

Untitled

Objective

Red Hat Customer Survey

The objective of this project is to classify customer potential based on a survey conducted by Red Hat. The database for this project has been downloaded from Kaggle. The main goal of this project is to predict the potential business value of a person based on an activity performed by that person over a given time frame. The business value outcome of each person is defined by Yes/No for each unique activity. This is used to indicate if a person has finished an activity within a time frame.

Library

```
##  
## Attaching package: 'dplyr'  
  
## The following objects are masked from 'package:plyr':  
##  
##     arrange, count, desc, failwith, id, mutate, rename, summarise,  
##     summarise  
  
## The following objects are masked from 'package:stats':  
##  
##     filter, lag  
  
## The following objects are masked from 'package:base':  
##  
##     intersect, setdiff, setequal, union  
  
## Loading required package: lattice  
  
##  
## Attaching package: 'kernlab'  
  
## The following object is masked from 'package:ggplot2':  
##  
##     alpha  
  
## Loading required package: gplots  
  
##  
## Attaching package: 'gplots'  
  
## The following object is masked from 'package:stats':  
##  
##     lowess
```

```

## 
## Attaching package: 'gridExtra'

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

```

Load data

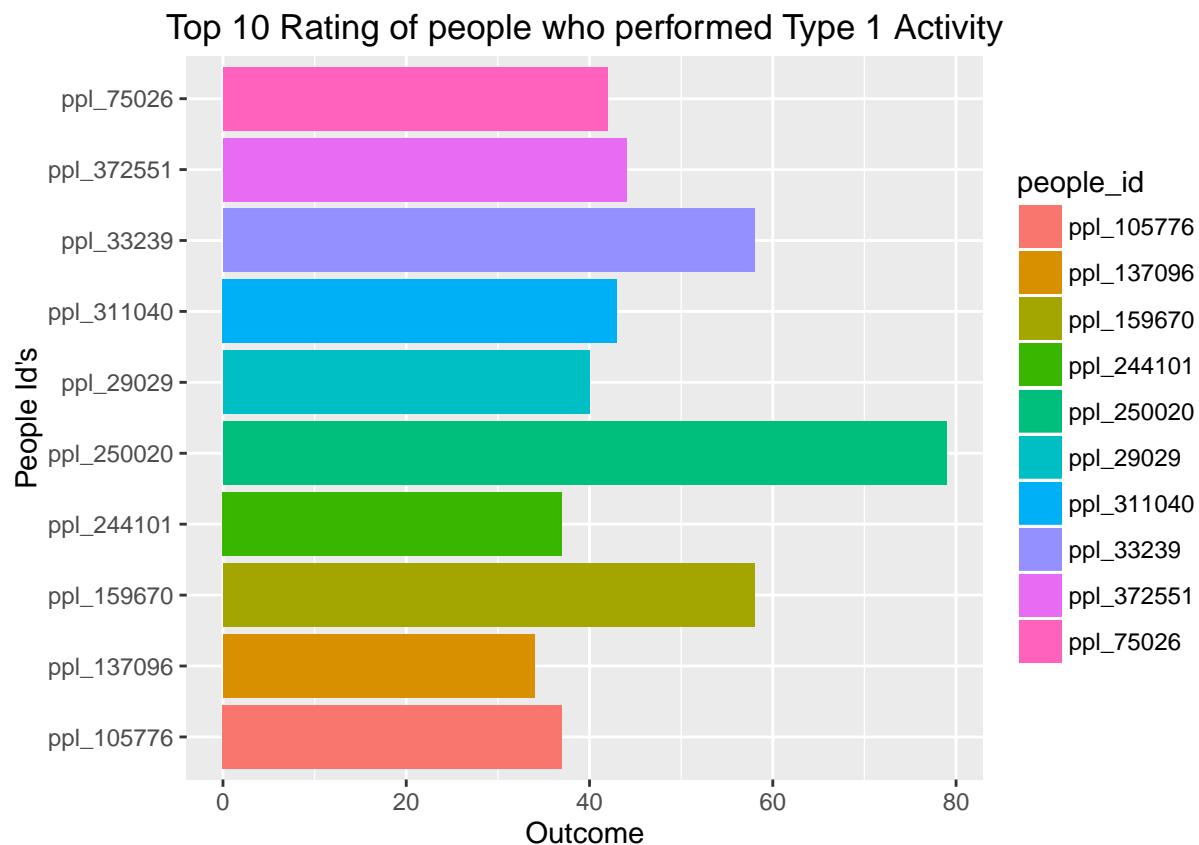
Cleaning and Prepare training data set

Analysis 1

Cleaning and Prepare testing data set

Split data in traing and testing sets for Type 1 activity

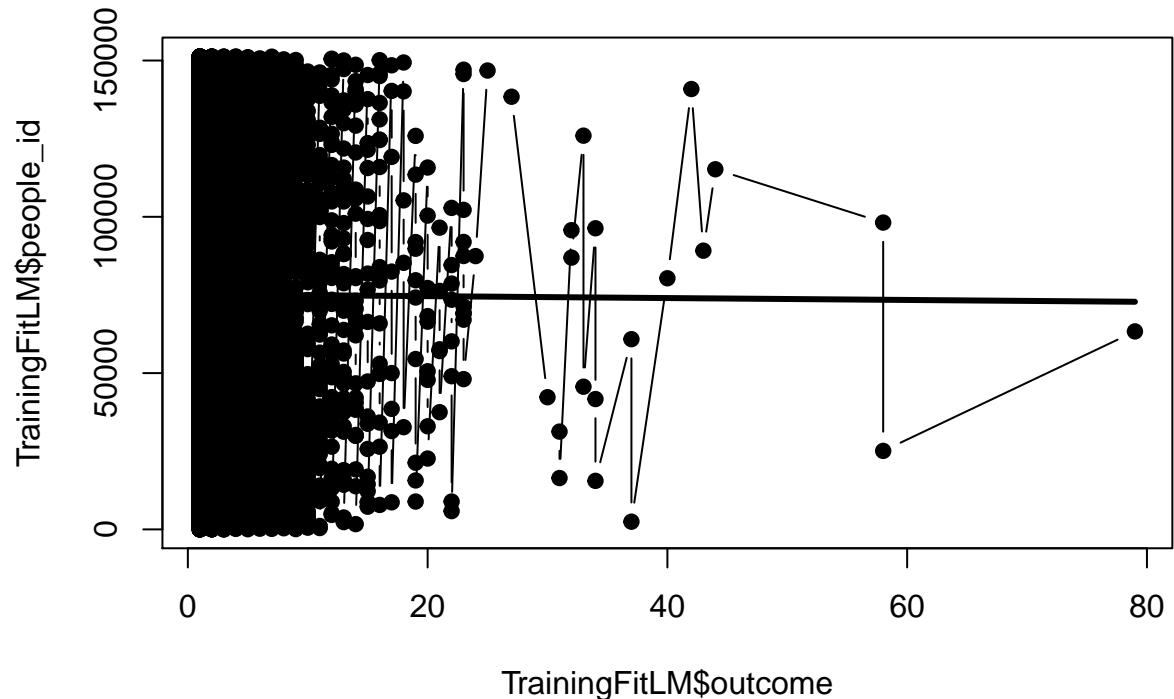
Prediction with Regression for Type 1 activity



