

# Assignment 2

Sri Chandana

2023-10-22

```
# It loads the class library  
# It loads the caret library
```

```
library(class)  
library(caret)
```

```
## Loading required package: ggplot2
```

```
## Loading required package: lattice
```

```
# It loads the e1071 library  
library(e1071)
```

```
# It imports the "UniversalBank" dataset from the specified file path  
# It displays the dimensions of the "Uni_Bank" dataset
```

```
Uni_Bank <- read.csv("C:/Users/srich/OneDrive/Desktop/R programming/UniversalBank.csv")  
dim(Uni_Bank)
```

```
## [1] 5000 14
```

```
# It displays the summary statistics for the "Uni_Bank" dataset  
summary(Uni_Bank)
```

```
##      ID      Age      Experience      Income      ZIP.Code  
## Min.   : 1    Min.   :23.00    Min.   :-3.0    Min.   : 8.00    Min.   : 9307  
## 1st Qu.:1251  1st Qu.:35.00    1st Qu.:10.0   1st Qu.: 39.00   1st Qu.:91911  
## Median :2500  Median :45.00    Median :20.0   Median : 64.00   Median :93437  
## Mean   :2500  Mean   :45.34    Mean   :20.1   Mean   : 73.77   Mean   :93153  
## 3rd Qu.:3750  3rd Qu.:55.00    3rd Qu.:30.0   3rd Qu.: 98.00   3rd Qu.:94608  
## Max.   :5000  Max.   :67.00    Max.   :43.0   Max.   :224.00   Max.   :96651  
##      Family      CCAvg      Education      Mortgage  
## Min.   :1.000    Min.   : 0.000    Min.   :1.000    Min.   : 0.0  
## 1st Qu.:1.000    1st Qu.: 0.700    1st Qu.:1.000    1st Qu.: 0.0  
## Median :2.000    Median : 1.500    Median :2.000    Median : 0.0  
## Mean   :2.396    Mean   : 1.938    Mean   :1.881    Mean   : 56.5  
## 3rd Qu.:3.000    3rd Qu.: 2.500    3rd Qu.:3.000    3rd Qu.:101.0  
## Max.   :4.000    Max.   :10.000    Max.   :3.000    Max.   :635.0  
## Personal.Loan  Securities.Account  CD.Account      Online
```

```
## Min. :0.000 Min. :0.0000 Min. :0.0000 Min. :0.0000
## 1st Qu.:0.000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000
## Median :0.000 Median :0.0000 Median :0.0000 Median :1.0000
## Mean :0.096 Mean :0.1044 Mean :0.0604 Mean :0.5968
## 3rd Qu.:0.000 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:1.0000
## Max. :1.000 Max. :1.0000 Max. :1.0000 Max. :1.0000
## CreditCard
## Min. :0.000
## 1st Qu.:0.000
## Median :0.000
## Mean :0.294
## 3rd Qu.:1.000
## Max. :1.000
```

```
# It removes columns 'ID' and 'ZIP.Code' from "Uni_Bank"
```

```
Uni_Bank$ID <- NULL
Uni_Bank$ZIP.Code <- NULL
```

```
# It shows the revised dataset summary after removing 'ID' and 'ZIP.Code' columns
summary(Uni_Bank)
```

```
##      Age      Experience      Income      Family
## Min. :23.00 Min. : -3.0 Min. : 8.00 Min. :1.000
## 1st Qu.:35.00 1st Qu.:10.0 1st Qu.: 39.00 1st Qu.:1.000
## Median :45.00 Median :20.0 Median : 64.00 Median :2.000
## Mean :45.34 Mean :20.1 Mean : 73.77 Mean :2.396
## 3rd Qu.:55.00 3rd Qu.:30.0 3rd Qu.: 98.00 3rd Qu.:3.000
## Max. :67.00 Max. :43.0 Max. :224.00 Max. :4.000
##      CCAvg      Education      Mortgage      Personal.Loan
## Min. : 0.000 Min. :1.000 Min. : 0.0 Min. :0.000
## 1st Qu.: 0.700 1st Qu.:1.000 1st Qu.: 0.0 1st Qu.:0.000
## Median : 1.500 Median :2.000 Median : 0.0 Median :0.000
## Mean : 1.938 Mean :1.881 Mean : 56.5 Mean :0.096
## 3rd Qu.: 2.500 3rd Qu.:3.000 3rd Qu.:101.0 3rd Qu.:0.000
## Max. :10.000 Max. :3.000 Max. :635.0 Max. :1.000
## Securities.Account CD.Account      Online      CreditCard
## Min. :0.0000 Min. :0.0000 Min. :0.0000 Min. :0.000
## 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.0000 1st Qu.:0.000
## Median :0.0000 Median :0.0000 Median :1.0000 Median :0.000
## Mean :0.1044 Mean :0.0604 Mean :0.5968 Mean :0.294
## 3rd Qu.:0.0000 3rd Qu.:0.0000 3rd Qu.:1.0000 3rd Qu.:1.000
## Max. :1.0000 Max. :1.0000 Max. :1.0000 Max. :1.000
```

