

FML Assignment 2

Priyanka Jonnala

2023-10-01

Required Libraries

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

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

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

```
library(e1071)
```

Read the data.

```
bankinfo <- read.csv("UniversalBank.csv")
dim(bankinfo)
```

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

```
t(t(names(bankinfo))) # The t function creates a transpose of the dataframe
```

```
##      [,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

```
bankinfo <- bankinfo[,-c(1,5)]
```

Split Data into 60% training and 40% validation. Before we split, let us transform categorical variables into dummy variables

```
# Only Education needs to be converted to factor
bankinfo$Education <- as.factor(bankinfo$Education)

# Now, convert Education to Dummy Variables

groups <- dummyVars(~., data = bankinfo) # This creates the dummy groups
m_bankinfo <- as.data.frame(predict(groups,bankinfo))

set.seed(1) # Important to ensure that we get the same sample if we rerun the code
train.index <- sample(row.names(m_bankinfo), 0.6*dim(m_bankinfo)[1])
val.index <- setdiff(row.names(m_bankinfo), train.index)
train.df <- m_bankinfo[train.index,]
val.df <- m_bankinfo[val.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"
```

Now, let us normalize the data

```
train.norm.df <- train.df[,-10] # Note that Personal Income is the 10th variable
val.norm.df <- val.df[,-10]

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

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?

```
# Calculate the accuracy for each value of k
# Set the range of k values to consider

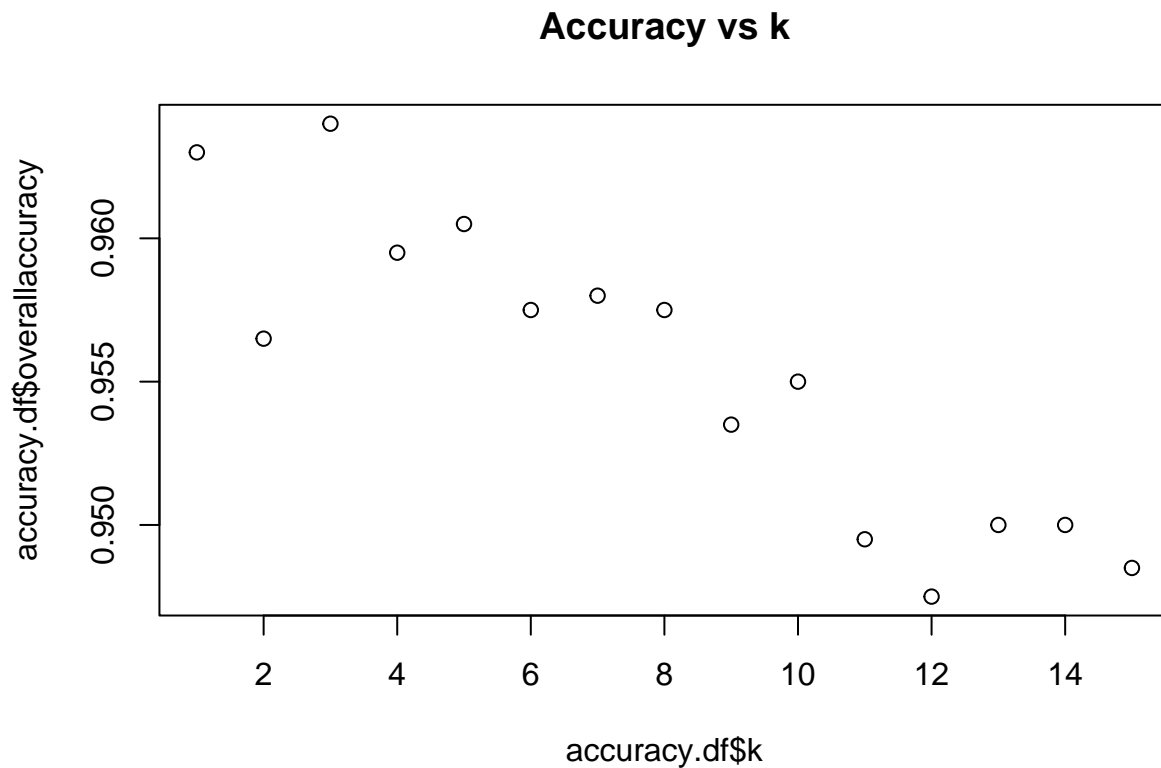
accuracy.df <- data.frame(k = seq(1, 15, 1), overallaccuracy = rep(0, 15))
for(i in 1:15) {
  knn.pred <- class::knn(train = train.norm.df,
    test = val.norm.df,
    cl = train.df$Personal.Loan, k = i)
  accuracy.df[i, 2] <- confusionMatrix(knn.pred,
    as.factor(val.df$Personal.Loan), positive =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

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

