

# Company-X Assessment - Monica Iyer

Monica Iyer  
02/01/2021

## Aim

To understand behaviours that are most predictive of a new user starting and staying active on Product-X. The Growth team at Company-X is interested in month-1ve retention, defined as whether a user remains active after signup. We will use the data sets to understand what factors are the best predictors of retention, and offer suggestions to operationalise these insights and help Company-X.

## Analyze the data

Company-X has an interactive design platform, where designers can make prototypes for online products that can be reviewed and accessed by different account holders. We are given two datasets - users data and their corresponding events. Our aim is to predict the number of users that will be retained in the first month of activity. For now, lets go ahead and convert the data into datasets.

	USER_ID	SIGNUP_DATE	SIGNUP_CHANNEL	EMAIL_TYPE	JOB_TITLE	COUNTRY	ML_RETAINED
1	7.61805e+17	2030-11-28 0:00	direct	Business	designer	Ukraine	1
2	7.75628e+17	2030-11-11 0:00	direct	Personal	other	China	1
3	7.85972e+17	2030-12-12 0:00	search	Business	other	Brazil	1
4	7.77574e+17	2030-11-16 0:00	direct	Personal	developer	Romania	0
5	7.91957e+17	2030-12-26 0:00	invites	Business	designer	Malaysia	1
6	7.74598e+17	2030-11-08 0:00	search	Personal	designer	Russia	1
6 rows							

RECEIVED_AT	USER_ID	FILE_KEY	EVENT_NAME
1 2030-11-17 14:51	7.76509e+17	4gyt8HJHwKdVrgB0u3G	file_edited
2 2030-12-13 19:52	7.86575e+17	XALINWvYQ1mFmhzY43kQ	file_edited
3 2030-12-16 19:21	7.88480e+17	u0P8GyNAGeNDHxuh7u8	file_edited
4 2030-11-19 13:53	7.70913e+17	lp0t85fD3SeJhuRjCp6	file_edited
5 2030-12-16 3:37	7.78333e+17	35dJhOUZvQp3oXoZEIY9	file_edited
6 2030-12-05 15:25	7.79314e+17	KH9WCC8GgGgTTrvRuud9P	file_edited
6 rows			

## Clean the data

The data doesn't seem to require extensive cleaning. I haven't found null values in key columns and the only redundant value would be 'FILE\_KEY' in events\_data, which will not be considered in this analysis, or prediction modelling for that matter.

```
#lets explore the data
#Result 1
sqlDF("SELECT JOB_TITLE, COUNT(JOB_TITLE) NUMBER_OF_USERS_RETAINED
FROM users_df
WHERE ML_RETAINED=1
GROUP BY JOB_TITLE
ORDER BY 2 DESC")
```

JOB_TITLE	NUMBER_OF_USERS_RETAINED
developer	2011
designer	1822
other	825
project-manager	498
marketer	353
5 rows	

```
#Result 2
sqlDF("SELECT EVENT_NAME, COUNT(EVENT_NAME) EVENTS_COUNT
FROM events_df
GROUP BY EVENT_NAME
ORDER BY 2 DESC")
```

EVENT_NAME	EVENTS_COUNT
file_opened	759631
file_edited	154825
prototype_viewed	53783
file_created	29772
comment_created	18037
5 rows	

```
#Result 3
sqlDF("SELECT U.JOB_TITLE, E.EVENT_NAME, HIGHEST_USER_EVENT
FROM users_df U INNER JOIN
(SELECT USER_ID, EVENT_NAME, COUNT(EVENT_NAME) EVENTS_COUNT
FROM events_df
GROUP BY EVENT_NAME) E
ON U.USER_ID=E.USER_ID
ON U.USER_ID=E.USER_ID
ORDER BY E.EVENTS_COUNT")
```

JOB_TITLE	HIGHEST_USED_EVENT
developer	comment_created
project-manager	file_created
designer	file_edited
other	file_opened
4 rows	

```
#Result 4
sqlDF("SELECT U.JOB_TITLE, SUM(E.EVENTS_COUNT) NUMBER_OF_EVENTS_BY_JOB
FROM users_df U INNER JOIN
(SELECT USER_ID, COUNT(EVENT_NAME) EVENTS_COUNT
FROM events_df
GROUP BY USER_ID) E
ON U.USER_ID=E.USER_ID
GROUP BY JOB_TITLE
ORDER BY NUMBER_OF_EVENTS_BY_JOB DESC")
```

