# Assingment_2_fml

## Chandu

## 2023-10-02

```r
#importing the requiored packages in r
library('caret')
```

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

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

```r
library('ISLR')
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
```

```r
library('class')

#Importing the dataset from local folders

sb.data <- read.csv("C:/Users/shash/Dropbox/PC/Downloads/UniversalBank (2).csv")
#Question_1
#conducting a k-NN classification
#predictors removed, i.e., removing ID and ZIP Code from each and every column from the data set
sb.data$ID <- NULL
sb.data$ZIP.Code <- NULL
summary(sb.data)
```

```
##       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
```

```r
#converting categorical variable "personal loan" into a factor that responses as "yes" or "no."

sb.data$Personal.Loan =  as.factor(sb.data$Personal.Loan)


#normalize the data by dividing
#training and validation, use preProcess() from the caret package.
M_norm <- preProcess(sb.data[, -8],method = c("center", "scale"))
sb.data_norm <- predict(M_norm,sb.data)
summary(sb.data_norm)
```

```
##       Age             Experience           Income          Family
## Min.   :-1.94871   Min.   :-2.014710   Min.   :-1.4288   Min.   :-1.2167
## 1st Qu.:-0.90188   1st Qu.:-0.881116   1st Qu.:-0.7554   1st Qu.:-1.2167
## Median :-0.02952   Median :-0.009121   Median :-0.2123   Median :-0.3454
## Mean   : 0.00000   Mean   : 0.000000   Mean   : 0.0000   Mean   : 0.0000
## 3rd Qu.: 0.84284   3rd Qu.: 0.862874   3rd Qu.: 0.5263   3rd Qu.: 0.5259
## Max.   : 1.88967   Max.   : 1.996468   Max.   : 3.2634   Max.   : 1.3973
##     CCAvg           Education         Mortgage       Personal.Loan
## Min.   :-1.1089   Min.   :-1.0490   Min.   :-0.5555   0:4520
## 1st Qu.:-0.7083   1st Qu.:-1.0490   1st Qu.:-0.5555   1: 480
## Median :-0.2506   Median : 0.1417   Median :-0.5555
## Mean   : 0.0000   Mean   : 0.0000   Mean   : 0.0000
## 3rd Qu.: 0.3216   3rd Qu.: 1.3324   3rd Qu.: 0.4375
## Max.   : 4.6131   Max.   : 1.3324   Max.   : 5.6875
## Securities.Account   CD.Account        Online          CreditCard
## Min.   :-0.3414   Min.   :-0.2535   Min.   :-1.2165   Min.   :-0.6452
## 1st Qu.:-0.3414   1st Qu.:-0.2535   1st Qu.:-1.2165   1st Qu.:-0.6452
## Median :-0.3414   Median :-0.2535   Median : 0.8219   Median :-0.6452
## Mean   : 0.0000   Mean   : 0.0000   Mean   : 0.0000   Mean   : 0.0000
## 3rd Qu.:-0.3414   3rd Qu.:-0.2535   3rd Qu.: 0.8219   3rd Qu.: 1.5495
## Max.   : 2.9286   Max.   : 3.9438   Max.   : 0.8219   Max.   : 1.5495
```

```r
#partition of the data into test and training sets as per the requirements
sb_train_index <- createDataPartition(sb.data$Personal.Loan, p = 0.6, list = FALSE)
my_train.df = sb.data_norm[sb_train_index,]
validate.sb.df = sb.data_norm[-sb_train_index,]
```

```r
print(head(my_train.df))
```

```
##            Age   Experience      Income     Family      CCAvg  Education
## 2  -0.02952064 -0.09632058 -0.8640230  0.5259383 -0.2505855 -1.0489730
## 6  -0.72740814 -0.61951767 -0.9726390  1.3972742 -0.8799989  0.1416887
## 7   0.66836686  0.60127554 -0.0385413 -0.3453975 -0.2505855  0.1416887
## 8   0.40665905  0.33967699 -1.1247014 -1.2167334 -0.9372183  1.3323505
## 10 -0.98911595 -0.96831574  2.3075645 -1.2167334  3.9836502  1.3323505
## 13  0.23218717  0.25247748  0.8738332 -0.3453975  1.0654607  1.3323505
##       Mortgage Personal.Loan Securities.Account  CD.Account     Online CreditCard
## 2  -0.5554684             0          2.9286223 -0.2535149 -1.2164961 -0.6452498
## 6   0.9684153             0         -0.3413892 -0.2535149  0.8218687 -0.6452498
## 7  -0.5554684             0         -0.3413892 -0.2535149  0.8218687 -0.6452498
## 8  -0.5554684             0         -0.3413892 -0.2535149 -1.2164961  1.5494774
## 10 -0.5554684             1         -0.3413892 -0.2535149 -1.2164961 -0.6452498
## 13 -0.5554684             0          2.9286223 -0.2535149 -1.2164961 -0.6452498
```

```r
#predict dataset from the above data given.
library(caret)
library(FNN)
```

```
##
## Attaching package: 'FNN'
```

```
## The following objects are masked from 'package:class':
##
##     knn, knn.cv
```

```r
sb.predict = data.frame(Age = 40, Experience = 10, Income = 84, Family = 2,
                    CCAvg = 2, Education = 1, Mortgage = 0, Securities.Account =
                       0, CD.Account = 0, Online = 1, CreditCard = 1)
print(sb.predict)
```

