
DECEMBER, 2021

THE LEAGUE OF LEGENDS RANKED MATCHES ANALYSIS REPORT

MGSC661 Final Project

Anqi Chen 261044081

Yulin Hong 260898713



Desautels Faculty of Management

McGill University

Montreal, Quebec, Canada

Contents

1. Introduction	1
2. Data Preprocessing and descriptive analysis	2
3. Model Selection & Results	3
3.1 Classification model	3
3.1.1 Principal component analysis	3
3.1.2 Model selection	3
3.1.3 Results and interpretations	4
3.2 Clustering	5
4. Managerial conclusion	7
5. Appendix	13

1. Introduction

The League of Legends (LOL) is a multiplayer online battle arena game developed by Riot in 2009. In each match, 10 players are divided into two teams, red and blue. In its most played mode, the Summoner's Rift, each player needs to control a character known as champion to defend their team's base and invade the other half of the map. A team wins by destroying the Nexus in front of the enemy team's base.

League of Legends has been the most played game in the world since 2012. Each year it will hold a ranked season where players play with or against others that share a similar level of gaming abilities. The annually ranked seasons first started in July 2010 and right now LOL is approaching the end of season 11. During a ranked season, the points will be added or deducted if a player succeeds or loses in one rank game. The final points will decide the annual ranking of a player. Currently, the League of Legends ranking system has 9 tires and 4 divisions within each tire. However, before season 9 there were only 7 tires and 5 divisions in each tire.

In this project, we used the League of Legend Ranked Matches dataset to explore the hidden information behind ranked games. Using classification models, we predicted the outcome of a game using game statistics such as the kill death assist ratio. We also conducted cluster analysis on the matches. By mapping the clustering result to the tires based on clusters' characteristics and size, we then analyzed the similarity and differences among tires.

Since the origin dataset comprised 184,070 ranked solo games across various seasons and platforms, we restrict our scope of work to analyzing North American platform Season 8 matches. However, the idea behind this project can be applied to the whole dataset.

2. Data Preprocessing and Descriptive Analysis

The original dataset comprised of 7 tables, but the datasets were manipulated into two final matches statistics tables. After extracting the match games data for games in the North American area and Season 8, we then combine the tables of matches, participants, general statistics, and team statistics by matching each player in the participants and general statistics tables using their match ID in matches and team statistics tables. Note that our project does not analyze the performance of individual players, but the performance of the whole team (comprised of 5 players) in a match. Hence, after extracting all the data, we then group the statistics by match ID and team ID (an indicator of the blue or red team) and calculate the average value of each continuous statistics variable within a team of each match. Binary variables were not included in the process of averaging because the binary variable (such as first blood) are already variables that indicate the team status.

The dataset obtained from the previous step still contains over 50 predictors, but most of them are irrelevant such as items or some of them are interdependent, such as total damage dealt and magical damage dealt. Therefore, we removed all the irrelevant predictors as well as the predictors that depend on others. The final dataset consists of 23 predictors.

Boxplots and histograms were used to visualize the distribution of these predictors (Figure 1). The predictor KDA represents the kill death assist ratio, and most of the teams have a KDA ratio of around 2. This means the number of kills or assists is twice as many as death. Then, total damage dealt, total damage to champion and total damage taken are bell-shaped continuous predictors with a little outlier. The gold earned and gold spent both have a mean of around \$12000, meaning players usually spend all the gold they earn. Next, the champion level is right-skewed, meaning most of the teams have a high champion level. In contrast, inhibitor kill, baron kill, dragon kill, are left-skewed, meaning most of the teams killed fewer inhibitors, barons, and dragons. Lastly, the duration is also bell-shaped, and the average matches lasted for 2000 seconds (about 33 and a half minutes).

A correlation matrix was created to examine the multicollinearity issue (Figure 2). Duration, champion level, gold spent, gold earned, the total damage to champions, and total damage dealt tend to have a relatively high correlation among each other. On the other hand, harry kill is uncorrelated to any other predictors or with the target variable. Therefore, in the classification model, gold earned, gold spent, duration, and harry kills are not considered.

