

Experiment 2.3

1. **Aim:** To perform the cluster analysis by k-means method using R.

2. **Objective:**

- To identify natural groupings or clusters within a dataset using the k-means clustering algorithm in R
- To apply the k-means clustering algorithm in R to a dataset with a known number of clusters, and to evaluate the effectiveness of the clustering method

3. **Script:**

K Means Clustering in R Programming is an Unsupervised Non-linear algorithm that cluster data based on similarity or similar groups. It seeks to partition the observations into a pre-specified number of clusters. Segmentation of data takes place to assign each training example to a segment called a cluster. In the unsupervised algorithm, high reliance on raw data is given with large expenditure on manual review for review of relevance is given. It is used in a variety of fields like Banking, healthcare, retail, Media, etc.

K-Means clustering groups the data on similar groups. The algorithm is as follows:

- Choose the number **K** clusters.
- Select at random K points, the centroids (Not necessarily from the given data).
- Assign each data point to closest centroid that forms K clusters.
- Compute and place the new centroid of each centroid.
- After final reassignment, name the cluster as Final cluster.

4. **Code:**

```
# Installing required packages
# ClusterR is an R package for cluster analysis and provides functions for k-means clustering,
hierarchical clustering, and more.
```

```
install.packages("ClusterR")
```

```
#The cluster package is an R package for cluster analysis, including functions for k-means
clustering, hierarchical clustering, and other algorithms.
```