Fit a linear model for Type 1 Activity

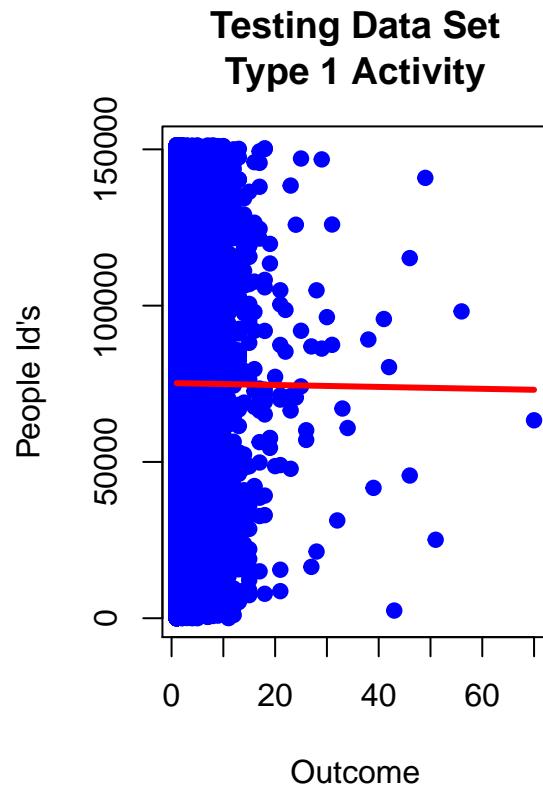
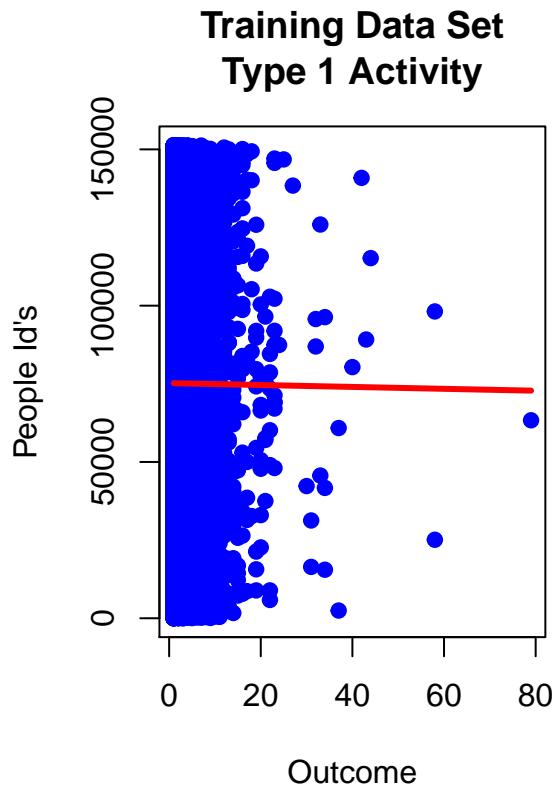
```
##  
## Call:  
## lm(formula = people_id ~ outcome, data = TrainingFitLM)  
##  
## Coefficients:  
## (Intercept)      outcome  
##       75226.6        -30.5
```

Model Fit for Type 1 Activity



Predict a new value for Type 1 Activity

Plot predictions - Testing and Training for Type 1 Activity



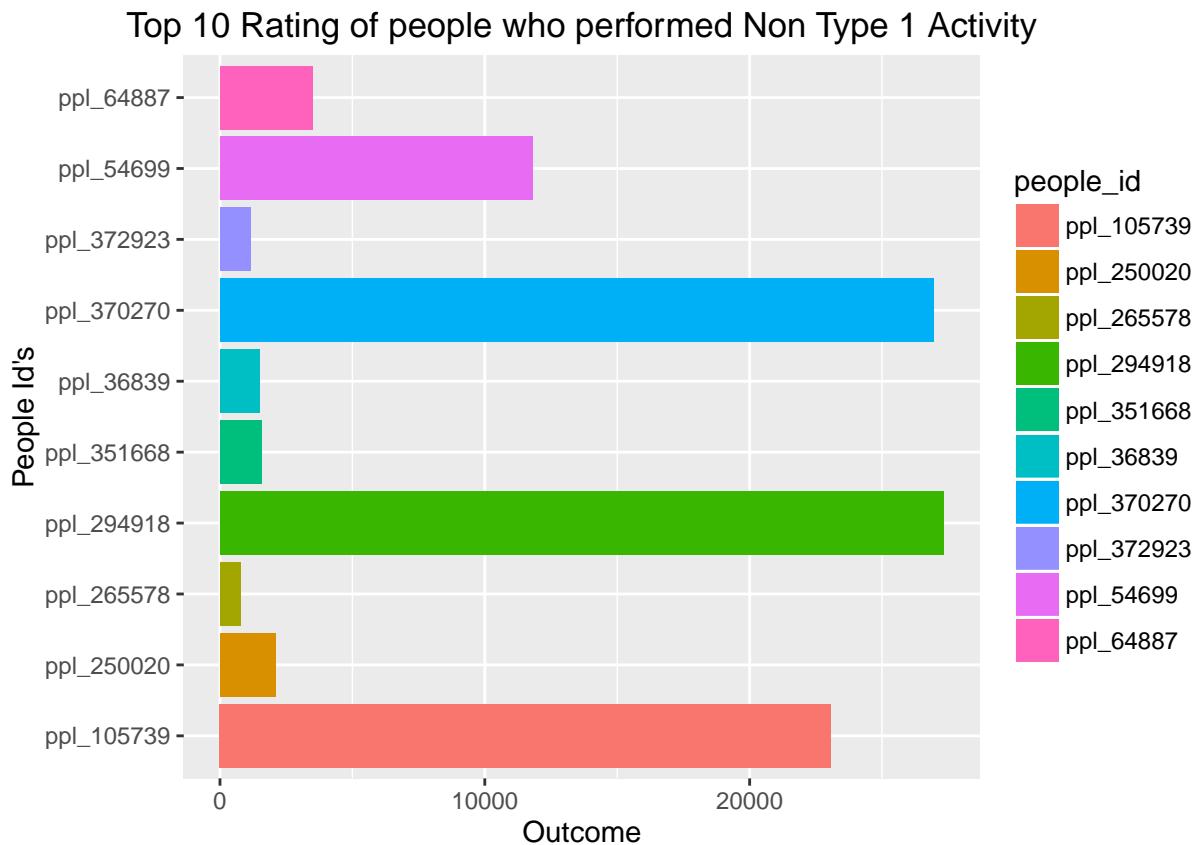
Calculate Traing and Testing data set errors for Type 1 Activity

```
## [1] 9503898
```

```
## [1] 9518819
```

