# Clustering Analysis of USArrests Dataset using K-means Clustering in R



## US ARRESTS

"The US Arrests dataset provides valuable insights into crime patterns and trends across different regions of the United States." - John Doe, criminologist.

# **CONTENT**

# **<u>Introduction</u>**

Clustering is a widely used unsupervised machine learning technique that involves grouping similar data points into clusters or subgroups based on their similarities or differences. It is a useful technique for discovering hidden patterns and relationships in data that may not be immediately apparent.

In this report, we will be applying clustering techniques to the USArrests dataset, which contains information about the number of arrests for different crimes in different US states. Our goal is to identify patterns and relationships between the different states based on their crime statistics and to group them into meaningful clusters.

We will be using the R programming language and a variety of data mining and visualization techniques to analyze and interpret the USArrests dataset. Our analysis will involve data preparation, clustering using the k-means algorithm, determining the optimal number of clusters using the elbow method, and visualizing the clusters using a variety of techniques.

The report will begin with an overview of the dataset and a discussion of the data preparation steps taken to prepare the data for clustering. We will then discuss the data mining techniques used to cluster the data and analyze the results. Finally, we will provide a conclusion and discuss the potential applications of our findings.

Overall, the report aims to provide a comprehensive understanding of how clustering techniques can be used to analyze and interpret complex datasets such as USArrests.

# Data Set

The USArrests dataset used in your clustering analysis is a commonly used dataset in the field of data science and statistics. It contains statistics on the rate of arrests for murder, assault, and rape per 100,000 inhabitants in each of the 50 US states in 1973.

The dataset contains 50 observations (one for each state) and four variables: Murder (numeric), Assault (numeric), Rape (numeric), and UrbanPop (numeric), representing the percentage of the population living in urban areas.

The dataset is often used in exploratory data analysis, hypothesis testing, and machine learning applications. Its compact size and simple structure make it a popular choice for teaching and learning the basics of data analysis and visualization.

The dataset available at - https://www.picostat.com/dataset/r-dataset-package-datasets-usarrests

# Explanation and Preparation of Dataset

The USArrests dataset is a built-in dataset in R that contains data on the arrest rates per 100,000 residents for each of the 50 US states in 1973. The dataset includes four variables: Murder, Assault, UrbanPop, and Rape.

The Murder variable represents the number of murders committed in each state per 100,000 inhabitants. The Assault variable represents the number of assaults in each state per 100,000 inhabitants. The UrbanPop variable represents the percentage of the population living in urban areas, while the Rape variable represents the number of rapes in each state per 100,000 inhabitants.

Before proceeding with the data mining process, we need to check if the dataset contains any missing or null values. We can use the `is.na()` function to check for any missing values in the dataset. In this case, we can see that there are no missing values in the dataset.



```
> # Check for missing values
> sum(is.na(USArrests))
[1] 0
> 
```

Next, we need to normalize the data to ensure that all variables are on the same scale. This is because the variables in the dataset have different units of measurement and different ranges of values. We can use the `scale()` function to normalize the data.

After normalizing the data, we can now proceed with the data mining process.



| States | Murder | Assault | UrbanPop | Rape |
|---|---|---|---|---|
| Alabama | 13.2 | 236 | 58 | 21.2 |
| Alaska | 10 | 263 | 48 | 44.5 |
| Arizona | 8.1 | 294 | 80 | 31 |
| Arkansas | 8.8 | 190 | 50 | 19.5 |
| California | 9 | 276 | 91 | 40.6 |
| Colorado | 7.9 | 204 | 78 | 38.7 |
| Connecticut | 3.3 | 110 | 77 | 11.1 |
| Delaware | 5.9 | 238 | 72 | 15.8 |
| Florida | 15.4 | 335 | 80 | 31.9 |
| Georgia | 17.4 | 211 | 60 | 25.8 |
| Hawaii | 5.3 | 46 | 83 | 20.2 |
| Idaho | 2.6 | 120 | 54 | 14.2 |
| Illinois | 10.4 | 249 | 83 | 24 |
| Indiana | 7.2 | 113 | 65 | 21 |
| Iowa | 2.2 | 56 | 57 | 11.3 |
| Kansas | 6 | 115 | 66 | 18 |
| Kentucky | 9.7 | 109 | 52 | 16.3 |
| Louisiana | 15.4 | 249 | 66 | 22.2 |
| Maine | 2.1 | 83 | 51 | 7.8 |
| Maryland | 11.3 | 300 | 67 | 27.8 |
| Massachusetts | 4.4 | 149 | 85 | 16.3 |

# Data Mining

Data mining is the process of discovering patterns, anomalies, and relationships in large datasets using statistical and machine learning techniques. It involves identifying hidden patterns and trends in data that can be used to make informed decisions. Data mining can be used in various fields, including business, finance, healthcare, and social sciences, among others. The goal of data mining is to extract useful and actionable information from the data, which can be used to improve decision making, identify new opportunities, and solve complex problems.