JOB_TITLE	NUMBER_OF_EVENTS_BY_JOB
designer	512608
developer	368173
other	182258
project-manager	87780
marketer	64008
5 rows	

```
temp_df <- data.frame(sqlDF("SELECT U.USER_ID, E.NUMBER_OF_EVENTS
FROM users_df U LEFT JOIN
(SELECT USER_ID, COUNT(EVENT_NAME) NUMBER_OF_EVENTS
FROM events_df
GROUP BY USER_ID) E
ON U.USER_ID=E.USER_ID"))
```

```
sqlDF("SELECT U.USER_ID, E1.FIRST_DATE - E2.SECOND_DATE DIFF
FROM users_df U
INNER JOIN (SELECT USER_ID, MIN(RECEIVED_AT) FIRST_DATE FROM events_df WHERE EVENT_NAME='file_opened'
GROUP BY USER_ID ORDER BY RECEIVED_AT) E1
ON U.USER_ID=E1.USER_ID
INNER JOIN (SELECT USER_ID, MIN(RECEIVED_AT) SECOND_DATE FROM events_df WHERE RECEIVED_AT NOT IN (SELECT MIN(RECEIVED_AT) FROM events_df WHERE EVENT_NAME='file_opened') GROUP BY USER_ID ORDER BY RECEIVED_AT) E2
ON U.USER_ID=E2.USER_ID")
```

	USER_ID	DIFF
	6.61881e+17	0
	6.61919e+17	0
	6.61922e+17	0
	6.61937e+17	0
	6.61976e+17	0
	6.61746e+17	0
	6.61851e+17	0
	6.62234e+17	0
	6.62239e+17	0
	6.62224e+17	0
1-10 of 10,000 rows	Previous	1 2 3 4 5 6 ...1000 Next

## Explore the data

- Explore the categories of the month-1 retained users. Notice in **Result #1** that the number of users retained is highest in developers, and the lowest with marketers. The number of customers retained after the first month are about 30% of those who initially signed up, which is an imbalance in the data that will be addressed later. We see that the 'other' category is the highest after developer and designer, however since there isn't much to be known about this category, we can analyze their events to understand their activity on Company-X, mainly to see if they differentiate themselves in any way from project-managers and marketers.
- Explore events. There's five types of actions and there's usability, in **Result #2**, it's key to the growth team to analyze the events most used by each category of users. **Result #3** shows that Designers edited files the most, while project managers created them, developers made comments and it seems like others mainly viewed the files (they could play the role of team members who view the design work and don't interact with it as much), which is expected intuitively and is even better to see that the data reflects that. Lastly, we access interactivity with the platform based on 'JOB\_TITLE' with **Result #4**, and the results are as expected. We see that designers use the platform the most followed by developers and other categories.
- Add features to the data. We can add the number of events per user, that indicates their extent of use of Company-X. Intuitively, the more a customer interacts with the platform, greater the chances of them coming back to it.
- Balance of the data. We notice that the proportion of retention is 30% and this represents an imbalance in the data. Just focusing on accuracy will not be the best choice in this case since it's a poor measure of imbalanced data, so it will be best to focus on both recall and accuracy.