Split data in traing and testing sets for Non Type 1 activity

Predict with regression for non Type 1 Activity

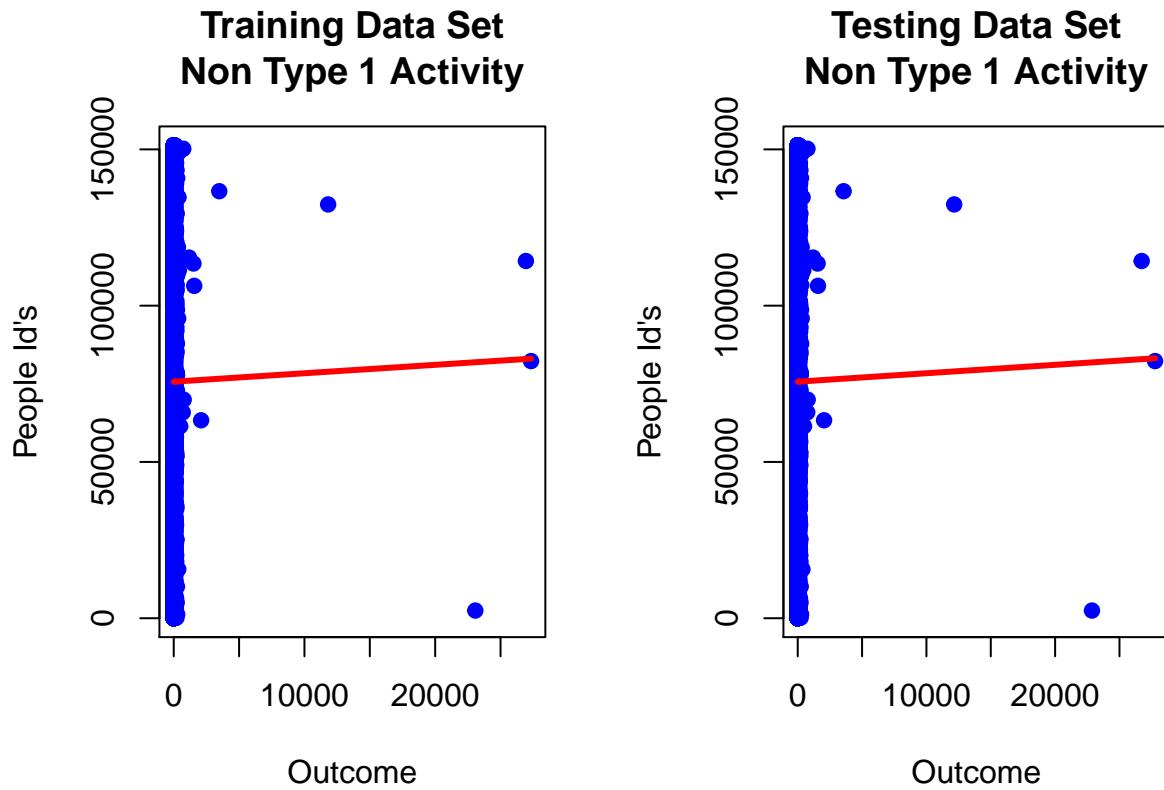


Fit a linear model for non Type 1 Activity

```
##  
## Call:  
## lm(formula = people_id ~ outcome, data = TrainingFitNonType1LM)  
##  
## Coefficients:  
## (Intercept)      outcome  
##   7.568e+04    2.695e-01
```

