

Analysis of Online Shoppers' Purchase Intention

Kaarthik Sundaramoorthy, Sahil Shah and Vidhi Shah

7/20/2020

Milestone 1: Project Proposal

Problem Definition

In the era of online shopping, known as e-shopping, people use online transactions to buy the items they need while exploring it online. This helps the buyers as well as sellers to understand the patterns, intentions, and behavior of various online customers. Thereby, helping businesses improve their revenue by focusing on customer experiences and marketing. Hence, the analysis of online shoppers' purchase intention has become an emerging field in data mining. Click-stream analysis refers to the online shoppers' behavior analysis as they invoke a sequence of web pages in a particular session. Therefore, analyzing this data is a primary goal for successful online businesses as they extract the clicks and behavior through web page requests. Our proposed solution is to provide a decisive and feasible recommendation algorithm that will allow us to predict the behavior of the shoppers'.

System Funtionality

Firstly, we will build the model and analyze the performance of various classification, Regression and Clustering Algorithms. Which includes Logistic Regression, Decision Tree, K-Means Clustering and K-NN Classification. Then, the values of different evaluation metrics like Accuracy, Precision, Recall, F-score will be calculated to compare the performance of each of the algorithms. Lastly, we also plan on using these models to predict the shopper's intentions.

Requirements and Benefits

In today's economy e-commerce is becoming more extensive and businesses within this sector need to understand, the factors which come into play when a shopper ventures into a website to make a purchase. The benefit this holds is that it will enable the websites to better target ads or other factors which may lead to an increase in sales. These findings support the feasibility of accurate and scalable purchasing intention prediction for virtual shopping environments. This also helps in knowing the market capabilities of the brand when released in the new market while finding out the problems in the existing market and helps in relevant marketing strategies that can help in conquering the market.

Dataset Details and Core Algorithm

The dataset which is used is based on the "Online Shoppers Purchasing Intention" UCI dataset. It consists of numerical as well as categorical data. There are a total of 12,330 records where each row corresponds to the session data of the particular user. The total no. of records for which the session ended without any purchase is 10,442 which contributes to 84.5%. The core algorithm which we will be implementing is a

Decision Tree. Being a classification algorithm, it creates a tree-like structure by creating rules for breaking the dataset into small subsets in each step. We create a training model that is used to predict the class of the variable by simply learning decision rules deduced from the training set. At each step, a decision is taken to classify the data in the beneath classes. The leaf node holds the final results.

Milestone 2: Data Summary/visualization

The dataset used in the project is based on “*online shoppers purchasing intention*” available on UCI Machine Learning dataset.

Importing Libraries

This are the important libraries that are to be installed for the execution of the file.

```
library(ggplot2)
library(tidyverse)
library(gmodels)
library(dplyr)
library(ggmosaic)
library(corrplot)
library(caret)
library(rpart)
library(rpart.plot)
library(cluster)
library(fpc)
library(data.table)
library(knitr)
library(kableExtra)
library(plyr)
library(caTools)
```

Importing the Dataset

The `read.csv()` command is used to import the dataset.

```
dataset <- read.csv("online_shoppers_intention.csv", header = TRUE)
attach(dataset)
```

Checking the number of columns and rows of the dataset.

```
ncol(dataset)
```

```
## [1] 18
```

```
nrow(dataset)
```

```
## [1] 12330
```

Looking at the dataset data structure.

```
str(dataset)
```

```
## 'data.frame': 12330 obs. of 18 variables:
## $ Administrative : int 0 0 0 0 0 0 0 1 0 0 ...
## $ Administrative_Duration: num 0 0 0 0 0 0 0 0 0 0 ...
## $ Informational : int 0 0 0 0 0 0 0 0 0 0 ...
## $ Informational_Duration : num 0 0 0 0 0 0 0 0 0 0 ...
## $ ProductRelated : int 1 2 1 2 10 19 1 0 2 3 ...
## $ ProductRelated_Duration: num 0 64 0 2.67 627.5 ...
## $ BounceRates : num 0.2 0 0.2 0.05 0.02 ...
## $ ExitRates : num 0.2 0.1 0.2 0.14 0.05 ...
## $ PageValues : num 0 0 0 0 0 0 0 0 0 0 ...
## $ SpecialDay : num 0 0 0 0 0 0 0.4 0 0.8 0.4 ...
## $ Month : Factor w/ 10 levels "Aug","Dec","Feb",...: 3 3 3 3 3 3 3 3 3 3 ...
## $ OperatingSystems : int 1 2 4 3 3 2 2 1 2 2 ...
## $ Browser : int 1 2 1 2 3 2 4 2 2 4 ...
## $ Region : int 1 1 9 2 1 1 3 1 2 1 ...
## $ TrafficType : int 1 2 3 4 4 3 3 5 3 2 ...
## $ VisitorType : Factor w/ 3 levels "New_Visitor",...: 3 3 3 3 3 3 3 3 3 3 ...
## $ Weekend : logi FALSE FALSE FALSE FALSE TRUE FALSE ...
## $ Revenue : logi FALSE FALSE FALSE FALSE FALSE FALSE ...
```

```
summary(dataset)
```

```
## Administrative Administrative_Duration Informational
## Min. : 0.000 Min. : 0.00 Min. : 0.0000
## 1st Qu.: 0.000 1st Qu.: 0.00 1st Qu.: 0.0000
## Median : 1.000 Median : 7.50 Median : 0.0000
## Mean : 2.315 Mean : 80.82 Mean : 0.5036
## 3rd Qu.: 4.000 3rd Qu.: 93.26 3rd Qu.: 0.0000
## Max. :27.000 Max. :3398.75 Max. :24.0000
##
## Informational_Duration ProductRelated ProductRelated_Duration
## Min. : 0.00 Min. : 0.00 Min. : 0.0
## 1st Qu.: 0.00 1st Qu.: 7.00 1st Qu.: 184.1
## Median : 0.00 Median : 18.00 Median : 598.9
## Mean : 34.47 Mean : 31.73 Mean : 1194.8
## 3rd Qu.: 0.00 3rd Qu.: 38.00 3rd Qu.: 1464.2
## Max. :2549.38 Max. :705.00 Max. :63973.5
##
## BounceRates ExitRates PageValues SpecialDay
## Min. :0.000000 Min. :0.00000 Min. : 0.000 Min. :0.00000
## 1st Qu.:0.000000 1st Qu.:0.01429 1st Qu.: 0.000 1st Qu.:0.00000
## Median :0.003112 Median :0.02516 Median : 0.000 Median :0.00000
## Mean :0.022191 Mean :0.04307 Mean : 5.889 Mean :0.06143
## 3rd Qu.:0.016813 3rd Qu.:0.05000 3rd Qu.: 0.000 3rd Qu.:0.00000
## Max. :0.200000 Max. :0.20000 Max. :361.764 Max. :1.00000
##
## Month OperatingSystems Browser Region
## May :3364 Min. :1.000 Min. : 1.000 Min. :1.000
## Nov :2998 1st Qu.:2.000 1st Qu.: 2.000 1st Qu.:1.000
## Mar :1907 Median :2.000 Median : 2.000 Median :3.000
## Dec :1727 Mean :2.124 Mean : 2.357 Mean :3.147
```

```
## Oct      : 549      3rd Qu.:3.000      3rd Qu.: 2.000      3rd Qu.:4.000
## Sep      : 448      Max.       :8.000      Max.       :13.000      Max.       :9.000
## (Other):1337
## TrafficType      VisitorType      Weekend      Revenue
## Min.      : 1.00      New_Visitor      : 1694      Mode :logical      Mode :logical
## 1st Qu.: 2.00      Other              :   85      FALSE:9462      FALSE:10422
## Median : 2.00      Returning_Visitor:10551      TRUE :2868      TRUE :1908
## Mean      : 4.07
## 3rd Qu.: 4.00
## Max.      :20.00
##
```

