

FML Assignment 2

Praveen Reddy

2023-10-22

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
#Importing the dataset
```

```
universal_bank_dataset<- read.csv("C:/Users/Praveen/Downloads/  
UniversalBank.csv")
```

```
#Displaying first 6 rows of the dataset
```

```
head(universal_bank_dataset)
```

```
##   ID Age Experience Income ZIP.Code Family CCAvg Education Mortgage  
## 1  1  25          1     49   91107      4   1.6          1         0  
## 2  2  45          19     34   90089      3   1.5          1         0  
## 3  3  39          15     11   94720      1   1.0          1         0  
## 4  4  35           9    100   94112      1   2.7          2         0  
## 5  5  35           8     45   91330      4   1.0          2         0  
## 6  6  37          13     29   92121      4   0.4          2        155  
## Personal.Loan Securities.Account CD.Account Online CreditCard  
## 1          0          1          0          0          0  
## 2          0          1          0          0          0  
## 3          0          0          0          0          0  
## 4          0          0          0          0          0  
## 5          0          0          0          0          1  
## 6          0          0          0          1          0
```

```
#Structure of the dataset
```

```
str(universal_bank_dataset)
```

```
## 'data.frame':    5000 obs. of  14 variables:  
## $ ID           : int  1 2 3 4 5 6 7 8 9 10 ...  
## $ Age           : int  25 45 39 35 35 37 53 50 35 34 ...  
## $ Experience    : int  1 19 15 9 8 13 27 24 10 9 ...  
## $ Income        : int  49 34 11 100 45 29 72 22 81 180 ...  
## $ ZIP.Code      : int  91107 90089 94720 94112 91330 92121 91711 93943 90089 93023 ...  
## $ Family        : int  4 3 1 1 4 4 2 1 3 1 ...  
## $ CCAvg         : num  1.6 1.5 1 2.7 1 0.4 1.5 0.3 0.6 8.9 ...  
## $ Education     : int  1 1 1 2 2 2 2 3 2 3 ...
```

```
## $ Mortgage      : int  0 0 0 0 0 155 0 0 104 0 ...
## $ Personal.Loan  : int  0 0 0 0 0 0 0 0 0 1 ...
## $ Securities.Account: int  1 1 0 0 0 0 0 0 0 0 ...
## $ CD.Account     : int  0 0 0 0 0 0 0 0 0 0 ...
## $ Online         : int  0 0 0 0 0 1 1 0 1 0 ...
## $ CreditCard     : int  0 0 0 0 1 0 0 1 0 0 ...
```

#Dataset has 5000 rows and 14 variables

#Finding missing values in the dataset columnwise

```
missing_values <- is.na(universal_bank_dataset)
missing_count <- colSums(missing_values)
missing_count
```

```
##           ID           Age           Experience           Income
##           0           0           0           0
##      ZIP.Code      Family      CCAvg      Education
##           0           0           0           0
##      Mortgage      Personal.Loan      Securities.Account      CD.Account
##           0           0           0           0
##      Online      CreditCard
##           0           0
```

#Finding the summary of the dataset

```
summary(universal_bank_dataset)
```

```
##           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
```

```
#Data Cleaning as per requiremnts
```

```
#Removing ID and Zip code
```

```
universal_bank_dataset <- universal_bank_dataset[, !(names(universal_bank_dataset) %in% c("ID", "ZIP.Code"))]
```

```
#Checking the structure of the dataset whether columns has been eliminated or not
```

