

# Assignment\_2\_FML

Anusha Banda

2023-10-01

## Summary

### Questions - Answers

1. How would this customer be classified? This new customer would be classified as 0, does not take the personal loan
2. The best K is 3

### Problem Statement

Universal bank is a young bank growing rapidly in terms of overall customer acquisition. The majority of these customers are liability customers (depositors) with varying sizes of relationship with the bank. The customer base of asset customers (borrowers) is quite small, and the bank is interested in expanding this base rapidly in more loan business. In particular, it wants to explore ways of converting its liability customers to personal loan customers.

A campaign that the bank ran last year for liability customers showed a healthy conversion rate of over 9% success. This has encouraged the retail marketing department to devise smarter campaigns with better target marketing. The goal is to use k-NN to predict whether a new customer will accept a loan offer. This will serve as the basis for the design of a new campaign.

The file UniversalBank.csv contains data on 5000 customers. The data include customer demographic information (age, income, etc.), the customer's relationship with the bank (mortgage, securities account, etc.), and the customer response to the last personal loan campaign (Personal Loan). Among these 5000 customers, only 480 (= 9.6%) accepted the personal loan that was offered to them in the earlier campaign.

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

---

First, load the required libraries

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

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

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

```
library(e1071)
```

Read the data.

```
universal.df <- read.csv("UniversalBank.csv")  
dim(universal.df)
```

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

```
t(t(names(universal.df)))
```

```
##      [,1]  
## [1,] "ID"  
## [2,] "Age"  
## [3,] "Experience"  
## [4,] "Income"  
## [5,] "ZIP.Code"  
## [6,] "Family"  
## [7,] "CCAvg"  
## [8,] "Education"  
## [9,] "Mortgage"  
## [10,] "Personal.Loan"  
## [11,] "Securities.Account"  
## [12,] "CD.Account"  
## [13,] "Online"  
## [14,] "CreditCard"
```

Drop ID and ZIP

```
universal.df <- universal.df[, -c(1,5)]
```

Split Data into 60% training and 40% validation. There are many ways to do this. We will look at 2 different ways. Before we split, let us transform categorical variables into dummy variables

```
universal.df$Education <- as.factor(universal.df$Education)
```

```
groups <- dummyVars(~., data = universal.df)  
universal_m.df <- as.data.frame(predict(groups, universal.df))
```

```
set.seed(1)  
train.index <- sample(row.names(universal_m.df), 0.6*dim(universal_m.df)[1])  
valid.index <- setdiff(row.names(universal_m.df), train.index)  
train.df <- universal_m.df[train.index,]  
valid.df <- universal_m.df[valid.index,]  
t(t(names(train.df)))
```

```
##      [,1]  
## [1,] "Age"  
## [2,] "Experience"
```

```
## [3,] "Income"
## [4,] "Family"
## [5,] "CCAvg"
## [6,] "Education.1"
## [7,] "Education.2"
## [8,] "Education.3"
## [9,] "Mortgage"
## [10,] "Personal.Loan"
## [11,] "Securities.Account"
## [12,] "CD.Account"
## [13,] "Online"
## [14,] "CreditCard"
```

```
library(caTools)
set.seed(1)
split <- sample.split(universal_m.df, SplitRatio = 0.6)
training_set <- subset(universal_m.df, split == TRUE)
validation_set <- subset(universal_m.df, split == FALSE)

print(paste("The size of the training set is:", nrow(training_set)))
```

```
## [1] "The size of the training set is: 2858"
```

```
print(paste("The size of the validation set is:", nrow(validation_set)))
```

```
## [1] "The size of the validation set is: 2142"
```

Now, let us normalize the data

```
train.norm.df <- train.df[, -10]
valid.norm.df <- valid.df[, -10]

norm.values <- preProcess(train.df[, -10], method=c("center", "scale"))
train.norm.df <- predict(norm.values, train.df[, -10])
valid.norm.df <- predict(norm.values, valid.df[, -10])
```

## Questions

Consider the following customer:

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?

```
# We have converted all categorical variables to dummy variables
# Let's create a new sample
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
)

# Normalize the new customer
new.cust.norm <- new_customer
new.cust.norm <- predict(norm.values, new_customer)

```

Now, let us predict using knn

```

knn.pred1 <- class::knn(train = train.norm.df,
                        test = new.cust.norm,
                        cl = train.df$Personal.Loan, k = 1)

knn.pred1

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

```

---

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

```

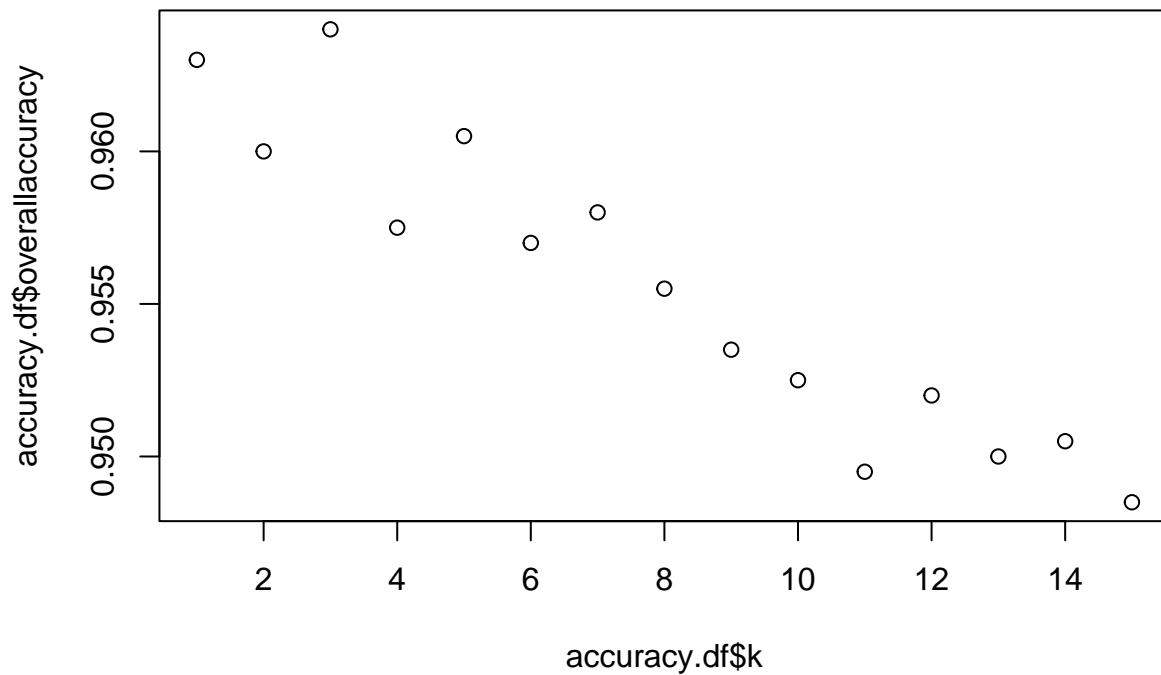
accuracy.df <- data.frame(k = seq(1, 15, 1), overallaccuracy = rep(0, 15))
for(i in 1:15) {
  knn.pred <- class::knn(train = train.norm.df,
                        test = valid.norm.df,
                        cl = train.df$Personal.Loan, k = i)
  accuracy.df[i, 2] <- confusionMatrix(knn.pred,
                                       as.factor(valid.df$Personal.Loan), positive = "1")$overall[1]
}

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

## [1] 3

plot(accuracy.df$k, accuracy.df$overallaccuracy)

```



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

```
best_k <- which(accuracy.df[, 2] == max(accuracy.df[, 2]))
knn.pred_valid <- class::knn(train = train.norm.df,
                             test = valid.norm.df,
                             cl = train.df$Personal.Loan, k = best_k)
confusionMatrix(knn.pred_valid,
                 as.factor(valid.df$Personal.Loan),
                 positive = "1")
```

```
## 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.