The purchasing intention model is designed as a classification problem which measures the purchasers' commitment to finalize purchase intent. Hence we have the session data of the users which has two categories : users who purchased the item and who didn't. The dataset consists of both numerical data and categorical data, and thus the target value is categorical. Table 1 refers to the numerical features and Table 2 refers to the categorical features used in the prediction model respectively. There are a total of 12,330 rows where each row represents session data of one particular user.

```
tab1 <- read.csv("table1.csv", header = TRUE)
kable(tab1) %>%
  kable_styling(full_width = T)
```

i.Feature.Name	Description	Min..value	Max..value	SD
Administrative	Number of pages visited by the visitor about account management	0	27.0	3.322e+00
Administrative duration	Total amount of time (in seconds) spent by the visitor on account management related pages	0	3399.0	1.768e+02
Informational	Number of pages visited by the visitor about Web site, communication and address information of the shopping site	0	24.0	1.270e+00
Informational duration	Total amount of time (in seconds) spent by the visitor on informational pages	0	2549.4	1.407e+02
Product related	Number of pages visited by visitor about product related pages	0	705.0	4.448e+01
Product related duration	Total amount of time (in seconds) spent by the visitor on product related pages	0	63974.0	1.914e+03
Bounce rates	Average bounce rate value of the pages visited by the visitor	0	0.2	4.849e-02
Exit rate	Average exit rate value of the pages visited by the visitor	0	0.2	4.860e-02
Page value	Average page value of the pages visited by the visitor	0	361.8	1.857e+01
Special day	Closeness of the site visiting time to a special day	0	1.0	1.989e-01

```
tab2 <- read.csv("table2.csv", header = TRUE)
kable(tab2) %>%
  kable_styling(full_width = T)
```

i.Name	Description	Values
OperatingSystems	Operating system of the visitor	8
Browser	Browser of the visitor	13
Region	Geographic region from which the session has been started by the visitor	9
TrafficType	Traffic source by which the visitor has arrived at the Web site (e.g., banner, SMS, direct)	20
VisitorType	Visitor type as New Visitor, Returning Visitor, and Other	3
Weekend	Boolean value indicating whether the date of the visit is weekend	2
Month	Month value of the visit date	10
Revenue	Class label indicating whether the visit has been finalized with a transaction	2

Taking the look at the **REVENUE** column which is the target column. The datatype of the REVENUE column is Logical which holds the value **TRUE** and **FALSE**.

```
library(gmodels)
summary(dataset$Revenue)
```

```
##      Mode   FALSE    TRUE
## logical  10422    1908
```

```
CrossTable(dataset$Revenue)
```

```
##
##
##      Cell Contents
## |-----|
## |                      N |
## |      N / Table Total |
## |-----|
##
##
## Total Observations in Table:  12330
##
##
##      |      FALSE |      TRUE |
##      |-----|-----|
##      |    10422 |    1908 |
##      |    0.845 |    0.155 |
##      |-----|-----|
##
##
##
##
```

Adding the new *Revenue_binary* column by using Logical Data of Shopper's Revenue into binary dependent variable that will helpful for potential regression models. The data will be converted with values 0 and 1, i.e. If it is false the value is 0 and if true it will be 1.

```
datasetbinary <- dataset %>%
  mutate(Revenue_binary = ifelse(dataset$Revenue == "TRUE", 1, 0))
```

Checking the dataset if it has any missing values. Even if it has we will remove it.

```
# colSums(is.na(dataset))
sum(is.na(dataset))
```

```
## [1] 0
```

```
sapply(dataset, function(x) sum(is.na(x)))
```

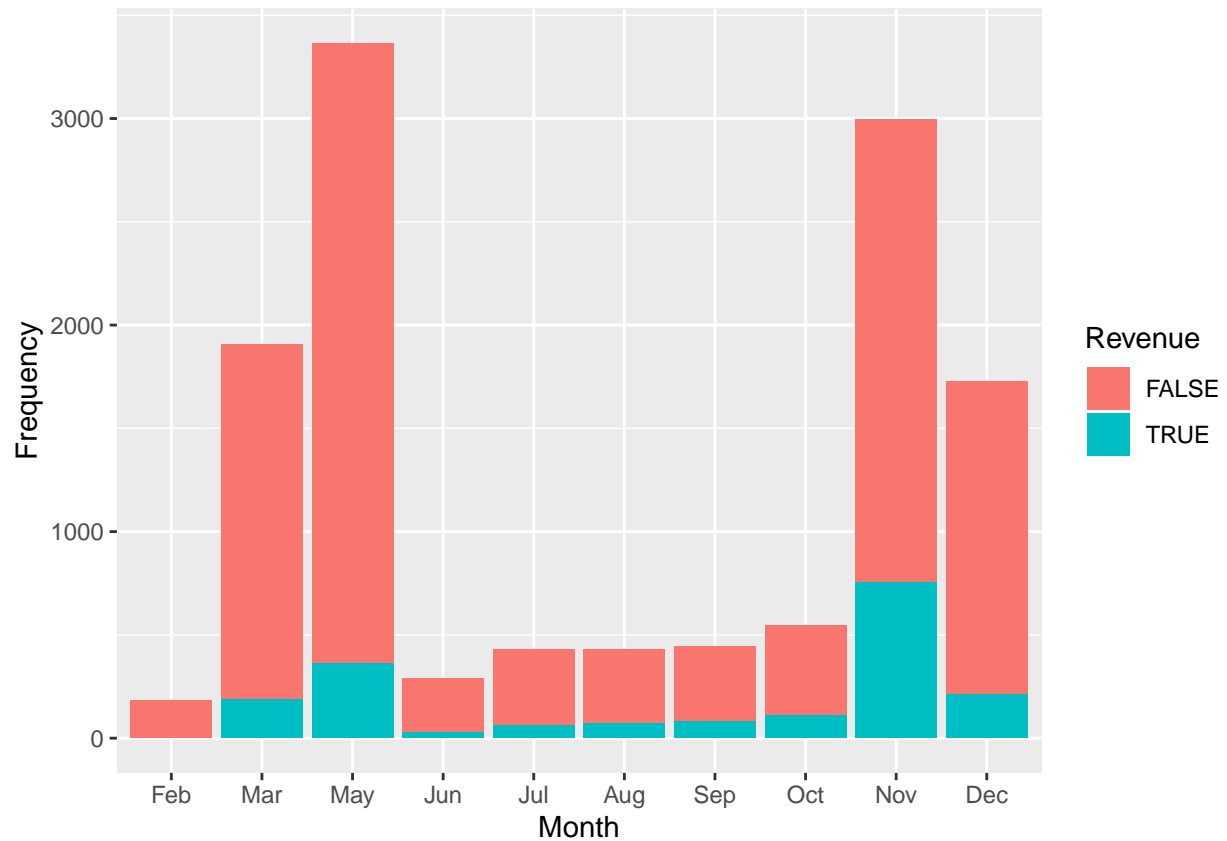
```
##      Administrative Administrative_Duration      Informational
##      0              0              0
## Informational_Duration      ProductRelated ProductRelated_Duration
##      0              0              0
##      BounceRates      ExitRates      PageValues
##      0              0              0
##      SpecialDay      Month      OperatingSystems
##      0              0              0
##      Browser      Region      TrafficType
##      0              0              0
##      VisitorType      Weekend      Revenue
##      0              0              0
```

```
dataset <- na.omit(dataset)
```

Visualizations

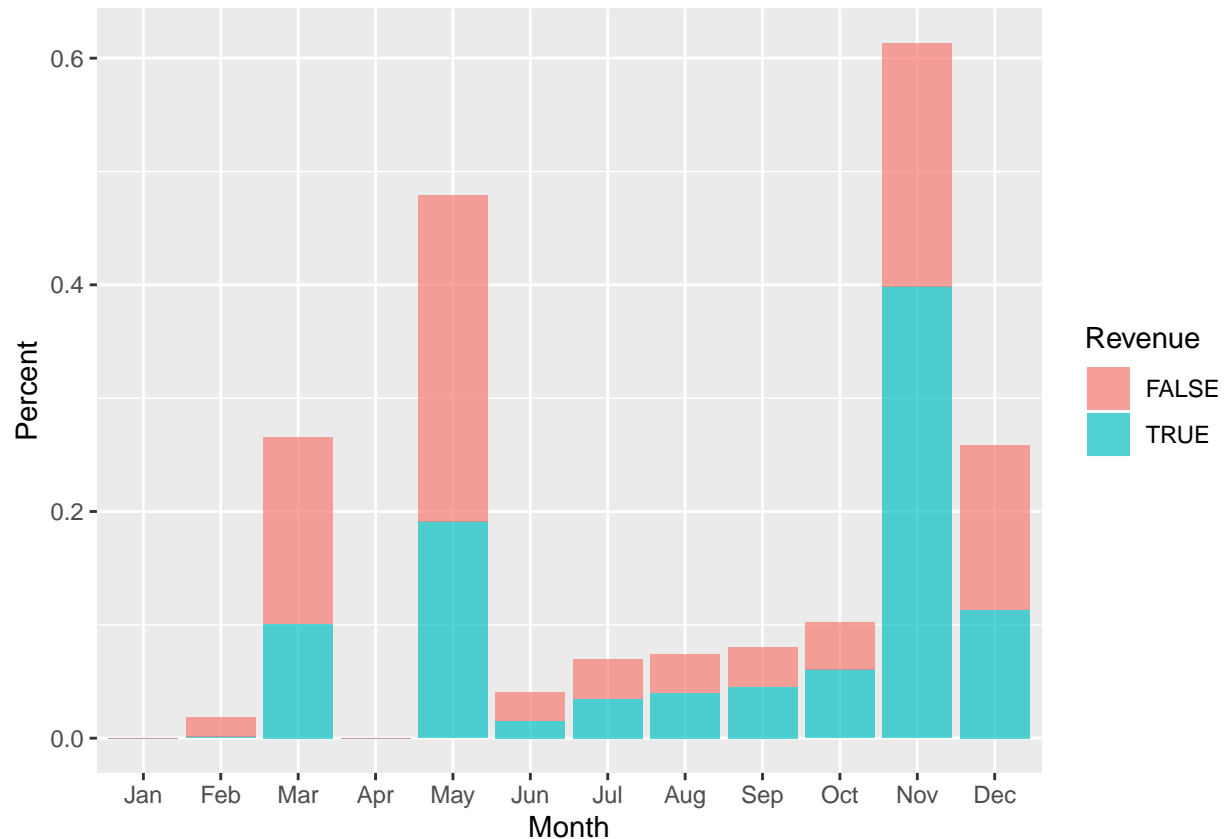
Month

```
dataset$Month = factor(dataset$Month, levels = month.abb)
dataset %>%
  ggplot() +
  aes(x = Month, Revenue = ..count../nrow(dataset), fill = Revenue) +
  geom_bar() +
  ylab("Frequency")
```