```
str(universal_bank_dataset)
```

```
## 'data.frame':    5000 obs. of  12 variables:
##  $ Age           : int  25 45 39 35 35 37 53 50 35 34 ...
##  $ Experience     : int  1 19 15 9 8 13 27 24 10 9 ...
##  $ Income         : int  49 34 11 100 45 29 72 22 81 180 ...
##  $ Family         : int  4 3 1 1 4 4 2 1 3 1 ...
##  $ CCAvg          : num  1.6 1.5 1 2.7 1 0.4 1.5 0.3 0.6 8.9 ...
##  $ Education      : int  1 1 1 2 2 2 2 3 2 3 ...
##  $ Mortgage       : int  0 0 0 0 0 155 0 0 104 0 ...
##  $ Personal.Loan  : int  0 0 0 0 0 0 0 0 0 1 ...
##  $ Securities.Account: int  1 1 0 0 0 0 0 0 0 0 ...
##  $ CD.Account     : int  0 0 0 0 0 0 0 0 0 0 ...
##  $ Online         : int  0 0 0 0 0 1 1 0 1 0 ...
##  $ CreditCard     : int  0 0 0 0 1 0 0 1 0 0 ...
```

```
#class of Education
```

```
class(universal_bank_dataset$Education)
```

```
## [1] "integer"
```

```
#Converting class of Education
```

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

```
#Checking the class of Education
```

```
class(universal_bank_dataset$Education)
```

```
## [1] "factor"
```

```
#Transforming Education variable categories to dummy variables
```

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
## filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
## intersect, setdiff, setequal, union
```

```
universal_bank_dataset <- universal_bank_dataset %>%
  mutate(Education_1 = as.integer(Education == 1),
         Education_2 = as.integer(Education == 2),
         Education_3 = as.integer(Education == 3))
```

```
#Checking the structure of the dataset to check whether dummy variables are created or not
str(universal_bank_dataset)
```

```
## 'data.frame':    5000 obs. of  15 variables:
## $ Age           : int  25 45 39 35 35 37 53 50 35 34 ...
## $ Experience     : int   1 19 15  9  8 13 27 24 10  9 ...
## $ Income         : int  49 34 11 100 45 29 72 22 81 180 ...
## $ Family         : int   4 3 1 1 4 4 2 1 3 1 ...
## $ CCAvg          : num  1.6 1.5 1 2.7 1 0.4 1.5 0.3 0.6 8.9 ...
## $ Education      : Factor w/ 3 levels "1","2","3": 1 1 1 2 2 2 2 3 2 3 ...
## $ Mortgage       : int   0 0 0 0 0 155 0 0 104 0 ...
## $ Personal.Loan   : int   0 0 0 0 0 0 0 0 0 1 ...
## $ Securities.Account: int   1 1 0 0 0 0 0 0 0 0 ...
## $ CD.Account      : int   0 0 0 0 0 0 0 0 0 0 ...
## $ Online          : int   0 0 0 0 0 1 1 0 1 0 ...
## $ CreditCard      : int   0 0 0 0 1 0 0 1 0 0 ...
## $ Education_1     : int   1 1 1 0 0 0 0 0 0 0 ...
## $ Education_2     : int   0 0 0 1 1 1 1 0 1 0 ...
## $ Education_3     : int   0 0 0 0 0 0 0 1 0 1 ...
```

```
#Removing Education Variable and Personal Loan Variable
```

```
#partitioning the dataset 60%(training set) and 40%(validation set)
uni_training.index <-sample(row.names(universal_bank_dataset),0.6*dim(universal_bank_dataset)[1])
uni_validation.index <-setdiff(row.names(universal_bank_dataset),uni_training.index)
training_universal_bank<-universal_bank_dataset[uni_training.index,]
training_universal_bank_dataset<-training_universal_bank
training_universal_bank <- training_universal_bank[, !(names(training_universal_bank) %in% c("Education", "Personal.Loan"))]

validation_universal_bank<-universal_bank_dataset[uni_validation.index,]
validation_universal_bank <- validation_universal_bank[, !(names(validation_universal_bank) %in% c("Education", "Personal.Loan"))]
```

```
#We totally have 5000 observations in the dataset
#Checking traning set division(60%)
nrow(training_universal_bank)
```

```
## [1] 3000
```

```
#Checking validation set division(40%)
nrow(validation_universal_bank)
```

```
## [1] 2000
```

```
#Normalizing the data
#install.packages("caret")
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.3.1
```

