# Programming For Data Science Lab Assignment 4

## L33-L34

## Tejas Rokade

## 20BDS0033

**Q1. Read gender_classification dataset from R and perform following model fitting techniques (Any other datasets also applicable)**

**a. Logistic Regression**

**b.Decision Tree**

**c. Naïve Bayes**

**d.SVM**

**e.Random forest**

**Q2. Compare each of the above models using the following parameters**

**a.Accuracy**

**b.Precision**

**c. Recall**

**d.Sensitivity**

**e. Specificity**

**Ans 1 & 2.**

**Code:**

```
#FUNCTION TO FIND accuracy, precision, recall, Sensitivity, Specificity

classification_report_for_model <- function(conf_matr){

print(paste0("Accuracy: ", (

(conf_matr[1,1]+conf_matr[2,2])/(conf_matr[1,1]+conf_matr[1,2]+conf_matr[2,1]+conf_matr[2,2]) )))

  print("")
```

```r
    print(paste0("Precision: ", (conf_matr[2,2]/(conf_matr[2,2]+conf_matr[1,2]))))

    print("")

    print(paste0("Recall: ", (conf_matr[2,2]/(conf_matr[2,2]+conf_matr[2,1])) ))

    print("")

    print(paste0("Senstivity: ", (conf_matr[2,2]/(conf_matr[2,2]+conf_matr[2,1])) ))

    print("")

    print(paste0("Specificity: ", (conf_matr[1,1]/(conf_matr[1,1]+conf_matr[1,2]))))


}
Gender_class=read.csv("C:\\Users\\Parshva Maniar\\Downloads\\archive\\Transformed
Data Set - Sheet1.csv")

library(superml)

label=LabelEncoder$new()

Gender_class$Favorite.Color=label$fit_transform(Gender_class$Favorite.Color)

label=LabelEncoder$new()

Gender_class$Favorite.Music.Genre=label$fit_transform(Gender_class$Favorite.Music.
Genre)

label=LabelEncoder$new()

Gender_class$Favorite.Beverage=label$fit_transform(Gender_class$Favorite.Beverage)

label=LabelEncoder$new()

Gender_class$Favorite.Soft.Drink    =label$fit_transform(Gender_class$Favorite.Soft.Drink
)

label=LabelEncoder$new()

Gender_class$Gender=label$fit_transform(Gender_class$Gender)

Gender_class$Gender = factor(Gender_class$Gender, levels = c(0, 1))

head(Gender_class)

library(caTools)

#splitting data into train and test
```

```r
split=sample.split(Gender_class$Gender,SplitRatio =0.7)

train=subset(Gender_class,split==TRUE)

test=subset(Gender_class,split==FALSE)

library(e1071)

train[-5]=scale(train[-5])

test[-5]=scale(test[-5])

#LOGISTIC

model_log_reg =glm(Gender ~ ., data = train,family = "binomial")

model_log_reg

#Prediction

pred_lr = predict(model_log_reg, newdata = test[-5],type="response")

pred_lr <- ifelse(pred_lr>0.5,1,0)

pred_lr

# Making the Confusion Matrix

conf_matr = table(test[,5], pred_lr)

print(conf_matr)

classification_report_for_model(conf_matr)

#Desicion tree

library(party)

Desc_tree = ctree(Gender~., data = train)

pred_desc_tree = predict(Desc_tree, newdata = test[-5])

pred_desc_tree

# Making the Confusion Matrix

conf_matr = table(test[,5], pred_desc_tree)

print(conf_matr)

classification_report_for_model(conf_matr)

#Naive Bayes

naive_bayes =naiveBayes(Gender ~ ., data = train)
```

```r
naive_bayes
#Prediction
pred_nb = predict(naive_bayes, newdata = test[-5])
# Making the Confusion Matrix
conf_matr = table(test[,5], pred_nb)
print(conf_matr)
classification_report_for_model(conf_matr)
#SVM
svm_model = svm(formula = Gender ~ ., data = train, type = 'C-classification', kernel
        = 'linear')
svm_model
#Prediction
pred_svm = predict(svm_model, newdata = test[-5])
#Confusion Matrix
conf_matr = table(test[,5], pred_svm)
print(conf_matr)
classification_report_for_model(conf_matr)
#Random Forest
