

User Engagement Analysis

Hans Peter Ndeffo

Showwcase

September 14, 2020

Abstract

In this project, we aim at **investigating criterias** and **defining metrics to measure user engagement**. In general, there are several components involve when assessing user involvement in a platform but according to the dataset that this project is based on they are: **session_id, customer_id, login_date, projects_added, likes_given, comment_given, inactive_status, bug_occured, session_projects_added, session_likes_given, session_comments_given, inactive_duration, bugs_in_session, session_duration**. However, not all of them are important as we will see in this analysis. The following process will be used in order to filter non-relevant variables when trying to identify highly or lowly active users using the following steps:

- Import showwcase sessions' dataset.
- Display an overview of the dataset to find missing values and check spelling mistakes. Also, because the dataset has about 300 rows, there is no need to use Apache Spark which is a big data computational tool.
- Replace TRUE and FALSE with 1s and 0s. This will help when counting and achieving other computations.

- Simplify the model. We will assume that session_id, login_date, session_likes_given are irrelevant for this analysis. Also, we will create a new variable called active_session_duration which is the difference between the total session duration and inactive_duration. Therefore, delete session duration and inactive_duration will be deleted.
- Determine thresholds for session duration and login frequency. These decision boundaries will help decide whether a user is engaged or not.
- Setup the random seed for sampling. The dataset will be split into 80% for training and 20% for testing.
- Use Random Forest, Logistic Regression, KNN to predict user engagement.
- Identify trends by determining important features
- Determine model effectiveness and predicting power by computing training error and testing error.
- Compare the different models

Introduction

In this age where technology dominates everything, understanding customers involvement in any given platform has become a real challenge. **What are the most frequent features by which we can help label highly and least active users?** Solving this problem can have significant economic benefits through better-targeted ads campaigns and positive feedbacks from users. Also, this might demonstrate that there are certain patterns through which we can determine which customers are more likely to be engaged. In this situation, we will first use association rules to evaluate what are the most frequent patterns. While looking for the right algorithm to use, there are certain procedures to follow to remain consistent and precise in any analysis. The outcome of those procedures will determine which machine learning algorithm suits best the problem.

Pre-requisites

Step 1: Read customers dataset' file and load the required libraries

```
# Load Libraries  
# dplyr is use to break complex formula into simpler piece of code  
# more understandable  
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 3.6.3
```

```
##  
## 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
```

```
# Spark to Leverage Large datasets  
library(sparklyr)
```

```
## Warning: package 'sparklyr' was built under R version 3.6.3
```

```
# ggplot to plot graphs  
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 3.6.3
```

```
library(cowplot)
```

```
## Warning: package 'cowplot' was built under R version 3.6.3
```

```
##  
## *****
```

```
## Note: As of version 1.0.0, cowplot does not change the
```

```
## default ggplot2 theme anymore. To recover the previous
```

```
## behavior, execute:  
## theme_set(theme_cowplot())
```

```
## *****
```

```
# Load random forest algorithm and other functions  
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 3.6.3
```

```
## randomForest 4.6-14
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##  
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:ggplot2':  
##  
## margin
```

```
## The following object is masked from 'package:dplyr':  
##  
##      combine
```

```
# Caret is used to predict the customers in the testing set  
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.6.3
```

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

```
# DataExplorer creates a preview of the dataset  
library(DataExplorer)
```

```
## Warning: package 'DataExplorer' was built under R version 3.6.3
```

```
# knitr is to display tables  
library(knitr)
```

```
## Warning: package 'knitr' was built under R version 3.6.3
```

```
library(flextable)
```

```
## Warning: package 'flextable' was built under R version 3.6.3
```

```
# Import functions that will help determine association rules  
library(arules)
```

```
## Warning: package 'arules' was built under R version 3.6.3
```

```
## Loading required package: Matrix
```

```
##  
## Attaching package: 'arules'
```

```
## The following object is masked from 'package:dplyr':  
##  
##      recode
```

```
## The following objects are masked from 'package:base':  
##  
##      abbreviate, write
```

```
# To compare strings  
library("pracma")
```

```
## Warning: package 'pracma' was built under R version 3.6.3
```

```
##  
## Attaching package: 'pracma'
```

```
## The following object is masked from 'package:arules':  
##  
##      size
```

```
## The following objects are masked from 'package:Matrix':  
##  
##      expm, lu, tril, triu
```

```
library("stringr")
```

```
## Warning: package 'stringr' was built under R version 3.6.3
```

```
# To import knn function  
library(class)  
  
# Read file
```