The plot describes the frequency of the revenue generated over the months.

```
table_month = table(dataset$Month, dataset$Revenue)
tab_mon = as.data.frame(prop.table(table_month,2))
colnames(tab_mon) = c("Month", "Revenue", "perc")
ggplot(data = tab_mon, aes(x = Month, y = perc, fill = Revenue)) +
  geom_bar(stat = 'identity', pdatasettion = 'dodge', alpha = 2/3) +
  xlab("Month")+
  ylab("Percent")
```

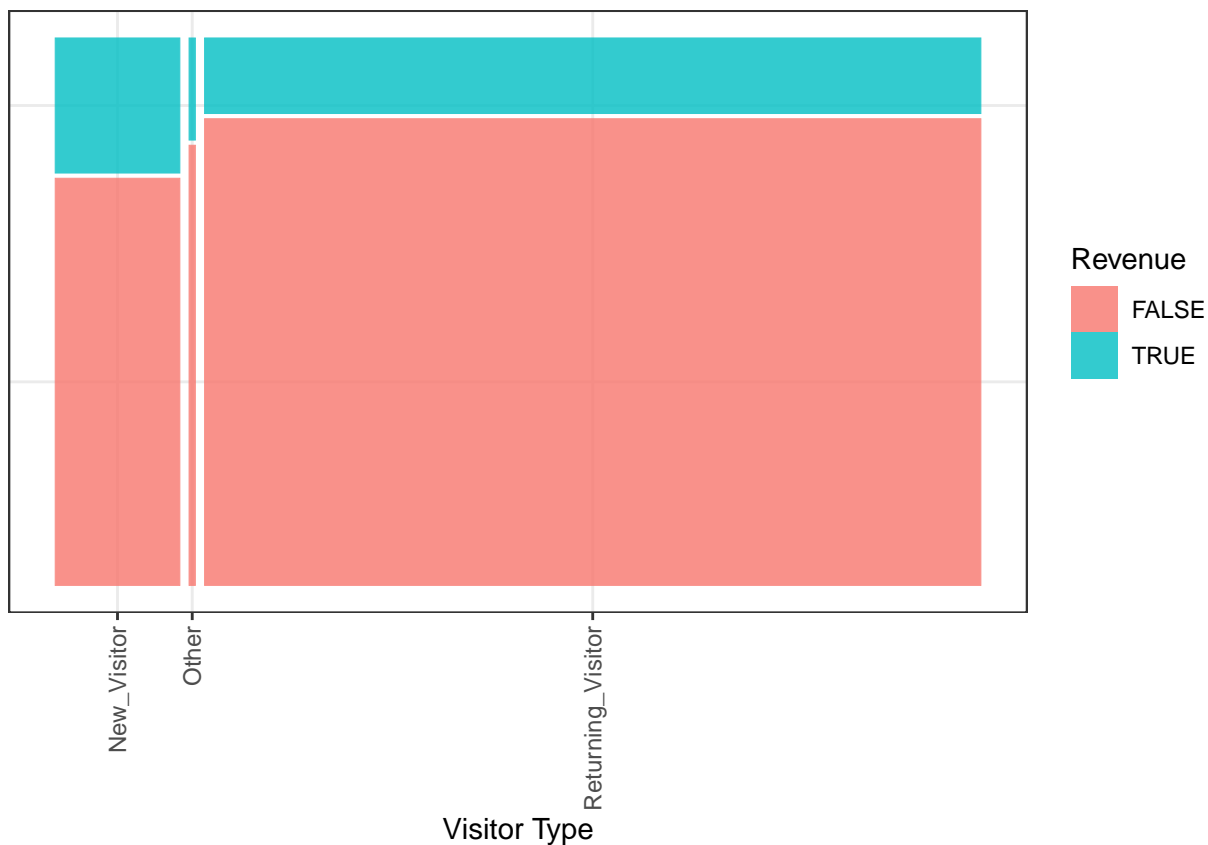
The plot portrays the high shopping rates in the months September, October and November with respect to the customers not buying the products. These months are comparatively considered as the *Holiday Season Months*. Also, there is high hits on the website with pdatasettive revenue in the month of may.

Visitor

```
theme_set(theme_bw())

## setting default parameters for mosaic plots
mosaic_theme = theme(axis.text.x = element_text(angle = 90,
                                                  hjust = 1,
                                                  vjust = 0.5),
                     axis.text.y = element_blank(),
                     axis.ticks.y = element_blank())

dataset %>%
  ggplot() +
  geom_mosaic(aes(x = product(Revenue, VisitorType), fill = Revenue)) +
  mosaic_theme +
  xlab("Visitor Type") +
  ylab(NULL)
```



The comparison of the VisitorType which are New_Visitors, Returning_Visitor and Others with Revenue generated. There are many returning visitors in the contrast to less new visitors. Although, the new visitors have high probability of purchasing the product and help the revenue than the returning visitors.

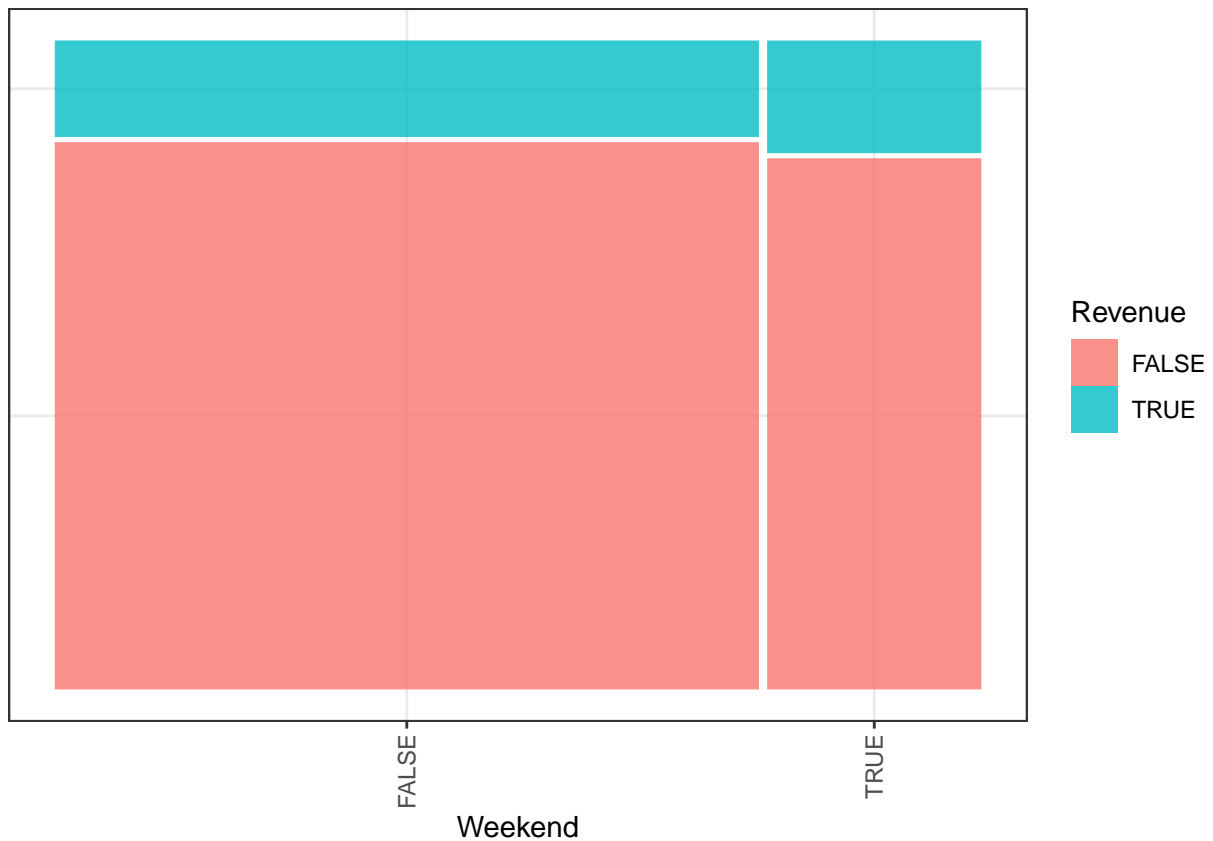
Weekend

```
CrossTable(dataset$Weekend, dataset$Revenue)
```

```
##
##
##   Cell Contents
## |-----|
## |                      N |
## | Chi-square contribution |
## |      N / Row Total |
## |      N / Col Total |
## |      N / Table Total |
## |-----|
##
##
## Total Observations in Table: 12330
##
##
##           | dataset$Revenue
```