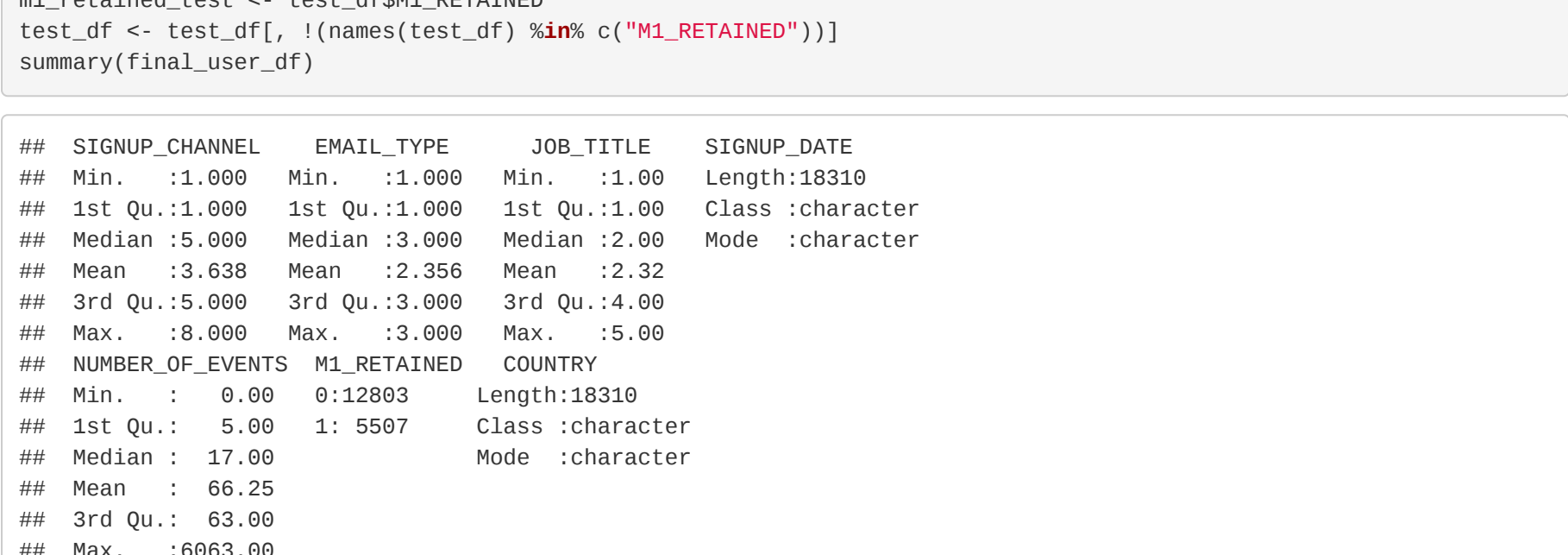
## Preprocess the data

Create a final dataframe, that can be split into train and test sets. Replace NAs (if any) with zero, change the 'SIGNUP\_DATE' to just the month since all the cases fall in 2030 and the day of the month is irrelevant, and remove 'JOB\_ID'.

```
#Create a final dataframe, that we can divide into training and test sets
#This dataframe contains the number of events corresponding to each user in users_df
final_user_df <- sqlDF("SELECT DISTINCT * FROM users_df U inner join temp_df T using (USER_ID) where U.USER_ID=T.
USER_ID")
final_user_df[is.na(final_user_df)] <- 0 #replace NA's in user activity with 0
final_user_df$SIGNUP_DATE <- format(as.Date(final_user_df$SIGNUP_DATE), "%m") #retain only month in date
#remove JOB_ID - does not add any useful information to the model
final_user_df <- final_user_df[, !names(final_user_df) %in% c("JOB_ID")] #remove JOB_ID from the data since it is irrelevant
```

## Visualization

We can Visualize the first month retention over the year 2030.



Separate the categorical and continuous variables from variables that don't require any further pre-processing. Convert categorical variables to numeric factors and join all the variable back into the 'final\_user\_df'. Once that's done, split into train and test sets and extract the test month 1 retain values ml\_retained\_test from the test set. We're ready to build our model!

```
final_user_df_cat <- final_user_df[,c("SIGNUP_CHANNEL", "EMAIL_TYPE", "JOB_TITLE")] #categorical variables
final_user_df_exclude <- final_user_df[, c("SIGNUP_DATE", "NUMBER_OF_EVENTS")] #variables to exclude
# From any conversion
final_user_df_cont <- final_user_df$CONTINUED # continuous variables
numeric_cat <- apply(final_user_df_cat, 2, FUN=function(x){ as.numeric(as.factor(x)) }) #convert categorical variables to numerical factors
numeric_cat_df <- as.data.frame(numeric_cat)
cont_df <- as.data.frame(final_user_df_exclude)
exclude_df <- as.data.frame(final_user_df_exclude)

#bind all dataframes after conversions to the final dataframe
final_user <- cbind(numeric_cat_df, exclude_df, cont_df)
colnames(final_user_df)[1] <- "CONTINUED"
final_user_df$ML_RETAINED <- as.factor(final_user_df$ML_RETAINED)

#Split final data into train and test data - 70%, 20%
trainIndex <- sample(nrow(final_user_df), size = round(0.8*nrow(final_user_df)), replace=FALSE)
train_df <- final_user_df[trainIndex,]
test_df <- final_user_df[-trainIndex,]

#remove ML_RETAINED from the test set
ml_retained_test <- test_df$ML_RETAINED
test_df <- test_df[, !names(test_df) %in% c("ML_RETAINED")]
summary(final_user_df)
```

