

Assignment 2

Harshini Balam

2024-02-25

#Summary The assignment's objective is to predict, using KNN(k-Nearest Neighbors)Classification, if the loan offer will be accepted by consumers of Universal Bank. The data set contains demographic information about the customers as well as other silent-related information. Following the reading of the data set and the installation of the required libraries, extra columns are removed, category categories are changed to dummy variables, and the data is eventually normalized. Following that, the data set was divided into two sets:training and validation, each of which contained 60% and 40% of the entire data. A new consumer was categorized as either accepting or rejecting a loan offer using k-NN with k=1.By assessing accuracy on the validation set, the optimal k value—which strikes a balance between over fitting and under fitting—was found, with k=1 being the result.

```
library(caret)
```

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

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

```
library(ISLR)
library(dplR)
```

```
## This is dplR version 1.7.6.
## dplR is part of openDendro https://opendendro.org.
## New users can visit https://opendendro.github.io/dplR-workshop/ to get started.
```

```
library('class')
library('gmodels')
library('FNN')
```

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

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

```
library("ggplot2")
```

Import dataset UniversalBank.csv

```
setwd("C:\\Users\\Harshini\\OneDrive - Kent State University\\Fundamentals of Machine Learning")
UniversalBank=read.csv("UniversalBank.csv")
head(UniversalBank)
```

```
##   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
```

```
colnames(UniversalBank)
```

```
## [1] "ID"           "Age"           "Experience"
## [4] "Income"       "ZIP.Code"      "Family"
## [7] "CCAvg"        "Education"     "Mortgage"
## [10] "Personal.Loan" "Securities.Account" "CD.Account"
## [13] "Online"       "CreditCard"
```

Summary of UniversalBank dataset

```
summary(UniversalBank)
```

```
##           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
```

Making columns ID and ZIP.Code as NULL

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

```
##      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
```

Making the Personal Loan column as factor

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

Normalization

```
Normal_Data <- preProcess(UniversalBank,method = "range")
UniversalBank_Norm <- predict(Normal_Data,UniversalBank)
summary(UniversalBank_Norm)
```

```
##      Age      Experience      Income      Family
## Min.   :0.0000   Min.   :0.0000   Min.    :0.0000   Min.    :0.0000
## 1st Qu.:0.2727   1st Qu.:0.2826   1st Qu.:0.1435   1st Qu.:0.0000
## Median :0.5000   Median :0.5000   Median :0.2593   Median :0.3333
## Mean    :0.5077   Mean    :0.5023   Mean    :0.3045   Mean    :0.4655
## 3rd Qu.:0.7273   3rd Qu.:0.7174   3rd Qu.:0.4167   3rd Qu.:0.6667
```

```
## Max. :1.0000 Max. :1.0000 Max. :1.0000 Max. :1.0000
## CCAvg Education Mortgage Personal.Loan
## Min. :0.0000 Min. :0.0000 Min. :0.00000 0:4520
## 1st Qu.:0.0700 1st Qu.:0.0000 1st Qu.:0.00000 1: 480
## Median :0.1500 Median :0.5000 Median :0.00000
## Mean :0.1938 Mean :0.4405 Mean :0.08897
## 3rd Qu.:0.2500 3rd Qu.:1.0000 3rd Qu.:0.15906
## Max. :1.0000 Max. :1.0000 Max. :1.00000
## 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
```

Partition the data into training 60% and validation 40% sets

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

Classifying the customer as per the data provided

```
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
```

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

```
## [1] 1
## attr(,"nn.index")
## [,1]
## [1,] 844
## attr(,"nn.dist")
## [,1]
## [1,] 92.3748
## Levels: 1
```

Customer is classified as 1.

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

```

set.seed(2808)
UniversalBank_control <- trainControl(method = "repeatedcv", number = 5, repeats = 2)
searchGrid = expand.grid(k=1:10)
knn.model = train(Personal.Loan~., data = train.df, method = 'knn', tuneGrid = searchGrid, trControl = UniversalBank_control)
knn.model