## dataset\$Weekend	FALSE	TRUE	Row Total
## -----	-----	-----	-----
## FALSE	8053	1409	9462
##	0.381	2.080	
##	0.851	0.149	0.767
##	0.773	0.738	
##	0.653	0.114	
## -----	-----	-----	-----
## TRUE	2369	499	2868
##	1.257	6.864	
##	0.826	0.174	0.233
##	0.227	0.262	
##	0.192	0.040	
## -----	-----	-----	-----
## Column Total	10422	1908	12330
##	0.845	0.155	
## -----	-----	-----	-----
##			
##			

```
dataset %>%
  ggplot() +
  mosaic_theme +
  geom_mosaic(aes(x = product(Revenue, Weekend), fill = Revenue)) +
  xlab("Weekend") +
  ylab(NULL)
```



The **Weekend** analysis shows that more than 70% of visitors are visiting the site on weekdays, with 15% chance of actually buying the products. The rest 30% visit on the weekend and there is 17% speculation of buying.

Milestone 3: Algorithm Testing

Data preparation for the algorithms testing, looking up for the features who has too many levels and transforming some feature with `as.factor()` function for easier execution.

```
dataset$Revenue <- as.factor(dataset$Revenue)
dataset$VisitorType <- as.factor(dataset$VisitorType)
dataset$TrafficType <- as.factor(dataset$TrafficType)
dataset$OperatingSystems <- as.factor(dataset$OperatingSystems)
dataset$Browser <- as.factor(dataset$Browser)
dataset$Region <- as.factor(dataset$Region)
dataset$Month <- as.factor(dataset$Month)
summary(dataset$Browser)
```

```
##      1      2      3      4      5      6      7      8      9     10     11     12     13
## 2462 7961  105  736  467  174   49  135    1  163    6   10   61
```

```
summary(dataset$TrafficType)
```

```
##      1      2      3      4      5      6      7      8      9     10     11     12     13     14     15     16
## 2451 3913 2052 1069  260  444   40  343  42  450  247    1  738   13   38    3
##      17     18     19     20
##      1     10     17    198
```

```
dataset$Browser = recode(dataset$Browser, '3' = '1')
dataset$Browser = recode(dataset$Browser, '4' = '1')
dataset$Browser = recode(dataset$Browser, '5' = '1')
dataset$Browser = recode(dataset$Browser, '6' = '1')
dataset$Browser = recode(dataset$Browser, '7' = '1')
dataset$Browser = recode(dataset$Browser, '8' = '1')
dataset$Browser = recode(dataset$Browser, '9' = '1')
dataset$Browser = recode(dataset$Browser, '10' = '1')
dataset$Browser = recode(dataset$Browser, '11' = '1')
dataset$Browser = recode(dataset$Browser, '12' = '1')
dataset$Browser = recode(dataset$Browser, '13' = '1')
dataset$TrafficType = recode(dataset$TrafficType, '6' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '7' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '8' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '9' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '10' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '11' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '12' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '13' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '14' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '15' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '16' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '17' = '5')
```

```
dataset$TrafficType = recode(dataset$TrafficType, '18' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '19' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '20' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '4' = '3')
dataset$TrafficType = recode(dataset$TrafficType, '5' = '4')
summary(dataset$Browser)
```

```
##      1      2
## 4369 7961
```

```
summary(dataset$TrafficType)
```

```
##      1      2      3      4
## 2451 3913 3121 2845
```

Logistic Regression

In this section we are implementing Logistic Regression, Data Preparing for the logistic regression and splitting the data in the ratio 75:25. Also preprocessing the data with `preProcess()` function and also feature scaling the data to get better results for the algorithm.

```
dataset_logistic <- dataset
set.seed(123)
split_log = createDataPartition(dataset_logistic$Revenue, p=0.75, list = FALSE, times = 1)
training_log = dataset_logistic[split_log,]
test_log = dataset_logistic[-split_log,]
preprocess_log <- preProcess(training_log, method = c("center", "scale"))
train_pre_log <- predict(preprocess_log, training_log)
test_pre_log <- predict(preprocess_log, test_log)
```

Building the Logistic Model with `glm()` function on the *Revenue* feature as the dependent variable.

```
set.seed(123)
log_reg <- glm(Revenue ~ ., data = train_pre_log, family = binomial)
summary(log_reg)
```

```
##
## Call:
## glm(formula = Revenue ~ ., family = binomial, data = train_pre_log)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -5.9831  -0.4654  -0.3292  -0.1586   3.3078
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -3.751072    0.773844  -4.847 1.25e-06 ***
## Administrative    0.034895    0.041685   0.837  0.40252
## Administrative_Duration 0.010610    0.038675   0.274  0.78383
## Informational    0.026663    0.039123   0.682  0.49554
```

```

## Informational_Duration      0.009812  0.035229  0.279  0.78061
## ProductRelated             0.034969  0.061687  0.567  0.57080
## ProductRelated_Duration    0.147627  0.063146  2.338  0.01939 *
## BounceRates                -0.079191  0.181232 -0.437  0.66214
## ExitRates                  -0.791956  0.136423 -5.805  6.43e-09 ***
## PageValues                 1.487787  0.051349 28.974 < 2e-16 ***
## SpecialDay                 -0.019035  0.053910 -0.353  0.72402
## MonthMar                   1.332588  0.767544  1.736  0.08253 .
## MonthMay                   1.376902  0.761498  1.808  0.07058 .
## MonthJun                   1.766948  0.800388  2.208  0.02727 *
## MonthJul                   1.953797  0.780886  2.502  0.01235 *
## MonthAug                   1.986493  0.780550  2.545  0.01093 *
## MonthSep                   1.869687  0.778831  2.401  0.01637 *
## MonthOct                   1.784202  0.775666  2.300  0.02144 *
## MonthNov                   2.454768  0.763003  3.217  0.00129 **
## MonthDec                   1.330743  0.767885  1.733  0.08310 .
## OperatingSystems2          0.105354  0.112706  0.935  0.34991
## OperatingSystems3          -0.202840  0.142393 -1.425  0.15430
## OperatingSystems4          -0.312258  0.209604 -1.490  0.13629
## OperatingSystems5           0.149420  1.236449  0.121  0.90381
## OperatingSystems6          -0.921330  0.939345 -0.981  0.32668
## OperatingSystems7           1.410105  1.191687  1.183  0.23670
## OperatingSystems8           0.380293  0.673965  0.564  0.57258
## Browser2                   -0.119399  0.094764 -1.260  0.20768
## Region2                    0.109857  0.127345  0.863  0.38832
## Region3                    -0.015605  0.099689 -0.157  0.87561
## Region4                    -0.047062  0.130217 -0.361  0.71779
## Region5                    -0.388573  0.248886 -1.561  0.11846
## Region6                    0.082921  0.152648  0.543  0.58698
## Region7                    0.046041  0.152237  0.302  0.76232
## Region8                    -0.080109  0.207849 -0.385  0.69993
## Region9                    -0.363039  0.197184 -1.841  0.06560 .
## TrafficType2               0.126742  0.109611  1.156  0.24756
## TrafficType3               -0.125890  0.125472 -1.003  0.31570
## TrafficType4               0.152666  0.119478  1.278  0.20133
## VisitorType0ther           -1.762535  0.971931 -1.813  0.06976 .
## VisitorTypeReturning_Visitor -0.252049  0.100949 -2.497  0.01253 *
## WeekendTRUE                0.066817  0.083298  0.802  0.42247
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 7968.8 on 9247 degrees of freedom
## Residual deviance: 5371.4 on 9206 degrees of freedom
## AIC: 5455.4
##
## Number of Fisher Scoring iterations: 7

```

The building of the logistic regression model, we can see that for the dataset, most significant features with respect to the Revenue are Exit Rates, Page Values, and Month of November.