```
# It converts "Education" column to a factor in Uni_Bank
```

```
Uni_Bank$Education <- as.factor(Uni_Bank$Education)
```

```
# It creates dummy variables for all columns in Uni_Bank
```

```
Dummy_Var <- dummyVars(~., data = Uni_Bank)
```

```
# It updates "Uni_Bank" with new dataset having dummy variables
```

```
Uni_updated <- as.data.frame(predict(Dummy_Var, Uni_Bank))
```

```
# Splits data into 60% training and 40% validation sets
# It creates training and validation datasets
```

```
set.seed(1)
train_data <- sample(row.names(Uni_updated), 0.6*dim(Uni_updated)[1])
valid_data <- setdiff(row.names(Uni_updated), train_data)
train_df <- Uni_updated[train_data,]
valid_df <- Uni_updated[valid_data,]
```

```
# It displays the summary statistics of the training dataset
summary(train_df)
```

```
##      Age      Experience      Income      Family
## Min.   :23.00  Min.   : -3.00  Min.    :  8.00  Min.    :1.000
## 1st Qu.:36.00  1st Qu.:10.00  1st Qu.: 39.00  1st Qu.:1.000
## Median :45.00  Median :20.00  Median : 63.00  Median :2.000
## Mean   :45.43  Mean   :20.19  Mean   : 73.08  Mean   :2.388
## 3rd Qu.:55.00  3rd Qu.:30.00  3rd Qu.: 98.00  3rd Qu.:3.000
## Max.   :67.00  Max.   :43.00  Max.   :224.00  Max.   :4.000
##      CCAvg      Education.1      Education.2      Education.3
## Min.    : 0.000  Min.    :0.0000  Min.    :0.000  Min.    :0.0000
## 1st Qu.: 0.700  1st Qu.:0.0000  1st Qu.:0.000  1st Qu.:0.0000
## Median : 1.500  Median :0.0000  Median :0.000  Median :0.0000
## Mean    : 1.915  Mean    :0.4173  Mean    :0.285  Mean    :0.2977
## 3rd Qu.: 2.500  3rd Qu.:1.0000  3rd Qu.:1.000  3rd Qu.:1.0000
## Max.    :10.000  Max.    :1.0000  Max.    :1.000  Max.    :1.0000
##      Mortgage      Personal.Loan      Securities.Account      CD.Account
## Min.    :  0.00  Min.    :0.00000  Min.    :0.0000  Min.    :0.00000
## 1st Qu.:  0.00  1st Qu.:0.00000  1st Qu.:0.0000  1st Qu.:0.00000
## Median :  0.00  Median :0.00000  Median :0.0000  Median :0.00000
## Mean    : 57.34  Mean    :0.09167  Mean    :0.1003  Mean    :0.05367
## 3rd Qu.:102.00  3rd Qu.:0.00000  3rd Qu.:0.0000  3rd Qu.:0.00000
## Max.    :635.00  Max.    :1.00000  Max.    :1.0000  Max.    :1.00000
##      Online      CreditCard
## Min.    :0.0000  Min.    :0.0000
## 1st Qu.:0.0000  1st Qu.:0.0000
## Median :1.0000  Median :0.0000
## Mean    :0.5847  Mean    :0.2927
## 3rd Qu.:1.0000  3rd Qu.:1.0000
## Max.    :1.0000  Max.    :1.0000
```

```
# It removes the 10th column which is 'Personal Income' from the training dataset
# Selecting only the 10th column from the validation dataset
```

```
train_normal_df <- train_df[,-10]
valid_normal_df <- valid_df[,10]
```

```
# It performs centering and scaling on the training data excluding 10th column
normal_values <- preProcess(train_df[,-10], method = c("center","scale"))
```

```
# It applies the centering and scaling transformation to the training and validation datasets
train_normal_df <- predict(normal_values, train_df[,-10])
valid_normal_df <- predict(normal_values, valid_df[,10])
```