library(randomForest)
randomforest_model = randomForest(x =train[-5], y =train$Gender, ntree = 500)
randomforest_model
#predict
pred_rft = predict(randomforest_model, newdata = test[-5])
#Confusion Matrix
conf_matr = table(test[,5], pred_rft)
conf_matr
classification_report_for_model(conf_matr)
```
**Output:**

```
> #FUNCTION TO FIND accuracy, precision, recall, Sensitivity, Specificity
> classification_report_for_model <- function(conf_matr){
+   print(paste0("Accuracy: ", (
+     (conf_matr[1,1]+conf_matr[2,2])/(conf_matr[1,1]+conf_matr[1,2]+conf_matr[2,1]+conf_matr[2,2]) )))
+   print("")
+   print(paste0("Precision: ", (conf_matr[2,2]/(conf_matr[2,2]+conf_matr[1,2]))))
+   print("")
+   print(paste0("Recall: ", (conf_matr[2,2]/(conf_matr[2,2]+conf_matr[2,1])) ))
+   print("")
+   print(paste0("Senstivity: ", (conf_matr[2,2]/(conf_matr[2,2]+conf_matr[2,1])) ))
+   print("")
+   print(paste0("Specificity: ", (conf_matr[1,1]/(conf_matr[1,1]+conf_matr[1,2]))))
+
+ }
> Gender_class=read.csv("C:\\Users\\Parshva Maniar\\Downloads\\archive\\Transformed Data Set - Sheet1.csv")
> library(superml)
> label=LabelEncoder$new()
> Gender_class$Favorite.Color=label$fit_transform(Gender_class$Favorite.Color)
> label=LabelEncoder$new()
> Gender_class$Favorite.Music.Genre=label$fit_transform(Gender_class$Favorite.Music.Genre)
> label=LabelEncoder$new()
> Gender_class$Favorite.Beverage=label$fit_transform(Gender_class$Favorite.Beverage)
> label=LabelEncoder$new()
> Gender_class$Favorite.Soft.Drink =label$fit_transform(Gender_class$Favorite.Soft.Drink
+ )
> label=LabelEncoder$new()
> Gender_class$Gender=label$fit_transform(Gender_class$Gender)
> Gender_class$Gender = factor(Gender_class$Gender, levels = c(0, 1))
> head(Gender_class)
  Favorite.Color Favorite.Music.Genre Favorite.Beverage Favorite.Soft.Drink Gender
1              0                    0                 0                   0      0
2              1                    1                 0                   1      0
3              2                    0                 1                   1      0
4              2                    2                 2                   2      0
5              0                    0                 0                   1      0
6              2                    3                 3                   2      0
```

```
> library(caTools)
> #splitting data into train and test
> split=sample.split(Gender_class$Gender,SplitRatio =0.7)
> train=subset(Gender_class,split==TRUE)
> test=subset(Gender_class,split==FALSE)
> library(e1071)
> train[-5]=scale(train[-5])
> test[-5]=scale(test[-5])
> #LOGISTIC
> model_log_reg =glm(Gender ~ ., data = train,family = "binomial")
> model_log_reg

Call:  glm(formula = Gender ~ ., family = "binomial", data = train)

Coefficients:
      (Intercept)      Favorite.Color  Favorite.Music.Genre      Favorite.Beverage    Favorite.Soft.Drink
          0.01176             -0.24580               -0.10460                0.06401                0.62048

Degrees of Freedom: 45 Total (i.e. Null);  41 Residual
Null Deviance:      63.77
Residual Deviance: 59.83        AIC: 69.83
> #Prediction
> pred_lr = predict(model_log_reg, newdata = test[-5],type="response")
> pred_lr <- ifelse(pred_lr>0.5,1,0)
> pred_lr
 4  9 12 14 15 21 27 28 29 33 38 39 41 42 43 51 53 55 64 65
 0  0  0  0  1  1  1  1  0  0  1  1  0  0  0  0  1  1  1  0
> # Making the Confusion Matrix
> conf_matr = table(test[,5], pred_lr)