Step	DF	Deviance	Resid. Df	Resid. Dev	AIC
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
	N/A	N/A	14476	15285.11	15629.11
-JOB_TITLE	1	0.06304798	14477	15285.18	15627.18
-COUNTRY	156	310.96258374	14633	15596.14	15626.14
3 rows					

```
backward <- stepAIC(full, direction="backward", trace=FALSE)
backward$anova
```

Step	DF	Deviance	Resid. Df	Resid. Dev	AIC
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
	N/A	N/A	14476	15285.11	15629.11
-JOB_TITLE	1	0.06304798	14477	15285.18	15627.18
-COUNTRY	156	310.96258374	14633	15596.14	15626.14
3 rows					

## Logistic Regression

```
#ML Model
model_glm <- glm(ML_RETAINED ~ SIGNUP_CHANNEL + EMAIL_TYPE + NUMBER_OF_EVENTS,
  data = train_df,
  family = binomial(link="logit"))

pred_glm <- predict(model_glm, newdata=test_df, type="response")
prob_glm <- as.factor(ifelse(pred_glm>0.5, 1, 0))

#Evaluation Metrics
result_glm <- confusionMatrix(data=pred_prob_glm, ml_retained_test)
precision_glm <- result_glm$byClass["Pos Pred Value"]
recall_glm <- result_glm$byClass["Sensitivity"]
f1_glm <- result_glm$byClass["F1"]
```

