

# PUBG Top 10% Placement Analysis

*Chance Robinson, Allison Roderick and William Arnost*

*Master of Science in Data Science, Southern Methodist University, USA*

## 1 Introduction

[Intro]

## 2 Data Description

The source data is available on Kaggle.com under the competition PUBG Finish Placement Prediction. The files used in our analysis were transformed to fit the requirements of a binomial logistic regression classifier. The data has also been pre-split into training and test files for consistency when comparing our different model types. Additionally, as the percentage-based nature of the top 10% of players is inherently unbalanced, we've also downsampled the higher frequency data to match that over the lower frequency outcome of the top 10% of players.

`pubg_solo_game_types.csv`

- Filtered for solo only game types

`pubg_solo_game_types_test_full.csv`

- Pre-split for test data

`pubg_solo_game_types_train_full.csv`

- Pre-split for train data without downsampling for the unbalanced response variable

`pubg_solo_game_types_train_downsampled.csv`

- Pre-split for train data with downsampling for the unbalanced response variable

### 2.1 Data Dictionary

Column Name	Type	Description
DBNOs		Number of enemy players knocked.
assists		Number of enemy players this player damaged that were killed by teammates.

Column Name	Type	Description
boosts		Number of boost items used.
damageDealt		Total damage dealt. Note: Self inflicted damage is subtracted.
headshotKills		Number of enemy players killed with headshots.
heals		Number of healing items used.
Id		Players Id
killPlace		Ranking in match of number of enemy players killed.
killPoints		Kills-based external ranking of player. (Think of this as an Elo ranking where only kills matter.) If there is a value other than -1 in rankPoints, then any 0 in killPoints should be treated as a None.
killStreaks		Max number of enemy players killed in a short amount of time.
kills		Number of enemy players killed.
longestKill		Longest distance between player and player killed at time of death.
matchDuration		Duration of match in seconds.
matchId		ID to identify match. There are no matches that are in both the training and testing set.
matchType		String identifying the game mode that the data comes from.
rankPoints		Elo-like ranking of player.
revives		Number of times this player revived teammates.
rideDistance		Total distance traveled in vehicles measured in meters.
roadKills		Number of kills while in a vehicle.
swimDistance		Total distance traveled by swimming measured in meters.
teamKills		Number of times this player killed a teammate.
vehicleDestroys		Number of vehicles destroyed.
walkDistance		Total distance traveled on foot measured in meters.
weaponsAcquired		Number of weapons picked up.
winPoints		Win-based external ranking of player. (Think of this as an Elo ranking where only winning matters.) If there is a value other than -1 in rankPoints, then any 0 in winPoints should be treated as a None.
groupId		ID to identify a group within a match. If the same group of players plays in different matches, they will have a different groupId each time.
numGroups		Number of groups we have data for in the match.
maxPlace		Worst placement we have data for in the match. This may not match with numGroups, as sometimes the data skips over placements.

Column Name	Type	Description
winPlacePerc		This is a percentile winning placement, where 1 corresponds to 1st place, and 0 corresponds to last place in the match. (to be removed from our binomial classifier so as not to influence our predictive results)
top.10		The target of prediction. This is a percentile winning placement, where 1 corresponds to a top 10% placement and 0 in the lower 90%.

## **2.2 Exploratory Data Analysis**

# **3 Objective I Analysis**

## **3.1 Question of Interest**

## **3.2 Model Selection**

[Logistic Regression, Ridge, Lasso and Elastic Net]

## **3.3 Comparing Competing Models**

- 1) AUC
- 2) Sen
- 3) Spec
- 4) Err Rate (1 - Accuracy)

## **3.4 Model Interpretation**

## **3.5 Conclusion**

## 4 Objective II Analysis

### 4.1 Question of Interest

### 4.2 Model Selection

[Logistic Regression with interactions, LDA and Random Forest]

#### 4.2.1 Linear Discriminant Analysis

Our next prediction tool is Linear Discriminant Analysis (LDA) for classifying the player as Top 10 or not. We have taken a subset of the continuous variables from our EDA to build the LDA off of. Before running the LDA, we will cover two things: a LASSO call to eliminate less important variables and assumption checking.