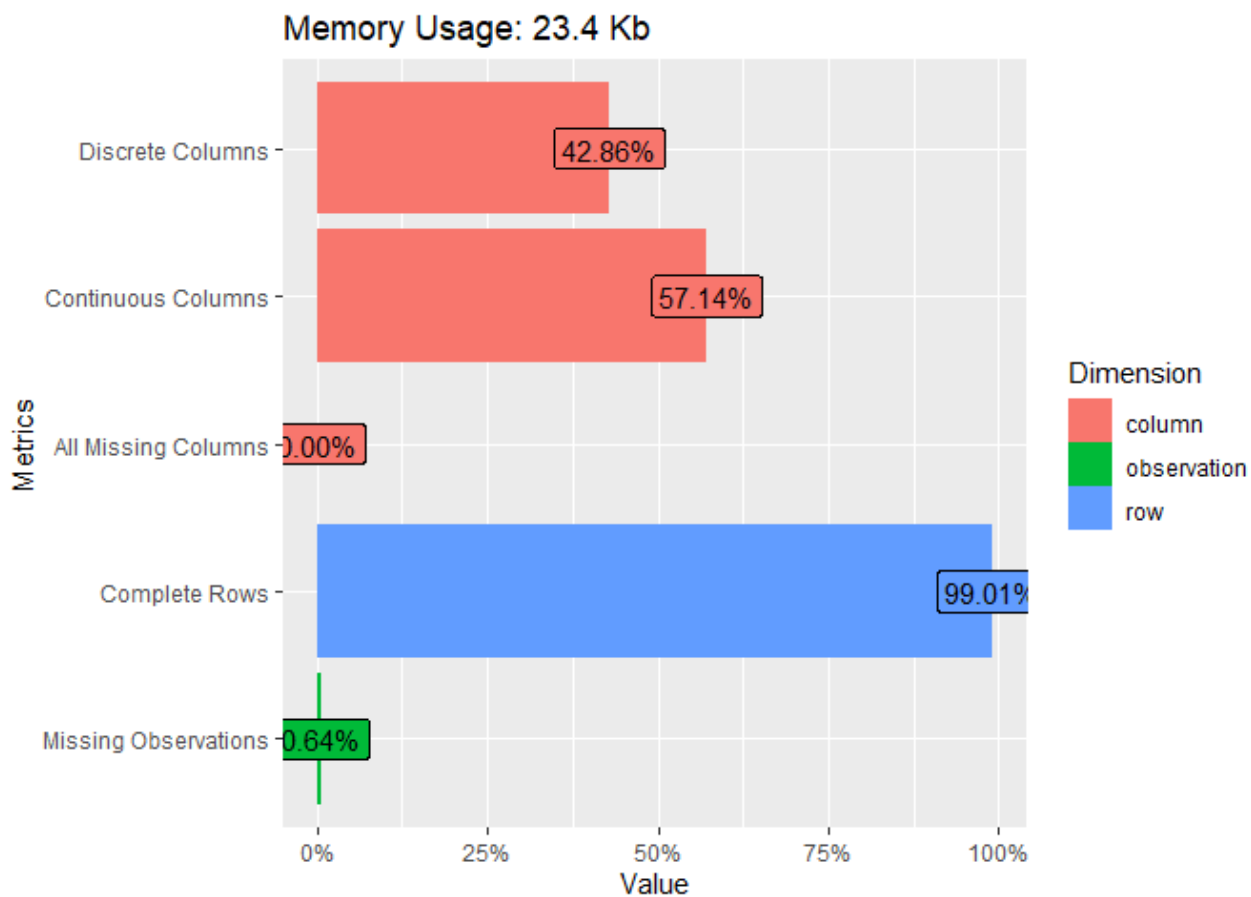
```
users_data <-read.csv("C:/Users/ndeff/OneDrive/Desktop/Research and online classes/showwcase_sessions.csv",  
                      sep=";",stringsAsFactors = FALSE, header=T,na.strings=c(  
"#NUM!","?","NA"," "))
```

