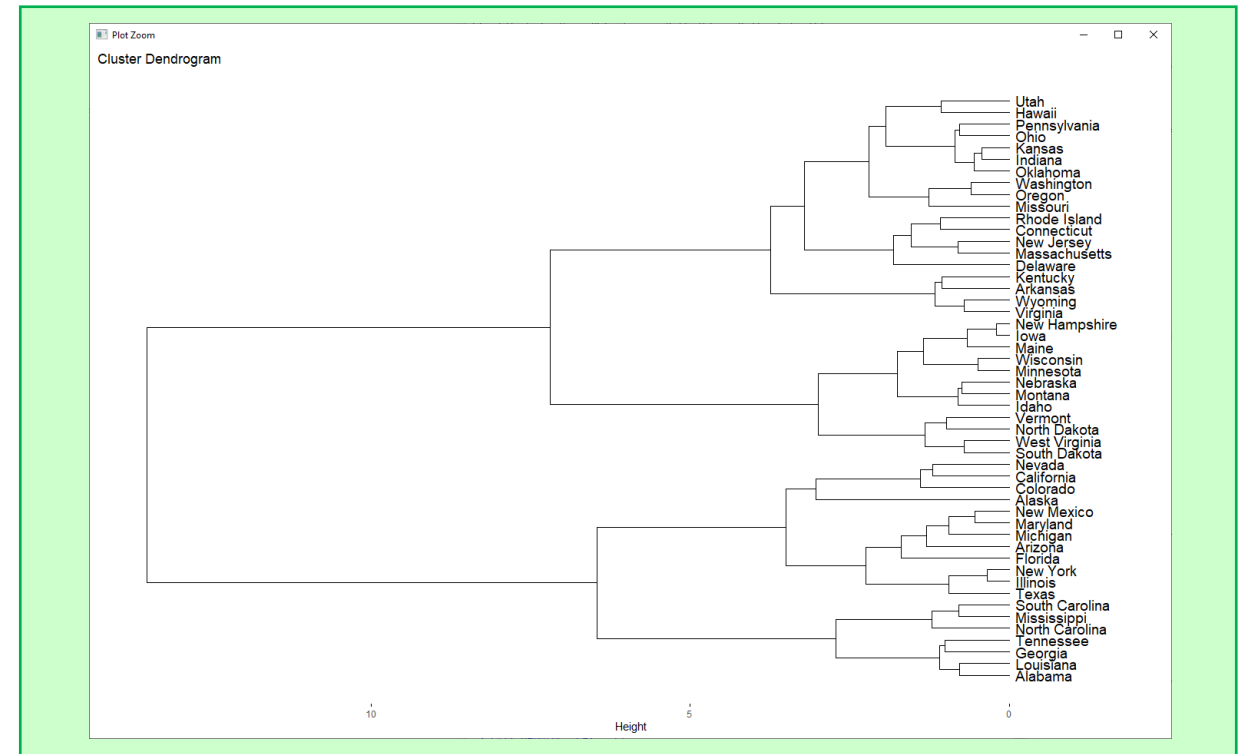
Phylogenetic tree plot

Heatmap

Interactive heatmap

We can rotate the orientation of the dendrogram by adding the *horiz* command

```
# Create a rotated dendrogram  
fviz_dend(hc, cex = 1, horiz = TRUE)
```