Predict a new value for non Type 1 Activity

Plot prediction - testing and traing for non Type 1 activity

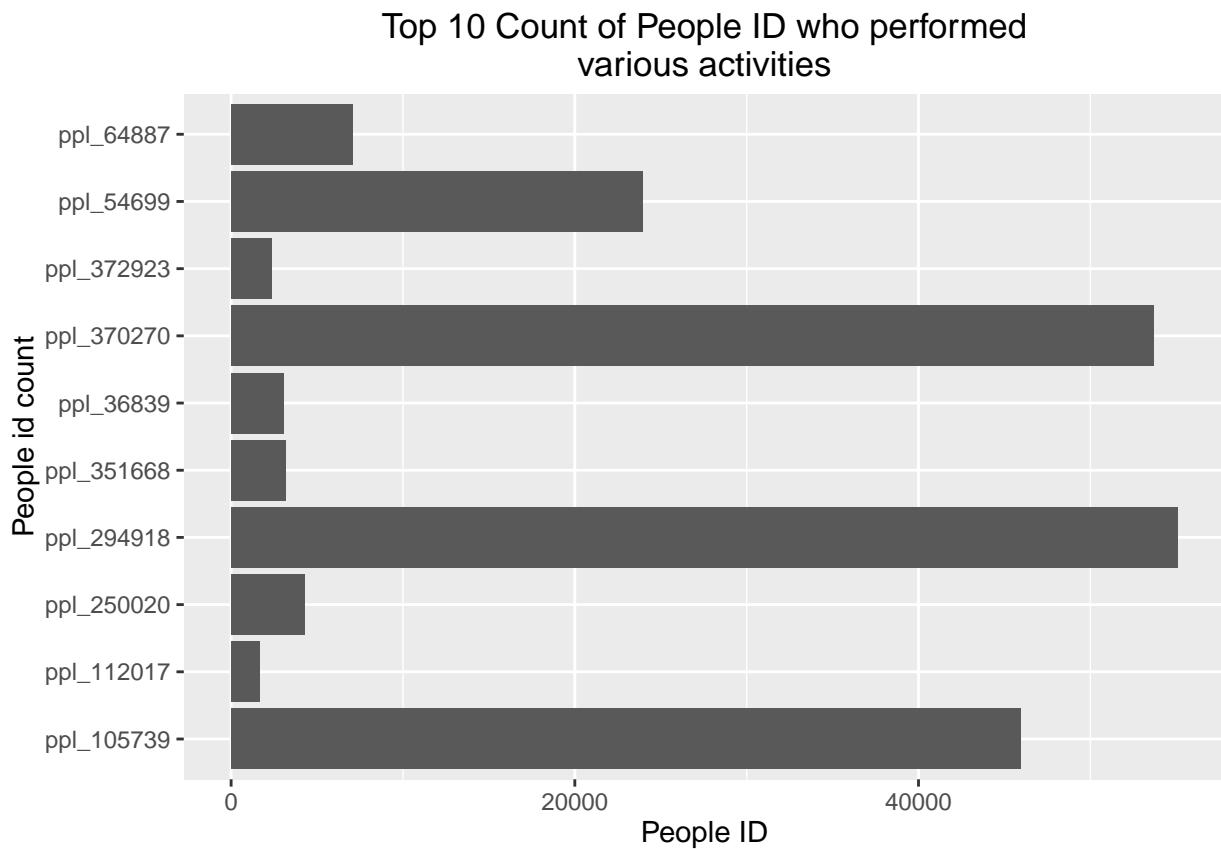


Calculate training and testing data set errors for non Type 1 activity

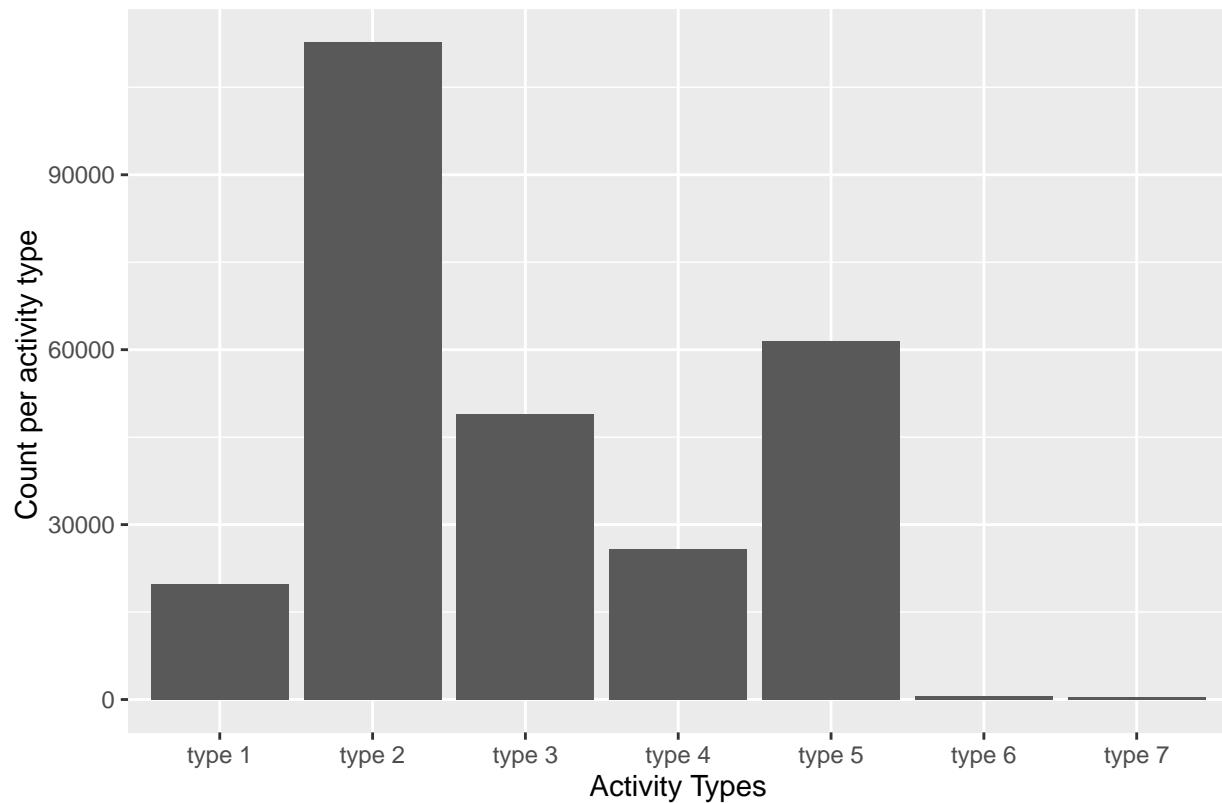
```
## [1] 15178340
```

```
## [1] 15173857
```