##### 4.2.1.1 LASSO

The LASSO call plus manual variable selection reduced the predictors considered for the LDA model to: boosts, heals, killPlace, killStreaks, longestKill, matchDuration, rideDistance, swimDistance, teamKills, walkDistance, weaponsAcquired.

##### 4.2.1.2 Assumption Checking

LDA performs optimally when the assumptions of MANOVA are met. That is,

1. The predictors are normally distributed for each class of the response.
2. The covariance matrices for each class of the response are homogeneous.

When we check the first assumption for the predictors that are to be included in the LDA model, we see that the assumption is not met. Most of the predictors are right skewed. The variable matchDuration is bimodal. To remedy this, we tried transforming the variables, but it did not help our overall prediction accuracy. However, because issues of normality exist, we will explore QDA as well as LDA.

We also checked the homogeneity of correlation matrices. we find that, overall, there are no major departures from homogeneity. The variables walkDistance and killPlace show the greatest deviances from homogenous correlations between bottom 90 and top 10 placements. Consequently, we tried removing those variables from the model. However, removing those variables reduced our prediction accuracy.

Thus, we will proceed with the variables selected and see if LDA or QDA performs better.

##### 4.2.1.3 LDA Results

LDA has a prediction accuracy of 0.9206, with a sensitivity of 0.57479 and a specificity of 0.96099. The area under the ROC curve is 0.941.

#### **4.2.1.4 QDA Results**

QDA has a prediction accuracy of 0.8779, with a sensitivity of 0.68262 and a specificity of 0.90071. The area under the ROC curve is 0.912.

#### **4.2.1.5 LDA Conclusion**

Although the QDA is better at predicting the top 10 placements that were true top 10 than LDA (QDA sensitivity of 0.68 > LDA sensitivity of 0.57), QDA predicts many more incorrect top 10 placements than LDA does (1907 LDA false positives < 4853 QDA false positives). Because LDA has a better overall accuracy (0.9206 for LDA > 0.8779 for QDA), we think the LDA is a stronger model for prediction than QDA.

### **4.3 Comparing Competing Models**

#### **4.4 Model Interpretation**

#### **4.5 Conclusion**

## 5 Appendix

### 5.1 Exploratory Data Analysis

### 5.2 Code

#### 5.2.1 LDA

##### 5.2.1.1 LASSO

```
# Reduce variables to relevant continuous variables
```

```
train1<-train[,c(5:6,8:10,12:15,19,21:27,29:30)]
```

```
test1<-test[,c(5:6,8:10,12:15,19,21:27,29:30)]
```

```
# train2 <- train1[,c(1,4:19)]
```

```
# test2 <- test1[,c(1,4:19)]
```

```
train2_alt <- train1[,c(1,4:17,19)]
```

```
test2_alt <- test1[,c(1,4:17,19)]
```

```
# Get data in format for LASSO
```

```
x=model.matrix(top.10~.,train2_alt)[,-1]
```

```
y=as.numeric(train2_alt$top.10)
```

```
xtest<-model.matrix(top.10~.,test2_alt)[,-1]
```

```
ytest<-as.numeric(test2_alt$top.10)
```