ENHANCING YOUR DENDROGRAMS

Circular plot

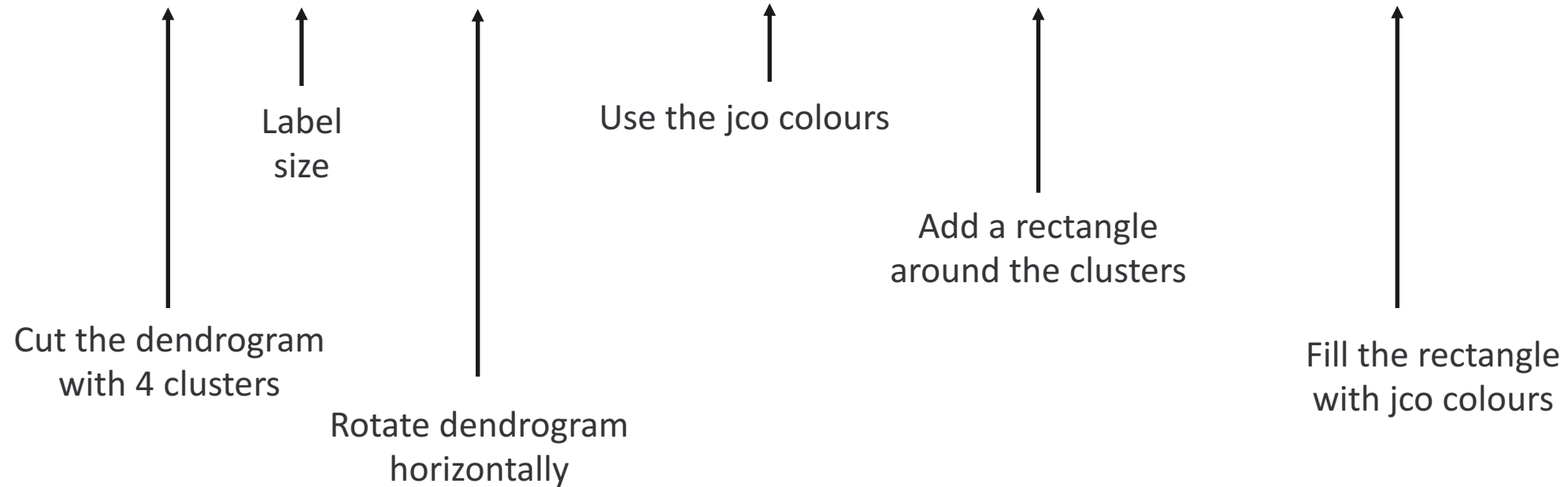
Phylogenetic tree plot

Heatmap

Interactive heatmap

Adding colour dendrogram

```
fviz_dend(hc, k = 4, cex = 0.8, horiz = TRUE, k_colors = "jco", rect = TRUE, rect_border = "jco", rect_fill = TRUE)
```