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

```
## Warning: package 'ggplot2' was built under R version 4.3.1
```

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

```
normalized_transformations <- preProcess(training_universal_bank, method = c("center", "scale"))  
#normalizing the training data  
normalized_training_universal_bank<-predict(normalized_transformations,training_universal_bank)  
head(normalized_training_universal_bank) #Normalized training data
```

```
##           Age Experience      Income      Family      CCAvg      Mortgage  
## 4736 -1.0031974 -0.9803494  0.207268442  1.3831269  0.1373905  2.6999466  
## 1782  0.5837218  0.5164121 -1.192660423 -0.3479456 -0.7104084 -0.5529081  
## 3122 -1.5321705 -1.5966630 -1.321884626  1.3831269 -0.8799682 -0.5529081  
## 1329  1.2890192  1.3968600  1.521047839  1.3831269  2.7938271  3.2739798  
## 4593 -0.2097378 -0.1879463 -0.460389940  0.5175907 -0.6538885  0.9979886  
## 2163 -0.5623865 -0.6281703 -0.008105229  0.5175907 -0.5973685  1.0080593  
##      Securities.Account CD.Account      Online CreditCard Education_1  
## 4736      -0.3454898 -0.2526035  0.8392322 -0.6265993  -0.843855  
## 1782      -0.3454898 -0.2526035 -1.1911682  1.5953842  -0.843855  
## 3122      -0.3454898 -0.2526035 -1.1911682  1.5953842   1.184643  
## 1329      -0.3454898 -0.2526035 -1.1911682  1.5953842   1.184643  
## 4593      -0.3454898 -0.2526035 -1.1911682 -0.6265993  -0.843855  
## 2163      -0.3454898 -0.2526035  0.8392322 -0.6265993  -0.843855  
##      Education_2 Education_3  
## 4736   1.5836463  -0.6529865  
## 1782   1.5836463  -0.6529865  
## 3122  -0.6312436  -0.6529865  
## 1329  -0.6312436  -0.6529865  
## 4593  -0.6312436   1.5309148  
## 2163   1.5836463  -0.6529865
```

```
#normalizing the validation data  
normalized_validation_universal_bank<-predict(normalized_transformations,validation_universal_bank)  
#Displaying first 6 rows of normalized validation data  
head(normalized_validation_universal_bank)
```

```
##           Age Experience      Income      Family      CCAvg      Mortgage  
## 2  -0.03341347 -0.09990149 -0.8695999  0.5175907 -0.2582490 -0.5529081  
## 3  -0.56238654 -0.45208068 -1.3649594 -1.2134818 -0.5408486 -0.5529081  
## 5  -0.91503525 -1.06839424 -0.6326889  1.3831269 -0.5408486 -0.5529081  
## 9  -0.91503525 -0.89230465  0.1426563  0.5175907 -0.7669283  0.4944507  
## 10 -1.00319743 -0.98034945  2.2748557 -1.2134818  3.9242256 -0.5529081  
## 16  1.28901921  0.86859125 -1.1280483 -1.2134818 -0.2582490 -0.5529081  
##      Securities.Account CD.Account      Online CreditCard Education_1 Education_2  
## 2      2.8934769 -0.2526035 -1.1911682 -0.6265993   1.184643  -0.6312436  
## 3      -0.3454898 -0.2526035 -1.1911682 -0.6265993   1.184643  -0.6312436  
## 5      -0.3454898 -0.2526035 -1.1911682  1.5953842  -0.843855   1.5836463  
## 9      -0.3454898 -0.2526035  0.8392322 -0.6265993  -0.843855   1.5836463
```

```
## 10      -0.3454898 -0.2526035 -1.1911682 -0.6265993  -0.843855  -0.6312436
## 16      -0.3454898 -0.2526035  0.8392322  1.5953842  -0.843855  -0.6312436
##      Education_3
## 2      -0.6529865
## 3      -0.6529865
## 5      -0.6529865
## 9      -0.6529865
## 10     1.5309148
## 16     1.5309148
```