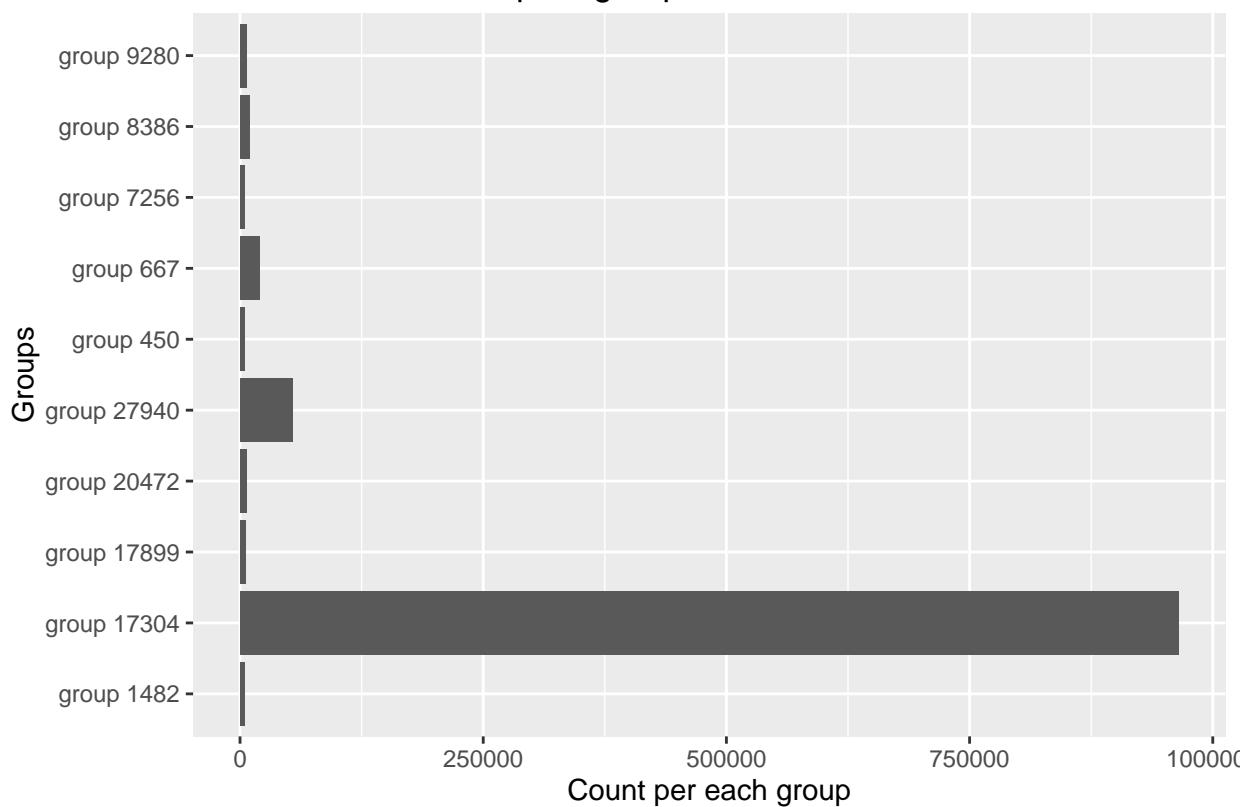
Analysis 2

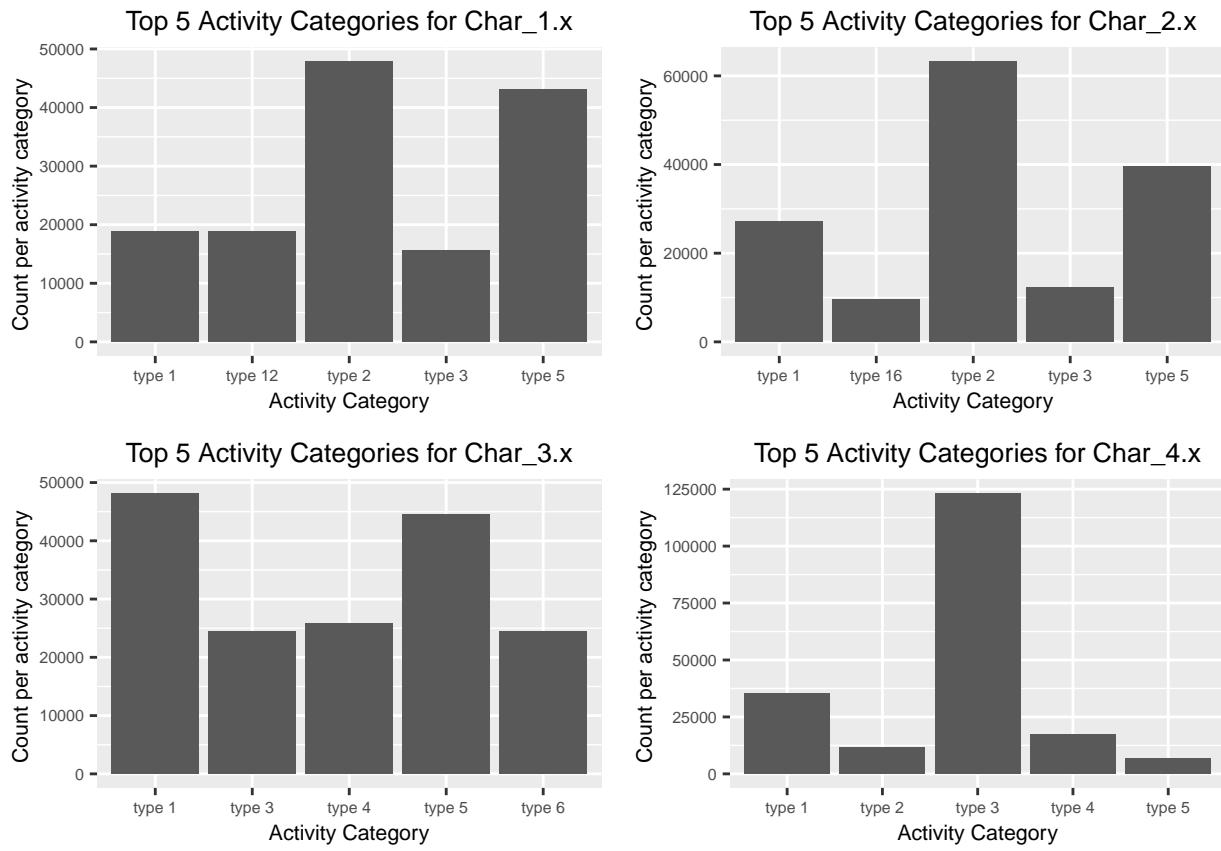


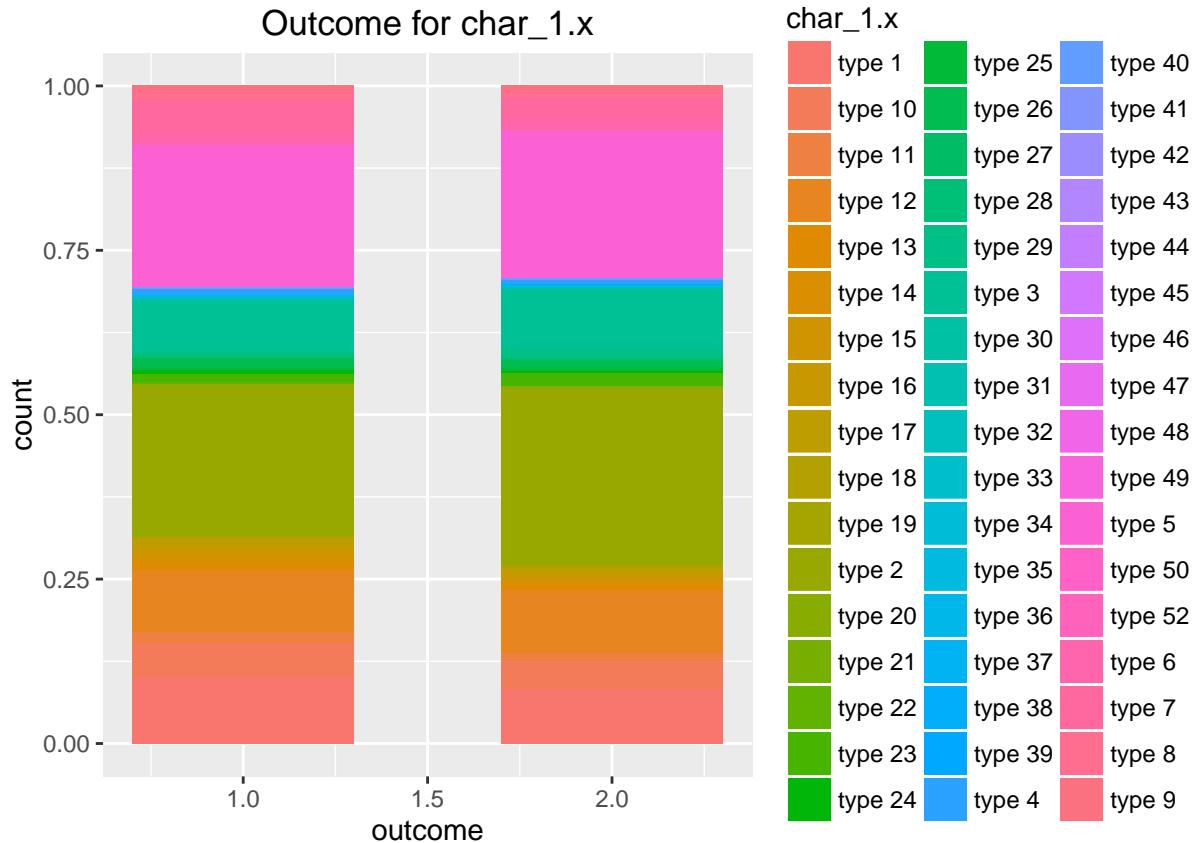
Activity Categories and their total count