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

```
# Loading packages
library(ClusterR) # Load ClusterR library
library(cluster) # Load cluster library
```

```
# Loading Seatbelts dataset
data(Seatbelts)
```

```
# Removing rows with missing values
Seatbelts_1 <- na.omit(Seatbelts[, -1]) # Remove rows with missing values in Seatbelts
dataset

# Fitting K-Means clustering Model to training dataset
set.seed(240) # Setting seed for reproducibility
kmeans.re <- kmeans(Seatbelts, centers = 3, nstart = 20) # Fit k-means clustering model to
Seatbelts dataset with 3 clusters and 20 starts

# Cluster identification for each observation
kmeans.re$cluster # Display the cluster identification for each observation

# Creating a confusion matrix
cm <- table(Seatbelts$front, kmeans.re$cluster) # Create a confusion matrix of Seatbelts
dataset and k-means clustering result
cm # Display the confusion matrix# Model Evaluation and visualization

# Plot drivers vs front for Seatbelts dataset
plot(Seatbelts[, c("drivers", "front")], ylim = c(0, max(Seatbelts_1$front)))

# Plot drivers vs front for Seatbelts dataset with cluster colors
plot(Seatbelts[, c("drivers", "front")], col = kmeans.re$cluster, ylim = c(0,
max(Seatbelts$front)))

# Plot drivers vs front for Seatbelts dataset with cluster colors and main title
plot(Seatbelts[, c("drivers", "front")], col = kmeans.re$cluster, main = "K-means with 3
clusters", ylim = c(0, max(Seatbelts_1$front)))

# Plotting cluster centers
# Display the cluster centers
kmeans.re$centers

# Display the cluster centers for drivers and front features
kmeans.re$centers[, c("drivers", "front")]

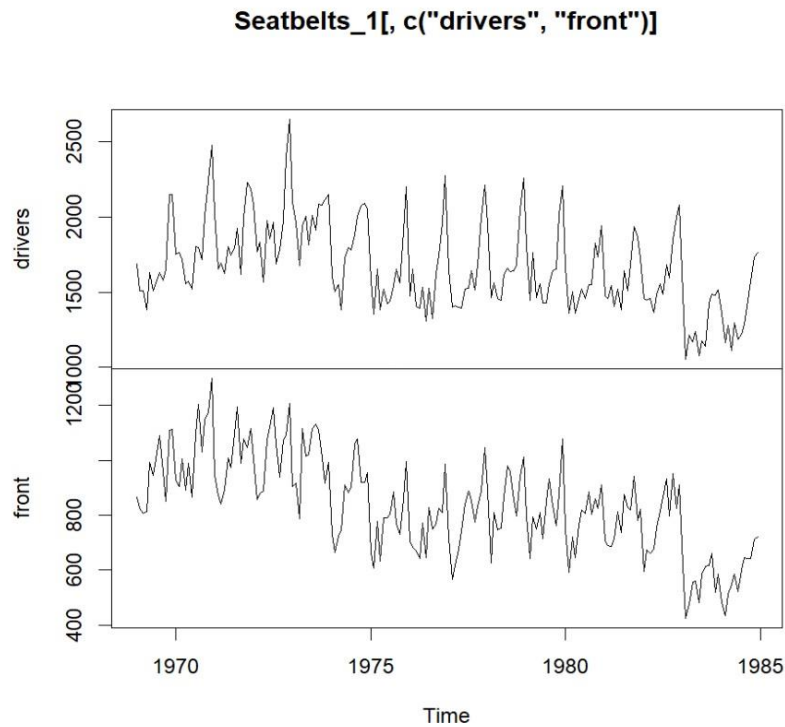
# Plot the cluster centers with different colors, shapes and size
points(kmeans.re$centers[, c("drivers", "front")], col = 1:3, pch = 8, cex = 3)

# Visualizing clusters
y_kmeans <- kmeans.re$cluster # Assign the cluster identification to y_kmeans variable
clusplot(Seatbelts[, c("drivers", "front")], # Plot a cluster plot of drivers vs front for
Seatbelts dataset

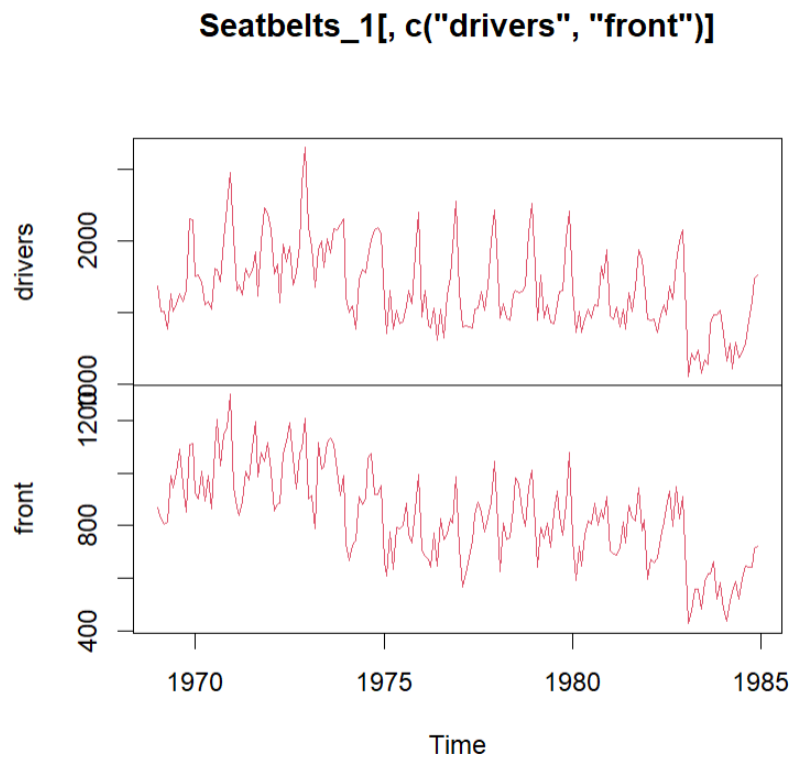
y_kmeans,
lines = 0,
shade = TRUE,
color = TRUE,
labels = 2,
plotchar = FALSE,
span = TRUE,
main = "Cluster Seatbelt",
xlab = 'drivers', ylab = 'front')
```

```
> cm
function (x)
2.54 * x
<bytecode: 0x0000017e18d2d878>
<environment: namespace:grDevices>
```