Now Implementing the model on the test dataset to test out the accuracy of the Logistic Model.

```
log_reg_pred <- predict(log_reg, newdata = test_pre_log, type = "response")
log_reg_pred <- ifelse(log_reg_pred>0.5,"TRUE","FALSE")
logcf <- table(log_reg_pred, test_pre_log$Revenue)
confusionMatrix(logcf)
```

```
## Confusion Matrix and Statistics
##
##
## log_reg_pred FALSE TRUE
##      FALSE  2553  288
##      TRUE    52  189
##
##              Accuracy : 0.8897
##              95% CI : (0.8781, 0.9005)
##      No Information Rate : 0.8452
##      P-Value [Acc > NIR] : 6.345e-13
##
##              Kappa : 0.4716
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##      Sensitivity : 0.9800
##      Specificity : 0.3962
##      Pos Pred Value : 0.8986
##      Neg Pred Value : 0.7842
##      Prevalence : 0.8452
##      Detection Rate : 0.8284
##      Detection Prevalence : 0.9218
##      Balanced Accuracy : 0.6881
##
##      'Positive' Class : FALSE
##
```

We got nearly 89% of accuracy in the Logistic regression model. We got good analysis running the Logistic Regression Model.

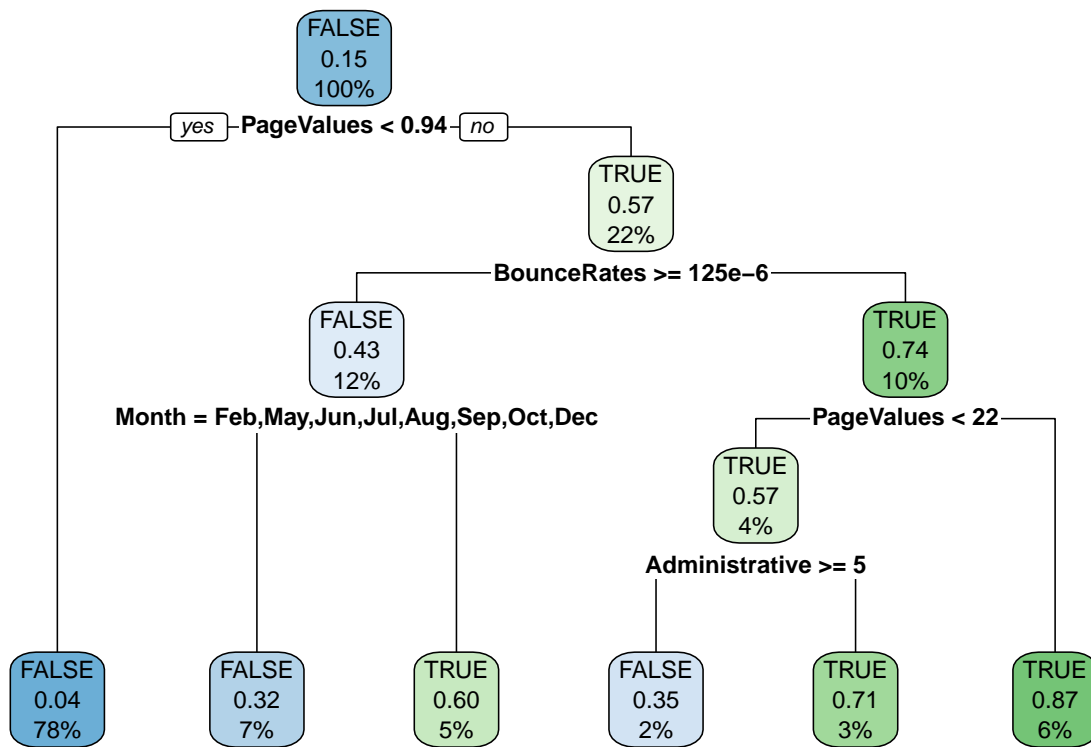
Decision Tree

Now preparing the dataset for the decision tree classification algorithm. Splitting the dataset into testing and training data for training the model. We will split the data in ratio 75:25.

```
dataset_classify = dataset
set.seed(123)
split = sample.split(dataset_classify$Revenue, SplitRatio = 0.75)
training_data = subset(dataset_classify, split == TRUE)
test_data = subset(dataset_classify, split == FALSE)
```

Running the decision tree from the `rpart()` library:

```
dt_model<- rpart(Revenue ~ . , data = training_data, method="class")
rpart.plot(dt_model)
```



The predictive model suggests that Page Values greater than 0.94 lead to a TRUE 57% of the time. On top of this, an effective Bounce Rate above 0 improves our TRUE to 74% and Administrative type '5' or below result in a TRUE 83% of the time. Also, we see that October and December are good months for shoppers' conversions.

```
dt.pred <- predict(dt_model, test_data, type = "class")
mean(dt.pred == test_data$Revenue)
```

```
## [1] 0.8965294
```

Confusion Matrix for Decision Tree

```
cfdt <- table(dt.pred, test_data$Revenue)
```

Accuracy measures for Decision Tree:

```
confusionMatrix(cfdt)
```

```
## Confusion Matrix and Statistics
##
##
## dt.pred FALSE TRUE
## FALSE 2479 192
## TRUE 127 285
```



```
##
##           Accuracy : 0.8965
##           95% CI : (0.8852, 0.9071)
##      No Information Rate : 0.8453
##      P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.5811
##
##      McNemar's Test P-Value : 0.0003393
##
##           Sensitivity : 0.9513
##           Specificity : 0.5975
##      Pos Pred Value : 0.9281
##      Neg Pred Value : 0.6917
##           Prevalence : 0.8453
##      Detection Rate : 0.8041
##      Detection Prevalence : 0.8664
##      Balanced Accuracy : 0.7744
##
##      'Positive' Class : FALSE
##
```

Achieved approximately 90% of accuracy in Decision Tree classification. Which is really good compared to the logistic regression as the classification techniques.

K - Means Clustering

Preparing the dataset for the K-Means Clustering Algorithm. Including the conversion of the Categorical values to numeric or factors.

```
dataset_kmeans <- dataset
dataset_kmeans$OperatingSystems <- factor(dataset$OperatingSystems,
                                           order = TRUE,
                                           levels = c(6,3,7,1,5,2,4,8))
dataset_kmeans$Browser <- factor(dataset$Browser,
                                  order = TRUE,
                                  levels = c(9,3,6,7,1,2,8,11,4,5,10,13,12))
dataset_kmeans$Region <- factor(dataset$Region,
                                 order = TRUE,
                                 levels = c(8,6,3,4,7,1,5,2,9))
dataset_kmeans$TrafficType <- factor(dataset$TrafficType,
                                     order = TRUE,
                                     levels = c(12,15,17,18,13,19,3,9,1,6,4,14,11,10,5,2,20,8,7,16))
dataset_kmeans$Month <- factor(dataset_kmeans$Month,
                               order = TRUE,
                               levels = c('Jan','Feb','Mar','Apr','May',
                                           'Jun','Jul','Aug','Sep','Oct','Nov','Dec'),
                               labels = c(1,2,3,4,5,6,7,8,9,10,11,12))
dataset_kmeans$VisitorType <- factor(dataset_kmeans$VisitorType,
                                     order = TRUE,
                                     levels=c('Returning_Visitor','Other','New_Visitor'),
                                     labels = c(1,2,3))
dataset_kmeans$Weekend <- mapvalues(dataset$Weekend,
```

```

        from = c('TRUE','FALSE'),
        to = c(1,0))
# dataset_kmeans$Revenue <- as.numeric(as.factor(dataset_kmeans$Revenue))
dataset_kmeans$Administrative <- as.numeric(dataset_kmeans$Administrative)
dataset_kmeans$ProductRelated <- as.numeric(dataset_kmeans$ProductRelated)
dataset_kmeans$Informational <- as.numeric(dataset_kmeans$Informational)