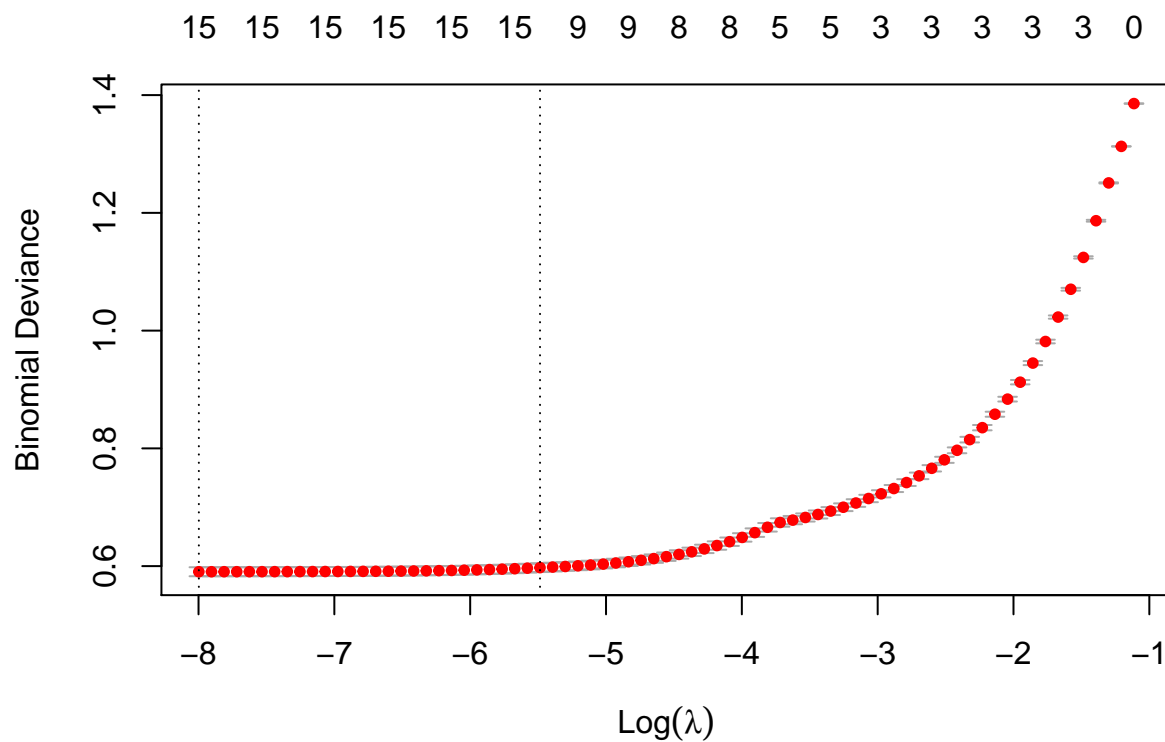
```
grid=10^seq(10,-2, length =100)
```

```
lasso.mod=glmnet(x,y,alpha=1, lambda =grid)
```

```
set.seed(23) #removes kills roadKills vehicleDestroys
```

```
cv.out=cv.glmnet(x,y,alpha=1,family="binomial") #alpha=1 performs LASSO
```

```
lda.lasso<-plot(cv.out)
```



```
# simplest model
lda.lasso.model.coef<-coef(cv.out, cv.out$lambda.1se)
lda.lasso.model.coef
```

```
## 16 x 1 sparse Matrix of class "dgCMatrix"
##              1
## (Intercept)  3.091783e+00
## boosts      2.672920e-01
## heals       -2.784586e-03
## killPlace    -8.033586e-02
## kills        .
## killStreaks  -1.117635e+00
## longestKill  -5.733176e-05
## matchDuration -1.609689e-03
## rankPoints    1.034136e-05
## rideDistance  1.948066e-04
## roadKills     .
## swimDistance  5.967117e-04
## teamKills     -1.509184e+00
## vehicleDestroys .
## walkDistance  8.511812e-04
```