> print(conf_matr)
   pred_lr
    0 1
  0 6 4
  1 5 5
> classification_report_for_model(conf_matr)
[1] "Accuracy: 0.55"
[1] ""
[1] "Precision: 0.555555555555556"
[1] ""
[1] "Recall: 0.5"
[1] ""
[1] "Senstivity: 0.5"
[1] ""
[1] "Specificity: 0.6"
> #Desicion tree
> library(party)
> Desc_tree = ctree(Gender~., data = train)
> pred_desc_tree = predict(Desc_tree, newdata = test[-5])
> pred_desc_tree
 [1] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
Levels: 0 1
```

```
> conf_matr = table(test[,5], pred_desc_tree)
> print(conf_matr)
   pred_desc_tree
    0  1
  0 10  0
  1 10  0
> classification_report_for_model(conf_matr)
[1] "Accuracy: 0.5"
[1] ""
[1] "Precision: NaN"
[1] ""
[1] "Recall: 0"
[1] ""
[1] "Senstivity: 0"
[1] ""
[1] "Specificity: 1"
> #Naive Bayes
> naive_bayes =naiveBayes(Gender ~ ., data = train)
> naive_bayes

Naive Bayes Classifier for Discrete Predictors

Call:
naiveBayes.default(x = X, y = Y, laplace = laplace)

A-priori probabilities:
Y
  0   1
0.5 0.5

Conditional probabilities:
   Favorite.Color
Y        [,1]      [,2]
  0  0.06922995 1.007715
  1 -0.06922995 1.009923

   Favorite.Music.Genre
Y        [,1]      [,2]
  0  0.02049915 0.950848
  1 -0.02049915 1.067925

   Favorite.Beverage
Y        [,1]      [,2]
  0 -0.07356222 1.0306934
  1  0.07356222 0.9858047

   Favorite.Soft.Drink
Y        [,1]      [,2]
  0 -0.2530263 0.7778618
  1  0.2530263 1.1430314
```

```
> #Prediction
> pred_nb = predict(naive_bayes, newdata = test[-5])
> # Making the Confusion Matrix
> conf_matr = table(test[,5], pred_nb)
> print(conf_matr)
   pred_nb
    0 1
  0 7 3
  1 5 5
> classification_report_for_model(conf_matr)
[1] "Accuracy: 0.6"
[1] ""
[1] "Precision: 0.625"
[1] ""
[1] "Recall: 0.5"
[1] ""
[1] "Senstivity: 0.5"
[1] ""
[1] "Specificity: 0.7"
> #SVM
> svm_model = svm(formula = Gender ~ ., data = train, type = 'C-classification', kernel
+                = 'linear')
> svm_model

Call:
svm(formula = Gender ~ ., data = train, type = "C-classification", kernel = "linear")


Parameters:
   SVM-Type:  C-classification
 SVM-Kernel:  linear
       cost:  1

Number of Support Vectors:  42

> #Prediction
> pred_svm = predict(svm_model, newdata = test[-5])
> #Confusion Matrix
> conf_matr = table(test[,5], pred_svm)
> print(conf_matr)
   pred_svm
    0 1
  0 7 3
  1 6 4
```

```
> classification_report_for_model(conf_matr)
[1] "Accuracy: 0.55"
[1] ""
[1] "Precision: 0.571428571428571"
[1] ""
[1] "Recall: 0.4"
[1] ""
[1] "Senstivity: 0.4"
[1] ""
[1] "Specificity: 0.7"
> #Random Forest
> library(randomForest)
randomForest 4.7-1.1
Type rfNews() to see new features/changes/bug fixes.
> randomforest_model = randomForest(x =train[-5], y =train$Gender, ntree = 500)
> randomforest_model

Call:
 randomForest(x = train[-5], y = train$Gender, ntree = 500)
               Type of random forest: classification
                     Number of trees: 500
No. of variables tried at each split: 2

        OOB estimate of  error rate: 41.3%
Confusion matrix:
   0  1 class.error
0 16  7   0.3043478
1 12 11   0.5217391
> #predict
> pred_rft = predict(randomforest_model, newdata = test[-5])
> #Confusion Matrix
> conf_matr = table(test[,5], pred_rft)
> conf_matr
   pred_rft
    0 1
  0 6 4
  1 6 4
> classification_report_for_model(conf_matr)
[1] "Accuracy: 0.5"
[1] ""
[1] "Precision: 0.5"
[1] ""
[1] "Recall: 0.4"
[1] ""
[1] "Senstivity: 0.4"
[1] ""
[1] "Specificity: 0.6"
>
```