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

```r
sb.predict_Norm <- predict(M_norm,sb.predict)

predictions <- knn(train= as.data.frame(my_train.df[,1:7,9:12]),
                test = as.data.frame(sb.predict_Norm[,1:7,9:12]),
                cl= my_train.df$Personal.Loan,
                k=1)
```

```
## Warning in drop && !has.j: 'length(x) = 4 > 1' in coercion to 'logical(1)'
```

```
## Warning in drop && length(y) == 1L: 'length(x) = 4 > 1' in coercion to
## 'logical(1)'
```

```
## Warning in drop && !mdrop: 'length(x) = 4 > 1' in coercion to 'logical(1)'

## Warning in drop && !has.j: 'length(x) = 4 > 1' in coercion to 'logical(1)'

## Warning in drop && length(y) == 1L: 'length(x) = 4 > 1' in coercion to
## 'logical(1)'

## Warning in drop && !mdrop: 'length(x) = 4 > 1' in coercion to 'logical(1)'
```

```r
print(predictions)
```

```
## [1] 0
## attr(,"nn.index")
##       [,1]
## [1,]   409
## attr(,"nn.dist")
##            [,1]
## [1,] 0.2986486
## Levels: 0
```

```r
#Question_2
#determining the K value that balances overfitting and underfitting from the data set

set.seed(123)
SB.Bank <- trainControl(method= "repeatedcv", number = 3, repeats = 2)
searchGrid = expand.grid(k=1:10)

knn.model = train(Personal.Loan~., data = my_train.df, method = 'knn', tuneGrid = searchGrid,trControl =

knn.model
```

```
## k-Nearest Neighbors
##
## 3000 samples
##   11 predictor
##    2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (3 fold, repeated 2 times)
## Summary of sample sizes: 2000, 2000, 2000, 2000, 2000, 2000, ...
## Resampling results across tuning parameters:
##
##   k  Accuracy   Kappa
##   1  0.9518333  0.6888533
##   2  0.9485000  0.6694436
##   3  0.9545000  0.6863764
##   4  0.9498333  0.6460216
##   5  0.9516667  0.6565383
##   6  0.9491667  0.6342150
##   7  0.9461667  0.5985264
##   8  0.9456667  0.5915846
##   9  0.9450000  0.5848058
```