#1 > Age = 40, Experience = 10, Income = 84, Family = 2, CCAvg = 2, Education\_1 = 0, Education\_2 = 1, Education\_3 = 0, Mortgage = 0, Securities Account = 0, CD Account = 0, Online = 1, and Credit Card = 1. Perform a k-NN classification with all predictors except ID and ZIP code using k = 1. Remember to transform categorical predictors with more than two categories into dummy variables first. Specify the success class as 1 (loan acceptance), and use the default cutoff value of 0.5. How would this customer be classified?

```
# It creates a data frame for a New_customer with these attributes
```

```
New_customer <- data.frame( Age = 40,
  Experience = 10,
  Income = 84,
  Family = 2,
  CCAvg = 2,
  Education.1 = 0,
  Education.2 = 1,
  Education.3 = 0,
  Mortgage = 0,
  Securities.Account = 0,
  CD.Account = 0,
  Online = 1,
  CreditCard = 1)
```

```
# It assigns the "New_customer" data to "New_customer_normal"
```

```
New_customer_normal <- New_customer
```

```
# Using "predict" function to transform the values in 'New_customer_normal' based on 'normal_values'
```

```
New_customer_normal <- predict(normal_values, New_customer_normal)
```

```
# It performs k-NN (k-Nearest Neighbors) classification with k=1.
```

```
# It displays the knn prediction1 result
```

```
knn.prediction1 <- class::knn(train = train_normal_df, test = New_customer_normal, cl = train_df$Personal.Loan)
knn.prediction1
```

```
## [1] 0
```

```
## Levels: 0 1
```

#2 > What is a choice of k that balances between overfitting and ignoring the predictor information?

```
# Creating "accuracy.df" data frame with k and overallaccuracy columns
```

```
# Performing k-NN (k-nearest neighbors) prediction using 'k' as the parameter.
```

```
# It stores the "overall accuracy" of the k-NN prediction in "accuracy.df".
```

```
# Finding the value of 'k' with the maximum overall accuracy
```

```
accuracy.df <- data.frame(k = seq(1, 15, 1), overallaccuracy = rep(0, 15))
```

```
for(i in 1:15)
```

```
{
```

```
  knn.prediction <- class::knn(train = train_normal_df,
    test = valid_normal_df,
    cl = train_df$Personal.Loan, k = i)
```

```
  accuracy.df[i, 2] <- confusionMatrix(knn.prediction,
    as.factor(valid_df$Personal.Loan), positive = "1")$overall[1]
```

```
}
```

```
which(accuracy.df[,2] == max(accuracy.df[,2]))
```



```
## [1555] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [1592] 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 1 1 0 0 0 0 0 0
## [1629] 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0
## [1666] 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0
## [1703] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0
## [1740] 0 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 1 1 0 0 0
## [1777] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [1814] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0
## [1851] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0
## [1888] 1 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [1925] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [1962] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [1999] 0 0
## Levels: 0 1
```