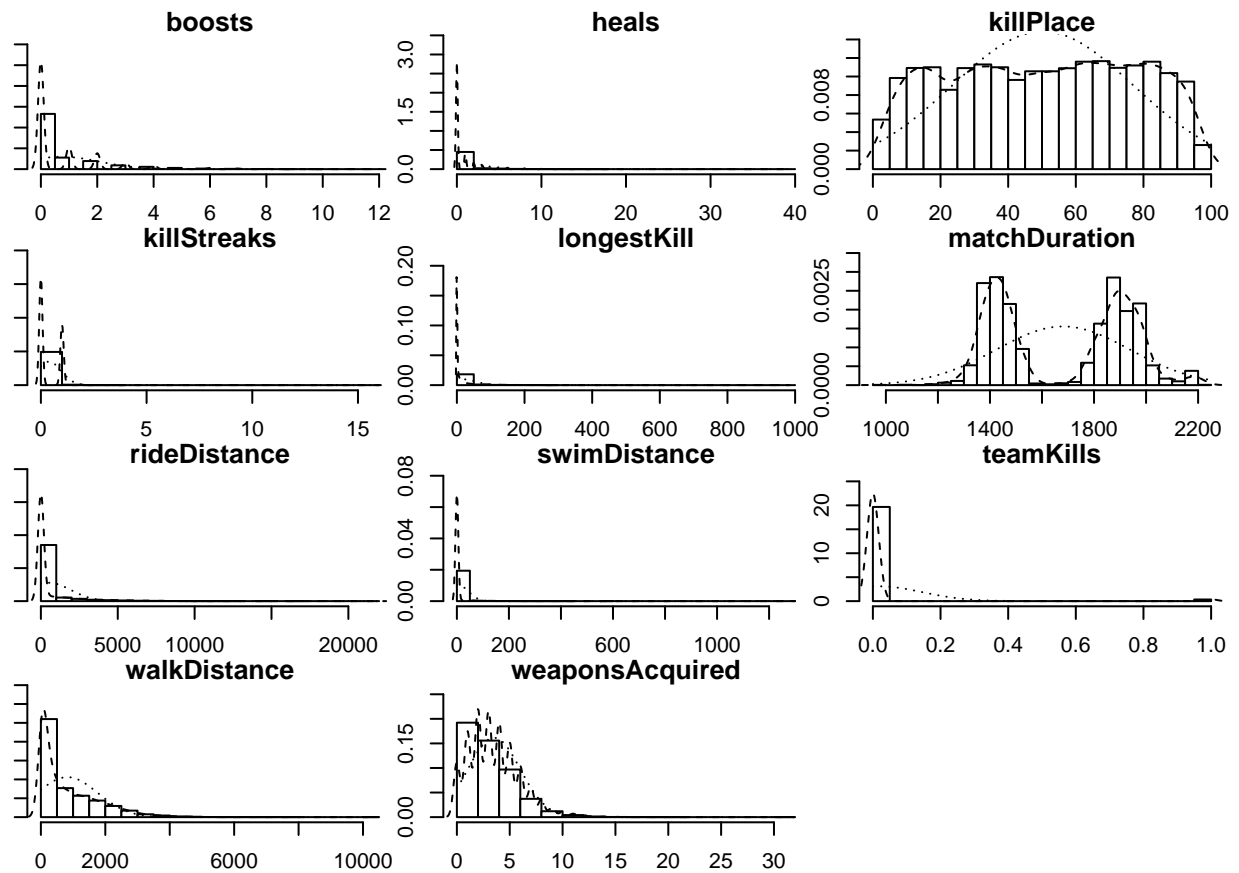


```
## weaponsAcquired 5.746021e-02
```