```

```

## k-Nearest Neighbors
##
## 3000 samples
## 11 predictor
## 2 classes: '0', '1'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold, repeated 2 times)
## Summary of sample sizes: 2400, 2400, 2399, 2400, 2401, 2400, ...
## Resampling results across tuning parameters:
##
##  k  Accuracy  Kappa
##  1  0.9586688  0.7293613
##  2  0.9521708  0.6824576
##  3  0.9536691  0.6796908
##  4  0.9513355  0.6628740
##  5  0.9483352  0.6313119
##  6  0.9475022  0.6214180
##  7  0.9441685  0.5891933
##  8  0.9438355  0.5840695
##  9  0.9430022  0.5752916
## 10  0.9423338  0.5649667
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was k = 1.

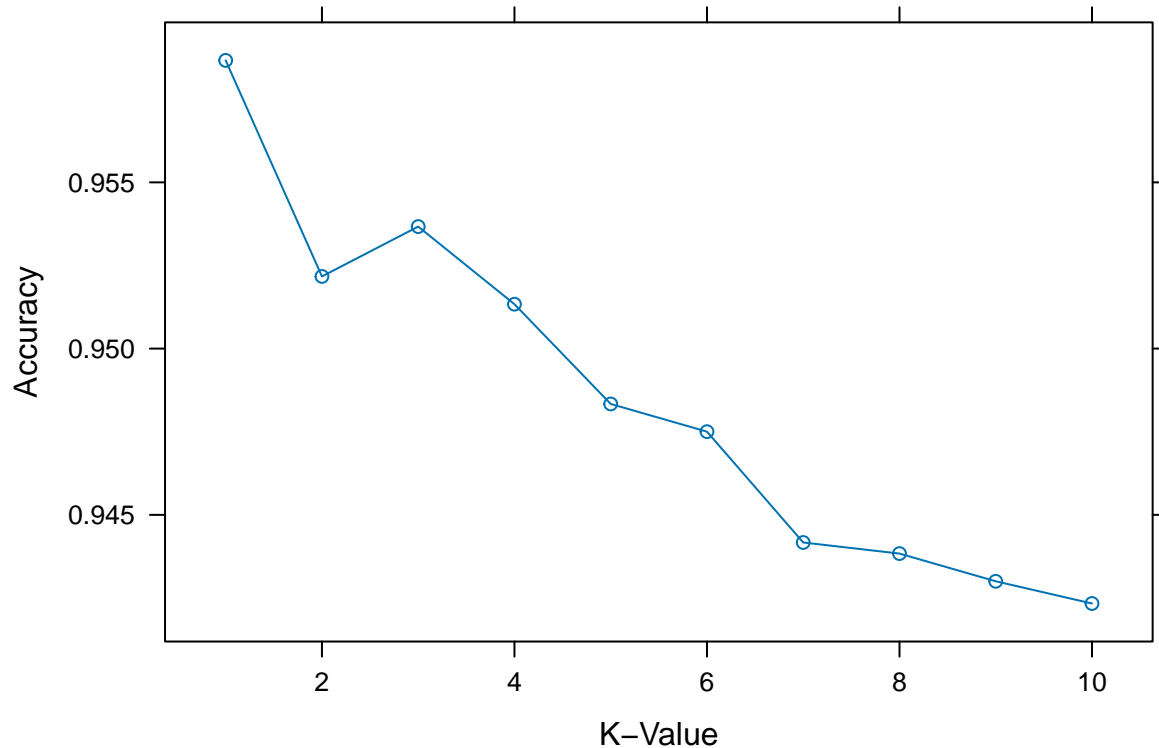
```

The choice of K that balances between overfitting and ignoring predictors $K=1$

```

plot(knn.model, type = "b", xlab = "K-Value", ylab = "Accuracy")

```



#finding the best K

```
best_k <- knn.model$bestTune[[1]]
best_k
```

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

3) Show the confusion matrix for the validation data that results from using the best k.

```
predictions <- predict(knn.model,validation.df)
confusionMatrix(predictions,validation.df$Personal.Loan)
```

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

```
##
```

```
##           Reference
```

```
## Prediction    0    1
```

```
##           0 1783   61
```

```
##           1   25  131
```

```
##
```

```
##           Accuracy : 0.957
```

```
##           95% CI : (0.9472, 0.9655)
```

```
## No Information Rate : 0.904
```

```
## P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##           Kappa : 0.7296
```

```
##
## McNemar's Test P-Value : 0.0001606
##
##          Sensitivity : 0.9862
##          Specificity : 0.6823
##          Pos Pred Value : 0.9669
##          Neg Pred Value : 0.8397
##          Prevalence : 0.9040
##          Detection Rate : 0.8915
##          Detection Prevalence : 0.9220
##          Balanced Accuracy : 0.8342
##
##          'Positive' Class : 0
##
```

4) Classify the customer using the best k

```
To_Predict_Normaliz = 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_Normaliz = predict(Normal_Data, To_Predict)
predict(knn.model, To_Predict_Normaliz)
```

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

5) Repartition the data into 50% for training ,30% for validation, 20% for test

```
train_size = 0.5
Train_index = createDataPartition(UniversalBank$Personal.Loan, p = 0.5,
list = FALSE)
train.df = UniversalBank_Norm[Train_index,]
test_size = 0.2
Test_index = createDataPartition(UniversalBank$Personal.Loan, p = 0.2,
list = FALSE)
Test.df = UniversalBank_Norm[Test_index,]
valid_size = 0.3
Validation_index = createDataPartition(UniversalBank$Personal.Loan, p = 0.3,
list = FALSE)
validation.df = UniversalBank_Norm[Validation_index,]
Testingknn <- 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)
Trainingknn <- knn(train = train.df[, -8],
test = train.df[, -8], cl = train.df[, 8], k = 3)
```

Comparing the confusion matrix of the test set with the training and validation sets.

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

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

```
##
##           Reference
## Prediction    0    1
##           0 900  38
##           1   4  58
##
##           Accuracy : 0.958
##           95% CI : (0.9436, 0.9696)
##           No Information Rate : 0.904
##           P-Value [Acc > NIR] : 9.200e-11
##
##           Kappa : 0.7125
##
## Mcnemar's Test P-Value : 3.543e-07
##
##           Sensitivity : 0.9956
##           Specificity : 0.6042
##           Pos Pred Value : 0.9595
##           Neg Pred Value : 0.9355
##           Prevalence : 0.9040
##           Detection Rate : 0.9000
##           Detection Prevalence : 0.9380
##           Balanced Accuracy : 0.7999
##
##           'Positive' Class : 0
##
```

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

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 2252  54
##           1   8 186
##
##           Accuracy : 0.9752
##           95% CI : (0.9683, 0.9809)
##           No Information Rate : 0.904
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.8437
##
## Mcnemar's Test P-Value : 1.097e-08
##
##           Sensitivity : 0.9965
##           Specificity : 0.7750
##           Pos Pred Value : 0.9766
##           Neg Pred Value : 0.9588
##           Prevalence : 0.9040
##           Detection Rate : 0.9008
##           Detection Prevalence : 0.9224
##           Balanced Accuracy : 0.8857
##
```



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