```
knn.pred_customer <- class::knn(train = train.norm.df,
                                test = new.cust.norm,
                                cl = train.df$Personal.Loan, k = best_k)

knn.pred_customer
```

```
## [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)
split1 <- sample.split(universal_m.df, SplitRatio = 0.5)
training_set1 <- subset(universal_m.df, split1 == TRUE)
temp <- subset(universal_m.df, split1 == FALSE)
split2 <- sample.split(temp, SplitRatio = 0.6)
validation_set1 <- subset(temp, split2 == TRUE)
test_set1 <- subset(temp, split2 == FALSE)

knn.pred_train <- class::knn(train = training_set1[, -10],
                             test = training_set1[, -10],
                             cl = training_set1$Personal.Loan, k = best_k)
knn.pred_valid <- class::knn(train = training_set1[, -10],
                              test = validation_set1[, -10],
                              cl = training_set1$Personal.Loan, k = best_k)
knn.pred_test <- class::knn(train = training_set1[, -10],
                             test = test_set1[, -10],
                             cl = training_set1$Personal.Loan, k = best_k)

confusion_train <- confusionMatrix(knn.pred_train,
                                   as.factor(training_set1$Personal.Loan),
                                   positive = "1")
confusion_valid <- confusionMatrix(knn.pred_valid,
```

```

                                as.factor(validation_set1$Personal.Loan),
                                positive = "1")
confusion_test <- confusionMatrix(knn.pred_test,
                                as.factor(test_set1$Personal.Loan),
                                positive = "1")

print("Confusion Matrix for Training Set:")

```

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

```
print(confusion_train)
```

```

## Confusion Matrix and Statistics
##
##              Reference
## Prediction    0    1
##              0 2236   97
##              1   31  137
##
##              Accuracy : 0.9488
##              95% CI : (0.9394, 0.9571)
##              No Information Rate : 0.9064
##              P-Value [Acc > NIR] : 2.074e-15
##
##              Kappa : 0.6546
##
##              McNemar's Test P-Value : 9.179e-09
##
##              Sensitivity : 0.58547
##              Specificity : 0.98633
##              Pos Pred Value : 0.81548
##              Neg Pred Value : 0.95842
##              Prevalence : 0.09356
##              Detection Rate : 0.05478
##              Detection Prevalence : 0.06717
##              Balanced Accuracy : 0.78590
##
##              'Positive' Class : 1
##

```

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

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

```
print(confusion_valid)
```

```

## Confusion Matrix and Statistics
##
##              Reference
## Prediction    0    1
##              0 1242   96

```

```
##          1   47   41
##
##          Accuracy : 0.8997
##          95% CI : (0.8829, 0.9148)
##    No Information Rate : 0.9039
##    P-Value [Acc > NIR] : 0.7231
##
##          Kappa : 0.3128
##
##    McNemar's Test P-Value : 5.971e-05
##
##          Sensitivity : 0.29927
##          Specificity : 0.96354
##    Pos Pred Value : 0.46591
##    Neg Pred Value : 0.92825
##    Prevalence : 0.09607
##    Detection Rate : 0.02875
##    Detection Prevalence : 0.06171
##    Balanced Accuracy : 0.63140
##
##    'Positive' Class : 1
##
```

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

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

```
print(confusion_test)
```

```
## Confusion Matrix and Statistics
##
##          Reference
## Prediction  0   1
##          0 934  73
##          1  30  36
##
##          Accuracy : 0.904
##          95% CI : (0.8848, 0.921)
##    No Information Rate : 0.8984
##    P-Value [Acc > NIR] : 0.2924
##
##          Kappa : 0.3626
##
##    McNemar's Test P-Value : 3.498e-05
##
##          Sensitivity : 0.33028
##          Specificity : 0.96888
##    Pos Pred Value : 0.54545
##    Neg Pred Value : 0.92751
##    Prevalence : 0.10158
##    Detection Rate : 0.03355
##    Detection Prevalence : 0.06151
##    Balanced Accuracy : 0.64958
```



```
##  
##      'Positive' Class : 1  
##
```