```
##    10  0.9413333  0.5499303
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was k = 3.
```

```
#perfect value of k is 3
#strikes a compromise between underfitting and overfitting of the data above.

#Question 3
#confusion Matrix is below
predictors_bank <- predict(knn.model,validate.sb.df)

confusionMatrix(predictors_bank,validate.sb.df$Personal.Loan)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 1799   75
##          1    9  117
##
##                Accuracy : 0.958
##                  95% CI : (0.9483, 0.9664)
##     No Information Rate : 0.904
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.7141
##
##  Mcnemar's Test P-Value : 1.321e-12
##
##             Sensitivity : 0.9950
##             Specificity : 0.6094
##          Pos Pred Value : 0.9600
##          Neg Pred Value : 0.9286
##              Prevalence : 0.9040
##          Detection Rate : 0.8995
##    Detection Prevalence : 0.9370
##       Balanced Accuracy : 0.8022
##
##        'Positive' Class : 0
##
```
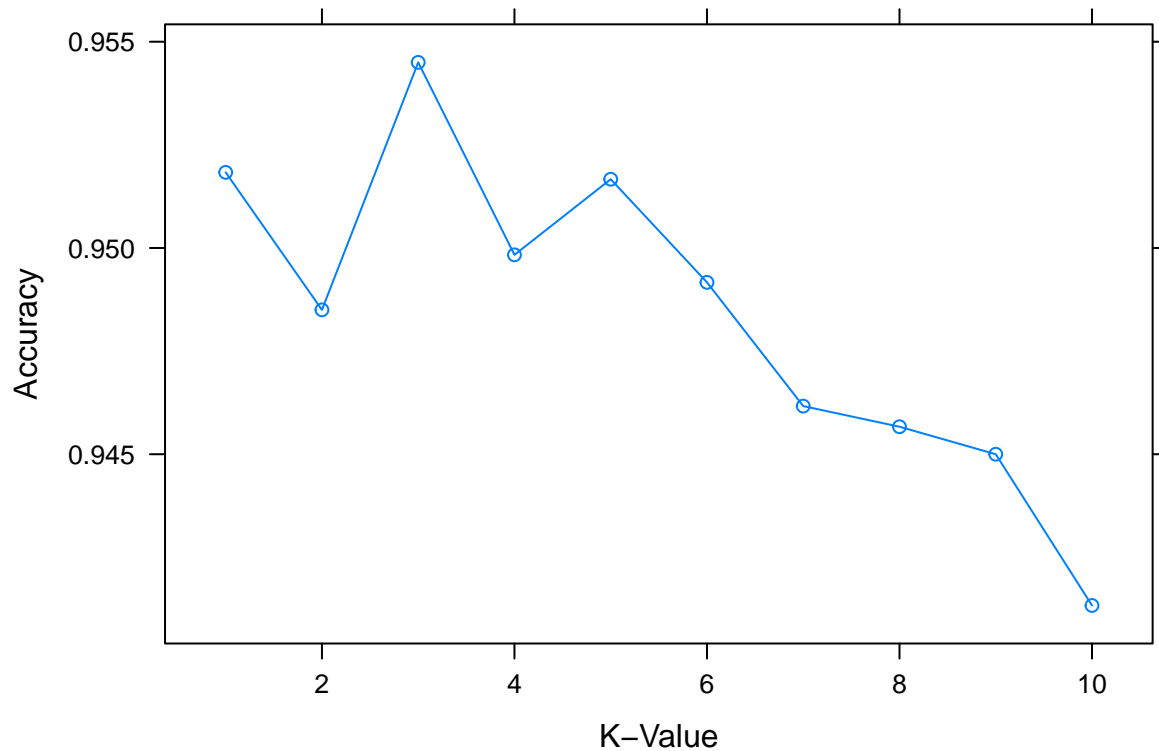
```
#The confustionmatrix has a 95.1% accuracy.
```

```
#Question 4
#Levels
#using the best K to classify the consumer.
sb.predict_Norm = data.frame(Age = 40, Experience = 10, Income = 84, Family = 2,
                             CCAvg = 2, Education = 1, Mortgage = 0,
                             Securities.Account =0, CD.Account = 0, Online = 1,
                             CreditCard = 1)
sb.predict_Norm = predict(M_norm, sb.predict)
predict(knn.model, sb.predict_Norm)
```