## Q3. Perform K Means Clustering on IRIS dataset

**Ans 3.**

**CODE:**

# Removing initial label of

x <- iris[, -5]

head(x)

# Fitting the clustering Model to the dataset

set.seed(240) # Setting seed

kmeans_clust <- kmeans(x, centers = 3, nstart = 20)

kmeans_clust

kmeans_clust$cluster

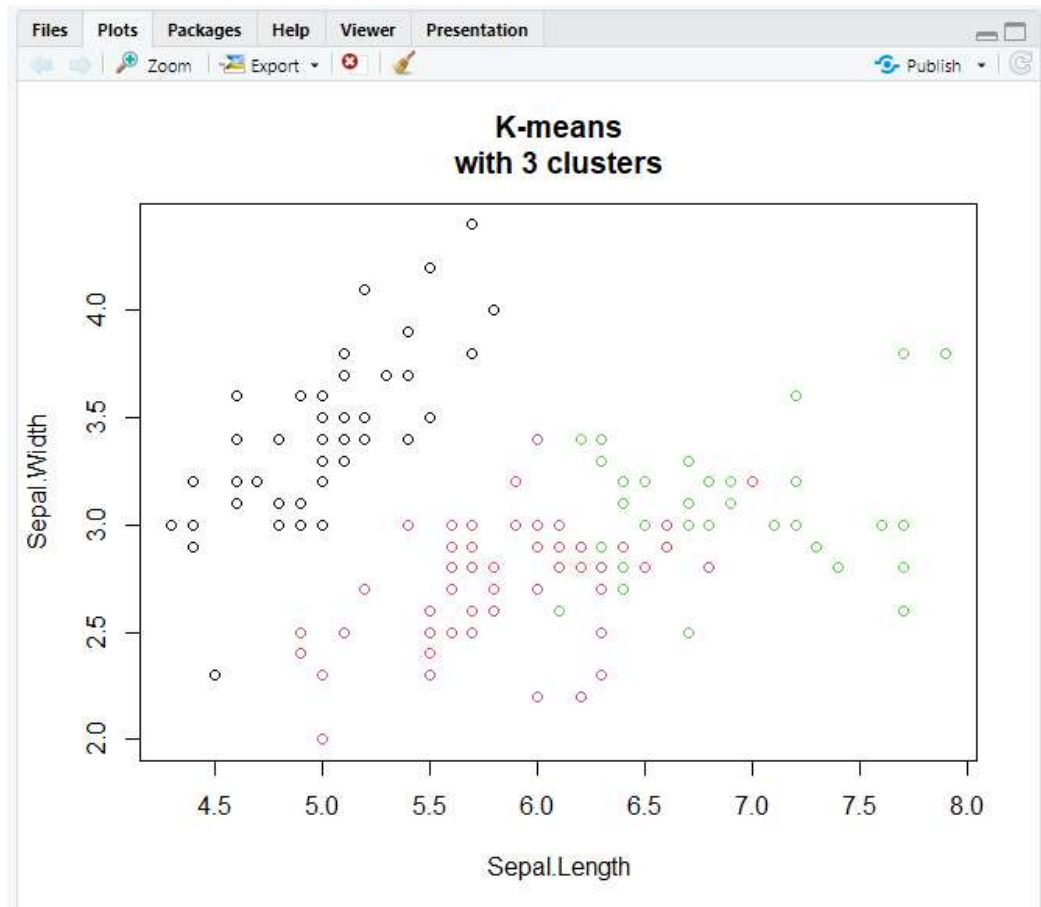# Confusion Matrix

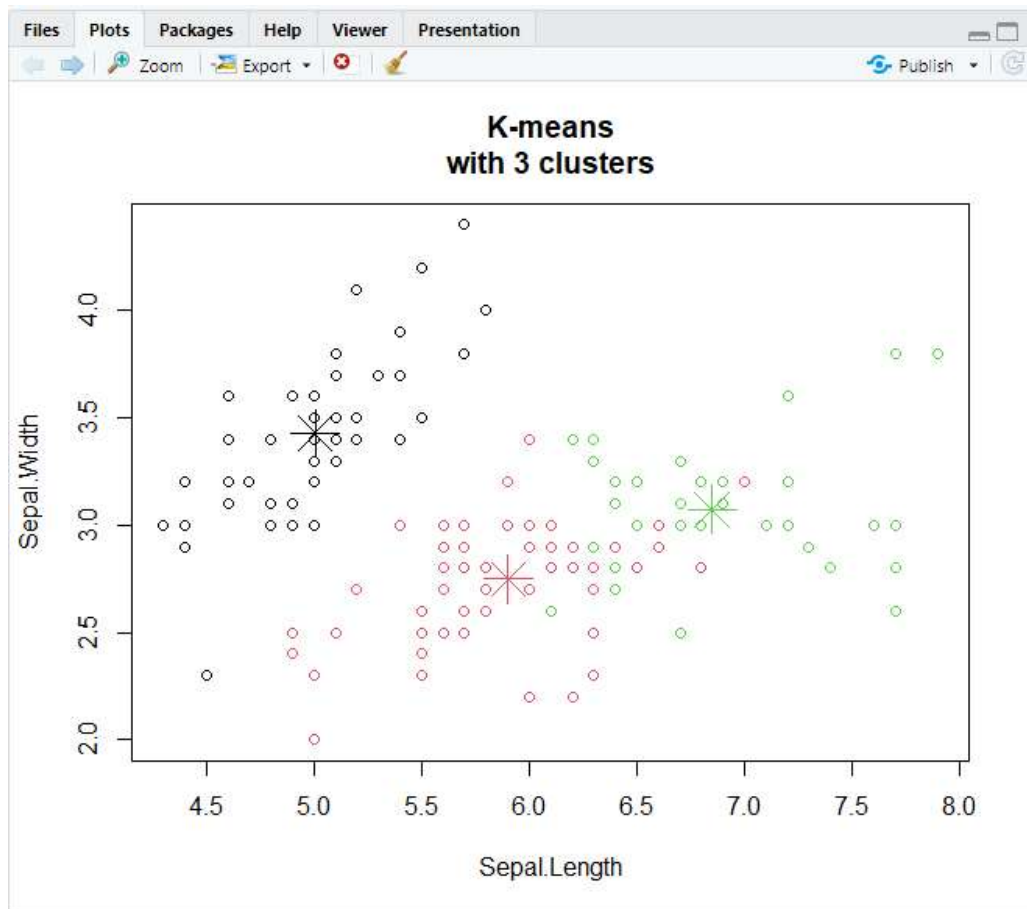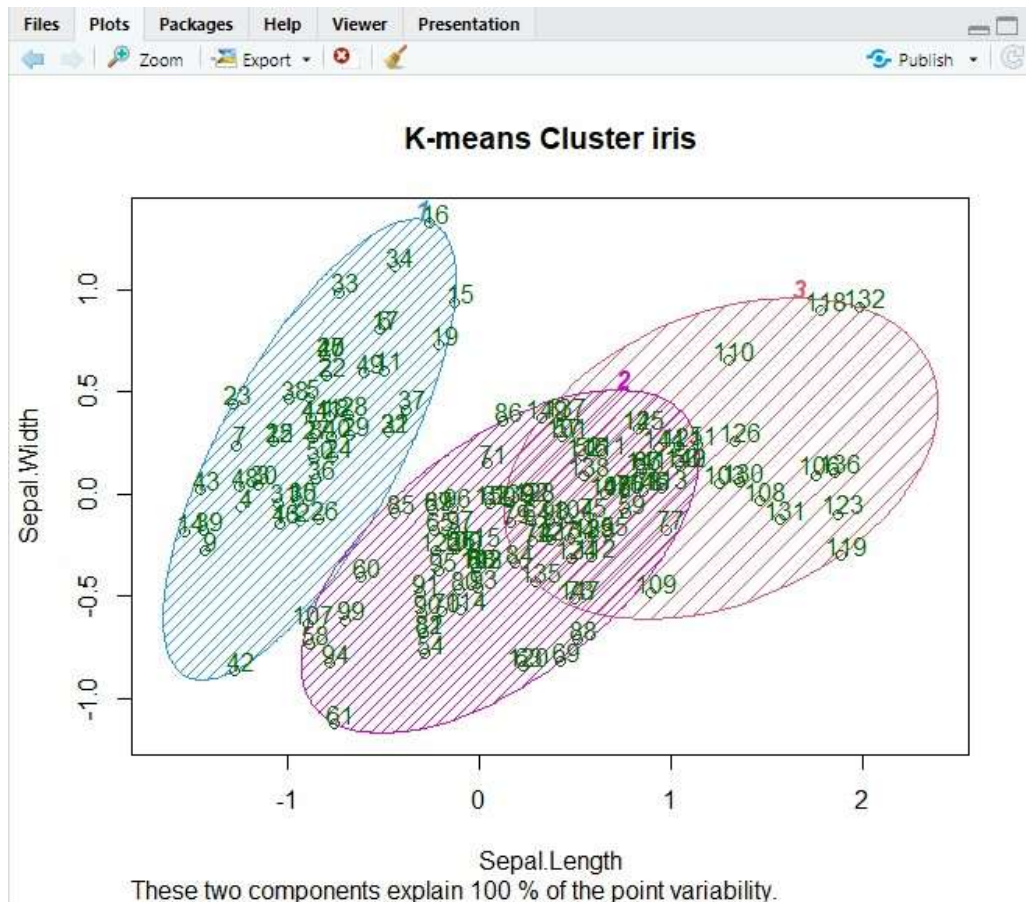cm <- table(iris$Species, kmeans_clust$cluster)