3. Model Selection & Results

3.1. Classification model

3.1.1. Principal Component Analysis

Principal component analysis (PCA) is used to find a low-dimensional representation of the data that captures maximum variability. Since there are 22 variables excluding the match ID and team ID and there are 3432 observations, K, in this case, will be 22.

From the principal component graph, two main directions are orthogonal to each other, where predictors are relatively correlated to each other in the two directions (Figure 3). For example, gold spent and champion level are highly correlated, and first blood and first dragon are highly correlated. In addition, the target variable win belongs to one of the directions, meaning predictors in the other direction might not be particularly useful to explain the variations.

Next, the percentage of variation in each component is calculated (Figure 4). With approximately 5 components, 80% accuracy can be recognized. In the first principal component, total damage dealt, the total damage to champions, gold earned, gold spent, champion level and duration are the main contributors. And in the second principal component, win, total damage taken, first inhibitor, tower kills, inhibitor kills are the main contributors. Similarly for other components, the greater the absolute value, the more it contributes to the component.

Overall, the result of the principal component analysis will be used as a tool to find the final list of predictors for the classification model. Overall, the predictors that will be used in the model are: KDA, the total damage to champions, total damage taken, champion level, tower kills, inhibitor kills, duration, first inhibitor, and baron kills.

3.1.2. Model selection

Both classification techniques and tree-based methods are used in the model selection process (Table 1). For the classification techniques, logistic regression and discriminant analysis are constructed. And for the tree-based methods, decision tree, random forest, and boosted forest are built. Among all the models that were built using the same set of predictors, the accuracy score of the boosted forest is the highest, which is 96.9%.

In the boosted forest model, the parameter of distribution is set to be Bernoulli because the target variable is binary. The model will develop 10000 trees, and the number of internal nodes for each tree is set to be 8. The means square error of the model is approximately 4858.

To test for accuracy, the dataset was first split into training and testing sets, with a ratio of 0.75 being the training set and the remaining being the testing set. Then, the boosted forest model was built with the same parameters as mentioned above. The model next predicts the testing set, and the value will round up to 1 if the result is greater than 0.5, and round down to 0 if the result is less or equal to 0.5. Lastly, the accuracy score is calculated by comparing the predicted result to the actual result. An accuracy score of 0.969 represents the boosted forest model that can correctly predict the outcome of the match in a 96.9% chance.

3.3.3 Results and interpretations

From the summary of the boosted forest model, the relative importance is shown. Tower kills has the highest score, which is approximately 39.7, and inhibitor kill is the second most important predictor with a score around 33.3. The remaining predictors have a smaller score compared to these two predictors.

To compare the current model with a model that includes all the predictors in the initial dataset, a second model with all available predictors is built. Nevertheless, this new model does not include identifiers such as match ID. The model shows that inhibitor kills and tower kills remain the top two most important predictors, while the rest of the predictors have a much lower importance score. The accuracy score of this new model is slightly lower, which is approximately 95.4%. However, the mean square error (MSE) is around 4608, which is 250 less than the MSE of the original model. Overall, the original model would perform better.

The boosted forest utilizes the boosting mechanism to improve the generalizability of the trees, where the trees are trained sequentially. Hence, the model performs the best among all the remaining models. The model could potentially help players predict the outcome using available information before the game and improve equipment or techniques if needed.

3.2 Clustering

The League of Legends ranking system assigns players to components that share the same level of gaming ability with them. This means that a match's participants should have around a similar rank. Through clustering models, we aimed to segment matches into different clusters, mapped each cluster to the possible tire, and examined the characteristic of each cluster and the tire it represents. However, most of the predictors in the dataset are highly dependent on the predictor duration. For example, if the duration of a match is long, the total damage dealt will accumulate as the match continues. Hence, we need to "normalize" the predictors to obtain a better clustering result. To normalize these predictors, we simply divide them by the duration so that all the match statistics become per second statistics. Moreover, because each match statistic relates to two teams' statistics and these two teams' statistics are complementary to each other's, we only keep all the winning teams' statistics to feed into the clustering model.