Step 2: Show an overview and a sample of the dataset

Having an overview of the dataset structure is extremely important because this could show the data distribution and missing values. For example, to evaluate whether the dataset is balanced or imbalanced. A dataset is balanced if it follows the following criterias: If there is n distinct preferably discrete output, Then the data should be evenly distributed among each output. In that case, the probability p of any output i for example to occur is approximatively $p = \frac{1}{n}$. However, in this case, we do not know how many users can be labeled as engaged and not. This scenario implies only 2 outcomes because we assume in general that users can be engaged or not. This assumption will be carried on throughout this analysis to simplify the model.

The plot below count the missing values and display how many discrete and continuous columns.

```
plot_intro(users_data)
```



As we can see there some missing values about 0.64% in data. Show all missing data

```
users_data[rowSums(is.na(users_data)) > 0,]
```

```
##      i..session_id customer_id login_date projects_added likes_given
## 36      862128      23404   10/26/19          TRUE        FALSE
## 301         NA         NA                NA            NA
## 302         NA         NA                NA            NA
##      comment_given inactive_status bug_occured session_projects_added
## 36          TRUE          TRUE        FALSE                2
## 301         NA          NA            NA                NA
## 302         NA          NA            NA                NA
##      session_likes_given session_comments_given inactive_duration
## 36              NA              2              1120
## 301              NA              NA              NA
## 302              NA              NA              NA
##      bugs_in_session session_duration
## 36              0              95
## 301             NA              NA
## 302             NA              NA
```


Since missing data at row 301 and 302 have all columns empty and are the end of the file, we can just delete those rows from the dataset. However, at row 36 we only see one missing value located at **session_likes_given**. That value can be replaced by 0 because any given user has likes or not.

```
# Remove row 301 and 302
users_data <- users_data[-c(301,302),]

# Replace NA with 0 at row 36 column at the session_likes_given
users_data$session_likes_given[36] <- 0
```

Check the spelling mistake in the column names

```
names(users_data)
```

```
## [1] "i..session_id"      "customer_id"        "login_date"
## [4] "projects_added"     "likes_given"        "comment_given"
## [7] "inactive_status"    "bug_occured"        "session_projects_added"
## [10] "session_likes_given" "session_comments_given" "inactive_duration"
## [13] "bugs_in_session"    "session_duration"
```

Fix session id name

```
names(users_data)[names(users_data)== names(users_data)[1]] <- "session_id"
```

After the data are cleaned we need to change the TRUE and FALSE to 1 and 0. This will help during computations because characters values such as "TRUE" can not be computed.

```
# number of rows of car dataset
n=nrow(users_data)
p=ncol(users_data)
```

```
# Go through the dataset and replace TRUE with 1 and FALSE with 0
for(i in 1:n){
  for(j in 1:p){
    if(strcmp(toString(users_data[i,j]),"TRUE"))
    {
      users_data[i,j]<- 1
    }
    if(strcmp(toString(users_data[i,j]),"FALSE"))
    {
      users_data[i,j]<- 0
    }
  }
}
```

Check if TRUE and FALSE have been replaced by 1s and 0s

```
table = flextable(users_data[1:5,],col_keys = names(users_data),
                  cwidth = 0.75, cheight = 1, theme_fun = theme_tron_legacy)
# Add color to text and the body of the table
table <- color(table,color = "black", part = "body")
table <- bg(table,bg = "white")
table
```

session_ id	customer_ id	login_date	projects_ added	likes_giv en	comment _given	inactive_ status	bug_occ ured	session_ projects_ added	session_ likes_giv en	session_ comment s_given	inactive_ duration	bugs_in_ session	session_ duration
624205	80746	10/30/19	0	1	1	1	0	0	24	3	1146	0	1564
624241	24520	10/30/19	1	1	1	1	0	2	3	5	133	0	1766
111002	32047	10/30/19	1	1	1	1	0	1	5	5	1571	0	2230
545113	23404	10/30/19	1	1	1	0	0	1	10	21	0	0	633
750269	40235	10/30/19	1	1	0	1	0	3	16	0	1405	0	1679

