# Predicting the Win Rate and Duration of Game LOL in Real Time
## Based on Classification and Linear Regression

Xiaohan Xu,  Yujie Ren, Directed by Dr. Zelenberg
Beijing Institute of Technology

ILLINOIS INSTITUTE
OF TECHNOLOGY

## Introduction

▪League of Legends (abbreviated LoL) is a multiplayer online battle arena video game. The goal is usually to destroy the opposing team's "nexus", a structure that lies at the heart of a base protected by defensive structures. There are many factors can influence the win rate and duration of game, such as tower, difference of gold, dragon ,inhibitor, etc. And the win rate and duration are what people are always concerned about.
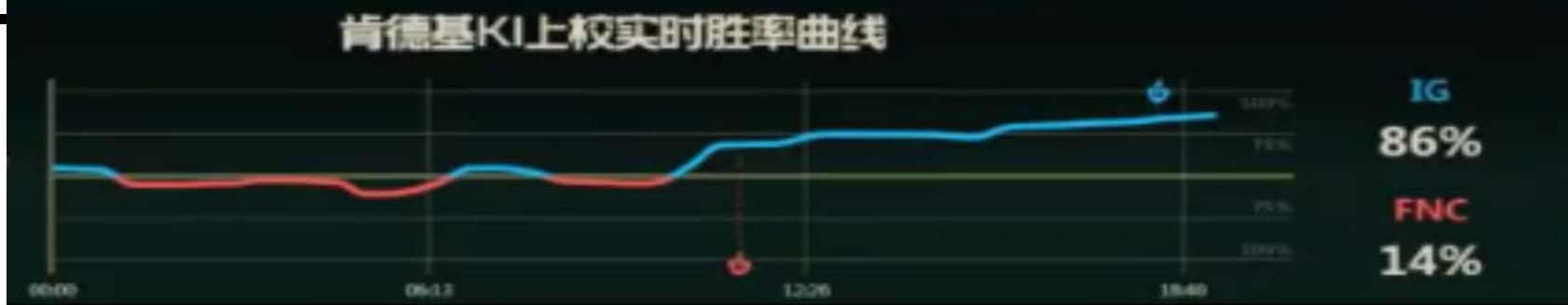
▪In our project, we will predict win rate and duration of game in real time by some predictors in the game based on the dataset "(LoL) League of Legends Ranked Games". And the techniques are classification and linear regression respectively. At the same time, we get our optimization model by different method of fitting, forward stepwise selection, lasso etc.

❑ Technique: Linear regression ,Logistic regression, LDA, KNN.

❑ Optimizing method: Forward stepwise selection, Lasso

## Theory and formula of linear regression

Multiple Linear Regression Model is used to regress and predict each match's game duration.

• **Linear Model:**
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon,$$
(Where $\beta_i$ means coefficient, X means predictors, and $\varepsilon$ means the error.)

• **Estimating coefficient:**
$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_p x_p.$$

• **Figuring out whether X affects Y:**
do the null hypothesis $H_0 : \beta_1 = 0$
calculate the t-value and p-value:
$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)},$$

• **F-statistic:**
null hypothesis: $H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$
$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)}.$$

• **Assessing the Accuracy of the Model:**
We define:
$$RSS = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2.$$
$$TSS = \sum (y_i - \bar{y})^2$$
$$RSE = \sqrt{\frac{1}{n-p-1}RSS} = \sqrt{\frac{1}{n-p-1}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}$$
$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$

RSE is used to evaluate the deviation of our model;
$R^2$ is used to evaluate our model's fitness to the data.

## Theory about Classification

• **Logistic Regression**
• Getting estimate β by get maximum of likelihood function
$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}.$$
$$\ell(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i':y_{i'}=0} (1 - p(x_{i'})).$$

• **LDA(Linear Discriminant Analysis)**
• Bayes' theorem
$$\Pr(Y = k | X = x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^{K} \pi_l f_l(x)}$$
• Compare discriminant function
$$\delta_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k)$$

▪**KNN(K-Nearest Neighbors)**
• By Calculating the number of the nearest points of every class respectively.
$$\Pr(Y = j | X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} I(y_i = j)$$