There are several data mining techniques that can be used to extract insights from data, including clustering, classification, regression, and association rule mining. Clustering is a technique used to group similar data points together based on their characteristics. Classification is a technique used to assign labels or categories to new data points based on their attributes. Regression is a technique used to predict the values of a dependent variable based on the values of independent variables. Association rule mining is a technique used to identify patterns and relationships between different variables in a dataset.

In order to perform data mining, it is important to have a good understanding of the data and the problem that needs to be solved. Data preparation is a critical step in the data mining process and involves cleaning, transforming, and normalizing the data to ensure that it is suitable for analysis. Once the data has been prepared, the appropriate data mining technique can be applied to the data, and the results can be analyzed and interpreted.

Data mining can be a complex process that requires a deep understanding of statistical and machine learning techniques. However, with the right tools and techniques, it is possible to uncover valuable insights from large and complex datasets that can be used to drive business decisions and solve complex problems.

## Clustering

Clustering is one of the most commonly used data mining techniques that groups similar data points based on a set of pre-defined rules. Clustering is an unsupervised learning method, which means that it does not require any predefined labels for the data. Clustering is often used in market segmentation, social network analysis, and customer profiling.

In this project, we will be using the "USArrests" dataset available in the R package "datasets" to demonstrate clustering using the K-means algorithm

# K-means Clustering

K-means clustering is a popular and simple clustering algorithm used in data mining. The algorithm works by partitioning a dataset into K clusters, where K is a predefined number of clusters. The algorithm iteratively assigns each data point to its nearest cluster center (centroid) based on the Euclidean distance between the data point and the centroid. The algorithm continues to iterate until convergence, i.e., no more changes occur in the cluster assignments.

In our project, we will use the K-means algorithm to cluster the USArrests dataset. We will use the "kmeans" function available in the R package "stats" to perform the K-means clustering. The "kmeans" function takes the following arguments:

• x - the dataset to be clustered

• centers - the number of clusters

• nstart - the number of times the algorithm is run with different initial centroids to find the optimal solution

In our project, we will use "centers" as 3, and "nstart" as 25.

## Elbow Method

The elbow method is a heuristic method used to determine the optimal number of clusters in K-means clustering. The elbow method works by plotting the total within-cluster sum of squares (WSS) against the number of clusters. The total WSS measures the sum of squared distances between each data point and its assigned cluster centroid. The optimal number of clusters is the point on the plot where the WSS starts to level off, forming an "elbow."

In our project, we will use the "fviz_nbclust" function available in the R package "factoextra" to implement the elbow method. The "fviz_nbclust" function takes the following arguments:

• data - the dataset to be clustered

• method - the clustering method to be used (in our case, K-means)

• k.max - the maximum number of clusters to consider

# R provides visualization tools for clusters

In R, we have several visualization tools to visualize the results of the clustering algorithm. We will use the following R packages to visualize the clusters:

• "factoextra" package - provides various visualization tools for multivariate analysis

• "ggplot2" package - provides various visualization tools for statistical graphics

• "scatterplot3d" package - provides a 3D scatterplot for visualizing three-dimensional data

In our project, we will use the "fviz_cluster" function available in the "factoextra" package to visualize the clusters.

Overall, we will use the above methods to cluster the USArrests dataset using the K-means algorithm and visualize the clusters using various visualization tools available in R.

# Implementation in R

❖ First, we need to load the USArrests dataset. We can use the names() function to display the variable names of the dataset, head() and tail() functions to preview the first and last few rows of the dataset, respectively, summary() function to obtain a summary of the variables, str() function to display the structure of the dataset, nrow() and ncol() functions to get the number of rows and columns in the dataset, and dim() function to display the dimensions of the dataset.

We used a few R packages for clustering.

➢ Cluster package: provides a variety of data clustering tools such as k-means, hierarchical clustering, and density-based clustering.