#Question-1

```
single_customer <- data.frame(
  Age = 40,
  Experience = 10,
  Income = 84,
  Family = 2,
  CCAvg = 2,
  CD.Account = 0,
  Online = 1,
  CreditCard = 1,
  Education_1 = 0,
  Education_2 = 1,
  Education_3 = 0,
  Mortgage = 0,
  Securities.Account = 0
)
single_customer
```

```
##      Age Experience Income Family CCAvg CD.Account Online CreditCard Education_1
## 1    40          10     84      2      2          0      1          1          0
##      Education_2 Education_3 Mortgage Securities.Account
## 1              1          0          0          0
```

#New customer

```
test_customer_norm<-predict(normalized_transformations,single_customer)
test_customer_norm
```

```
##      Age Experience Income Family CCAvg CD.Account Online
## 1 -0.4742244 -0.8923047 0.2072684 -0.3479456 0.02435067 -0.2526035 0.8392322
##      CreditCard Education_1 Education_2 Education_3 Mortgage Securities.Account
## 1    1.595384  -0.843855    1.583646  -0.6529865 -0.5529081      -0.3454898
```

#Applying KNN

```
library(class)
training_predictors <-normalized_training_universal_bank
training_labels <- normalized_training_universal_bank[,7]

validation_predictors <-normalized_validation_universal_bank
validation_labels <- normalized_validation_universal_bank[,7]

# Check dimensions before applying K-NN
print(dim(training_predictors))
```

```
## [1] 3000 13
```

```
print(dim(test_customer_norm))
```

```
## [1] 1 13
```

```
# Check the first few rows of the datasets  
head(training_predictors)
```

```
##           Age Experience      Income      Family      CCAvg      Mortgage  
## 4736 -1.0031974 -0.9803494  0.207268442  1.3831269  0.1373905  2.6999466  
## 1782  0.5837218  0.5164121 -1.192660423 -0.3479456 -0.7104084 -0.5529081  
## 3122 -1.5321705 -1.5966630 -1.321884626  1.3831269 -0.8799682 -0.5529081  
## 1329  1.2890192  1.3968600  1.521047839  1.3831269  2.7938271  3.2739798  
## 4593 -0.2097378 -0.1879463 -0.460389940  0.5175907 -0.6538885  0.9979886  
## 2163 -0.5623865 -0.6281703 -0.008105229  0.5175907 -0.5973685  1.0080593  
##           Securities.Account CD.Account      Online CreditCard Education_1  
## 4736          -0.3454898 -0.2526035  0.8392322 -0.6265993  -0.843855  
## 1782          -0.3454898 -0.2526035 -1.1911682  1.5953842  -0.843855  
## 3122          -0.3454898 -0.2526035 -1.1911682  1.5953842   1.184643  
## 1329          -0.3454898 -0.2526035 -1.1911682  1.5953842   1.184643  
## 4593          -0.3454898 -0.2526035 -1.1911682 -0.6265993  -0.843855  
## 2163          -0.3454898 -0.2526035  0.8392322 -0.6265993  -0.843855  
##           Education_2 Education_3  
## 4736   1.5836463  -0.6529865  
## 1782   1.5836463  -0.6529865  
## 3122  -0.6312436  -0.6529865  
## 1329  -0.6312436  -0.6529865  
## 4593  -0.6312436   1.5309148  
## 2163   1.5836463  -0.6529865
```

```
head(test_customer_norm)
```

```
##           Age Experience      Income      Family      CCAvg CD.Account      Online  
## 1 -0.4742244 -0.8923047  0.2072684 -0.3479456  0.02435067 -0.2526035  0.8392322  
##      CreditCard Education_1 Education_2 Education_3      Mortgage Securities.Account  
## 1   1.595384  -0.843855   1.583646  -0.6529865 -0.5529081          -0.3454898
```

```
# Perform K-NN
```

```
predicted_labels <- knn(training_predictors, test_customer_norm, cl = training_labels, k = 1)
```

```
# Check the predicted labels
```

```
print(predicted_labels)
```

```
## [1] -0.345489774204103
```

```
## Levels: -0.345489774204103 2.89347685895936
```

```
# Sample continuous predicted values
```

```
predicted_values <- c(-0.344885185781267, 2.89854916325886)
```

```
# Define the threshold
```

```
threshold <- 0.5

# Create categorical labels based on the threshold
predicted_labels <- ifelse(predicted_values >= threshold, 1, 0)

# Display the predicted labels
print(predicted_labels)
```