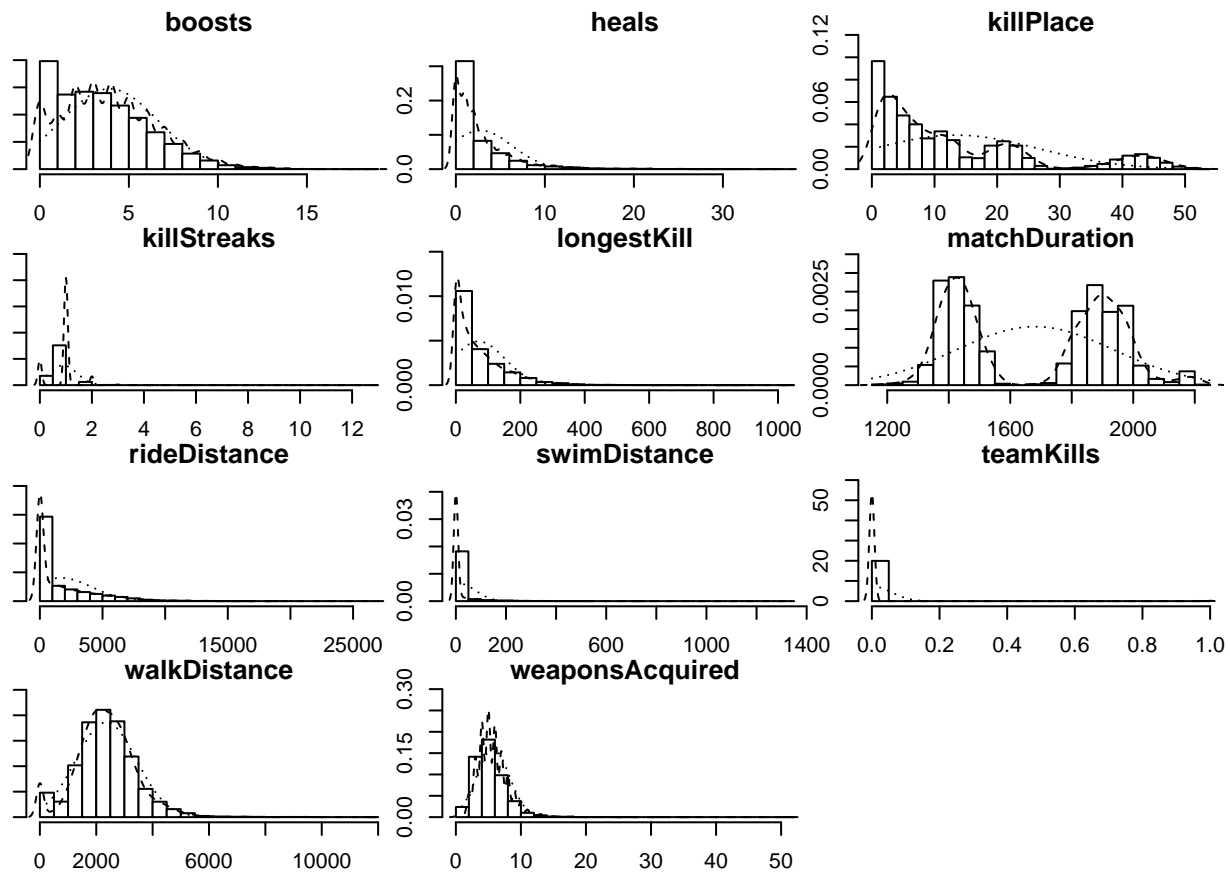
```
# Removing the kills, roadKills, and vehicleDestroys as shown above  
# Also removing rankPoints because  
# (a) the majority of the players in these data aren't ranked and  
# (b) the Kaggle description says "This ranking is inconsistent and is being deprecated in the API's next version"  
train_final <- train2_alt[,c(-4,-8,-10,-13)]
```

#### 5.2.1.2 Assumption Checking

```
train_final_no <- train_final[which(train_final$top.10==0),]  
train_final_yes <- train_final[which(train_final$top.10==1),]  
  
# nrow(train_final)  
# nrow(train_final_no)  
# nrow(train_final_yes)  
  
max_cols<-ncol(train_final)  
no_matrix<-as.matrix(train_final_no[, -max_cols])  
yes_matrix<-as.matrix(train_final_yes[, -max_cols])  
  
no_hist<-multi.hist(no_matrix)
```



```
yes_hist<-multi.hist(yes_matrix)
```



```
# par(mfrow=c(1,1))
```

```
#http://www.sthda.com/english/wiki/ggplot2-quick-correlation-matrix-heatmap-r-software-and-data-visuali
```

```
get_upper_tri <- function(cormat){
  cormat[lower.tri(cormat)]<- NA
  return(cormat)
}
```

```
custom_corr_plot <- function(cormat){
```

```
  upper_tri <- get_upper_tri(cormat)
```

```
  # Melt the correlation matrix
```

```
  melted_cormat <- melt(upper_tri, na.rm = TRUE)
```

```
  # Create a ggheatmap
```

```
  ggheatmap <- ggplot(melted_cormat, aes(Var2, Var1, fill = value))+
```

```
    geom_tile(color = "white")+
```

```
    scale_fill_gradient2(low = "blue", high = "red", mid = "white",
                        midpoint = 0, limit = c(-1,1), space = "Lab",
                        name="Pearson\nCorrelation") +
```

```
    theme_minimal()+ # minimal theme
```

```

    theme(axis.text.x = element_text(angle = 45, vjust = 1,
                                      size = 12, hjust = 1))+

    coord_fixed()

p<-ggheatmap +
  geom_text(aes(Var2, Var1, label = value), color = "black", size = 4) +
  theme(
    axis.title.x = element_blank(),
    axis.title.y = element_blank(),
    panel.grid.major = element_blank(),
    panel.border = element_blank(),
    panel.background = element_blank(),
    axis.ticks = element_blank(),
    legend.justification = c(1, 0),
    legend.position = c(0.6, 0.7),
    legend.direction = "horizontal")+
  guides(fill = guide_colorbar(barwidth = 7, barheight = 1,
                              title.position = "top", title.hjust = 0.5))

  return(p)
}

cormat <- round(cor(no_matrix),2)
no <- custom_corr_plot(cormat)

cormat <- round(cor(yes_matrix),2)
yes <- custom_corr_plot(cormat)