cm

# Visualization

```r
plot(x[c("Sepal.Length", "Sepal.Width")], col = kmeans_clust$cluster, main = "K-means
with 3 clusters")
# Plotting clusters with centres
kmeans_clust$centers
kmeans_clust$centers[, c("Sepal.Length", "Sepal.Width")]
points(kmeans_clust$centers[, c("Sepal.Length", "Sepal.Width")], col = 1:3, pch = 8,
    cex = 3)
# Visualizing clusters
y <- kmeans_clust$cluster
library(cluster)
clusplot(x[, c("Sepal.Length", "Sepal.Width")],y,lines = 0,shade = TRUE, color = TRUE,
    labels = 2, plotchar = FALSE, span = TRUE, main = paste("K-means Cluster iris"),
xlab
    = 'Sepal.Length', ylab = 'Sepal.Width')
```

**Output:**

```
> # Removing initial label of
> x <- iris[, -5]
> head(x)
  Sepal.Length Sepal.Width Petal.Length Petal.Width
1          5.1         3.5          1.4         0.2
2          4.9         3.0          1.4         0.2
3          4.7         3.2          1.3         0.2
4          4.6         3.1          1.5         0.2
5          5.0         3.6          1.4         0.2
6          5.4         3.9          1.7         0.4
> # Fitting the clustering Model to the dataset
> set.seed(240) # Setting seed
> kmeans_clust <- kmeans(x, centers = 3, nstart = 20)
> kmeans_clust
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

Clustering vector:
  [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 [71] 2 2 2 2 2 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 3 2 3 3 3 3 2 3 3 3 3 3 3 2 2 3 3 3 3 2 3 2 3 2 3 3 2 2 3 3 3 3 3 2 3 3 3 3 3 2
[141] 3 3 2 3 3 3 2 3 3 2

within cluster sum of squares by cluster:
[1] 15.15100 39.82097 23.87947
 (between_SS / total_SS =  88.4 %)

Available components:

[1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss" "betweenss"    "size"         "iter"         "ifault"
> kmeans_clust$cluster
  [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 [71] 2 2 2 2 2 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 3 2 3 3 3 3 2 3 3 3 3 3 3 2 2 3 3 3 3 2 3 2 3 2 3 3 2 2 3 3 3 3 3 2 3 3 3 3 3 2
[141] 3 3 2 3 3 3 2 3 3 2
> kmeans_clust$cluster
  [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 [71] 2 2 2 2 2 2 2 3 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 3 2 3 3 3 3 2 3 3 3 3 3 3 2 2 3 3 3 3 2 3 2 3 2 3 3 2 2 3 3 3 3 3 2 3 3 3 3 3 2
[141] 3 3 2 3 3 3 2 3 3 2
> # Confusion Matrix
> cm <- table(iris$Species, kmeans_clust$cluster)
> cm

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


> # Visualization
> plot(x[c("Sepal.Length", "Sepal.Width")], col = kmeans_clust$cluster, main = "K-means
+ with 3 clusters")
> # Plotting clusters with centres
> kmeans_clust$centers
  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
> kmeans_clust$centers[, c("Sepal.Length", "Sepal.Width")]
  Sepal.Length Sepal.Width
1     5.006000    3.428000
2     5.901613    2.748387
3     6.850000    3.073684
> points(kmeans_clust$centers[, c("Sepal.Length", "Sepal.Width")], col = 1:3, pch = 8,
+        cex = 3)
> # Visualizing clusters
> y <- kmeans_clust$cluster
> library(cluster)
> clusplot(x[, c("Sepal.Length", "Sepal.Width")],y,lines = 0,shade = TRUE, color = TRUE,
+          labels = 2, plotchar = FALSE, span = TRUE, main = paste("K-means Cluster iris"), xlab
+          = 'Sepal.Length', ylab = 'Sepal.Width')
> |
```