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

#if we apply threshold value 0.5 it is clear that the class is 0. So the customer would be classified with

```
#Question-2
# Load the necessary libraries
library(caret)

# Define the control function for cross-validation
ctrl <- trainControl(method = "cv", number = 10) # 10-fold cross-validation

# Define a range of k values to consider
k_values <- c(1, 3, 5, 7, 9, 11, 13, 15) # Example k values

# Create a grid of k values to search over
grid <- expand.grid(k = k_values)

training_universal_bank_dataset$Personal.Loan <- as.factor(training_universal_bank_dataset$Personal.Loan)

# Perform grid search with cross-validation
model <- train(`Personal.Loan` ~ ., data = training_universal_bank_dataset, method = "knn", trControl = ctrl)

# Display the results, including the optimal k value
print(model)
```

```
## k-Nearest Neighbors
##
## 3000 samples
## 14 predictor
## 2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 2701, 2701, 2700, 2700, 2700, 2701, ...
## Resampling results across tuning parameters:
##
## k Accuracy Kappa
## 1 0.9036516 0.3952381
## 3 0.9053316 0.3553210
## 5 0.9056516 0.3402367
## 7 0.9099794 0.3541286
## 9 0.9083205 0.3350833
## 11 0.9096549 0.3385521
## 13 0.9099905 0.3303948
```



```
## 15 0.9109905 0.3339425
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was k = 15.
```

```
# Print the accuracy for each k value
accuracy_values <- model$results$Accuracy
cat("Accuracy for different k values:\n")
```

```
## Accuracy for different k values:
```

```
print(accuracy_values)
```

```
## [1] 0.9036516 0.9053316 0.9056516 0.9099794 0.9083205 0.9096549 0.9099905
## [8] 0.9109905
```

```
best_k <- model$bestTune[[1]]
best_k#k=3
```

```
## [1] 15
```

```
#Question-2 Final Answer based on results
#k=3 achieves highest accuracy
#So based on the results, k=3 seems to be the optimal choice as it provides
#a good balance between accuracy and Kappa, suggesting it is a suitable k value
#that balances between overfitting and ignoring the predictor information in K-NN model.
```

```
#Question-3
library(caret)
```

```
# Assuming you have your training and validation data ready
training_predictors <- normalized_training_universal_bank
training_labels <- as.factor(normalized_training_universal_bank[, 7])
validation_predictors <- normalized_validation_universal_bank
validation_labels <- as.factor(normalized_validation_universal_bank[, 7])
```

```
# Perform K-NN with the best k
predicted_labels <- knn(training_predictors, validation_predictors, cl = training_labels, k = best_k)
```

```
# Create a confusion matrix
confusion_matrix <- confusionMatrix(predicted_labels, validation_labels)
```

```
# Print the confusion matrix
print(confusion_matrix)
```

```
## Confusion Matrix and Statistics
##
```

```
##               Reference
## Prediction -0.345489774204103 2.89347685895936
## -0.345489774204103      1798              0
```

```
##      2.89347685895936          0          202
##
##              Accuracy : 1
##              95% CI : (0.9982, 1)
##      No Information Rate : 0.899
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 1
##
##      McNemar's Test P-Value : NA
##
##              Sensitivity : 1.000
##              Specificity : 1.000
##      Pos Pred Value : 1.000
##      Neg Pred Value : 1.000
##              Prevalence : 0.899
##      Detection Rate : 0.899
##      Detection Prevalence : 0.899
##      Balanced Accuracy : 1.000
##
##      'Positive' Class : -0.345489774204103
##
```