```
## [1] 3
```

```
#To find the best k value
```

```
plot(accuracy.df$k,accuracy.df$overallaccuracy,main = "Accuracy vs k")
```



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

```
knn.pred2 <- class::knn(train = train.norm.df,
  test = val.norm.df,
  cl = train.df$Personal.Loan, k = 3)
knn.pred2
```

```
## [1] 0 0 0 0 1 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 1 0 0
## [38] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0
## [75] 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0
## [112] 0 0 1 0 0 1 0 0 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [149] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0
## [186] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [223] 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [260] 0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 0 0 0 0 1
## [297] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0
## [334] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 0 0 0 0 0 0 1 0 0 0 0 0 0 0
## [371] 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0
## [408] 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0
## [445] 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 0 0 0
## [482] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 0
```



```
##      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 CreditCard = 1. Classify the customer using the best k.

```
new_customer1 <- 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 second customer
new.cust.norm1 <- new_customer1
new.cust.norm1 <- predict(norm.values, new.cust.norm1)
```

Using the best k value to predict the second customer

```
knn.pred3 <- class::knn(train = train.norm.df,
  test = new.cust.norm1,
  cl = train.df$Personal.Loan, k = 3)

knn.pred3
```

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

- ```
set.seed(2)
train.index1 <- sample(row.names(m_bankinfo), 0.5*dim(m_bankinfo)[1])
train.df1 <- m_bankinfo[train.index1,]
val.index1 <- setdiff(row.names(m_bankinfo), train.index1)
val.df1 <- m_bankinfo[val.index1,]
val.index2 <- sample(row.names(val.df1), 0.6*dim(val.df1)[1])
val.df2 <- val.df1[val.index2,]
test.index1 <- setdiff(row.names(val.df1), val.index2)
test.df1 <- val.df1[test.index1,]
```

```
train.norm.df1 <- train.df1[, -10]
val.norm.df2 <- val.df2[, -10]
test.norm.df1 <- test.df1[, -10]
norm.values1 <- preprocess(train.df1[, -10], method = c("center", "scale"))
train.norm.df1 <- predict(norm.values1, train.df1[, -10])
val.norm.df2 <- predict(norm.values1, val.df2[, -10])
test.norm.df1 <- predict(norm.values1, test.df1[, -10])
```

```
knn.pred4 <- class::knn(train = train.norm.df1,
 test = train.norm.df1,
 cl = train.df1$Personal.Loan, k = 3)
knn.pred4
```

7





### KNN prediction of Validation data set

[illegible]

[illegible]

```
confusion_matrix1 <-confusionMatrix(knn.pred5,as.factor(val.df2$Personal.Loan))
confusion_matrix1
```



```

1 4 46
##
Accuracy : 0.968
95% CI : (0.9551, 0.978)
No Information Rate : 0.926
P-Value [Acc > NIR] : 1.208e-08
##
Kappa : 0.7256
##
McNemar's Test P-Value : 4.785e-05
##
Sensitivity : 0.9957
Specificity : 0.6216
Pos Pred Value : 0.9705
Neg Pred Value : 0.9200
Prevalence : 0.9260
Detection Rate : 0.9220
Detection Prevalence : 0.9500
Balanced Accuracy : 0.8087
##
'Positive' Class : 0
##

```

#### Comparing and commenting on the Confusion Matrices

1. Training set confusion matrix The training set typically gets the best results because the model is already trained on this data. It gets more true positives and true negatives, and less false positives and false negatives. This is because the model is already trained well and sometimes it even memorizes it.
2. Validation set confusion matrix The validation set gives a realistic view of model's performance as it wasn't part of the training. It gets more balanced results with moderate values for true positives, true negatives, false positives, false negatives.
3. Test set confusion matrix This confusion matrix represents how the model performs on entirely new and unseen data. It shows lower performance compared to the training and validation sets. It gets more false positives and false negatives.