▪**Forward Stepwise Selection**
• Aim: reduce the number of predictors

▪**The Lasso**
• Aim: reduce the variance.
• By minimizing the quality of right equation.
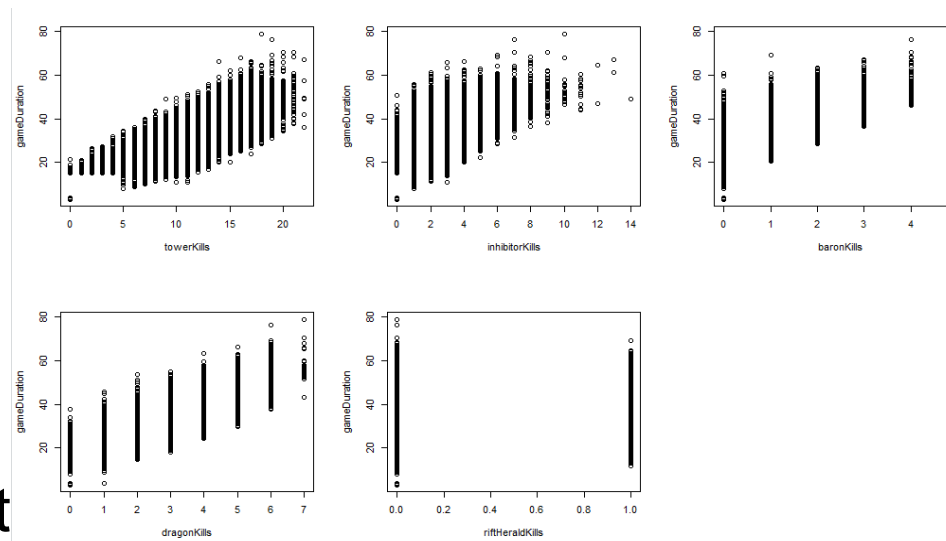$$\sum_{i=1}^{n}\left(y_i - \beta_0 - \sum_{j=1}^{p}\beta_j x_{ij}\right)^2 + \lambda\sum_{j=1}^{p}|\beta_j| = RSS + \lambda\sum_{j=1}^{p}|\beta_j|.$$

## Implementing Linear Regression Model to regress and predict the game duration

**1. Process the Dataset**
There are several works to do to deal with the dataset:
① Delete the useless or missing observations.
② Select predictors roughly.
③ Change yes-no question into 0-1 value
④ Combine two-teams' variables into one.
⑤ Change Response's unit from seconds to minut

**2. Select predictors using stepwise**
There are 5 variables left, and we will select the best ones to predict, using the method stepwise. (In R, stepwise can be done by a function called "stepAIC()", which is from MASS package.)
The result of this step is on the right side. We can conclude that all of the variables are chosen.
**3. Optimize this model using non-linear relationships.**

It is possible that these variables has more complicated relationships with the response, so we try polynomial relationship here, and evaluate the new model using the method called "k-fold cross validation".

There are four figures on the right side that indicate the relationship between mean test error and highest power of the polynomial of each variable

When choosing the power with smallest mean test error, we can get a better model:

gameDuration ~ poly(towerKills, 13) + poly(inhibitorKills, 9) + poly(baronKills, 4) + poly(dragonKills, 2) + riftHeraldKills