Before we started to apply the clustering model, we first examined the outliers of each predictor by checking the data points that lie outside of the 0.01% to 99.9% quantile of each predictor. Since the K-Means clustering method clusters data based on the distance between each data point, outliers in the dataset can affect the clustering result. Hence, removing the outliers can greatly improve the quality of clustering. We chose to remove data points that are outliers for two or more predictors.

After preprocessing the dataset, we used Silhouette Score and Elbow Method to find the optimal number of clusters for the K-Means method. The results both indicate that the optimal cluster number is around 4 (Figure 5 and Figure 6). After setting the K value and running the K-Means function, we obtained the clustering result and mean values of variables in each cluster.

There are in total 1712 observations that were fed into the K-Means model. The resulting 4 clusters have the size 28,797,162 and 725 respectfully. The first cluster contains 1.63% of the total matches. The average game statistics in the first cluster are the worst. The total damage and total damage to champion statistics are very low, meaning that these players did not land their skills on the enemy's champions at all. Moreover, they are not interested in push towers or killing dragons, which are two curtail factors that lead to the success of one match. Hence, the corresponding tire of this cluster should be the lowest rank in the League, Bronze V, in which the actual percentage in Season 8 is 1.70% (Table 2). The fourth cluster has the second-lowest performance. The amount of damage they dealt as well as gold earned per second is lower than

Cluster 3 and Cluster 2. This cluster includes around 42.35% of rank matches and it matches the distribution of Bronze IV to Silver I in Seasons 8.

Cluster 2 and Cluster 3 did not vary much by game statistics. In general, Cluster 3 has higher average statistics in most variables, but Cluster 2 has higher KDA, tower Kills, and Baron Kills. However, we still consider Cluster 3 to be the highest rank cluster because one of the crucial differences between entry-level players and skillful players is the ability to accumulate gold. In fact, in the League of Legends World Championship gold earn is one of the key elements when we measure the performance of a team. Given that Cluster 3 has the highest gold earn per second, the players in this Cluster are top-ranked players. By mapping the cluster size to the Season 8 rank distribution again, we concluded that players in this cluster are from Diamond V to Challengers. Finally for Cluster 2, the players are from Gold V to Diamond V.

4. Managerial Conclusion

Overall, League of Legends is the most popular game in the world since 2012, and it held the annual World Championship hosted by Riot Games. The goals of this analysis are: classifying the League of Legends rank matches' outcomes using game statistics related to the matches, and clustering teams to examine their similarities and differences. Throughout the analysis, data collection, preprocessing, descriptive analysis, model selection, and result analysis were conducted.

For the classification task, nine predictors were used in the classification models, and five models were built. The best model with the highest accuracy score is the boosted forest model. Tower kill and inhibitor kill are the most important predictors, meaning the number of towers and inhibitors that a team killed are significant predictors to forecast the outcome of the match for that team. This result explains the main difference between League of Legends with other games such as PUBG. The League of Legends is not a game that individual players' ability can dominate the outcome of a match. Conversely, a team needs to work together to take down enemies' towers and inhibitors to win a game.

The clustering result shares the same information as discussed above, but it also provides more insights on how high-level players differ from entry-level players. They might not be the group with the highest KDA, but they focus more on the farm. After all, farming is the most stable way to accumulate gold and finish your final build as soon as possible. Moreover, they did not have the highest number of tower kills but they have the highest number of inhibitor kills. This might be counter-intuitive because every inhibitor is guarded by a tower. However, if you destroy an inhibitor it will respawn after 5 minutes and every time you destroy it, you can send 8 waves of super minions, the strongest minions of all types, to help you push and fight in the lane. Thus, when players are considering winning a game as efficiently as possible, they should target more on inhibitors and utilize the super minion waves to help them secure success.