ENHANCING YOUR DENDROGRAMS

Circular plot

Phylogenetic tree plot

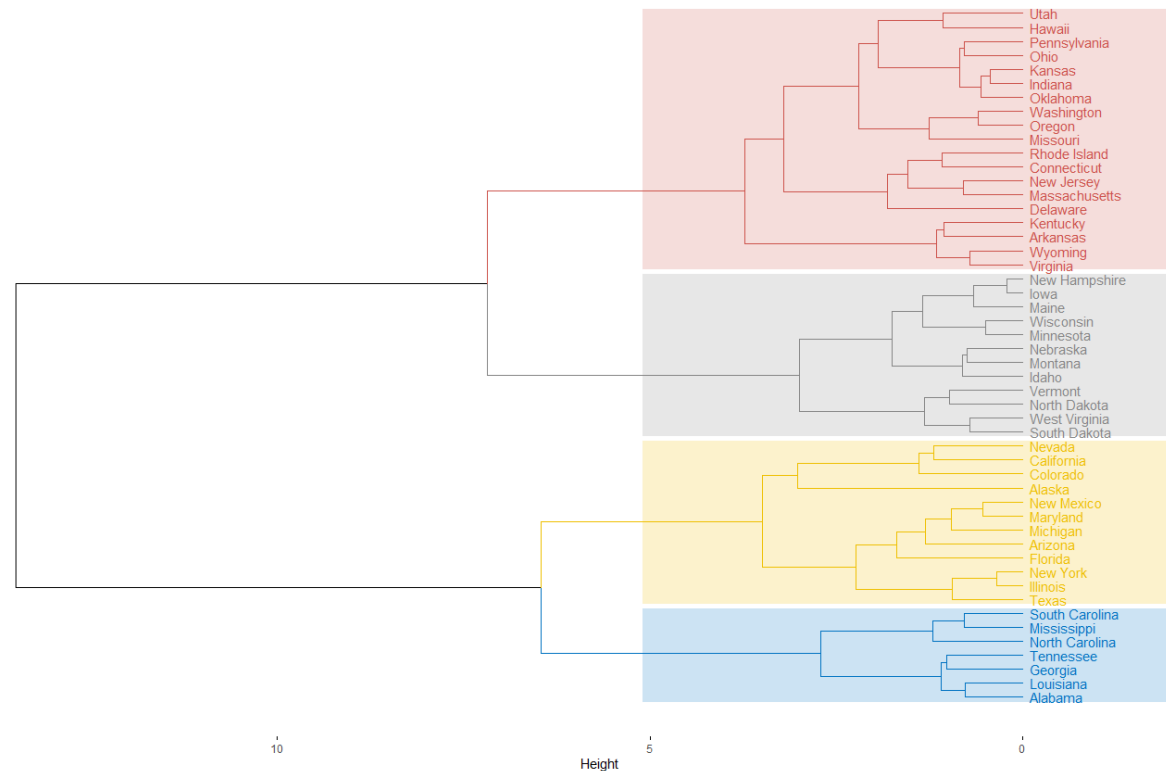
Heatmap

Interactive heatmap

Adding colour dendrogram

```
fviz_dend(hc, k = 4, cex = 0.8, horiz = TRUE, k_colors = "jco", rect = TRUE, rect_border = "jco", rect_fill = TRUE)
```

Cluster Dendrogram



Source: Kassambara (2017). Practical Guide To Cluster Analysis in R

ENHANCING YOUR DENDROGRAMS

Circular plot

Phylogenetic tree plot

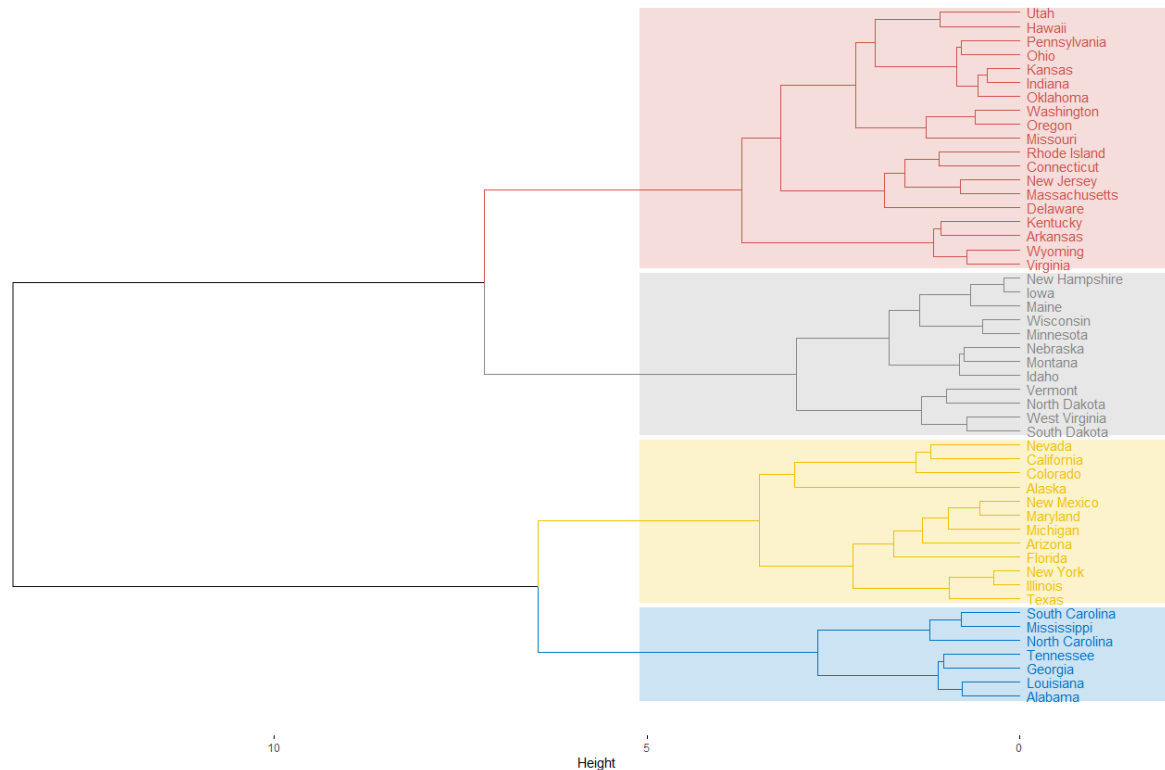
Heatmap

Interactive heatmap

Adding colour dendrogram

```
fviz_dend(hc, k = 4, cex = 0.8, horiz = TRUE, k_colors = "jco", rect = TRUE, rect_border = "jco", rect_fill = TRUE)
```