## Decision Trees

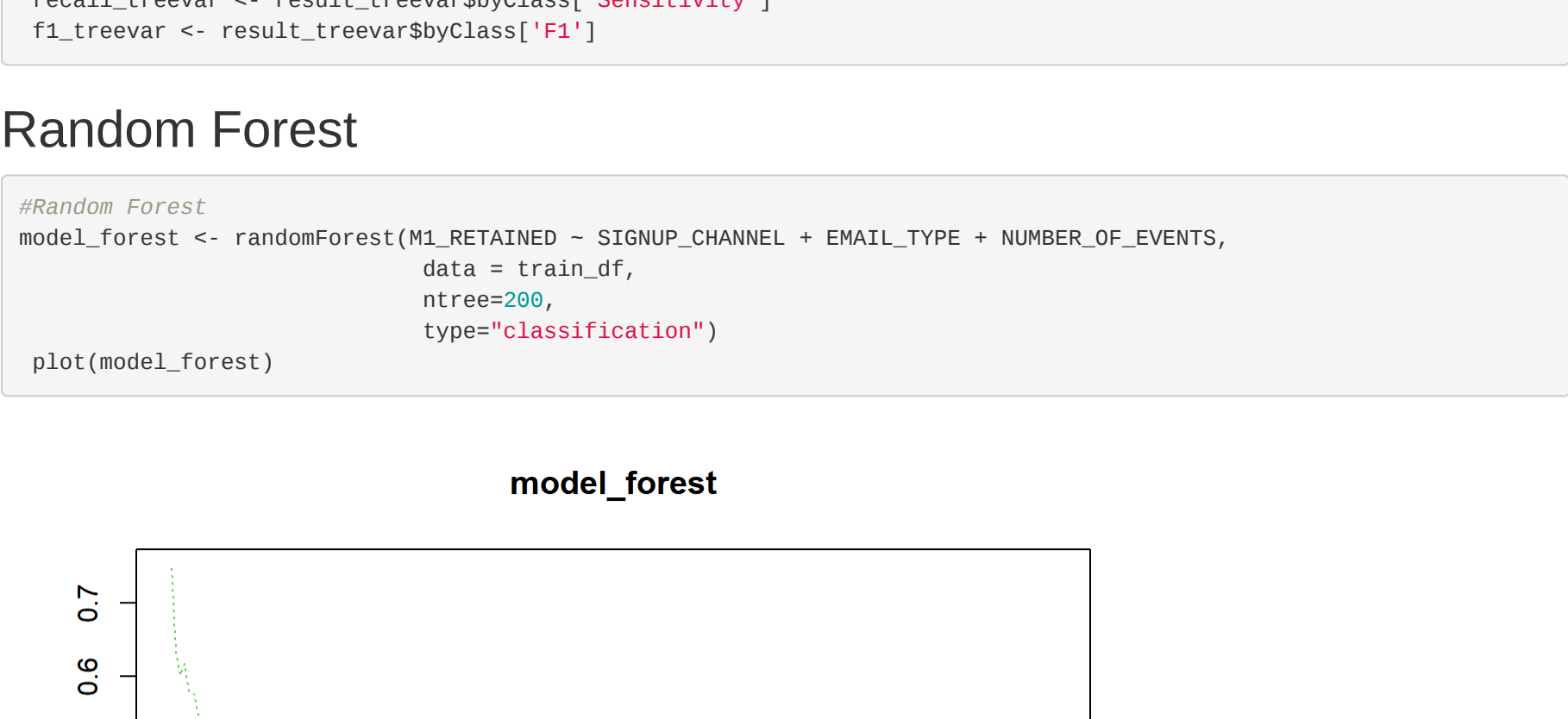
```
#Decision Tree
model_tree <- rpart(ML_RETAINED ~ SIGNUP_CHANNEL + EMAIL_TYPE + NUMBER_OF_EVENTS,
  data=train_df,
  method="class",
  control = rpart.control(xval=10))

rpart.plot(model_tree)

#Evaluation Metrics
pred_tree <- predict(model_tree, newdata=test_df, type="class")
result_tree <- confusionMatrix(data=pred_tree, ml_retained_test)
precision_tree <- result_tree$byClass["Pos Pred Value"]
recall_tree <- result_tree$byClass["Sensitivity"]
f1_tree <- result_tree$byClass["F1"]

#Decision tree - variant with job_title
model_treear <- rpart(ML_RETAINED ~ SIGNUP_CHANNEL + EMAIL_TYPE + NUMBER_OF_EVENTS + JOB_TITLE,
  data=train_df,
  method="class",
  control = rpart.control(xval=10))

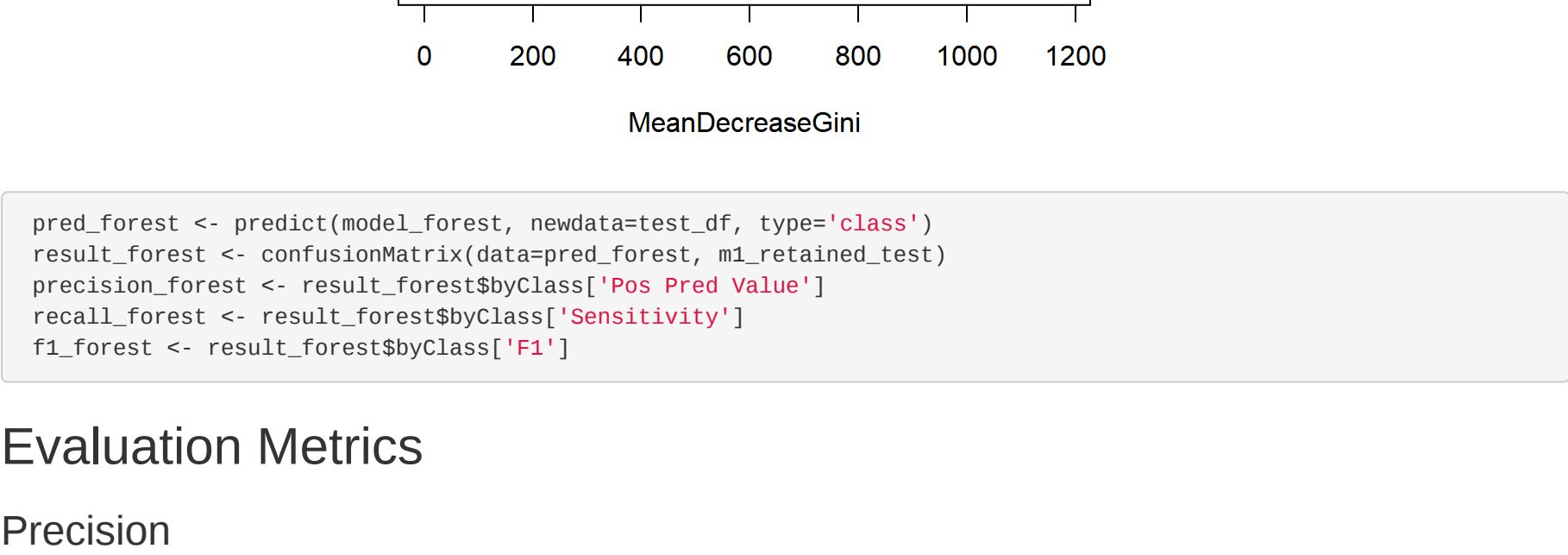
rpart.plot(model_treear)
```



```
#Evaluation Metrics
pred_treear <- predict(model_treear, newdata=test_df, type="class")
result_treear <- confusionMatrix(data=pred_treear, ml_retained_test)
precision_treear <- result_treear$byClass["Pos Pred Value"]
recall_treear <- result_treear$byClass["Sensitivity"]
f1_treear <- result_treear$byClass["F1"]
```