```

Creating the normalizing function for the 10 variables that needs feature scaling for the better implementation of the Euclidian distance in the algorithm.

```

normalize <- function(x) {
  return((x-min(x))/(max(x) - min(x)))
}
dataset_kmeans$Administrative <- normalize(dataset_kmeans$Administrative)
dataset_kmeans$Administrative_Duration <- normalize(dataset_kmeans$Administrative_Duration)
dataset_kmeans$Informational <- normalize(dataset_kmeans$Informational_Duration)
dataset_kmeans$Informational_Duration <- normalize(dataset_kmeans$Administrative)
dataset_kmeans$ProductRelated <- normalize(dataset_kmeans$ProductRelated)
dataset_kmeans$ProductRelated_Duration <- normalize(dataset_kmeans$ProductRelated_Duration)
dataset_kmeans$BounceRates <- normalize(dataset_kmeans$BounceRates)
dataset_kmeans$ExitRates <- normalize(dataset_kmeans$ExitRates)
dataset_kmeans$PageValues <- normalize(dataset_kmeans$PageValues)
dataset_kmeans$SpecialDay <- normalize(dataset_kmeans$SpecialDay)

```

Building the K-Means model with `kmeans()` function and assigning K as 2, as we know there are only two clusters in Revenue feature that is either *True* or *False*. In other words to mention whether the customer helped to generate the revenue of not.

```

kmeans_clust <- kmeans(dataset_kmeans[-18], centers = 2, iter.max = 150)
kmeans_clust$size

```

```
## [1] 6155 6175
```

```

cm_km <- table(kmeans_clust$cluster, dataset_kmeans$Revenue)
cm_km

```

```

##
##      FALSE TRUE
##  1  4902 1253
##  2  5520  655

```

```

kmeans_accuracy = sum(diag(cm_km))/sum(cm_km)
kmeans_precision<- cm_km[1,1]/(sum(cm_km[1,]))
kmeans_recall<- cm_km[1,1]/(sum(cm_km[,1]))
kmeans_fscore <- 2*kmeans_precision*kmeans_recall/(kmeans_precision+kmeans_recall)
kmeans_accuracy

```

```
## [1] 0.4506894
```


Milestone 4: Core Algorithm Fine Tuning

KNN

Accuracy tells us the percentage of test cases that have been correctly classified, i.e, the number of test cases among all for which we could correctly identify if the “Revenue” is positive or negative. Precision gives us the ratio of correctly positive test cases among all the positive test cases predicted by the algorithm. The sensitivity represents the ratio of the correctly predicted positive test cases among all the actually positive test cases. Specificity represents the ratio of the correctly predicted negative test cases among all the actually negative test cases. These parameters are used as a comparison measure for various algorithms.

Again, we have to prepare the data for analysis. The data preprocessing has to be done before we run any machine learning algorithm. The data preprocessing includes the preparation of the data converting from categorical data to ordinal factors. Creating Binary variables for the *Weekend* column for False as ‘0’ and True as ‘1’.

```
dataset_knn <- dataset
dataset_knn$OperatingSystems <- factor(dataset$OperatingSystems,
                                       order = TRUE,
                                       levels = c(6,3,7,1,5,2,4,8))
dataset_knn$Browser <- factor(dataset$Browser,
                              order = TRUE,
                              levels = c(9,3,6,7,1,2,8,11,4,5,10,13,12))
dataset_knn$Region <- factor(dataset$Region,
                              order = TRUE,
                              levels = c(8,6,3,4,7,1,5,2,9))
dataset_knn$TrafficType <- factor(dataset$TrafficType,
                                  order = TRUE,
                                  levels = c(12,15,17,18,13,19,3,9,1,6,4,14,11,10,5,2,20,8,7,16))
dataset_knn$Month <- factor(dataset_knn$Month,
                            order = TRUE,
                            levels = c('Jan', 'Feb', 'Mar', 'Apr', 'May',
                                         'Jun', 'Jul', 'Aug', 'Sep', 'Oct', 'Nov', 'Dec'),
                            labels = c(1,2,3,4,5,6,7,8,9,10,11,12))
dataset_knn$VisitorType <- factor(dataset_knn$VisitorType,
                                  order = TRUE,
                                  levels=c('Returning_Visitor', 'Other', 'New_Visitor'),
                                  labels = c(1,2,3))
dataset_knn$Weekend <- mapvalues(dataset$Weekend,
                                 from = c('TRUE', 'FALSE'),
                                 to = c(1,0))
# dataset_knn$Revenue <- as.numeric(as.factor(dataset_knn$Revenue))
dataset_knn$Administrative <- as.numeric(dataset_knn$Administrative)
dataset_knn$ProductRelated <- as.numeric(dataset_knn$ProductRelated)
dataset_knn$Informational <- as.numeric(dataset_knn$Informational)
```

Splitting the data for the testing and training data in the ratio 75:25 for KNN. Moreover, feature scaling the features except the *Revenue* as we will need the factor for the prediction of the model on the dataset.

```
set.seed(123)
split_knn = createDataPartition(dataset_knn$Revenue, p=0.75, list = FALSE, times = 1)
training_knn = dataset_knn[split_knn,]
test_knn = dataset_knn[-split_knn,]
preprocess_knn <- preProcess(training_knn[-18], method = c("center", "scale"))
```

```
train_pre_knn <- predict(preprocess_knn, training_knn)
test_pre_knn <- predict(preprocess_knn, test_knn)
```

Running the KNN algorithm from the “class” package and bulding the KNN classifier as the y_pred.

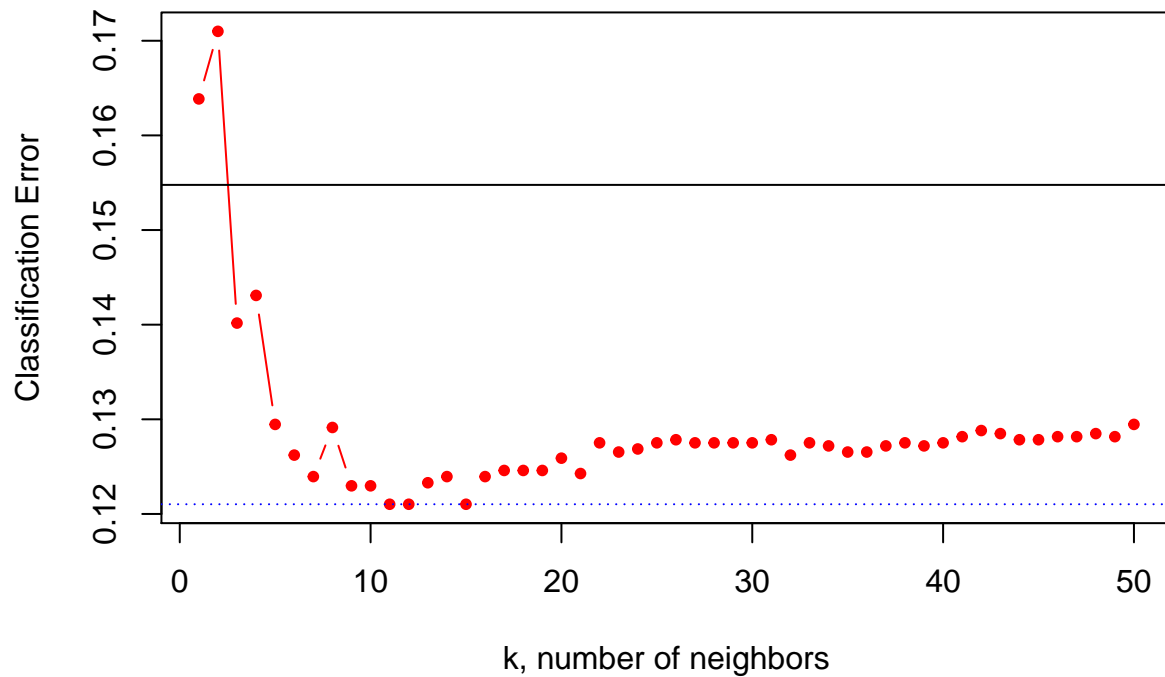
```
calc_class_err = function(actual, predicted) {
  mean(actual != predicted)
}
```

```
library(class)
set.seed(24)
k_to_try = 1:50
err_k = rep(x = 0, times = length(k_to_try))
for (i in seq_along(k_to_try)) {
  k.mod = knn(train = train_pre_knn[, -18],
              test  = test_pre_knn[, -18],
              cl     = train_pre_knn[, 18],
              k      = k_to_try[i])
  err_k[i] = calc_class_err(test_pre_knn[, 18], k.mod)
}
```