Cluster Dendrogram



Source: Kassambara (2017). Practical Guide To Cluster Analysis in R

ENHANCING YOUR DENDROGRAMS

Circular plot

Phylogenetic tree plot

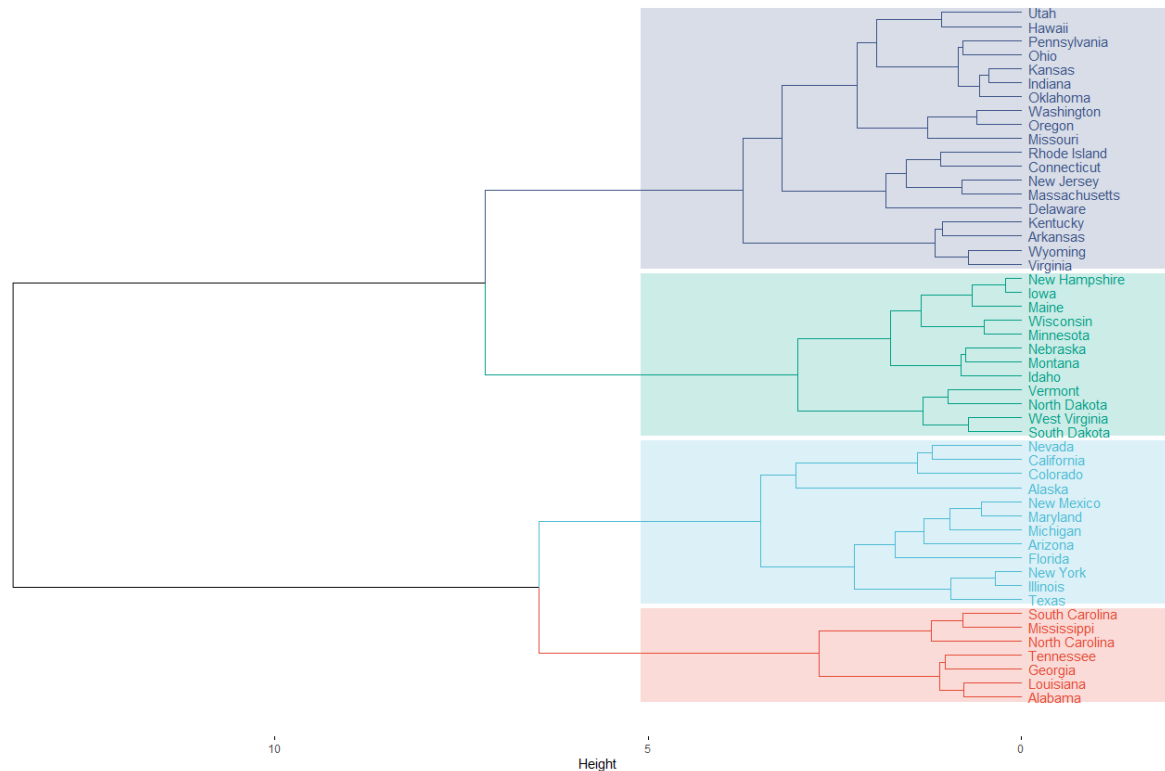
Heatmap

Interactive heatmap

Adding colour dendrogram

```
fviz_dend(hc, k = 4, cex = 0.8, horiz = TRUE, k_colors = "npg", rect = TRUE, rect_border = "npg", rect_fill = TRUE)
```