Step 3: Simplify the model

If the inactive duration is greater than session duration for some cases that implies that inactive duration included a measure of log out + login time. Which might be caused by

bug or connexion issues. On the other hand session duration only measure the login time. So the difference measure the active time.

```
active_session_duration <- users_data$session_duration - users_data$inactive_duration
# Replace all negative active session duration by 0
active_session_duration[active_session_duration< 0]<- 0
```

Add active session duration to the dataset

```
users_data <- cbind(users_data,active_session_duration)
```

Remove session duration and inactive duration

```
users_data <- users_data[,!(names(users_data) %in% c("session_duration","inactive_duration","login_date","session_id"))]
```

Show a sample after binding and deletion of certain columns

```
table = flextable(users_data[1:3,],col_keys = names(users_data),
  cwidth = 0.75, cheight = 1, theme_fun = theme_tron_legacy)
# Add color to text and the body of the table
table <- color(table,color = "black", part = "body")
table <- bg(table,bg = "white")
table
```

customer_id	projects_added	likes_gIVEN	comment_given	inactive_status	bug_occured	session_projects_added	session_likes_gIVEN	session_comments_gIVEN	bugs_in_session	active_session_duration
80746	0	1	1	1	0	0	24	3	0	418
24520	1	1	1	1	0	2	3	5	0	1633
32047	1	1	1	1	0	1	5	5	0	659

Step 4: Determining a threshold for active users

We need now decide a threshold for session duration that will define is the user engaged or not. We will use the **average duration** as a threshold for active and not engaged users. Lets declare a variable called `user_engagement`. That variable reflects whether or not an user is engaged if his active session duration is above the threshold. **This approach assumes** for a given session that the **active session duration is the most determinant factor to evaluate an user engagement**.

```
# This variable has TRUE or FALSE values which will be represented by 1s and 0s
user_engagement <- c()
```

```
threshold <- mean(users_data$active_session_duration)

# Go through all rows of the dataset
for(i in 1:dim(users_data)[1]){

  if(users_data$active_session_duration[i] > threshold){
    user_engagement <- c(user_engagement, 1)
  }
  else{
    user_engagement <- c(user_engagement, 0)
  }
}
user_engagement <- as.factor(user_engagement)
```

Add user engagement values to the dataset

```
users_data <- cbind(users_data, user_engagement)
```

Step 5: Determining other relevant factors to predict user engagement

Here we will **ignore the active session duration as a variable** and **try to find other important factors**. In this dataset, since there is no obvious relationship between data, we will use several machine learning algorithm and plot their performance. The dependent variable to predict will be user_engagement and all the other variables will be considered as independent.

Split the dataset into 80% training and 20% testing

```
set.seed(1)
train_sample = sample(1:dim(users_data)[1], floor(dim(users_data)[1]*0.80))
train_data = users_data[train_sample, ]
test_data = users_data[-train_sample, ]
```

Step 6: Predict test and train data to test model effectiveness

Random Forest

```
rf <- randomForest(user_engagement ~ customer_id + projects_added+ likes_given+
  session_projects_added+ inactive_status+ comment_given+ bug_occured+ se
  ssion_likes_given
  + session_comments_given, data = train_data,
  type = classification, importance = TRUE, proximity = TRUE)
```

```
test_predict_rf <- predict(rf, test_data)
# Compute the error
rdmFor_testing_error_rf <- sum(test_predict_rf!= test_data$user_engagement)/nrow
(test_data)
```

The model effectiveness is 65% accuracy

Get important features from the random forest model

```
importance(rf,2)
```

```
##                MeanDecreaseGini
## customer_id      23.518306
## projects_added   2.743130
## likes_given      3.231945
## session_projects_added 9.410683
## inactive_status  25.107729
## comment_given    2.987632
## bug_occured      5.084175
## session_likes_given 19.031958
## session_comments_given 12.630227
```

The important features are:

- **customer_id** counts for about **23.51%**
- **inactive_status** use by the customer counts for **25.10%**
- **session_likes_given** influence the model for about **19.03%**
- **session_comment_given** influence the model for about **12.63%**
- **session_project_added** influence the model for about **9.41%**
- **bug_occured** influence the model for about **5.08%**

Any feature that influences the model by less than 5% will be rejected because p value = 5%.