Plotting the error classification with the K.

```
plot(err_k, type = "b", col = "red", cex = 1, pch = 20,
     xlab = "k, number of neighbors", ylab = "Classification Error",
     main = "(Test) Error Rate vs Neighbors")
abline(h = min(err_k), col = "blue", lty = 3)
abline(h = mean(test_pre_knn[, 18] == "TRUE", col = "grey", lty = 3))
```

(Test) Error Rate vs Neighbors



Looking up the best K value possible from the range 1 to 50. Here the best K values with less error rate are 11, 12 and 15. We are gonna select the maximum K, as it has the least possible chance of overfitting.

```
which(err_k == min(err_k))
```

```
## [1] 11 12 15
```

```
max(which(err_k == min(err_k)))
```

```
## [1] 15
```

```
library(class)
y_pred = knn(train = train_pre_knn[, -18],
              test = test_pre_knn[, -18],
              cl = train_pre_knn[, 18],
              k = 15,
              prob = TRUE)
```

Building the Confusion Matrix for the K-Nearest Neighbor.

```
cm = table(test_knn[, 18], y_pred)
cm
```

```
##          y_pred
```

```
##          FALSE TRUE
##  FALSE  2565   40
##   TRUE   333  144
```

Accuracy measures for KNN.

```
confusionMatrix(cm)
```

```
## Confusion Matrix and Statistics
##
##          y_pred
##          FALSE TRUE
##  FALSE  2565   40
##   TRUE   333  144
##
##              Accuracy : 0.879
##              95% CI : (0.8669, 0.8903)
##   No Information Rate : 0.9403
##   P-Value [Acc > NIR] : 1
##
##              Kappa : 0.3825
##
##  Mcnemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.8851
##              Specificity : 0.7826
##              Pos Pred Value : 0.9846
##              Neg Pred Value : 0.3019
##              Prevalence : 0.9403
##              Detection Rate : 0.8323
##   Detection Prevalence : 0.8452
##   Balanced Accuracy : 0.8339
##
##              'Positive' Class : FALSE
##
```

K-NN results in giving an accuracy of 88%. In this experiment we took $K = 5$, which is by far the perfect fit for the prediction.

It shows that, Decision Tree algorithm performs best among these four mostly used classification, clustering, and regression algorithms in terms of accuracy. It gives us 89.65% accuracy for the used dataset. Even for the specificity, K-means performs best, i.e, it is more accurate than other algorithms to predict the customers who give revenue. For sensitivity, we see that Decision tree algorithm works better than others. It is because, our dataset has only 15.5% positive test cases and the relationship among the attributes for these test cases can be derived according to decision tree more accurately. But, if we look at specificity, we will see Logistic Regression works better to find the customers who don't generate revenue. So, we can confidently say, Decision Tree is the best algorithm to predict the purchase intention of customers from empirical data.

Conclusion

We have implemented Logistic Regression, Decision Tree, K-Means Algorithm and K-Nearest Neighbors. From these, we got nearly of 89% accuracy in the Logistic regression model. For Decision Tree, we got 90% accuracy, making it slightly better than the Regression model. For K-Means Algorithm, the accuracy is

nearly 45%, it was discovered that clustering is not suitable for this dataset. Finally for K-Nearest Neighbors, the results showed an accuracy of 88%, therefore making Decision tree the best model for this dataset with an accuracy of 90%.

In the future to be able to better use Clustering models, we need more variables and observations. Apart from this, having more observations would also have enabled us to better train our model. Moreover, for the given dataset it can be seen that classification techniques results better clustering techniques.

To conclude, we would like to recommend more additional variables and collection of more observations so that we would be better able to analyze and predict shoppers' intentions.

Appendix—Code

```
knitr::opts_chunk$set(echo= TRUE, warning=FALSE, message=FALSE)
library(ggplot2)
library(tidyverse)
library(gmodels)
library(dplyr)
library(ggmosaic)
library(corrplot)
library(caret)
library(rpart)
library(rpart.plot)
library(cluster)
library(fpc)
library(data.table)
library(knitr)
library(kableExtra)
library(plyr)
library(caTools)
dataset <- read.csv("online_shoppers_intention.csv", header = TRUE)
attach(dataset)
ncol(dataset)
nrow(dataset)
str(dataset)
summary(dataset)
tab1 <- read.csv("table1.csv", header = TRUE)
kable(tab1) %>%
  kable_styling(full_width = T)
tab2 <- read.csv("table2.csv", header = TRUE)
kable(tab2) %>%
  kable_styling(full_width = T)
library(gmodels)
summary(dataset$Revenue)
CrossTable(dataset$Revenue)
datasetbinary <- dataset %>%
  mutate(Revenue_binary = ifelse(dataset$Revenue == "TRUE", 1, 0))
# colSums(is.na(dataset))
sum(is.na(dataset))
sapply(dataset, function(x) sum(is.na(x)))
dataset <- na.omit(dataset)
dataset$Month = factor(dataset$Month, levels = month.abb)
```



```

dataset %>%
  ggplot() +
    aes(x = Month, Revenue = ..count../nrow(dataset), fill = Revenue) +
    geom_bar() +
    ylab("Frequency")
table_month = table(dataset$Month, dataset$Revenue)
tab_mon = as.data.frame(prop.table(table_month,2))
colnames(tab_mon) = c("Month", "Revenue", "perc")
ggplot(data = tab_mon, aes(x = Month, y = perc, fill = Revenue)) +
  geom_bar(stat = 'identity', pdatasettion = 'dodge', alpha = 2/3) +
  xlab("Month")+
  ylab("Percent")
theme_set(theme_bw())

## setting default parameters for mosaic plots
mosaic_theme = theme(axis.text.x = element_text(angle = 90,
                                                  hjust = 1,
                                                  vjust = 0.5),
                     axis.text.y = element_blank(),
                     axis.ticks.y = element_blank())

dataset %>%
  ggplot() +
  geom_mosaic(aes(x = product(Revenue, VisitorType), fill = Revenue)) +
  mosaic_theme +
  xlab("Visitor Type") +
  ylab(NULL)
CrossTable(dataset$Weekend, dataset$Revenue)
dataset %>%
  ggplot() +
  mosaic_theme +
  geom_mosaic(aes(x = product(Revenue,Weekend), fill = Revenue)) +
  xlab("Weekend") +
  ylab(NULL)
dataset$Revenue <- as.factor(dataset$Revenue)
dataset$VisitorType <- as.factor(dataset$VisitorType)
dataset$TrafficType <- as.factor(dataset$TrafficType)
dataset$OperatingSystems <- as.factor(dataset$OperatingSystems)
dataset$Browser <- as.factor(dataset$Browser)
dataset$Region <- as.factor(dataset$Region)
dataset$Month <- as.factor(dataset$Month)
summary(dataset$Browser)
summary(dataset$TrafficType)
dataset$Browser = recode(dataset$Browser, '3' = '1')
dataset$Browser = recode(dataset$Browser, '4' = '1')
dataset$Browser = recode(dataset$Browser, '5' = '1')
dataset$Browser = recode(dataset$Browser, '6' = '1')
dataset$Browser = recode(dataset$Browser, '7' = '1')
dataset$Browser = recode(dataset$Browser, '8' = '1')
dataset$Browser = recode(dataset$Browser, '9' = '1')
dataset$Browser = recode(dataset$Browser, '10' = '1')
dataset$Browser = recode(dataset$Browser, '11' = '1')
dataset$Browser = recode(dataset$Browser, '12' = '1')
dataset$Browser = recode(dataset$Browser, '13' = '1')