Cluster Dendrogram



ENHANCING YOUR DENDROGRAMS

Circular plot

Phylogenetic tree plot

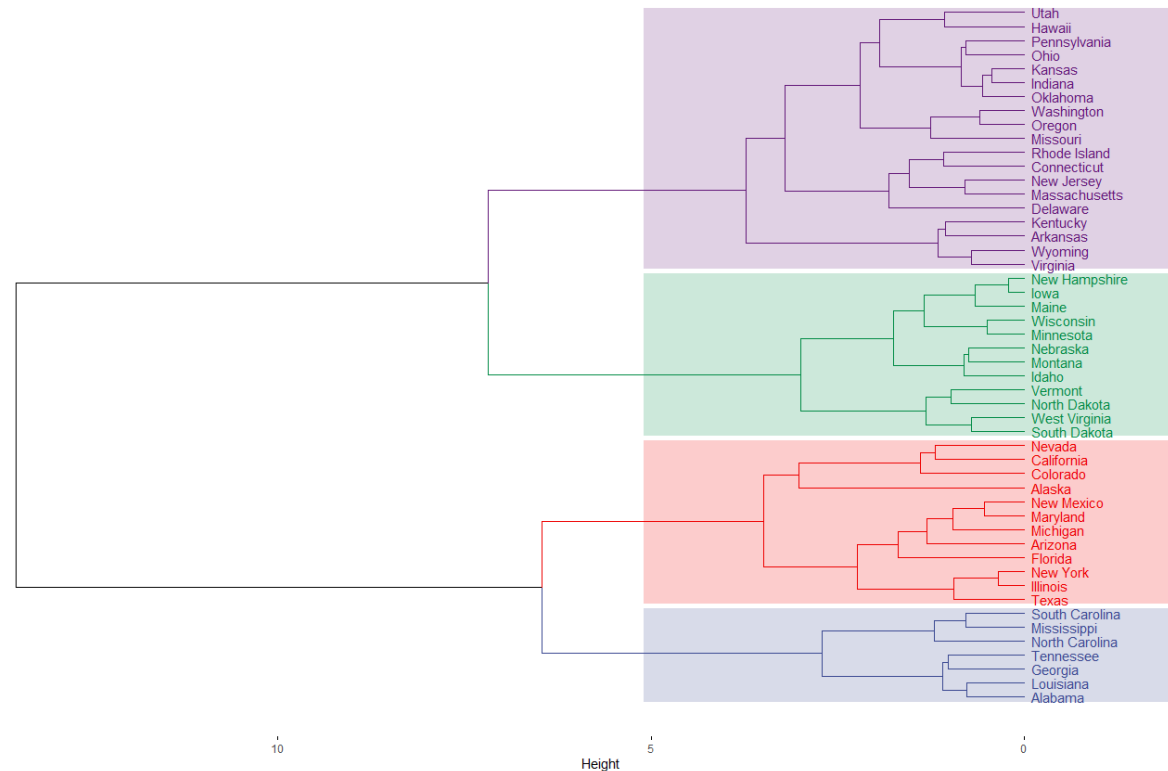
Heatmap

Interactive heatmap

Adding colour dendrogram

```
fviz_dend(hc, k = 4, cex = 0.8, horiz = TRUE, k_colors = "aaas", rect = TRUE, rect_border = "aaas", rect_fill = TRUE)
```