Determining important features using a logistic regression model

The assumption using this model is the user engagement is linearly depend on all the other independent variables such customer_id, inactive_status. This model might suit this problem because we face only 2 outcomes(an user can be engaged or not). However, because we assume we do not know exactly which feature is important that we will have to add more variables(higher dimensionality) to the model. This might reduce the model accuracy. The logistic regression will be define as follows

```
lr <- glm(user_engagement ~ customer_id + projects_added+ likes_given+ session_projects_added+ inactive_status+ comment_given+ bug_occured+ session_likes_given
          + session_comments_given , family = binomial(link="logit"), data = users_data) # fit lm() model
```

**** Fixing a threshold for logistic regression.****

We expect the logistic regression model to perform potentially poorly because of there are too many variables involves and also due the fact we have to choose manually a random threshold. In this case, we will use the mean.

```
# Predict user engagement using Logistic regression
test_predict_lr <- predict(lr, test_data)

# Compute threshold using mean function
threshold_lr <- mean(test_predict_lr)

# Go through all rows of the dataset
for(i in 1:length(test_predict_lr)){
  if(test_predict_lr[i] > threshold_lr){
    test_predict_lr[i] <- 1
  }
  else{
    test_predict_lr[i] <- 0
  }
}
```

Computing testing error for logistic regression

```
testing_error_lr <- sum(test_predict_lr!= test_data$user_engagement)/nrow(test_data)
```

The **testing error for the logistic regression** is **30%**

Using the knn neighbors

Knn can be tested because this analysis uses supervised learning. We will compute knn with different parameters k and select the one that gives the highest accuracy. The maximum value of k has to be less than the number of observations or rows in this dataset which is about 300.

The code below defines different values of k to be tested

```
k <- seq(from = 1, to = 240, by = 5)
```

This function computes the accuracy of each model

```
accuracy = function(actual, predicted) {  
  mean(actual == predicted)  
}
```

use a matrix store the model by saving its accuracy and k

```
models_accuracy <- c()  
#models_k
```

For all k, compute the predicted value

```
# Values to be predicted  
actual <- test_data$user_engagement  
  
for(i in k){  
  predicted = knn(train = train_data[, -c(12)],  
                  test = test_data[, -c(12)],  
                  cl = train_data$user_engagement, k = i)  
  
  # Store the accuracy of each model  
  models_accuracy <- c(models_accuracy, accuracy(actual, predicted))  
}
```

Now we find the model with highest accuracy

```
# Get the matching index of the best model accuracy
matching_index <- match(max(models_accuracy), models_accuracy)

# Get the best accuracy and best k
knn_accuracy <- models_accuracy[matching_index]
Knn_k <- k[matching_index]
```

With **Knn** we get an accuracy of **96.66%** with **k = 1**. Therefore, the error is:

```
knn_error <- 1 - knn_accuracy
```