```

### 5.2.1.3 LDA Results

```

# Run LDA
lda <- lda(top.10 ~ . , data=train_final, prior=c(.9,.1))

lda

```

```

## Call:
## lda(top.10 ~ ., data = train_final, prior = c(0.9, 0.1))
##
## Prior probabilities of groups:
##    0    1
## 0.9 0.1
##

```

```
## Group means:
##      boosts      heals killPlace killStreaks longestKill matchDuration
## 0 0.7381078 0.8007815 50.25445 0.3915984 15.00272 1681.215
## 1 3.9244758 2.7710228 13.42301 0.8962200 70.77070 1680.553
##      rideDistance swimDistance teamKills walkDistance weaponsAcquired
## 0      515.6142      5.365372 0.016908394      823.7333      3.517848
## 1     1666.6767     16.136750 0.003005937     2351.6425     5.710528
##
## Coefficients of linear discriminants:
##                                LD1
## boosts          0.1636559657
## heals           -0.0071779101
## killPlace       -0.0313509817
## killStreaks     -0.4645084746
## longestKill      0.0003963294
## matchDuration   -0.0007443145
## rideDistance     0.0001033605
## swimDistance     0.0006812726
## teamKills        -1.0147862026
## walkDistance     0.0004217550
## weaponsAcquired  0.0155413057
```

```
predict <- predict(lda,newdata=test)

test_final <- test
test_final$predcited_place <- as.vector(predict$class)
test_final$top.10<-as.factor(test_final$top.10)
test_final$predcited_place<-as.factor(test_final$predcited_place)

xtab<-table(test_final$predcited_place,test_final$top.10)
confusionMatrix(xtab, positive="1")
```

```
## Confusion Matrix and Statistics
##
##
##      0      1
## 0 46972 2425
## 1 1907 3278
##
##              Accuracy : 0.9206
##              95% CI : (0.9183, 0.9229)
##      No Information Rate : 0.8955
##      P-Value [Acc > NIR] : < 2.2e-16
```

