

Assignment2

Ashvitha M

20/02/2022

```
#required packages
library('caret')
```

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

```
## Warning in register(): Can't find generic 'scale_type' in package ggplot2 to
## register S3 method.
```

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

```
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
```

```
library('class')
```

```
Bankdata <- read.csv("C:/Users/mashv/Downloads/UniversalBank (2).csv", sep = ',') )
```

```
Bankdata$ID <- NULL
Bankdata$ZIP.Code <- NULL
summary(Bankdata)
```

```
##      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
```

```
Bankdata$Personal.Loan = as.factor(Bankdata$Personal.Loan)
```

```
#Normalisation of the data
```

```
Model_normalized <- preProcess(Bankdata[, -8],method = c("center", "scale"))
Bank_normalized <- predict(Model_normalized,Bankdata)
summary(Bank_normalized)
```

```
## 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
```

```
Train_index <- createDataPartition(Bankdata$Personal.Loan, p = 0.6, list = FALSE)
train.df = Bank_normalized[Train_index,]
validation.df = Bank_normalized[-Train_index,]
```

```
#Prediction the sample test data with best model.
```

```
To_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(To_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
```

```
To_Predict_Normalized <- predict(Model_normalized, To_Predict)

Prediction <- knn(train= train.df[,1:7,9:12],
                  test = To_Predict_Normalized[,1:7,9:12],
                  cl= train.df$Personal.Loan,
                  k=1)
print(Prediction)
```

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

#Task 2

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

knn.model = train(Personal.Loan~., data = train.df, method = 'knn', tuneGrid = searchGrid, trControl = Bankcontrol)
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  0.9560000  0.7206138
##  2  0.9466667  0.6669685
##  3  0.9561667  0.7034998
##  4  0.9548333  0.6894926
##  5  0.9558333  0.6922267
##  6  0.9546667  0.6843400
##  7  0.9538333  0.6725590
##  8  0.9523333  0.6589127
##  9  0.9496667  0.6372399
```

```
## 10 0.9468333 0.6068365
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was k = 3.
```

#Task 3

```
predictions <- predict(knn.model, validation.df)

confusionMatrix(predictions, validation.df$Personal.Loan)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 1794   75
##           1   14  117
##
##           Accuracy : 0.9555
##           95% CI : (0.9455, 0.9641)
##       No Information Rate : 0.904
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.7012
##
##  McNemar's Test P-Value : 2.018e-10
##
##           Sensitivity : 0.9923
##           Specificity : 0.6094
##           Pos Pred Value : 0.9599
##           Neg Pred Value : 0.8931
##           Prevalence : 0.9040
##           Detection Rate : 0.8970
##       Detection Prevalence : 0.9345
##           Balanced Accuracy : 0.8008
##
##           'Positive' Class : 0
##
```

#Task 4

```
To_Predict_Normalized = 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)

To_Predict_Normalized = predict(Model_normalized, To_Predict)
predict(knn.model, To_Predict_Normalized)
```

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

#Task5

*#Splitting dataset into training, validation, and test sets (50% : 30% : 20%)**

```
train_size = 0.5
Train_index = createDataPartition(Bankdata$Personal.Loan, p = 0.5, list = FALSE)
train.df = Bank_normalized[Train_index,]

valid_size = 0.3
Validation_index = createDataPartition(Bankdata$Personal.Loan, p = 0.3, list = FALSE)
validation.df = Bank_normalized[Validation_index,]

test_size = 0.2
Test_index = createDataPartition(Bankdata$Personal.Loan, p = 0.2, list = FALSE)
Test.df = Bank_normalized[Test_index,]

Testknn <- knn(train = train.df[, -8], test = Test.df[, -8], cl = train.df[, 8], k = 3)
Validationknn <- knn(train = train.df[, -8], test = validation.df[, -8], cl = train.df[, 8], k = 3)
Trainknn <- knn(train = train.df[, -8], test = train.df[, -8], cl = train.df[, 8], k = 3)

confusionMatrix(Testknn, Test.df[, 8])
```

Confusion Matrix and Statistics

```
##
##           Reference
## Prediction    0    1
##           0 899  22
##           1   5  74
##
##           Accuracy : 0.973
##           95% CI : (0.961, 0.9821)
##       No Information Rate : 0.904
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.8311
##
##  McNemar's Test P-Value : 0.002076
##
##           Sensitivity : 0.9945
##           Specificity : 0.7708
##       Pos Pred Value : 0.9761
##       Neg Pred Value : 0.9367
##           Prevalence : 0.9040
##       Detection Rate : 0.8990
##   Detection Prevalence : 0.9210
##       Balanced Accuracy : 0.8827
##
##       'Positive' Class : 0
##
```

```
confusionMatrix(Trainknn, train.df[, 8])
```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 2255   58
##           1    5  182
##
##           Accuracy : 0.9748
##           95% CI : (0.9679, 0.9806)
##           No Information Rate : 0.904
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.8389
##
## Mcnemar's Test P-Value : 5.701e-11
##
##           Sensitivity : 0.9978
##           Specificity : 0.7583
##           Pos Pred Value : 0.9749
##           Neg Pred Value : 0.9733
##           Prevalence : 0.9040
##           Detection Rate : 0.9020
##           Detection Prevalence : 0.9252
##           Balanced Accuracy : 0.8781
##
##           'Positive' Class : 0
##

```

```
confusionMatrix(Validationknn, validation.df[,8])
```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 1352   44
##           1    4  100
##
##           Accuracy : 0.968
##           95% CI : (0.9578, 0.9763)
##           No Information Rate : 0.904
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.7895
##
## Mcnemar's Test P-Value : 1.811e-08
##
##           Sensitivity : 0.9971
##           Specificity : 0.6944
##           Pos Pred Value : 0.9685
##           Neg Pred Value : 0.9615
##           Prevalence : 0.9040
##           Detection Rate : 0.9013
##           Detection Prevalence : 0.9307
##           Balanced Accuracy : 0.8457
##

```

```
##  
##      'Positive' Class : 0  
##
```

##Conclusion:

*#From both performance metrics, we can infer that the amount of sample data considered for training can be a significant factor in achieving higher accuracy.
#Hence, more the training data there are more chances of achieving higher accuracy.*