Plot drivers vs front for Seatbelts dataset

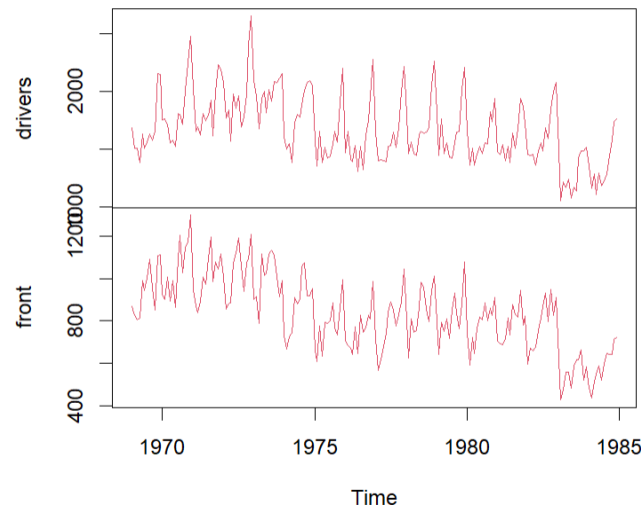


Plot drivers vs front for Seatbelts dataset with cluster colors



Plot drivers vs front for Seatbelts dataset with cluster colors and main title

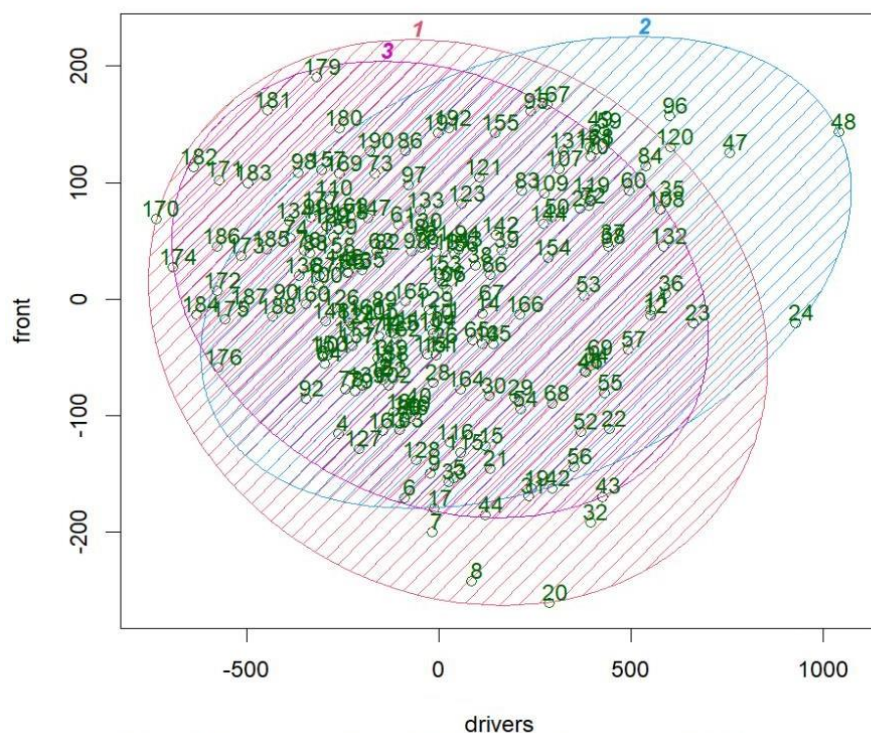
K-means with 3 clusters



```
> # Plotting cluster centers
> kmeans.re$centers
  drivers    front    rear    kms PetrolPrice VanKilled    law
1 1692.463 858.2388 401.2388 14798.22  0.1001847  9.059701 0.02985075
2 1815.172 896.1724 367.9483 11512.90  0.1003265 11.586207 0.00000000
3 1522.746 765.1642 429.9701 18202.13  0.1099178  6.865672 0.31343284
> kmeans.re$centers[, c("drivers", "front")]
  drivers    front
1 1692.463 858.2388
2 1815.172 896.1724
3 1522.746 765.1642
```

Plot a cluster plot of drivers vs front for Seatbelts dataset

Cluster Seatbelt



These two components explain 100 % of the point variability.



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

Discover. Learn. Empower.