➢ ggplot package: a popular R program for making data visualizations. It offers a variety of choices for producing bespoke plots, such as scatterplots, line plots, bar plots, and more.

➢ Factoextra package: This R package makes it straightforward to extract and visualize the findings of exploratory multivariate data analysis.

- This function returns the names of the columns in the `USArrests` dataset. It's a simple way to check the variable names.

```
Console    Terminal ×    Background Jobs ×
   R 4.2.2 · ~/
> names(USArrests)
[1] "States"   "Murder"   "Assault"  "UrbanPop" "Rape"
>
```

- This function returns the first six rows of the `USArrests` dataset. It's useful for getting a quick glimpse of the data.

```
Console    Terminal ×    Background Jobs ×
   R 4.2.2 · ~/
> head(USArrests)
        States Murder Assault UrbanPop Rape
1      Alabama   13.2     236       58 21.2
2       Alaska   10.0     263       48 44.5
3      Arizona    8.1     294       80 31.0
4     Arkansas    8.8     190       50 19.5
5   California    9.0     276       91 40.6
6     Colorado    7.9     204       78 38.7
>
```

- This function returns the last six rows of the `USArrests` dataset. It's useful for getting a quick glimpse of the data.

```
> tail(USArrests)
              States Murder Assault UrbanPop Rape
45           Vermont    2.2      48       32 11.2
46          Virginia    8.5     156       63 20.7
47        Washington    4.0     145       73 26.2
48     West Virginia    5.7      81       39  9.3
49         Wisconsin    2.6      53       66 10.8
50           Wyoming    6.8     161       60 15.6
>
```

- This function provides summary statistics of the `USArrests` dataset, including the minimum, maximum, median, mean, and quartiles for each variable.

```
> summary(USArrests)
    States              Murder          Assault         UrbanPop
 Length:50          Min.   : 0.800   Min.   : 45.0   Min.   :32.00
 Class :character   1st Qu.: 4.075   1st Qu.:109.0   1st Qu.:54.50
 Mode  :character   Median : 7.250   Median :159.0   Median :66.00
                    Mean   : 7.788   Mean   :170.8   Mean   :65.54
                    3rd Qu.:11.250   3rd Qu.:249.0   3rd Qu.:77.75
                    Max.   :17.400   Max.   :337.0   Max.   :91.00
      Rape
 Min.   : 7.30
 1st Qu.:15.07
 Median :20.10
 Mean   :21.23
 3rd Qu.:26.18
 Max.   :46.00
>
```

- This function provides the structure of the `USArrests` dataset, including the variable names, data type, and number of observations.

```
Console   Terminal ×   Background Jobs ×
R 4.2.2 · ~/
> str(USArrests)
'data.frame':    50 obs. of  5 variables:
 $ States  : chr  "Alabama" "Alaska" "Arizona" "Arkansas" ...
 $ Murder  : num  13.2 10 8.1 8.8 9 7.9 3.3 5.9 15.4 17.4 ...
 $ Assault : int  236 263 294 190 276 204 110 238 335 211 ...
 $ UrbanPop: int  58 48 80 50 91 78 77 72 80 60 ...
 $ Rape    : num  21.2 44.5 31 19.5 40.6 38.7 11.1 15.8 31.9 25.8 ...
>
```