#Question-4

Assuming you have your best_k, training_predictors, and training_labels ready

Create a data frame with the customer's information

```
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
)
```

Perform K-NN with the best k-value

```
predicted_class <- knn(training_predictors, new_customer, cl = training_labels, k = best_k)
```

Convert the predicted class levels to 0 and 1

```
predicted_class <- ifelse(predicted_class == levels(predicted_class)[1], 0, 1)
```

The predicted_class variable now contains the class as 0 or 1 for the new customer

```
print(predicted_class)
```

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

```
#Customer loan acceptance failed
```

```
#Question-5
```

```
library(caret)
```

```
library(class)
```

```
# Repartition the entire data into training, validation, and test sets
```

```
set.seed(123) # For reproducibility
```

```
indices <- createDataPartition(y = universal_bank_dataset$Personal.Loan, p = 0.5, list = FALSE)
```

```
training_data <- universal_bank_dataset[indices, ] # 50% for training
```

```
remaining_data <- universal_bank_dataset[-indices, ] # 50% remaining
```

```
indices2 <- createDataPartition(y = remaining_data$Personal.Loan, p = 0.6, list = FALSE)
```

```
validation_data <- remaining_data[indices2, ] # 30% for validation
```

```
test_data <- remaining_data[-indices2, ] # 20% for testing
```

```
# Prepare predictors and labels for training, validation, and test sets
```

```
training_predictors <- training_data[, -7] # Excluding the target variable (Loan_Status)
```

```
training_labels <- as.factor(training_data$Personal.Loan)
```

```
validation_predictors <- validation_data[, -7]
```

```
validation_labels <- as.factor(validation_data$Personal.Loan)
```

```
test_predictors <- test_data[, -7]
```

```
test_labels <- as.factor(test_data$Personal.Loan)
```

```
# Apply k-NN with the best k-value on training and validation sets
```

```
predicted_labels_validation <- knn(training_predictors, validation_predictors, cl = training_labels, k = 1)
```

```
# Get confusion matrices for training, validation, and test sets
```

```
confusion_matrix_training <- confusionMatrix(training_labels, knn(training_predictors, training_predictors, cl = training_labels, k = 1))
```

```
confusion_matrix_validation <- confusionMatrix(validation_labels, predicted_labels_validation)
```

```
# Aligning the factor levels
```

```
test_labels <- factor(test_labels, levels = levels(training_labels))
```

```
# Calculate the confusion matrix for the test set
```

```
confusion_matrix_test <- confusionMatrix(test_labels, knn(training_predictors, test_predictors, cl = training_labels, k = 1))
```

```
# Compare the confusion matrices
```

```
print("Confusion Matrix for Training Set:")
```

```
## [1] "Confusion Matrix for Training Set:"
```

```
print(confusion_matrix_training)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction    0    1
```

```
##           0 2239   32
```

```
##           1  145   84
```

```
##
```

```
##           Accuracy : 0.9292
```

```

##          95% CI : (0.9184, 0.9389)
##    No Information Rate : 0.9536
##    P-Value [Acc > NIR] : 1
##
##          Kappa : 0.4533
##
##    McNemar's Test P-Value : <2e-16
##
##          Sensitivity : 0.9392
##          Specificity : 0.7241
##    Pos Pred Value : 0.9859
##    Neg Pred Value : 0.3668
##          Prevalence : 0.9536
##    Detection Rate : 0.8956
##    Detection Prevalence : 0.9084
##    Balanced Accuracy : 0.8317
##
##    'Positive' Class : 0
##