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

```r
#A plot that shows the best value of K (3), the one with the highest accuracy, is also present.
plot(knn.model, type = "b", xlab = "K-Value", ylab = "Accuracy")
```



```r
#Question 5
#creating training, test, and validation sets from the data collection.
t_size = 0.5 #training(50%)
sb_train_index = createDataPartition(sb.data$Personal.Loan, p = 0.5, list = FALSE)
my_train.df = sb.data_norm[sb_train_index,]


t.data_size = 0.2 #Test Data(20%)
Test.data_index = createDataPartition(sb.data$Personal.Loan, p = 0.2, list = FALSE)
t.data.df = sb.data_norm[Test.data_index,]


validation_size = 0.3 #validation(30%)
Validation.sb_index = createDataPartition(sb.data$Personal.Loan, p = 0.3, list = FALSE)
validate.sb.df = sb.data_norm[Validation.sb_index,]


Test.data.knn <- knn(train = my_train.df[,-8], test = t.data.df[,-8], cl = my_train.df[,8], k =3)
```

```
Validation.knn <- knn(train = my_train.df[,-8], test = validate.sb.df[,-8], cl = my_train.df[,8], k =3)
Training.knn <- knn(train = my_train.df[,-8], test = my_train.df[,-8], cl = my_train.df[,8], k =3)

confusionMatrix(Test.data.knn, t.data.df[,8])
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 901  30
##          1   3  66
##
##                Accuracy : 0.967
##                  95% CI : (0.954, 0.9772)
##     No Information Rate : 0.904
##     P-Value [Acc > NIR] : 1.058e-14
##
##                   Kappa : 0.7825
##
##  Mcnemar's Test P-Value : 6.011e-06
##
##             Sensitivity : 0.9967
##             Specificity : 0.6875
##          Pos Pred Value : 0.9678
##          Neg Pred Value : 0.9565
##              Prevalence : 0.9040
##          Detection Rate : 0.9010
##    Detection Prevalence : 0.9310
##       Balanced Accuracy : 0.8421
##
##        'Positive' Class : 0
##
```

```
confusionMatrix(Validation.knn, validate.sb.df[,8])
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 1349   37
##          1    7  107
##
##                Accuracy : 0.9707
##                  95% CI : (0.9608, 0.9786)
##     No Information Rate : 0.904
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.8136
##
##  Mcnemar's Test P-Value : 1.232e-05
##
##             Sensitivity : 0.9948
##             Specificity : 0.7431
```

```
##          Pos Pred Value : 0.9733
##          Neg Pred Value : 0.9386
##             Prevalence : 0.9040
##         Detection Rate : 0.8993
##   Detection Prevalence : 0.9240
##       Balanced Accuracy : 0.8689
##
##         'Positive' Class : 0
##
```

```
confusionMatrix(Training.knn, my_train.df[,8])
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##          0 2254   55
##          1    6  185
##
##                Accuracy : 0.9756
##                  95% CI : (0.9688, 0.9813)
##     No Information Rate : 0.904
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.8453
##
##  Mcnemar's Test P-Value : 7.958e-10
##
##             Sensitivity : 0.9973
##             Specificity : 0.7708
##          Pos Pred Value : 0.9762
##          Neg Pred Value : 0.9686
##             Prevalence : 0.9040
##         Detection Rate : 0.9016
##   Detection Prevalence : 0.9236
##       Balanced Accuracy : 0.8841
##
##         'Positive' Class : 0
##
```

*#Final Verdict: The training data have improved accuracy and sensitivity. According to the aforemention*
*#matrices, the values for the Test, Training, and Validation sets are 96.3%, 97.32%, and 96.73%, respec*