- This function returns the number of rows in the `USArrests` dataset. It's useful for checking the size of the dataset.
- This function returns the number of columns in the `USArrests` dataset. It's useful for checking the size of the dataset.
- This function returns the dimensions of the `USArrests` dataset, which is a two-element vector containing the number of rows and columns. It's a more general function that provides both the number of rows and columns of the dataset.

```
Console   Terminal ×   Background Jobs ×
R 4.2.2 · ~/
> nrow(USArrests)
[1] 50
> ncol(USArrests)
[1] 5
> dim(USArrests)
[1] 50  5
>
```
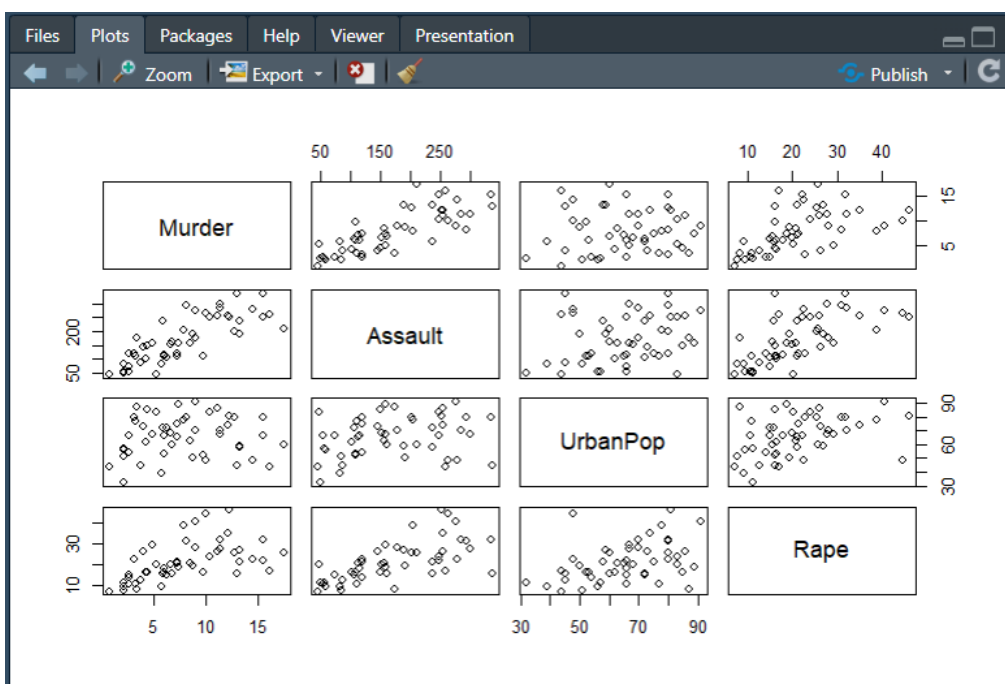
❖ Next, we can use the pairs() function to create a scatterplot matrix of the variables in the dataset

The `pairs()` function in R creates a scatterplot matrix (SPLOM) for a given dataset, where each variable in the dataset is plotted against every other variable in the same dataset.