*# Calculating the confusion matrix for K-Nearest Neighbors prediction and then viewing the confusion matrix*

```
confusion.matrix <- confusionMatrix(knn.prediction2, as.factor(valid_df$Personal.Loan), positive = "1")
confusion.matrix
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 1786   63
##           1    9  142
##
##           Accuracy : 0.964
##           95% CI : (0.9549, 0.9717)
##           No Information Rate : 0.8975
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.7785
##
##           Mcnemar's Test P-Value : 4.208e-10
##
##           Sensitivity : 0.6927
##           Specificity : 0.9950
##           Pos Pred Value : 0.9404
##           Neg Pred Value : 0.9659
##           Prevalence : 0.1025
##           Detection Rate : 0.0710
##           Detection Prevalence : 0.0755
##           Balanced Accuracy : 0.8438
##
##           'Positive' Class : 1
##
```

#4 > Consider the following customer: Age = 40, Experience = 10, Income = 84, Family = 2, CCAvg = 2, Education\_1 = 0, Education\_2 = 1, Education\_3 = 0, Mortgage = 0, Securities Account = 0, CD Account = 0, Online = 1 and Credit Card = 1. Classify the customer using the best k.

```

# Creating a data frame for "New_customer1" with these attributes
New_customer1 <- data.frame(
  Age = 40,
  Experience = 10,
  Income = 84,
  Family = 2,
  CCAvg = 2,
  Education.1 = 0,
  Education.2 = 1,
  Education.3 = 0,
  Mortgage = 0,
  Securities.Account = 0,
  CD.Account = 0,
  Online = 1,
  CreditCard = 1
)

# Creating a new variable 'New_customer_normal1' and assigning the values of 'New_customer1' to it.
New_customer_normal1 <- New_customer1
New_customer_normal1 <- New_customer1

# Using the "predict" function to normalize the 'New_customer_normal1' data using normal_values.
New_customer_normal1 <- predict(normal_values, New_customer_normal1)

# Performing k-NN (k-Nearest Neighbors) classification on the normalized test data.
knn.prediction3 <- class::knn(train = train_normal_df,
                             test = New_customer_normal1,
                             cl= train_df$Personal.Loan, k= 3)

knn.prediction3

```

```

## [1] 0
## Levels: 0 1

```

#5 > Repartition the data, this time into training, validation, and test sets (50% : 30% : 20%). Apply the k-NN method with the k chosen above. Compare the confusion matrix of the test set with that of the training and validation sets. Comment on the differences and their reason.

```

set.seed(1)

# Splitting the data into 50% for training set, 30% for validation set and 20% for testing set.

# Sample training data

train_index1 <- sample(row.names(Uni_updated), 0.5*dim(Uni_updated)[1])
train_df1 <-Uni_updated[train_index1,]

# Creating a validation set by excluding the training data
valid_index1 <- setdiff(row.names(Uni_updated), train_index1)
valid_df1 <- Uni_updated[valid_index1, ]

# Splitting the validation set into a second validation set
valid_index2 <- sample(row.names(valid_df1), 0.6*dim(valid_df1)[1])
valid_df2 <- valid_df1[valid_index2, ]

```











```
cl= train_df1$Personal.Loan, k= 3)
```

[illegible]

```
# Calculating the confusion matrix for K-Nearest Neighbors prediction and then viewing the confusion ma
```

```
confusion_matrix3 <- confusionMatrix(knn_prediction6, as.factor(test_df1$Personal.Loan))
confusion_matrix3
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 884  35
##           1   4  77
##
##           Accuracy : 0.961
##           95% CI : (0.9471, 0.9721)
##           No Information Rate : 0.888
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.777
##
##           McNemar's Test P-Value : 1.556e-06
```

```
##
##      Sensitivity : 0.9955
##      Specificity : 0.6875
##      Pos Pred Value : 0.9619
##      Neg Pred Value : 0.9506
##      Prevalence : 0.8880
##      Detection Rate : 0.8840
##      Detection Prevalence : 0.9190
##      Balanced Accuracy : 0.8415
##
##      'Positive' Class : 0
##
```