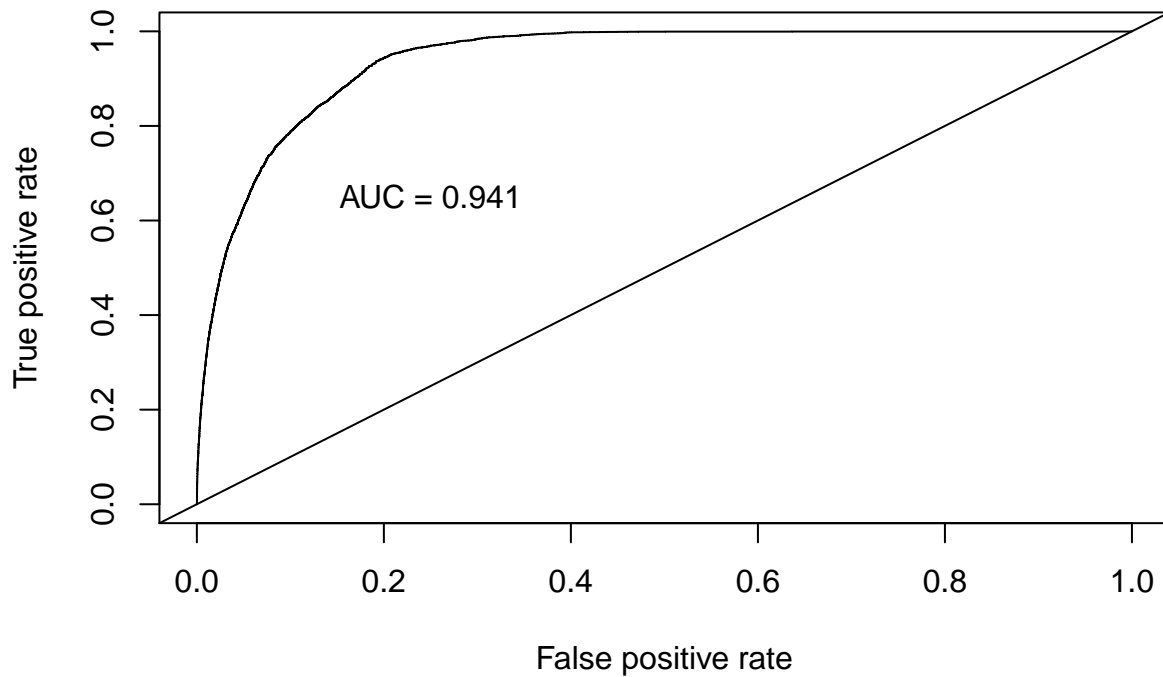
```
##
##           Kappa : 0.5582
##
## McNemar's Test P-Value : 3.998e-15
##
##           Sensitivity : 0.57479
##           Specificity : 0.96099
##           Pos Pred Value : 0.63221
##           Neg Pred Value : 0.95091
##           Prevalence : 0.10448
##           Detection Rate : 0.06006
##           Detection Prevalence : 0.09499
##           Balanced Accuracy : 0.76789
##
##           'Positive' Class : 1
##
```

```
predict.posterioriors <- as.data.frame(predict$posterior)

# Evaluate the model
pred <- prediction(predict.posterioriors[,2], test$top.10)
roc.perf = performance(pred, measure = "tpr", x.measure = "fpr")
auc.train <- performance(pred, measure = "auc")
auc.train <- auc.train@y.values

# Plot
plot(roc.perf, main="ROC Curve - LDA")
abline(a=0, b= 1)
text(x = .25, y = .65 ,paste("AUC = ", round(auc.train[[1]],3), sep = ""))
```

## ROC Curve – LDA



### 5.2.1.4 QDA Results

```
qda <- qda(top.10 ~ . , data=train_final, prior=c(.9,.1))
```

```
qda
```

```
## Call:
```

```
## qda(top.10 ~ . , data = train_final, prior = c(0.9, 0.1))
```

```
##
```

```
## Prior probabilities of groups:
```

```
## 0 1
```

```
## 0.9 0.1
```

```
##
```

```
## Group means:
```

```
##      boosts      heals killPlace killStreaks longestKill matchDuration
```

```
## 0 0.7381078 0.8007815 50.25445 0.3915984 15.00272 1681.215
```

```
## 1 3.9244758 2.7710228 13.42301 0.8962200 70.77070 1680.553
```

```
## rideDistance swimDistance teamKills walkDistance weaponsAcquired
```

```
## 0 515.6142 5.365372 0.016908394 823.7333 3.517848
```

```
## 1 1666.6767 16.136750 0.003005937 2351.6425 5.710528
```

```

predict <- predict(qda,newdata=test)

test_final <- test
test_final$predcited_place <- as.vector(predict$class)
test_final$top.10<-as.factor(test_final$top.10)
test_final$predcited_place<-as.factor(test_final$predcited_place)
# test_final$predcited_place<-as.factor(test_final$predcited_place)
str(test_final)