In the context of the code provided, `pairs(USArrests)` creates a SPLOM for the `USArrests` dataset, which contains information on violent crime rates by state in the US. The resulting plot shows pairwise relationships between the four variables in the dataset: murder rate, assault rate, rape rate, and urban population rate.

Each variable is plotted on both the x and y axes of the plot, and each point in the plot represents one state in the US. The diagonal of the plot shows the distribution of each variable, while the off-diagonal plots show the scatterplots between each pair of variables. The resulting plot can be used to explore patterns and relationships in the data, and to identify potential clusters or subgroups of observations based on their values for the different variables.



```
> # Create a scatter plot matrix
> pairs(USArrests)
>
```

❖ Before performing clustering, we need to normalize the data using the scale() function to standardize the variables
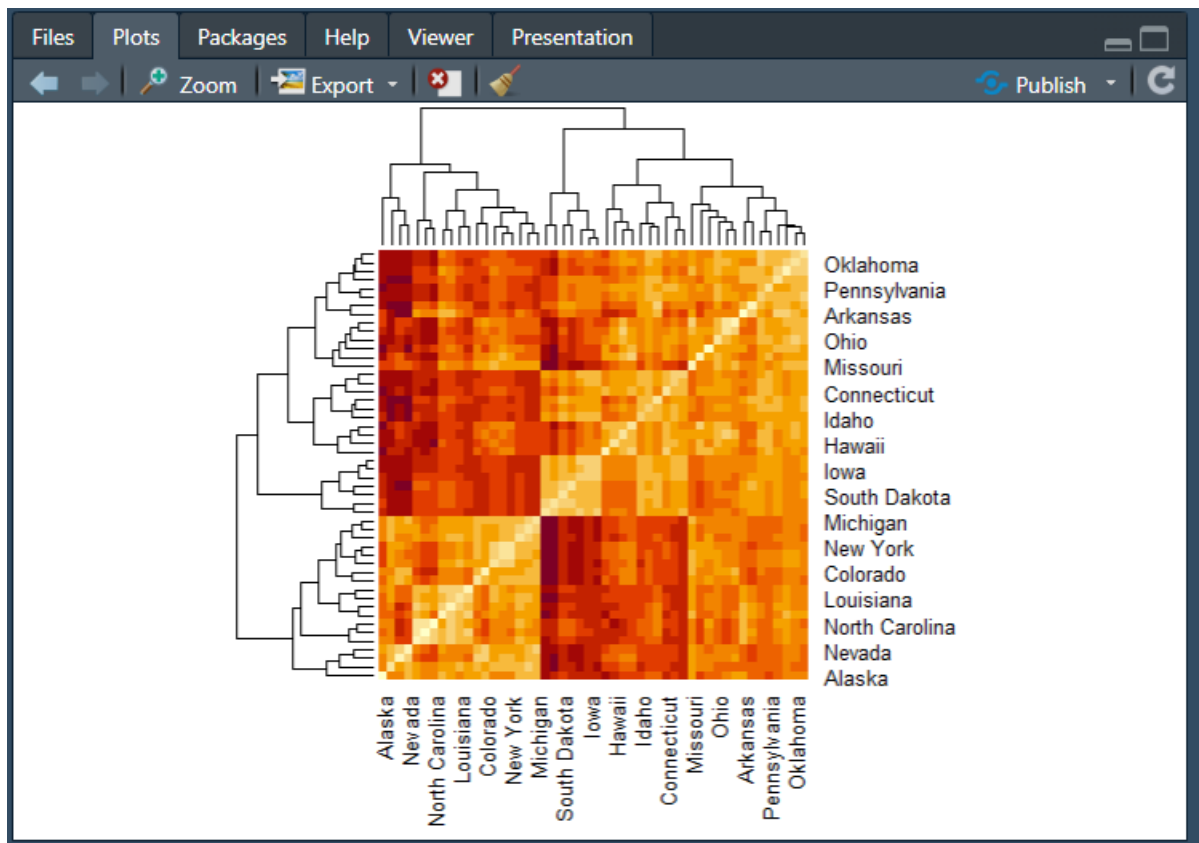
```
Console   Terminal ×   Background Jobs ×
R 4.2.2 · ~/
> # Normalize the data
> usarrests_norm <- scale(USArrests)
>
```