## Conclusion and Model Evaluation

**1.Conclusion:**
We can get quantitative relationships between predictors and the response, here are some of them below:
• With the number of broken towers increasing, the game would probably last longer.
• With the number of baron and dragon killed increasing, the game will absolutely last longer. This makes sense, because baron and dragon need time to rebirth.
• With the number of Rift Herald killed increasing, the game will last shorter.
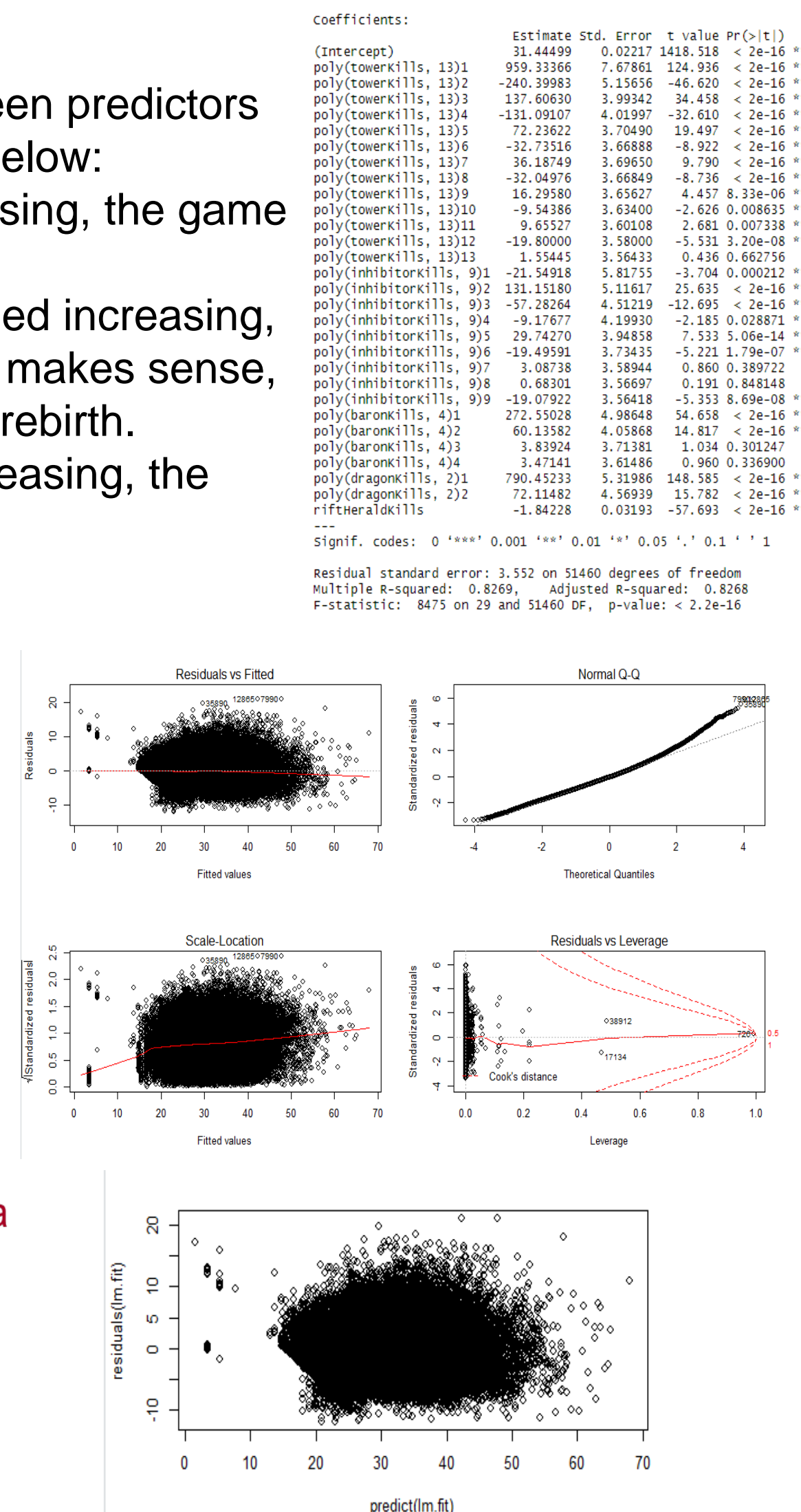
**2.Model Evaluation:**
• **Accuracy:**
① In this model, adjusted $R^2$ = 0.8268 which means this model fits the dataset very well;