```

```

## 'data.frame':    54582 obs. of  31 variables:
##  $ Id           : chr  "5fd62798396ca8" "d08ce24e7a7973" "50820adcefc866" "029a693a75ef57" ...
##  $ groupId      : chr  "bb19a05801d30d" "d57ed9de010a4e" "212fc68772a25a" "11d78ae0ca90b6" ...
##  $ matchId      : chr  "9e3c46f8acde82" "1eda9747e31f1f" "c9259d5a2a4f19" "bfccecf5202cc8" ...
##  $ assists      : int   0 0 0 1 0 0 0 0 0 0 ...
##  $ boosts       : int   0 0 9 1 2 1 6 1 0 1 ...
##  $ damageDealt   : num   36 0 167.4 568.7 29.1 ...
##  $ DBNOs        : int   0 0 0 0 0 0 0 0 0 0 ...
##  $ headshotKills : int   0 0 1 1 0 0 0 0 0 0 ...
##  $ heals        : int   0 0 5 0 3 5 3 2 0 0 ...
##  $ killPlace     : int   84 65 12 3 53 60 16 69 52 57 ...
##  $ killPoints    : int   0 0 0 1339 1180 0 1297 977 0 0 ...
##  $ kills        : int   0 0 2 5 0 0 2 0 0 0 ...
##  $ killStreaks   : int   0 0 1 1 0 0 1 0 0 0 ...
##  $ longestKill   : num   0 0 82.6 92.8 0 ...
##  $ matchDuration : int  1999 1471 1801 1964 1381 1892 1439 1976 1436 2190 ...
##  $ matchType     : chr   "solo" "solo" "solo" "solo" ...
##  $ maxPlace      : int   94 99 91 98 95 94 95 91 96 93 ...
##  $ numGroups     : int   92 94 85 92 93 90 93 87 92 91 ...
##  $ rankPoints    : int  1507 1500 1500 -1 -1 1500 -1 0 1492 1507 ...
##  $ revives       : int   0 0 0 0 0 0 0 0 0 0 ...
##  $ rideDistance  : num   0 0 4146 5187 1093 ...
##  $ roadKills     : int   0 0 0 1 0 0 0 0 0 0 ...
##  $ swimDistance  : num   0 0 0 0 0 0 0 0 0 0 ...
##  $ teamKills     : int   0 0 0 0 0 0 0 0 0 0 ...
##  $ vehicleDestroys: int   0 0 0 0 0 0 0 0 0 0 ...
##  $ walkDistance  : num   293 871 3592 827 1135 ...
##  $ weaponsAcquired: int   1 3 7 4 2 5 2 2 3 4 ...
##  $ winPoints     : int   0 0 0 1537 1546 0 1596 1485 0 0 ...
##  $ winPlacePerc  : num   0.107 0.388 0.922 0.598 0.692 ...
##  $ top.10        : Factor w/ 2 levels "0","1": 1 1 2 1 1 1 1 1 1 1 ...
##  $ predcited_place: Factor w/ 2 levels "0","1": 1 1 2 2 1 1 2 1 1 1 ...

```



```

xtab<-table(test_final$predcited_place,test_final$top.10)
confusionMatrix(xtab, positive="1")

```

```

## Confusion Matrix and Statistics
##
##
##           0      1
##  0 44026  1810
##  1  4853  3893
##
##              Accuracy : 0.8779
##              95% CI : (0.8752, 0.8807)
##      No Information Rate : 0.8955
##      P-Value [Acc > NIR] : 1
##
##              Kappa : 0.4721
##
##  Mcnemar's Test P-Value : <2e-16
##
##      Sensitivity : 0.68262
##      Specificity : 0.90071
##      Pos Pred Value : 0.44512
##      Neg Pred Value : 0.96051
##      Prevalence : 0.10448
##      Detection Rate : 0.07132
##      Detection Prevalence : 0.16024
##      Balanced Accuracy : 0.79167
##
##      'Positive' Class : 1
##

```

```

# confusionMatrix(test_final$predcited_place, test_final$top.10)

```

```

predict.posterior <- as.data.frame(predict$posterior)

```

```

# Evaluate the model

```

```

pred <- prediction(predict.posterior[,2], test$top.10)
roc.perf = performance(pred, measure = "tpr", x.measure = "fpr")
auc.train <- performance(pred, measure = "auc")
auc.train <- auc.train@y.values

```

```

# Plot

```

```
plot(roc.perf, main="ROC Curve - QDA")
abline(a=0, b= 1)
text(x = .25, y = .65 ,paste("AUC = ", round(auc.train[[1]],3), sep = ""))
```