```

```

dataset$TrafficType = recode(dataset$TrafficType, '6' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '7' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '8' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '9' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '10' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '11' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '12' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '13' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '14' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '15' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '16' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '17' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '18' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '19' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '20' = '5')
dataset$TrafficType = recode(dataset$TrafficType, '4' = '3')
dataset$TrafficType = recode(dataset$TrafficType, '5' = '4')
summary(dataset$Browser)
summary(dataset$TrafficType)
dataset_logistic <- dataset
set.seed(123)
split_log = createDataPartition(dataset_logistic$Revenue, p=0.75, list = FALSE, times = 1)
training_log = dataset_logistic[split_log,]
test_log = dataset_logistic[-split_log,]
preprocess_log <- preprocess(training_log, method = c("center", "scale"))
train_pre_log <- predict(preprocess_log, training_log)
test_pre_log <- predict(preprocess_log, test_log)
set.seed(123)
log_reg <- glm(Revenue ~ ., data = train_pre_log, family = binomial)
summary(log_reg)
log_reg_pred <- predict(log_reg, newdata = test_pre_log, type = "response")
log_reg_pred <- ifelse(log_reg_pred > 0.5, "TRUE", "FALSE")
logcf <- table(log_reg_pred, test_pre_log$Revenue)
confusionMatrix(logcf)
dataset_classify = dataset
set.seed(123)
split = sample.split(dataset_classify$Revenue, SplitRatio = 0.75)
training_data = subset(dataset_classify, split == TRUE)
test_data = subset(dataset_classify, split == FALSE)
dt_model <- rpart(Revenue ~ ., data = training_data, method = "class")
rpart.plot(dt_model)
dt.pred <- predict(dt_model, test_data, type = "class")
mean(dt.pred == test_data$Revenue)
cfdt <- table(dt.pred, test_data$Revenue)
confusionMatrix(cfdt)
dataset_kmeans <- dataset
dataset_kmeans$OperatingSystems <- factor(dataset$OperatingSystems,
                                           order = TRUE,
                                           levels = c(6,3,7,1,5,2,4,8))
dataset_kmeans$Browser <- factor(dataset$Browser,
                                  order = TRUE,
                                  levels = c(9,3,6,7,1,2,8,11,4,5,10,13,12))
dataset_kmeans$Region <- factor(dataset$Region,

```

```

        order = TRUE,
        levels = c(8,6,3,4,7,1,5,2,9))
dataset_kmeans$TrafficType <- factor(dataset$TrafficType,
        order = TRUE,
        levels = c(12,15,17,18,13,19,3,9,1,6,4,14,11,10,5,2,20,8,7,16))
dataset_kmeans$Month <- factor(dataset_kmeans$Month,
        order = TRUE,
        levels = c('Jan','Feb','Mar','Apr','May',
        'Jun','Jul','Aug','Sep','Oct','Nov','Dec'),
        labels = c(1,2,3,4,5,6,7,8,9,10,11,12))
dataset_kmeans$VisitorType <- factor(dataset_kmeans$VisitorType,
        order = TRUE,
        levels=c('Returning_Visitor','Other','New_Visitor'),
        labels = c(1,2,3))
dataset_kmeans$Weekend <- mapvalues(dataset$Weekend,
        from = c('TRUE','FALSE'),
        to = c(1,0))
# dataset_kmeans$Revenue <- as.numeric(as.factor(dataset_kmeans$Revenue))
dataset_kmeans$Administrative <- as.numeric(dataset_kmeans$Administrative)
dataset_kmeans$ProductRelated <- as.numeric(dataset_kmeans$ProductRelated)
dataset_kmeans$Informational <- as.numeric(dataset_kmeans$Informational)
normalize <- function(x) {
  return((x-min(x))/(max(x) - min(x)))
}
dataset_kmeans$Administrative <- normalize(dataset_kmeans$Administrative)
dataset_kmeans$Administrative_Duration <- normalize(dataset_kmeans$Administrative_Duration)
dataset_kmeans$Informational <- normalize(dataset_kmeans$Informational_Duration)
dataset_kmeans$Informational_Duration <- normalize(dataset_kmeans$Administrative)
dataset_kmeans$ProductRelated <- normalize(dataset_kmeans$ProductRelated)
dataset_kmeans$ProductRelated_Duration <- normalize(dataset_kmeans$ProductRelated_Duration)
dataset_kmeans$BounceRates <- normalize(dataset_kmeans$BounceRates)
dataset_kmeans$ExitRates <- normalize(dataset_kmeans$ExitRates)
dataset_kmeans$PageValues <- normalize(dataset_kmeans$PageValues)
dataset_kmeans$SpecialDay <- normalize(dataset_kmeans$SpecialDay)
kmeans_clust <- kmeans(dataset_kmeans[-18], centers = 2, iter.max = 150)
kmeans_clust$size
cm_km <- table(kmeans_clust$cluster, dataset_kmeans$Revenue)
cm_km
kmeans_accuracy = sum(diag(cm_km))/sum(cm_km)
kmeans_precision<- cm_km[1,1]/(sum(cm_km[1,]))
kmeans_recall<- cm_km[1,1]/(sum(cm_km[,1]))
kmeans_fscore <- 2*kmeans_precision*kmeans_recall/(kmeans_precision+kmeans_recall)
kmeans_accuracy
kmeans_precision
kmeans_recall
kmeans_fscore
library(cluster)
clusplot(dataset_kmeans,
        kmeans_clust$cluster,
        lines = 0,
        shade = TRUE,
        color = TRUE,
        labels = 2,

```

```

    plotchar = FALSE,
    span = TRUE,
    main = paste('Cluster of Revenues'))
dataset_knn <- dataset
dataset_knn$OperatingSystems <- factor(dataset$OperatingSystems,
                                       order = TRUE,
                                       levels = c(6,3,7,1,5,2,4,8))
dataset_knn$Browser <- factor(dataset$Browser,
                              order = TRUE,
                              levels = c(9,3,6,7,1,2,8,11,4,5,10,13,12))
dataset_knn$Region <- factor(dataset$Region,
                             order = TRUE,
                             levels = c(8,6,3,4,7,1,5,2,9))
dataset_knn$TrafficType <- factor(dataset$TrafficType,
                                  order = TRUE,
                                  levels = c(12,15,17,18,13,19,3,9,1,6,4,14,11,10,5,2,20,8,7,16))
dataset_knn$Month <- factor(dataset_knn$Month,
                            order = TRUE,
                            levels = c('Jan', 'Feb', 'Mar', 'Apr', 'May',
                                         'Jun', 'Jul', 'Aug', 'Sep', 'Oct', 'Nov', 'Dec'),
                            labels = c(1,2,3,4,5,6,7,8,9,10,11,12))
dataset_knn$VisitorType <- factor(dataset_knn$VisitorType,
                                  order = TRUE,
                                  levels=c('Returning_Visitor', 'Other', 'New_Visitor'),
                                  labels = c(1,2,3))
dataset_knn$Weekend <- mapvalues(dataset$Weekend,
                                 from = c('TRUE', 'FALSE'),
                                 to = c(1,0))
# dataset_knn$Revenue <- as.numeric(as.factor(dataset_knn$Revenue))
dataset_knn$Administrative <- as.numeric(dataset_knn$Administrative)
dataset_knn$ProductRelated <- as.numeric(dataset_knn$ProductRelated)
dataset_knn$Informational <- as.numeric(dataset_knn$Informational)
set.seed(123)
split_knn = createDataPartition(dataset_knn$Revenue, p=0.75, list = FALSE, times = 1)
training_knn = dataset_knn[split_knn,]
test_knn = dataset_knn[-split_knn,]
preprocess_knn <- preProcess(training_knn[-18], method = c("center", "scale"))
train_pre_knn <- predict(preprocess_knn, training_knn)
test_pre_knn <- predict(preprocess_knn, test_knn)
calc_class_err = function(actual, predicted) {
  mean(actual != predicted)
}
library(class)
set.seed(24)
k_to_try = 1:50
err_k = rep(x = 0, times = length(k_to_try))
for (i in seq_along(k_to_try)) {
  k.mod = knn(train = train_pre_knn[-18],
              test = test_pre_knn[-18],
              cl = train_pre_knn[, 18],
              k = k_to_try[i])
  err_k[i] = calc_class_err(test_pre_knn[,18], k.mod)
}

```

```

plot(err_k, type = "b", col = "red", cex = 1, pch = 20,
     xlab = "k, number of neighbors", ylab = "Classification Error",
     main = "(Test) Error Rate vs Neighbors")
abline(h = min(err_k), col = "blue", lty = 3)
abline(h = mean(test_pre_knn[,18] == "TRUE", col = "grey", lty = 3))
which(err_k == min(err_k))
max(which(err_k == min(err_k)))
library(class)
y_pred = knn(train = train_pre_knn[, -18],
             test = test_pre_knn[, -18],
             cl = train_pre_knn[, 18],
             k = 15,
             prob = TRUE)
cm = table(test_knn[, 18], y_pred)
cm
confusionMatrix(cm)

```