## Random Forest

```
#Random Forest
model_forest <- randomForest(ML_RETAINED ~ SIGNUP_CHANNEL + EMAIL_TYPE + NUMBER_OF_EVENTS,
  data = train_df,
  ntree=200,
  type="classification")
plot(model_forest)
```



```
#Evaluation Metrics
varImpPlot(model_forest, sort=T, main="Variable Importance")
```

	F1
forest	0.8434258
tree	0.8457467
treear	0.8457467
forest	0.8461192

## Evaluation Metrics

### Precision

precision_glm
## Pos Pred Value
## 0.7572634
precision_tree
## Pos Pred Value
## 0.8155394
precision_treear
## Pos Pred Value
## 0.8155394

### Recall

recall_glm
## Sensitivity
## 0.8787879
recall_tree
## Sensitivity
## 0.8782882
recall_treear
## Sensitivity
## 0.8782882
recall_forest
## Sensitivity
## 0.8778829

### F1 Statistic

f1_glm
## F1
## 0.8434258
f1_tree
## F1
## 0.8457467
f1_treear
## F1
## 0.8457467
f1_forest
## F1
## 0.8461192

## Conclusions

Assessing from the Precision and Recall metrics, its best to choose the Decision Tree Model although the Random Forest model is a close second. With more time, I would tweak the models with different hyperparameters to see which one is fact performs better. For now, I will choose the Decision Trees Model.

Additionally, I made a model **model\_treear** that used 'JOB\_TITLE' as one of the variables to consider even though it wasn't a feature variable through stepAIC model selection. Intuitively, it made more sense that a designer is more likely to continue to use Company-X after the first month vs a marketer or project designer. This model has slightly better recall than **model\_tree** and similar precision.

## Insights

- It's interesting to note that the 'SIGNUP\_CHANNEL' plays an important role in whether the customer continues on Company-X after the first month. We can continue to build on the model to find which channels are likely to create customer 'stickiness'. The growth team can determine which marketing channels are valuable investments on financial and creative resources, and invest in analytics through CTRs and monitoring behaviour using Google Analytics for the chosen streams.
- EMAIL\_TYPE plays a key role in understanding client accounts. Intuitively, Businesses and Personal users are more likely to continue on the platform. Investing in marketing streams that target these specific users and combining that with a focus on building features that these streams are more likely to use is advantageous in retaining them and diversifying the product line!
- By enhancing the functionality for key and subsidiary users of Company-X accounts, we can better retain them. For instance, when including 'JOB\_TITLE' as a variable in 'model\_treear', we were able to deduce that it's significant to know what role the user plays in their corporation. By creating functionality for project-managers and developers alongside the key user (a designer), Company-X as a platform is enables more collaborative work and retains better in the long run.
- Using 'NUMBER\_OF\_EVENTS' turned out to be very useful. Intuitively, a customer who uses Company-X more is more likely to be retained. As mentioned previously, understanding the most used event by 'JOB\_TITLE' helps with adding features beneficial to specific use cases. Hence, by improving events/activities that keep customers coming back and fixing those that have not been used as much could help.

## Errors and Assumptions

When adding 'NUMBER\_OF\_EVENTS' to 'final\_user\_df', I encountered an issue with duplication that was caused during the use of both an inner join and left join. I solved this issue with by using a different SQL query when creating the dataframe, however instead of 18323 users, I was left with 18310 users. I believe this error occurred since R detected duplicate entries of 'USER\_ID' in one of the tables that maps to all duplicate values in the other table through 'inner join' or 'left join'. However, the 'users\_data' has no duplicate entries on inspecting the CSV. Since this accounts for less than 5% of the total users, I continued with building the predictive model using 'final\_user\_df'.