Overall, these two models serve as great tools to analyze the League of Legends matches. In terms of potential improvement, the classification models can incorporate additional predictors such as type of the dragons killed in a game if the data is available. By including more detailed data, this model can be considered a powerful tool for coaches of the professional gaming teams to analyze the current situation of the teams during important circumstances such as the World Championship. The coaches can input their strategies into the model, for example, take the fire dragon instead of baron, and simulate the win rate based on these strategies. By doing so, coaches can communicate with players in time and help them to seize more opportunities to win.

5. Appendix

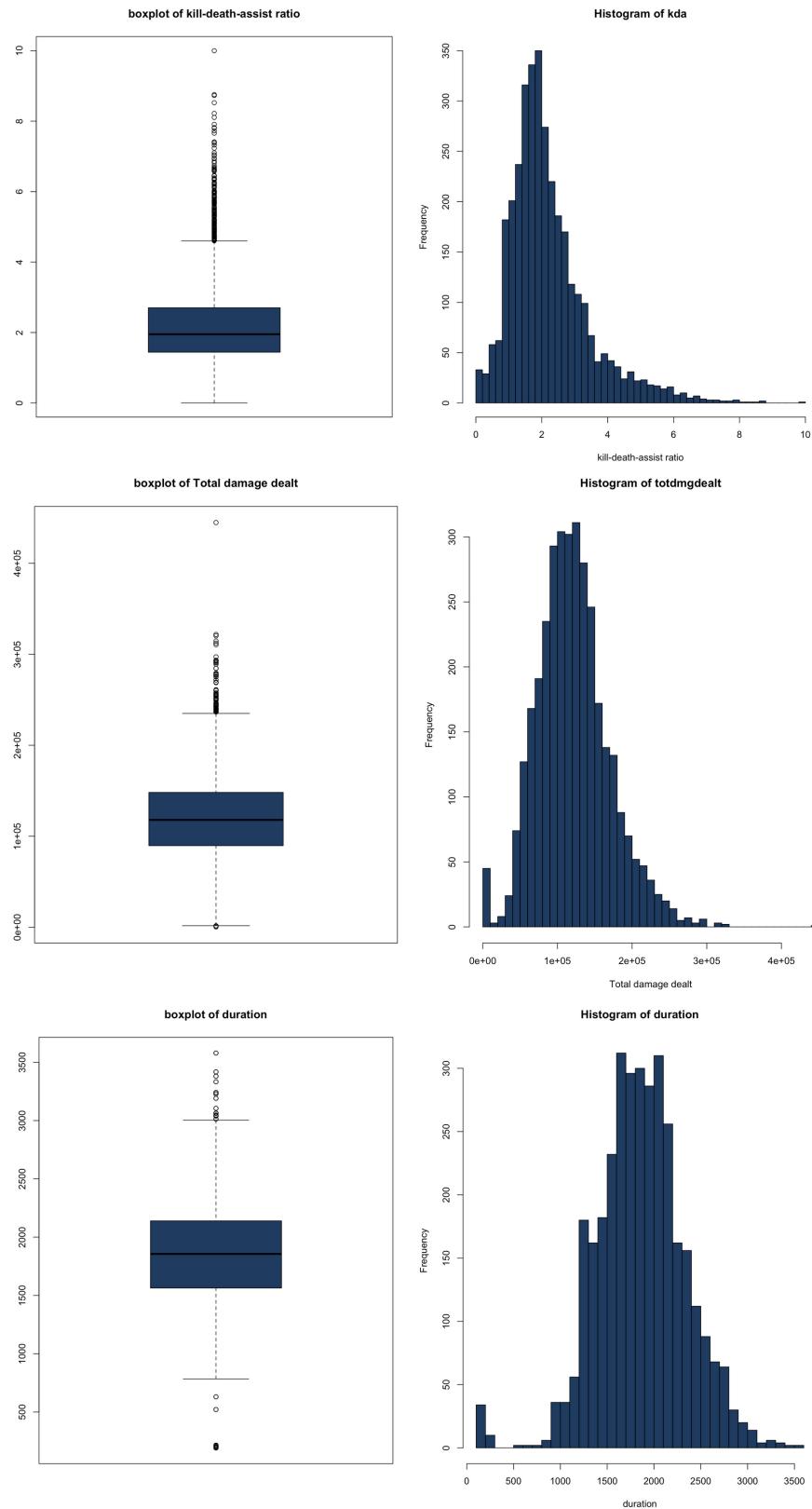


Figure 1. Boxplots and histograms for kda, total damage dealt, and duration