② RSE = 3.552, which means the error of this model is 3.552 minutes, regarding that the mean of the duration of the games is approximately 31.22minutes, this error seems acceptable.
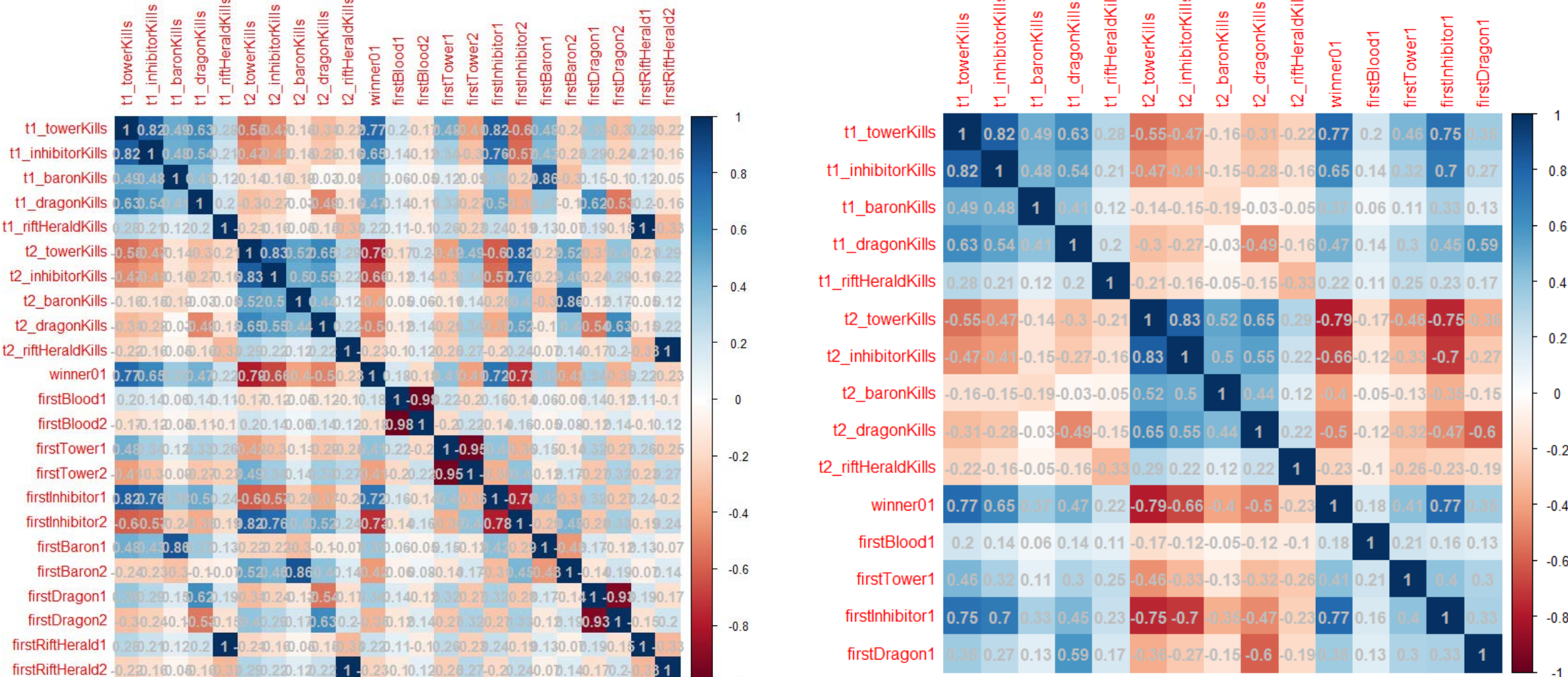③ F-statistic = 8475
• **Evaluation from the plots:**
① The "Residual vs fitted" indicates that there are some outliers, and the current model fits well.
② The Q-Q plot indicates that the skew of our data is not so strong.
③ In the "Scale Location" plot the read line is pretty flattish, which means data spreading out.
④ The leverage plot suggests that there is an extremely high leverage point in the data.
⑤ The last one indicate that predicting residual is also normal.