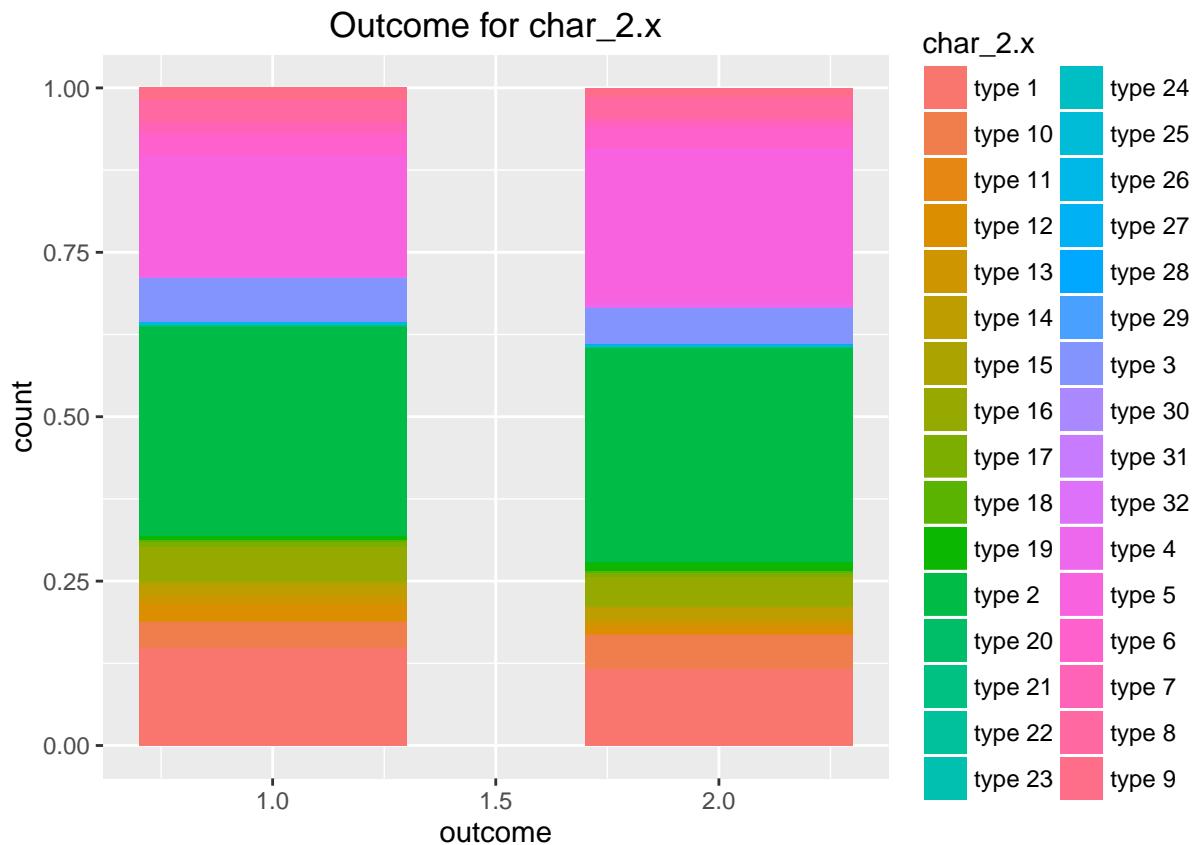


Top 10 groups and their count

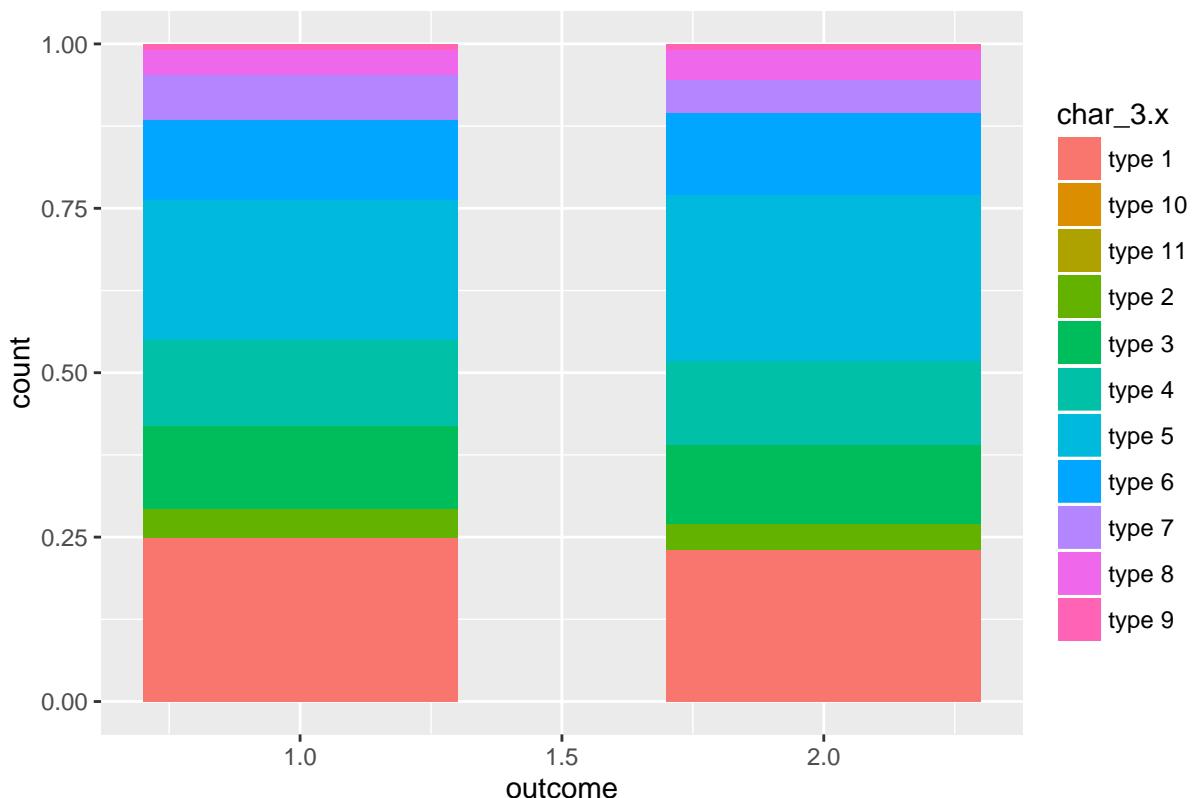




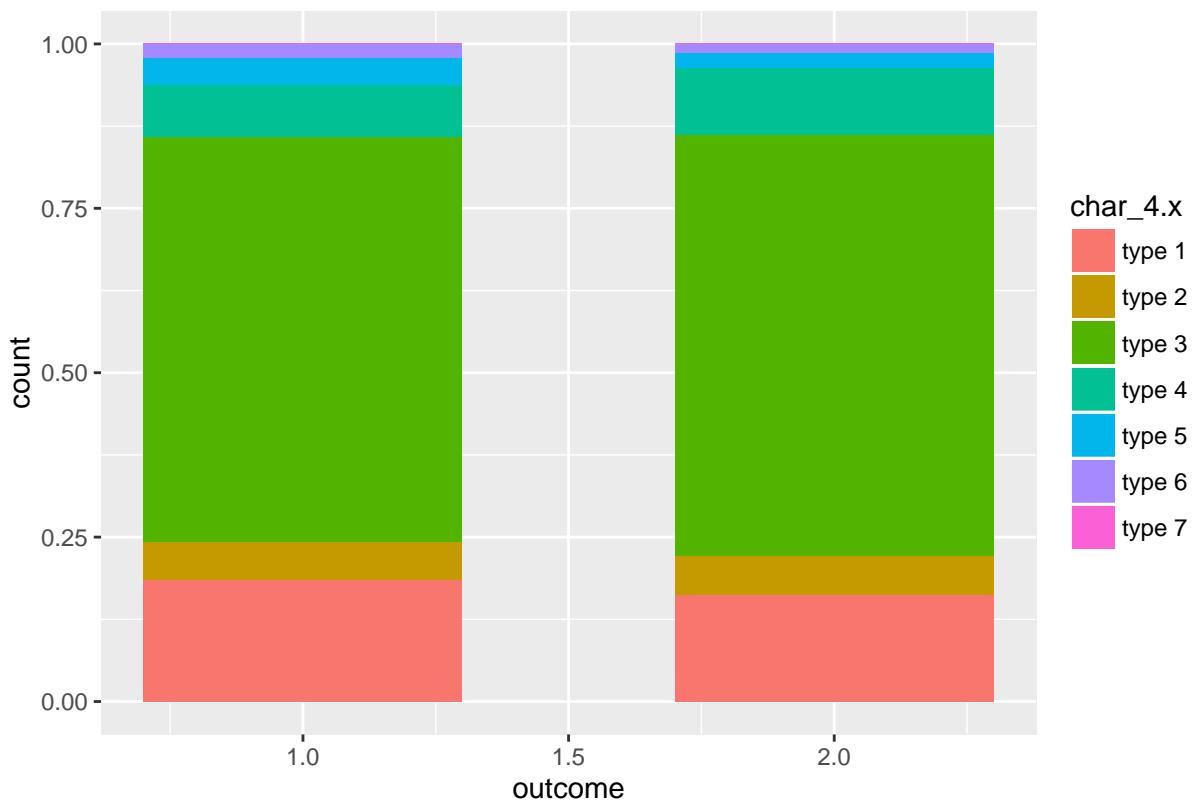




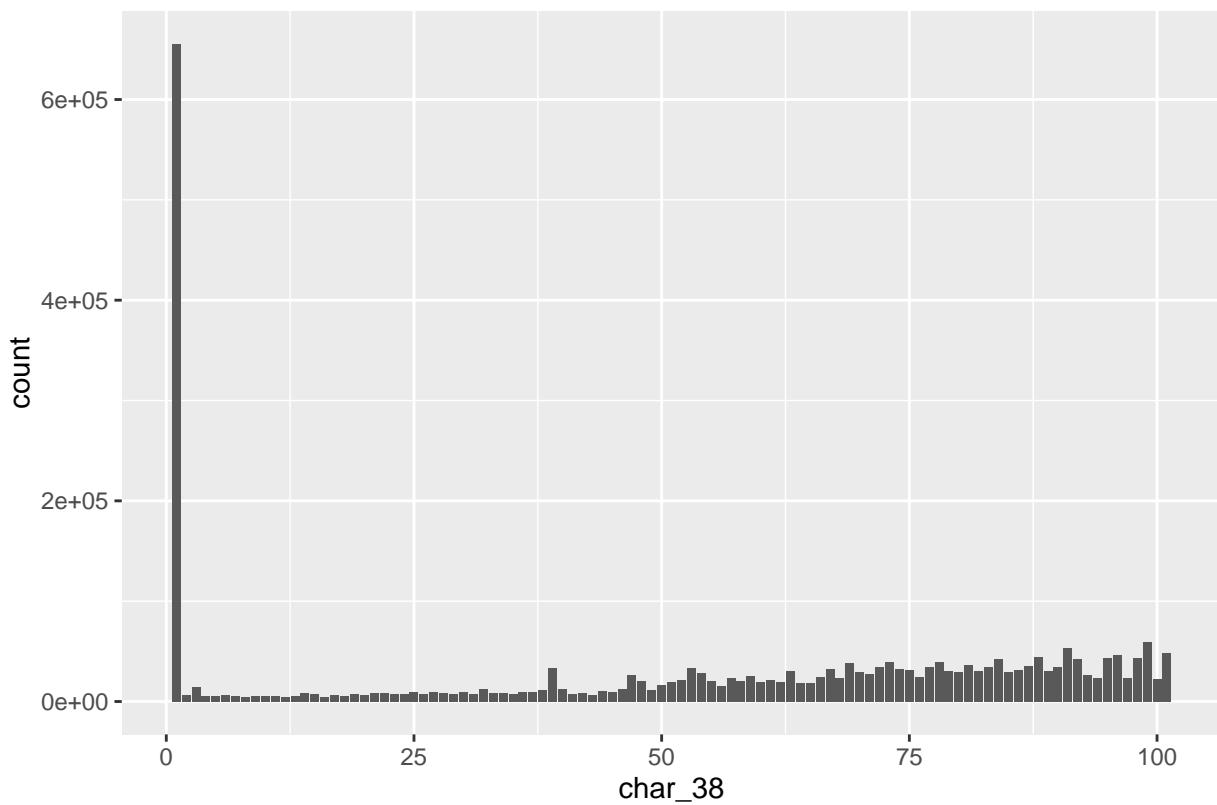
Outcome for char_3.x



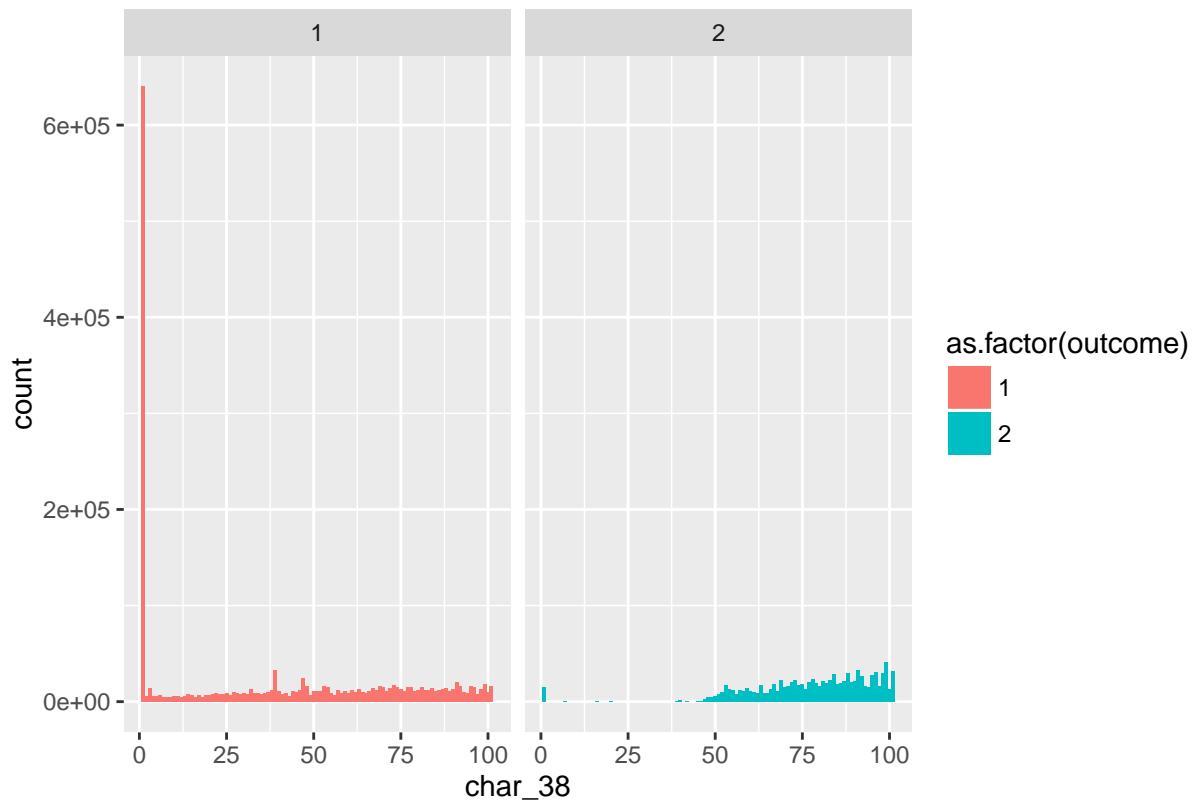