THE LEAGUE OF LEGENDS RANKED MATCHES ANALYSIS

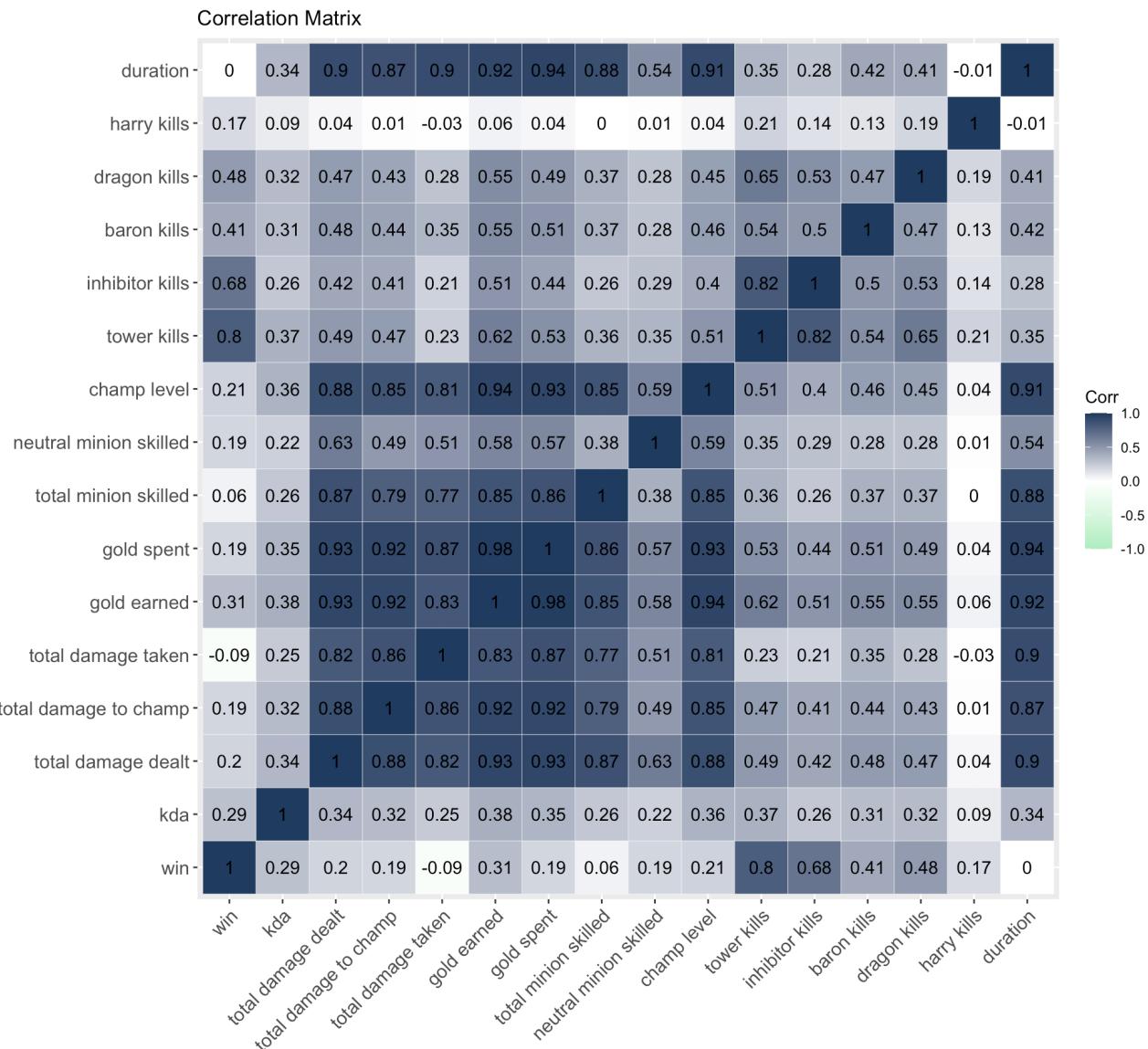


Figure 2. Correlation Matrix of numerical predictors

THE LEAGUE OF LEGENDS RANKED MATCHES ANALYSIS



Figure 3. Principal Component Analysis Plot

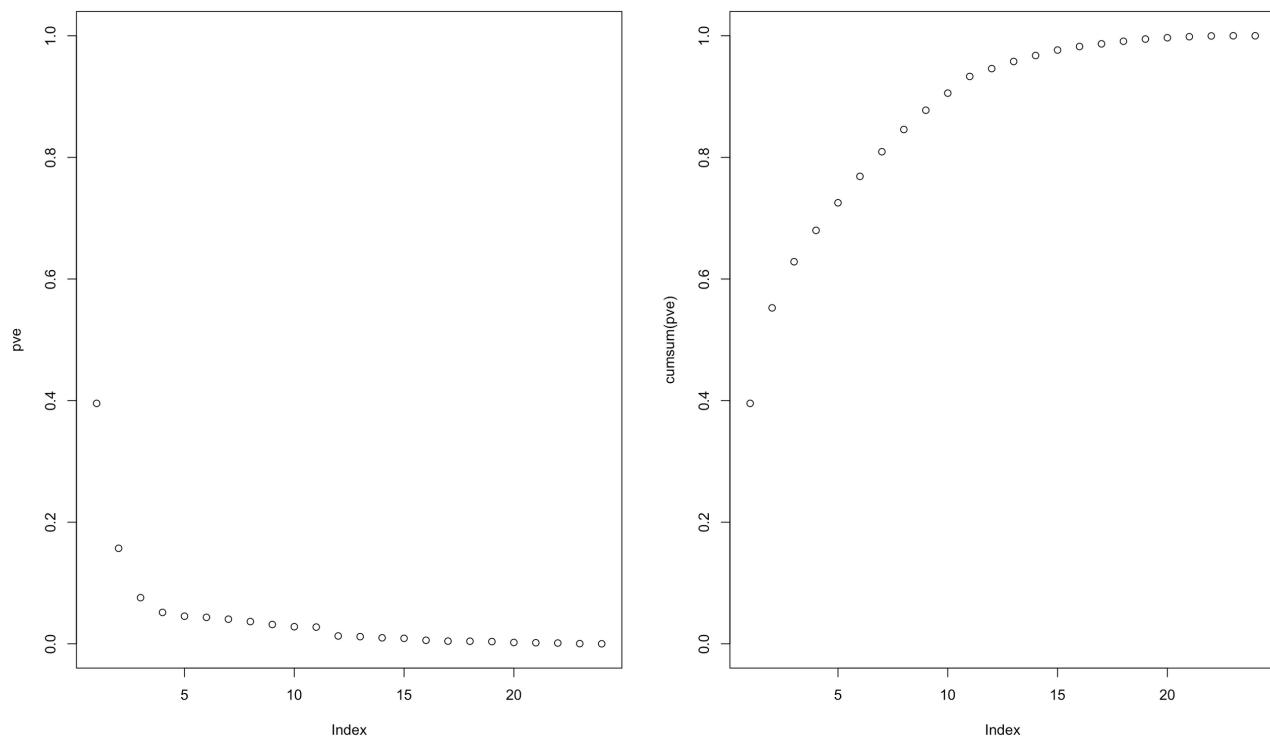


Figure 4. Percentage of Variance explained plot

	1	2	3	4	5
Model	Logistic regression	LDA	Decision Tree	Random Forest	Boosted Forest
Accuracy Score	0.948	0.943	0.909	0.950	0.969