Cluster Dendrogram



ENHANCING YOUR DENDROGRAMS

Circular plot

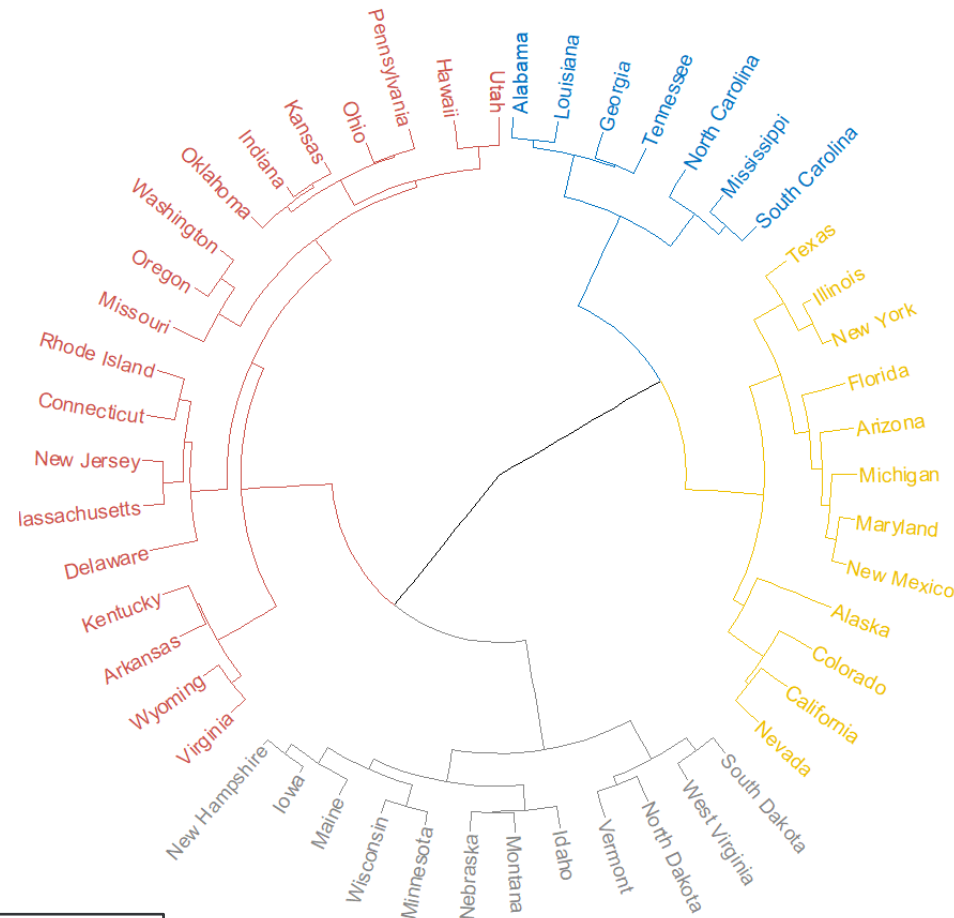
Phylogenetic tree plot

Heatmap

Interactive heatmap

Creating a circular dendrogram chart

```
fviz_dend(hc, cex = 1, k = 4, k_colors = "jco", type = "circular")
```