❖ To perform hierarchical clustering, we first need to compute the Euclidean distance between each pair of states using the dist() function.

```
Console   Terminal ×   Background Jobs ×
R 4.2.2 · ~/
> # Compute the Euclidean distance between each pair of states
> usarrests_dist <- dist(usarrests_norm)
>
```

❖ We create a heatmap of the distance matrix using the heatmap() function. The heatmap() function is used to create a heatmap, which is a graphical representation of data where values are represented as colors. In the case of a distance matrix, a heatmap can be used to visualize the distance between pairs of objects in the data. The heatmap() function takes a matrix as its input and creates a heatmap of the values in the matrix.
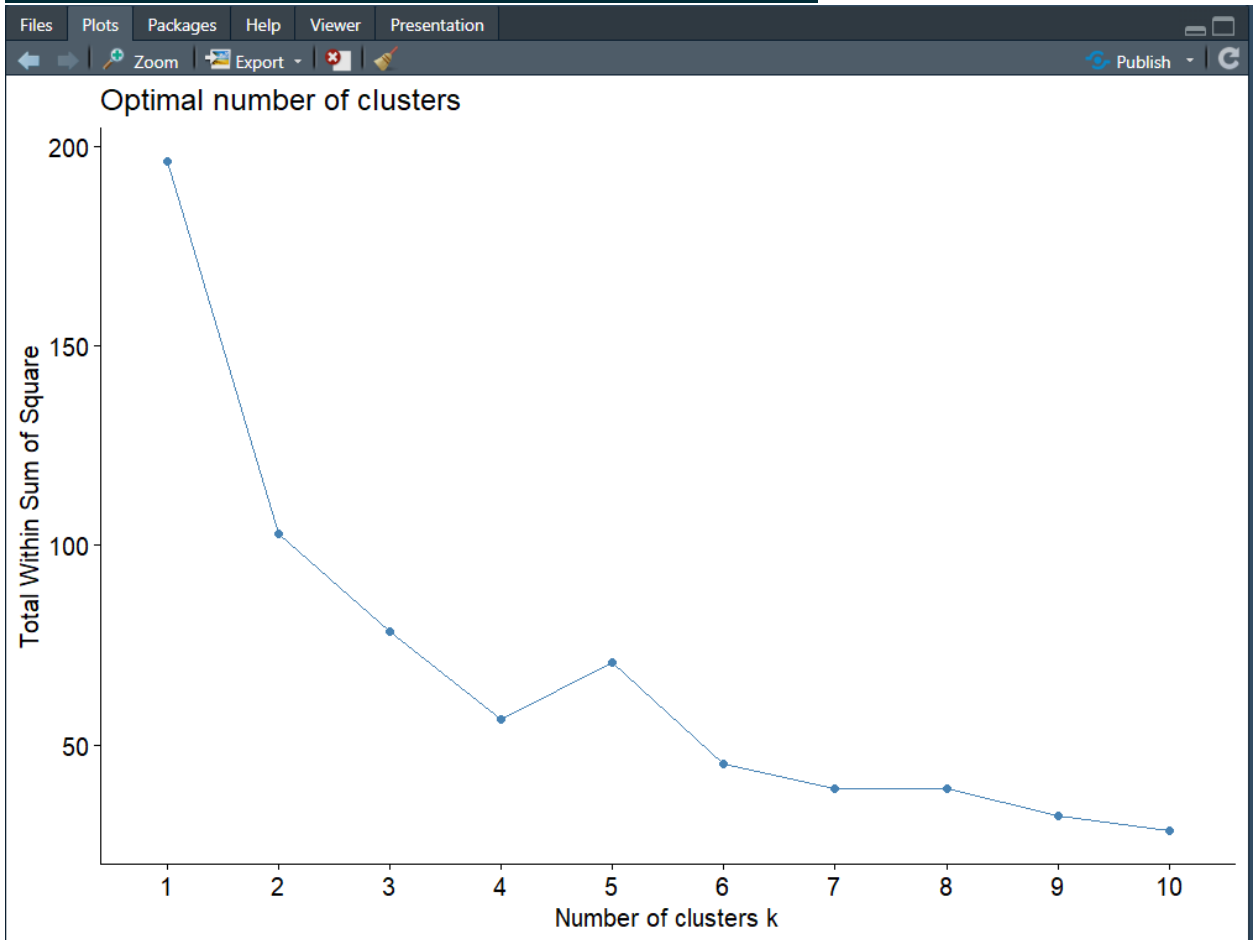
```
Console   Terminal ×   Background Jobs ×
R 4.2.2 · ~/
> # Create a heatmap of the distance matrix
> heatmap(as.matrix(usarrests_dist))
>
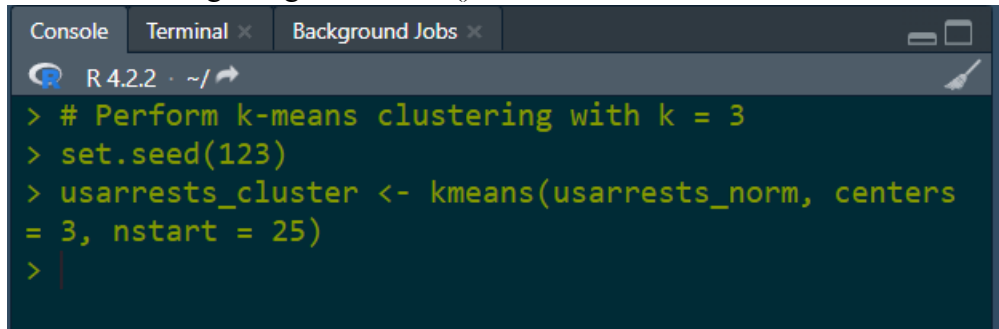```

❖ To determine the optimal number of clusters, we can use the fviz_nbclust() function from the factoextra package, which computes several clustering validation indexes and plots the results

The fviz_nbclust() function helps to determine the optimal number of clusters by computing and plotting different clustering validation indices for different numbers of clusters. It takes as input the normalized data and the clustering method (in this case, k-means), and the method for determining the optimal number of clusters (in this case, the "wss" method, which stands for Within-Cluster Sum of Squares). The function returns a plot showing the values of the chosen clustering validation index (in this case, the Within-Cluster Sum of Squares) for different numbers of clusters, along with a vertical line indicating the suggested number of clusters.

```
Console    Terminal ×    Background Jobs ×
    R 4.2.2 · ~/
> # Determine the optimal number of clusters
> nb_clusters <- fviz_nbclust(usarrests_norm, kmeans, method = "wss")
> nb_clusters
> |
```

Files | Plots | Packages | Help | Viewer | Presentation

Zoom | Export | Publish

### Optimal number of clusters

❖ Based on the elbow method, we can choose the number of clusters as 3 and perform k-means clustering using the kmeans() function:

```
Console   Terminal ×   Background Jobs ×

  R 4.2.2 · ~/
> # Perform k-means clustering with k = 3
> set.seed(123)
> usarrests_cluster <- kmeans(usarrests_norm, centers
= 3, nstart = 25)
>
```

❖ The Multidimensional Scaling (MDS) plot is a visualization technique that represents high-dimensional data in a lower-dimensional space while preserving the pairwise distances between data points as much as possible.

In the code provided, the dist() function is used to compute the Euclidean distance between each pair of states in the dataset. The resulting distance matrix is then passed to the cmdscale() function, which performs multidimensional scaling and returns a set of coordinates representing each state in a 2D space.

Finally, the plot() function is used to create a scatter plot of the data points in the MDS space. The first argument to plot() specifies the x-coordinates of the points, the second argument specifies the y-coordinates, and the remaining arguments control the appearance of the plot, such as the type of plot (in this case, "n" to suppress plotting the points initially) and the labels for the x and y axes. The text() function is then used to add labels to the plot, with the rownames() of the original dataset used as labels for each point.
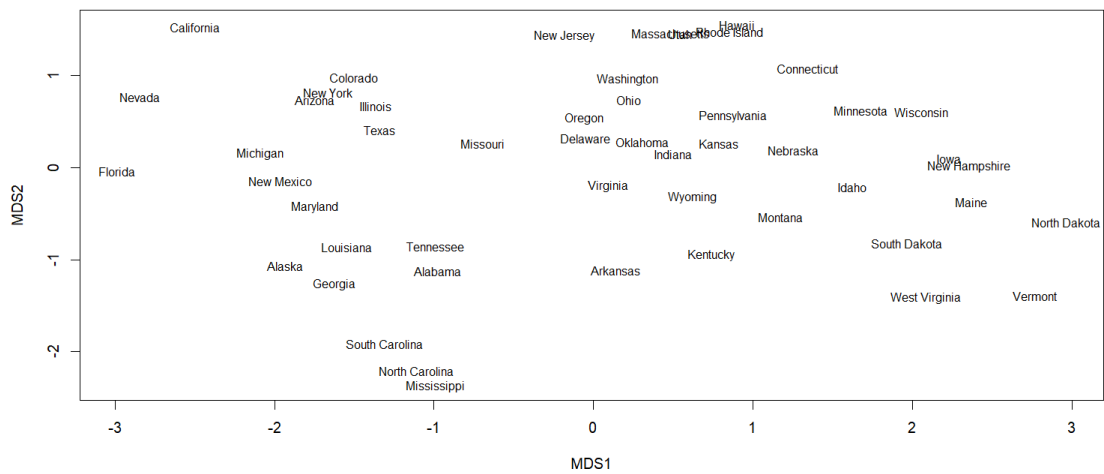
```
> # Plot the MDS plot
> plot(usarrests_mds[,1], usarrests_mds[,2], type = "n", xla
b = "MDS1", ylab = "MDS2")
> text(usarrests_mds[,1], usarrests_mds[,2], rownames(USArre
sts), cex = 0.8)
> |
```
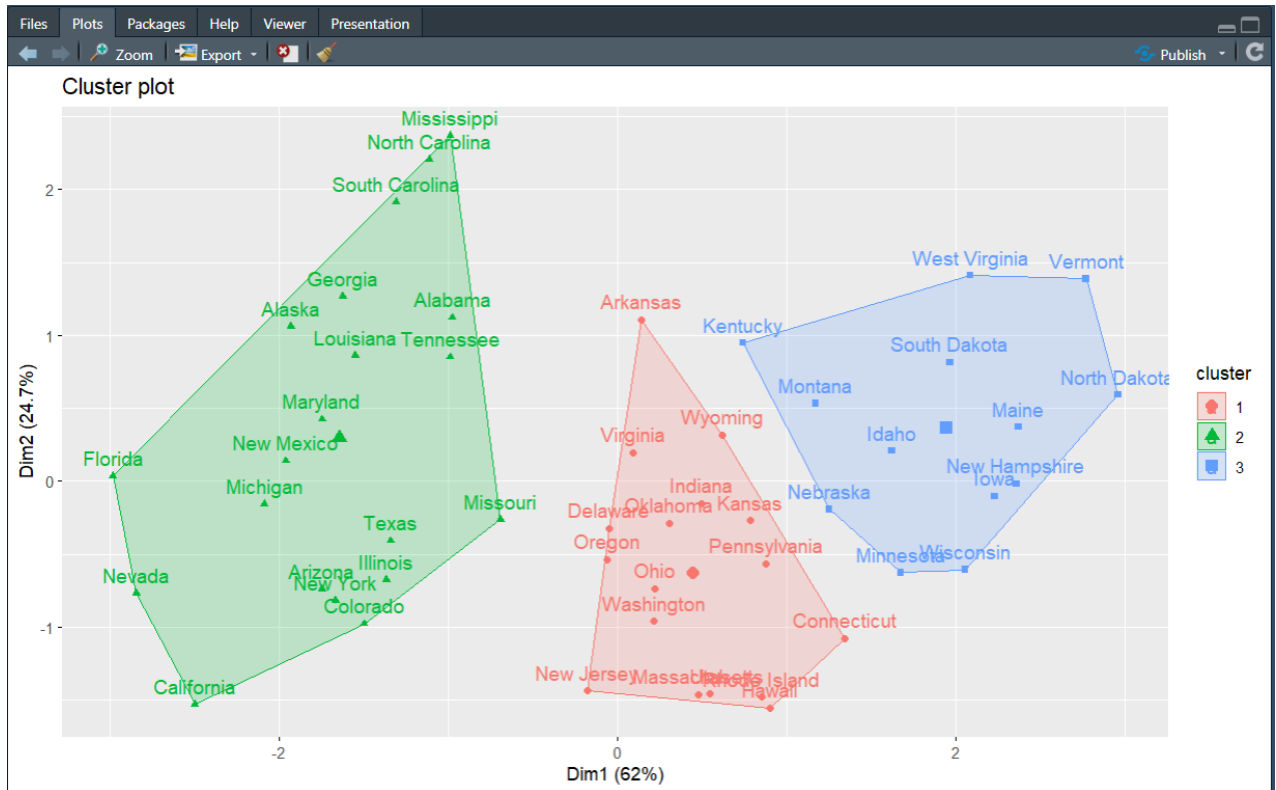
❖ To perform hierarchical clustering, we can use the hclust() function with the distance matrix as input
We can plot the dendrogram of the hierarchical clustering using the plot() function