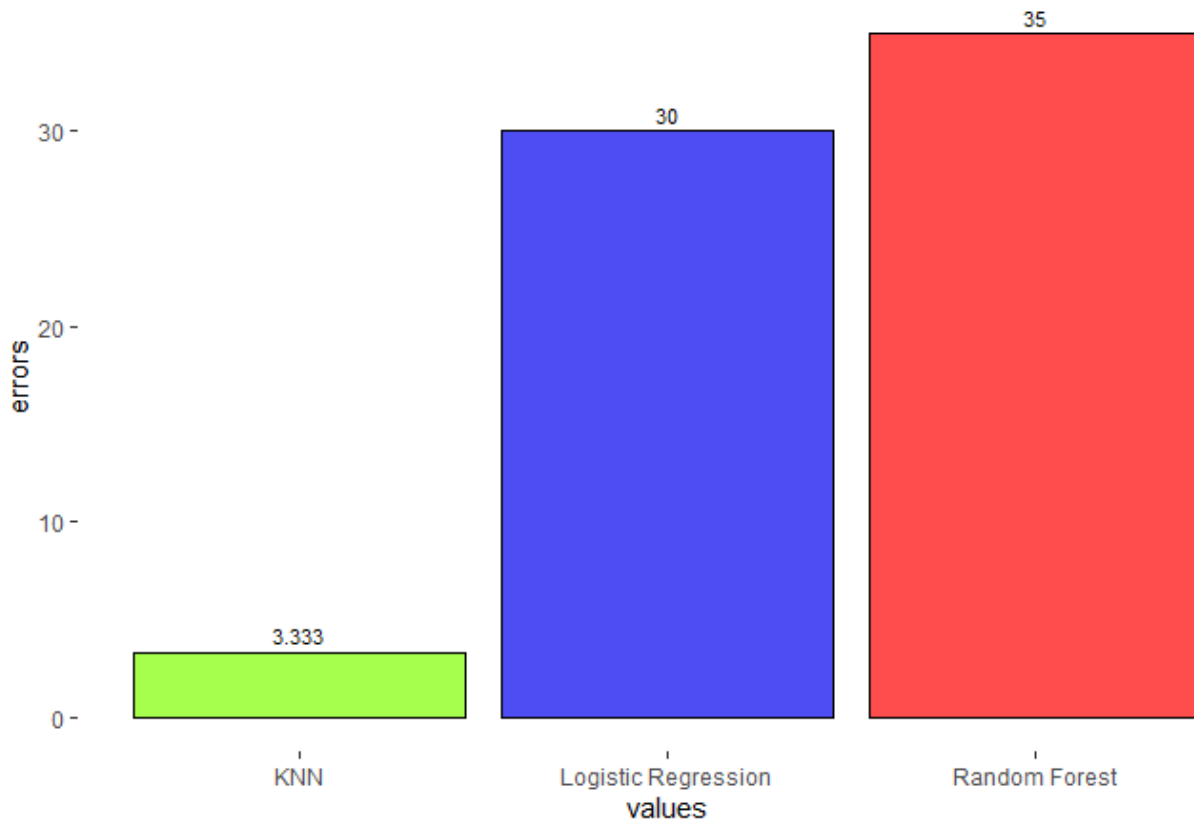
Step 7: Compare the different models

```
values <- c("KNN", "Random Forest", "Logistic Regression")
errors <- (c(knn_error, rdmFor_testing_error_rf, testing_error_lr))*100
```

```
ggplot()+ aes(x= values, y= errors)+
  geom_bar(colour = "black", fill= c("chartreuse1", "blue2", "red"),
    alpha= 0.7, stat = "summary", show.legend = TRUE)+
  theme(plot.title = element_text(size = 20, hjust = 0.5), panel.background = element_rect("white"))+
  geom_text(aes(label= round(errors, digits = 3)),position = position_dodge(width = 0.9),vjust= -0.5, size= 3)+
  ggtitle("Models Errors Comparaison in %")
```

```
## No summary function supplied, defaulting to `mean_se()`
```

Models Errors Comparison in %



Conclusion

In this project, **KNN was the best model in this analysis with an accuracy of 96.66%**. However, due to time constraints when using KNN, we are not able to determine exactly what are the most significant features and the correlations between them. Nonetheless, **we are sure about 67.5% based on the average accuracy from the random forest and logistic models that the most relevant features are customer_id, inactive_status, session_likes_given, session_comment_given, session_project_added, bug_occured**. Those variables define “engaged customers”. For instance, with customer_id, we can track the frequency at which a customer login overtime during different sessions. This implies we can increase the predicting power of the random forest and logistic models given a set of input features such as **customer_id, session_likes_given** by 32.5% in the ideal situation by building a much more complex model. However, the tradeoffs can be huge for such an approach which might not only cause **overfitting** but raise the computational cost and training time. In addition, a such model is likely to have less predicting power on new datasets.