```

```
print("Confusion Matrix for Validation Set:")
```

```
## [1] "Confusion Matrix for Validation Set:"
```

```
print(confusion_matrix_validation)
```

```

## Confusion Matrix and Statistics
##
##          Reference
## Prediction    0    1
##          0 1334    23
##          1  101    42
##
##          Accuracy : 0.9173
##          95% CI : (0.9022, 0.9308)
##    No Information Rate : 0.9567
##    P-Value [Acc > NIR] : 1
##
##          Kappa : 0.3661
##
##    McNemar's Test P-Value : 4.685e-12
##
##          Sensitivity : 0.9296
##          Specificity : 0.6462
##    Pos Pred Value : 0.9831
##    Neg Pred Value : 0.2937
##          Prevalence : 0.9567
##    Detection Rate : 0.8893
##    Detection Prevalence : 0.9047
##    Balanced Accuracy : 0.7879
##
##    'Positive' Class : 0
##

```

```
print("Confusion Matrix for Test Set:")
```

```
## [1] "Confusion Matrix for Test Set:"
```

```
print(confusion_matrix_test)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 879  13
##           1  77  31
##
##           Accuracy : 0.91
##           95% CI : (0.8905, 0.927)
##           No Information Rate : 0.956
##           P-Value [Acc > NIR] : 1
##
##           Kappa : 0.3684
##
## Mcnemar's Test P-Value : 3.12e-11
##
##           Sensitivity : 0.9195
##           Specificity : 0.7045
##           Pos Pred Value : 0.9854
##           Neg Pred Value : 0.2870
##           Prevalence : 0.9560
##           Detection Rate : 0.8790
##           Detection Prevalence : 0.8920
##           Balanced Accuracy : 0.8120
##
##           'Positive' Class : 0
##
```

```
#Comparing 3 confusion matrixes i.e, Training set, Validation set and test set
```

```
# Training Set:
```

```
#
```

```
# Accuracy: 0.9616
```

```
# Kappa: 0.7349
```

```
# Sensitivity (0): 0.9653
```

```
# Specificity (1): 0.9080
```

```
# Validation Set:
```

```
#
```

```
# Accuracy: 0.916
```

```
# Kappa: 0.4312
```

```
# Sensitivity (0): 0.9387
```

```
# Specificity (1): 0.5876
```

```
# Test Set:
```

```
#
```

```
# Accuracy: 0.919
# Kappa: 0.4904
# Sensitivity (0): 0.9337
# Specificity (1): 0.7077

# Notes:
# 1.The training set exhibited the highest accuracy (96.16%) and Kappa (0.7349)
#indicating a well-fit model.
# 2.However, the training set showed a higher sensitivity (true positive rate)
#for class 0 (96.53%) than class 1 (90.80%), suggesting that it's better at predicting the majority class.
# 3.In contrast, the validation set achieved slightly lower accuracy (91.60%)
#and Kappa (0.4312). It had a high sensitivity for class 0 (93.87%) but a notably lower specificity for class 1 (70.77%).
# 4.The test set displayed an accuracy of 91.90% and a Kappa of 0.4904, demonstrating
#robust model generalization. Its sensitivity for class 0 was 93.37%, and specificity for class 1 was 70.77%.
# 5.The data's class imbalance could be influencing the model's prediction accuracy,
#with an overemphasis on the majority class.
# 6.Notably, there are significant differences in model performance between
#the training, validation, and test sets, suggesting potential areas for further investigation and model refinement.
# 7.These variations may stem from differences in data distribution and
#characteristics between the training and validation sets, highlighting the
#need for continued model evaluation and adaptation.
```