Outcome for char_4.x



Distribution of char_38



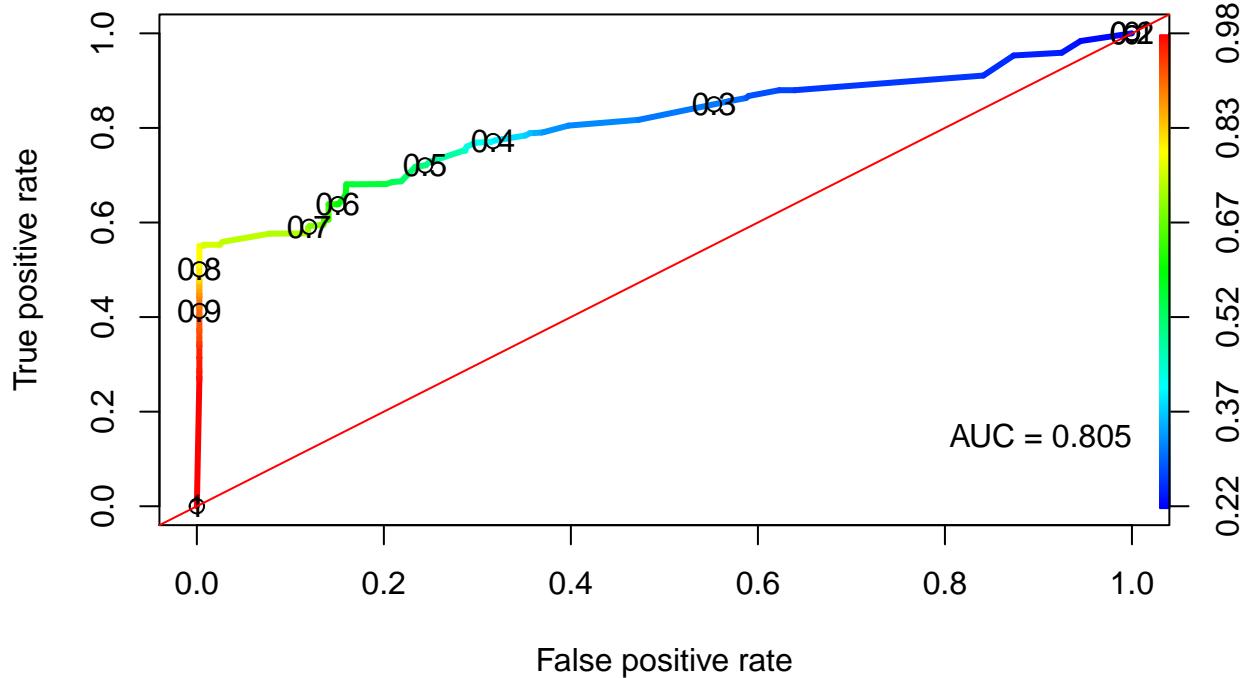
Distribution of char_38



Receiver Operating Characteristics

```
## numeric(0)
```

char_38 ROC curve for people with only 0/1



Conclusion

The ROC curve for 'char_38' was plotted with a sample size of 5000 records from the combined data set comprising of people, act_train and act_test. The AUC value comes to 0.805 which shows that the prediction is correct as it is very near to AUC value of 1.