## Predicting the Win Rate by Classification

**1. Judge Correlation and Remove Factor with High  Correlation.**



(FIGURE 1)                    (FIGURE 2)

• Use corrplot function to visualize the correlation and remove 8 predictors with large correlation.
• The Figure 1 is about original correlation and the Figure 2 is the ultimate correlation of left predictors.

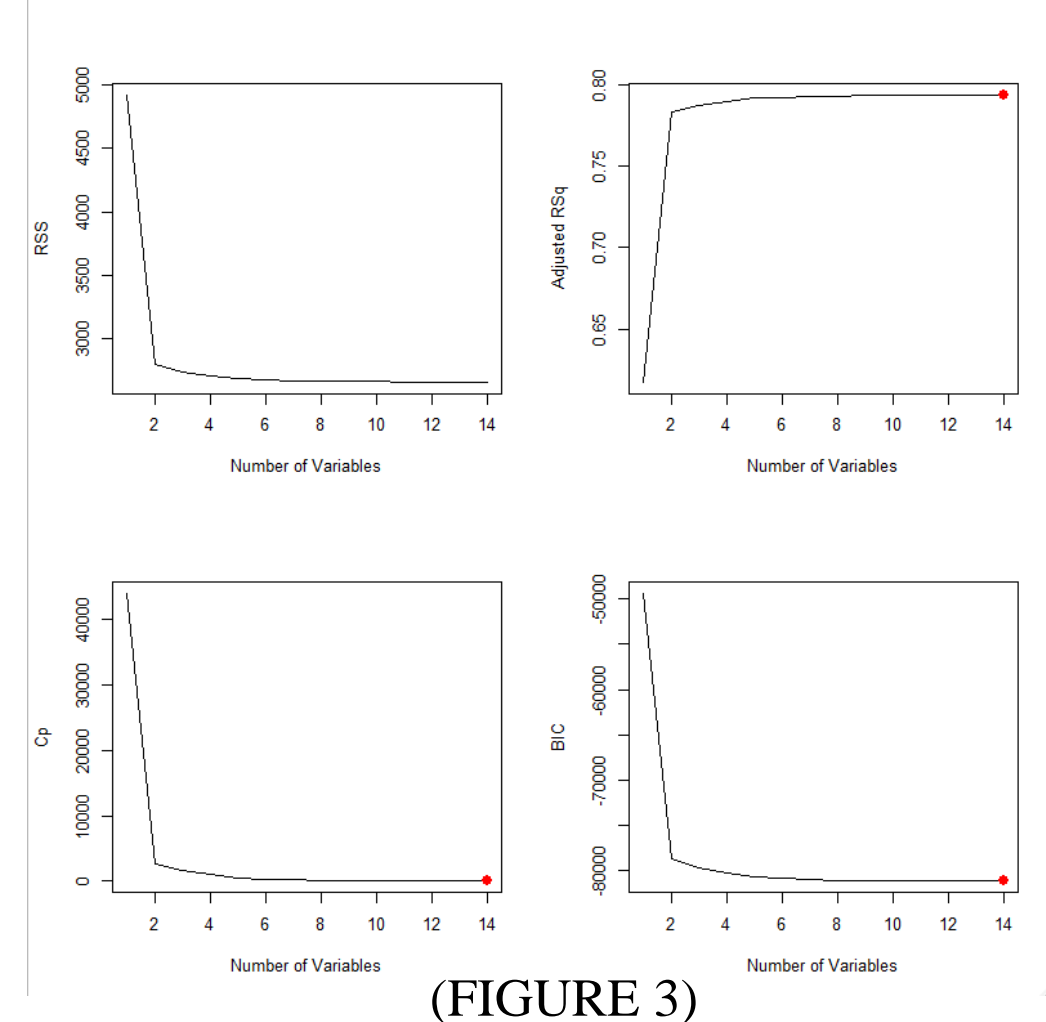**2. Classifying by Logistic Regression, LDA and KNN.**

| | Confusion matrix | | Test error |
|---|---|---|---|
| Logistic Regression | 0    1 <br> 0   7328   189 <br> 1   370   7560 | | 3.62% |
| LDA | 0    1 <br> 0   7348   245 <br> 1   350   7504 | | 3.85% |
| KNN(K = 10) | 0    1 <br> 0   7471   321 <br> 1   227   7429 | | 3.34% |

• From the result we can get the conclusion that the KNN model is the best, but their difference is so small. We will optimize the model further at the next step.