Source: Kassambara (2017). Practical Guide To Cluster Analysis in R

ENHANCING YOUR DENDROGRAMS

Circular plot

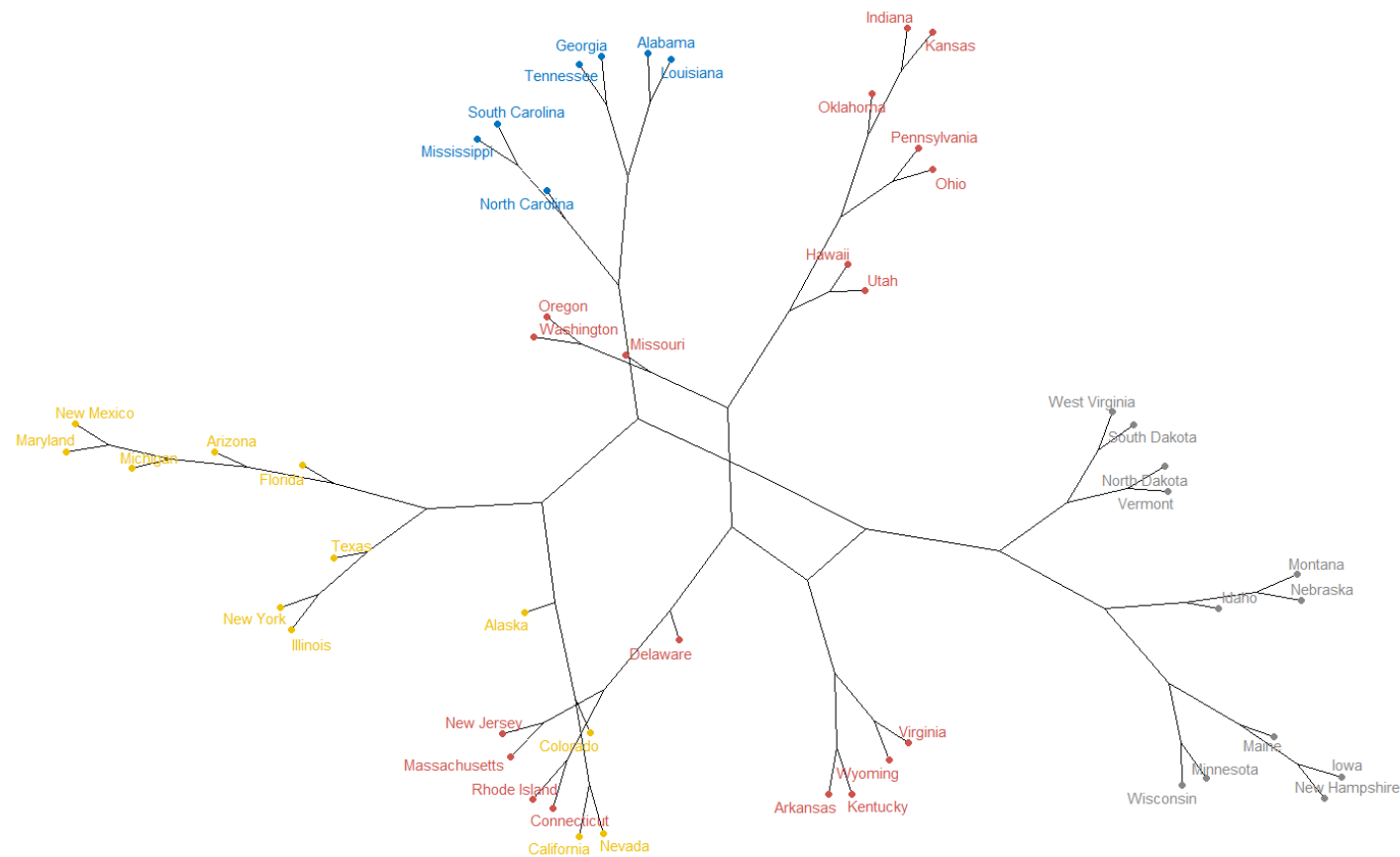
Phylogenetic tree plot

Heatmap

Interactive heatmap

Creating a circular dendrogram chart

```
fviz_dend(hc, cex = 1, k = 4, k_colors = "jco", type = "phylogenetic", repel = TRUE)
```



Source: Kassambara (2017). Practical Guide To Cluster Analysis in R

ENHANCING YOUR DENDROGRAMS

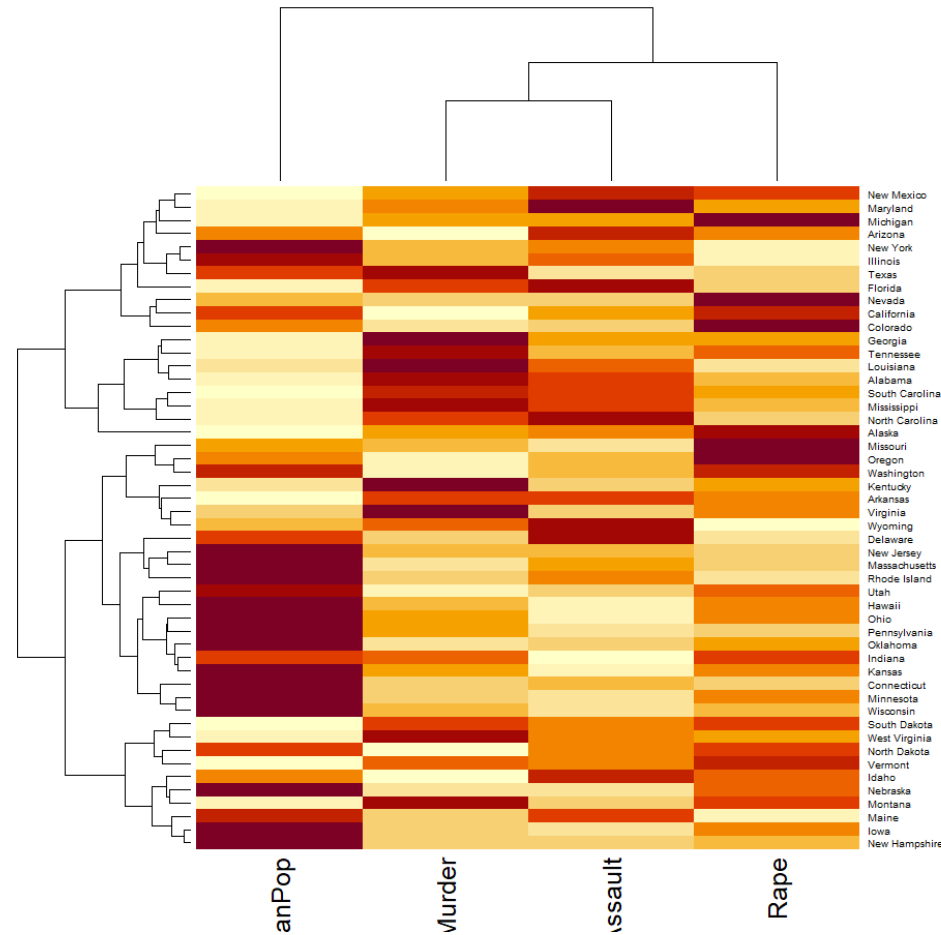
Circular plot

Phylogenetic tree plot

Heatmap

Interactive heatmap

```
# Create a heat map  
df<-scale(USArrests)  
heatmap(df)
```