**Q4. Perform Hierarchical clustering in mtcars dataset**

**Ans 4**

**Code:**

```
mtcars_temp <- dist(mtcars, method = 'euclidean')

mtcars_temp

set.seed(240)

hierarchial_cluster_model <- hclust(mtcars_temp, method = "average")

hierarchial_cluster_model

# Plotting dendrogram

plot(hierarchial_cluster_model)

# Cutting tree by height
```

abline(h = 110, col = "green")

# Cutting by no. of clusters

final_fit <- cutree(hierarchial_cluster_model, k = 3 )

final_fit

table(final_fit)

rect.hclust(hierarchial_cluster_model, k = 3, border = "yellow")

**Output:**

```
Merc 280C
Merc 450SE
Merc 450SL
Merc 450SLC
Cadillac Fleetwood
Lincoln Continental
Chrysler Imperial
Fiat 128
Honda Civic
Toyota Corolla
Toyota Corona
Dodge Challenger
AMC Javelin
Camaro Z28
Pontiac Firebird
Fiat X1-9
Porsche 914-2
Lotus Europa          33.7678653
Ford Pantera L       288.5852993  297.5376920
Ferrari Dino          87.9105966   80.4553451   224.4587490
Maserati Bora        303.9222549  303.2796468    86.9383253  223.5342175
Volvo 142E            18.7555858   27.8104457   277.4803312   70.4751034  289.1157363
> set.seed(240)
> hierarchial_cluster_model <- hclust(mtcars_temp, method = "average")
> hierarchial_cluster_model

Call:
hclust(d = mtcars_temp, method = "average")

Cluster method   : average
Distance         : euclidean
Number of objects: 32

> # Plotting dendrogram
> plot(hierarchial_cluster_model)
> # Cutting tree by height
> abline(h = 110, col = "green")
> # Cutting by no. of clusters
> final_fit <- cutree(hierarchial_cluster_model, k = 3 )
> final_fit
          Mazda RX4        Mazda RX4 Wag          Datsun 710      Hornet 4 Drive   Hornet Sportabout              Valiant          Duster 360
                  1                    1                   1                   2                   2                    2                   2
          Merc 240D             Merc 230            Merc 280            Merc 280C           Merc 450SE           Merc 450SL          Merc 450SLC
                  1                    1                   1                   1                   2                    2                   2
 Cadillac Fleetwood  Lincoln Continental   Chrysler Imperial            Fiat 128          Honda Civic        Toyota Corolla        Toyota Corona
                  2                    2                   2                   1                   1                    1                   1
    Dodge Challenger          AMC Javelin           Camaro Z28     Pontiac Firebird            Fiat X1-9         Porsche 914-2         Lotus Europa
                  2                    2                   2                   2                   1                    1                   1
      Ford Pantera L         Ferrari Dino        Maserati Bora           Volvo 142E
                  2                    1                   3                   1
> table(final_fit)
final_fit
 1  2  3
16 15  1
> rect.hclust(hierarchial_cluster_model, k = 3, border = "yellow")
> |
```

Cluster Dendrogram