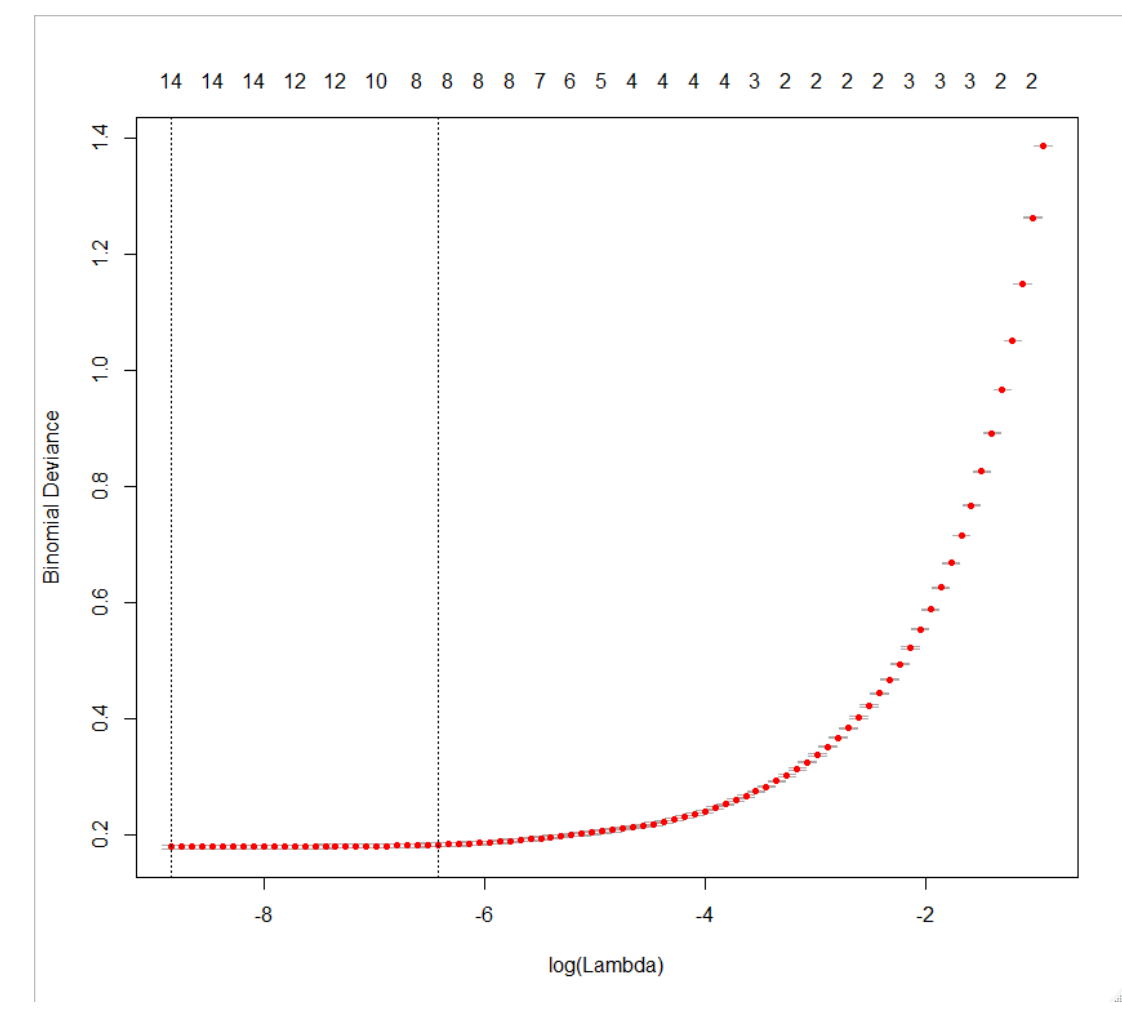**3. Optimize the model by forward selection and the lasso**

• **Forward Selection:**
• get the most appopriate number of predictors by the plots.
• The result is N = 14 is the best situation

• **The Lasso:** enhance model' interpretability.
• When λ = 0.0001441945, the model is best.



(FIGURE 3)
(R2, RSS, adjusted R2, Cp, and BIC change with different number of predictors

(FIGURE 4)
(Binomial deviance versus log(Lambda))

• Under this condition , all coefficients of predictors are 0 except 4  predictors, which make model more interpretable.

**4. Analysis of ultimate model**
• The final model only correlate to 4 predictors and the test error of  logistic regression with these 4 predictors is 4.12% which is a little bigger than original model's , but this model has a stronger interpretation and we can predict the result just by 4 predictors. So the final model is much better.
• The four predictors are the destroyed tower number of team1, the killing baron number of team 1, destroyed tower number of team2,  the killing baron number of team 2. Their coefficients are 0.68, 0.2523, -1.623 and -0.7029.