Source: Kassambara (2017). Practical Guide To Cluster Analysis in R

ENHANCING YOUR DENDROGRAMS

Circular plot

Phylogenetic tree plot

Heatmap

Interactive heatmap

Install / Load the required packages

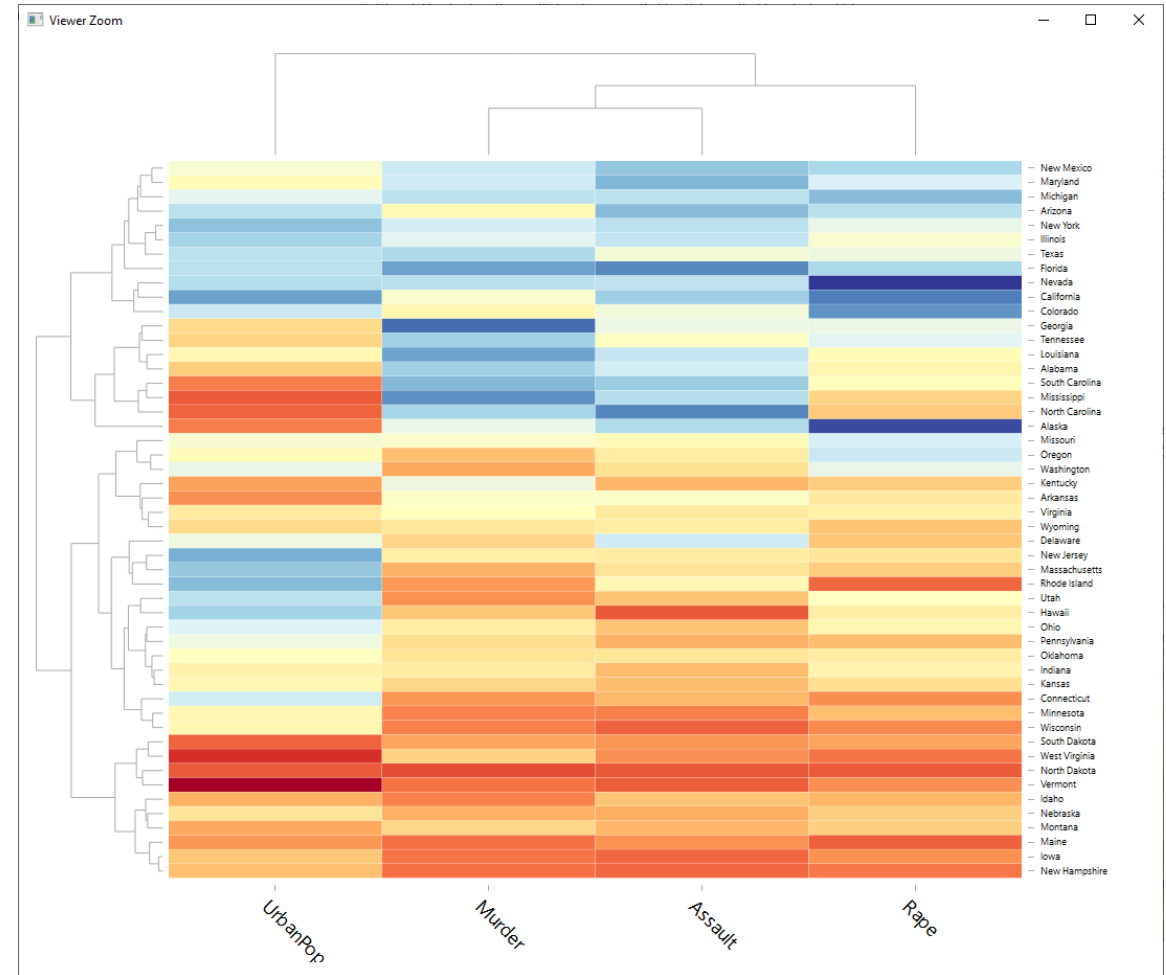
```
install.packages("d3heatmap")
```

```
library(d3heatmap)
```

Create the interactive heatmap

```
d3heatmap(df)
```

You can click on the states and / or variables which will highlight parts of the heat map and provide you with cluster values. You can also hover your mouse over parts of the heat map as well for more information. This is what makes this an interactive chart



TASKS

Now repeat these three tasks using the mtcars data file (this is also a base R data set)

Task A: K-means Clustering

Using the base R data file, USArrests, conduct a K-means cluster analysis and produce an appropriate visualisation to present your findings. Your visualisation should clearly display which US states belong to which cluster.

Task B: Hierarchical Clustering

Using the base R data file, USArrests, conduct a hierarchical (agglomerative) cluster analysis and produce a dendrogram to present your findings. You should compare different distance and linkage methods to find the optimal dendrogram for this scenario.

Task C: Enhancing your dendrograms

Recreate your dendrograms (from Task B) as a:

- Circular plot
- Phylogenetic tree plot
- Heatmap
- Interactive heatmap