Table 1. Summary of classification models and corresponding accuracy scores

*LDA: Linear Discriminant analysis

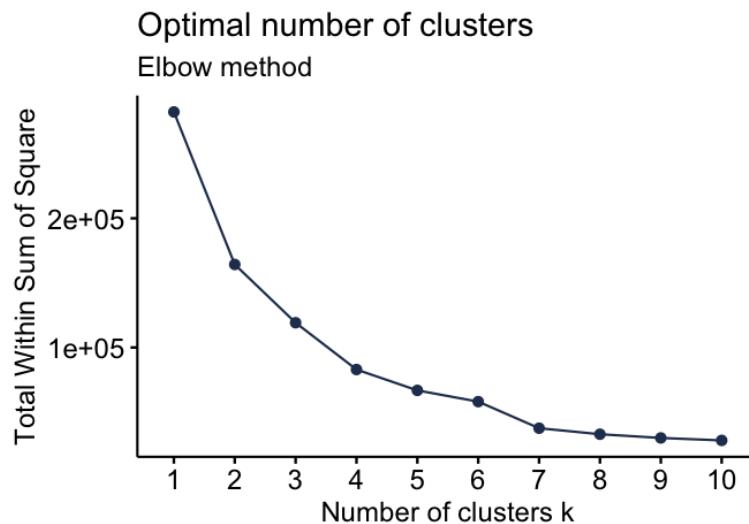


Figure 5. Finding Optimal Cluster Numbers, Elbow Method

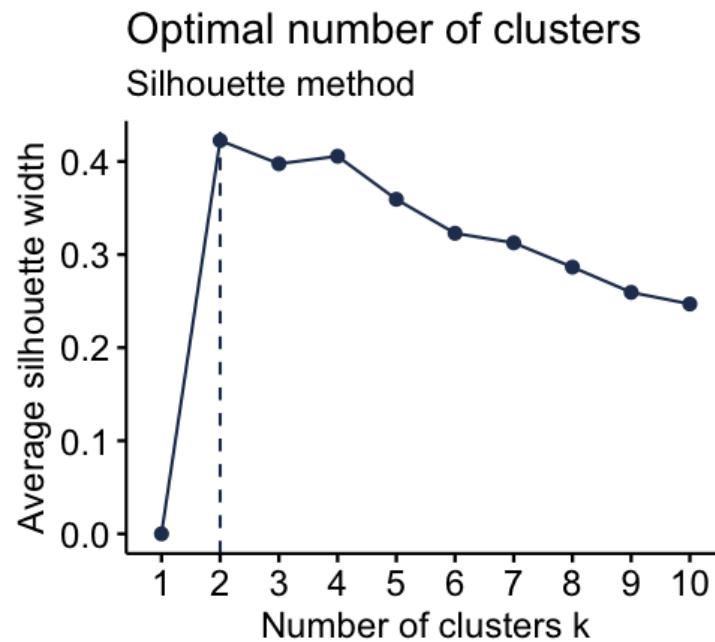


Figure 6. Finding Optimal Cluster Numbers, Silhouette Method

THE LEAGUE OF LEGENDS RANKED MATCHES ANALYSIS

Rank	Percentage	Rank	Percentage	Rank	Percentage
Bronze V	1.70	Gold V	14.94	Diamond V	1.67
Bronze IV	1.68	Gold IV	5.94	Diamond IV	0.35
Bronze III	2.27	Gold III	4.10	Diamond III	0.19
Bronze II	2.68	Gold II	2.59	Diamond II	0.10
Bronze I	3.26	Gold I	3.90	Diamond I	0.05
Silver V	10.77	Platinum V	4.48	Master	0.05
Silver IV	10.49	Platinum IV	1.70	Challenger	1.67
Silver III	10.34	Platinum III	1.45		
Silver II	8.65	Platinum II	1.12		
Silver I	5.03	Platinum I	0.47		

Table 2. Season 8 Rank Distribution

Source: <https://www.esportstales.com/league-of-legends/rank-distribution-season-8>