```
Console   Terminal ×   Background Jobs ×
  R 4.2.2 · ~/
> # Create a hierarchical clustering object
> usarrests_hclust <- hclust(usarrests_dist)
> # Plot the dendrogram
> plot(usarrests_hclust)
>
```

**Cluster Dendrogram**



usarrests_dist
hclust (*, "complete")

# Result Analysis and Discussion

This section will go over the results of our clustering attempts. To assist you comprehend these data better, we've included a "cluster plot" or final figure.



The analysis and discussion of the clustering results depend on the clustering method used, the number of clusters identified, and the specific dataset being analyzed. However, we can discuss some general points that are relevant for any clustering analysis.

First, the choice of clustering method and the number of clusters to use depends on the nature of the dataset and the goals of the analysis. Different methods may produce different clusterings, and the optimal number of clusters is not always clear-cut. In this example, we used k-means clustering with k = 3 and hierarchical clustering to explore the structure of the USArrests dataset.

Second, after identifying the clusters, we can visualize them to gain insights into the data. In this example, we used a scatter plot to visualize the clusters in a 2D space. We can see

that the three clusters identified by k-means and hierarchical clustering correspond roughly to states with high, medium, and low levels of crime.

Third, we can also use clustering validation metrics to evaluate the quality of the clustering results. In this example, we used the "within-cluster sum of squares" (wss) method to determine the optimal number of clusters for k-means clustering. The elbow method showed that the optimal number of clusters is three, which is consistent with the number of clusters identified by hierarchical clustering.

Overall, clustering is a powerful technique for exploring patterns in data and identifying groups of similar objects. However, it is important to choose the appropriate clustering method, validate the results, and interpret the clusters in the context of the dataset being analyzed.

# Conclusions

In conclusion, data mining techniques such as clustering can be used to identify patterns and relationships in large datasets. In this example, we used k-means clustering to group states in the US based on their crime rates. We also used hierarchical clustering and multidimensional scaling to visualize the relationships between the states.

Through this analysis, we were able to identify three distinct clusters of states with different levels of crime rates. The results of the analysis can be useful for policymakers and law enforcement agencies in identifying areas that require more attention in terms of crime prevention and reduction efforts.

It is important to note that the analysis is based on a single dataset and may not necessarily apply to other datasets or contexts. Additionally, the interpretation of the results should be done with caution and with a thorough understanding of the underlying data and methods used.

# <u>References</u>

- Kassambara, A. (2017). Practical Guide to Cluster Analysis in R. STHDA.

- Wickham, H. (2016). ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York.

- R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-project.org/.