STATS 101C Final Report
Prediction of NBA Game Outcomes

Group Members: Fanchu Xu (UID: 405800796), Katherine Jin (505650329), Qianping Wu (906230284)

## 1 Introduction

Predicting game outcomes has become an important part of sports analytics, providing valuable insights for teams, coaches and fans. This project uses historical game data from the 2023-2024 NBA season to predict whether a team will win or lose.

Feature engineering is essential for transforming raw game statistics into meaningful inputs. Thus, key features like home court advantage, team performance consistency, previous matchup results between teams and average scores before the game are created to improve the accuracy of predictions. To further refine the model, variable selection techniques like lasso regression and random forest permutation are implemented to select the most significant variables.

A variety of modeling approaches, including logistic regression, QDA, decision tree with XGBoost, ridge classifier and Support Vector Machine will be tested to determine the best approach for predicting game outcomes. Train-test split will also be applied to ensure the models perform well on unseen data.

By combining feature engineering, variable selection, and a range of prediction models, we hope to achieve better accuracy than the baseline of 65% and provide deeper insights into the factors driving NBA game outcomes.

## 2 Descriptive Analysis

The data set consists of 2460 matches in total, with 82 games for each of the 30 teams. The summary of each team's winning rate is presented below.

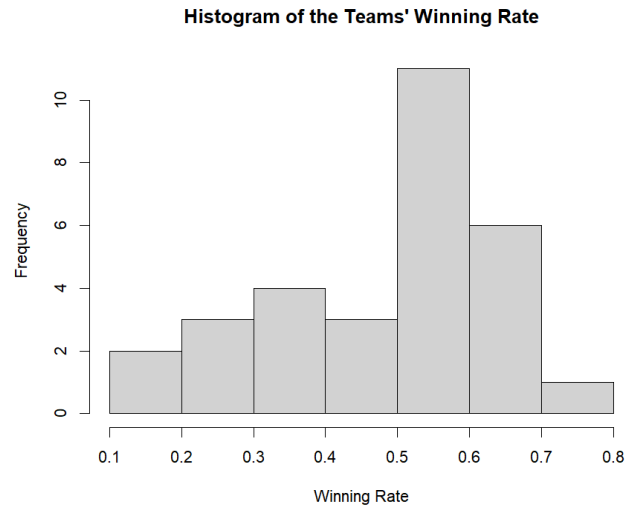**Histogram of the Teams' Winning Rate**

Figure 1: The Distribution of Winning Rate

The distribution of the team's average winning rate is slightly right-skewed, with the most frequent winning rate centered around 0.5-0.6.
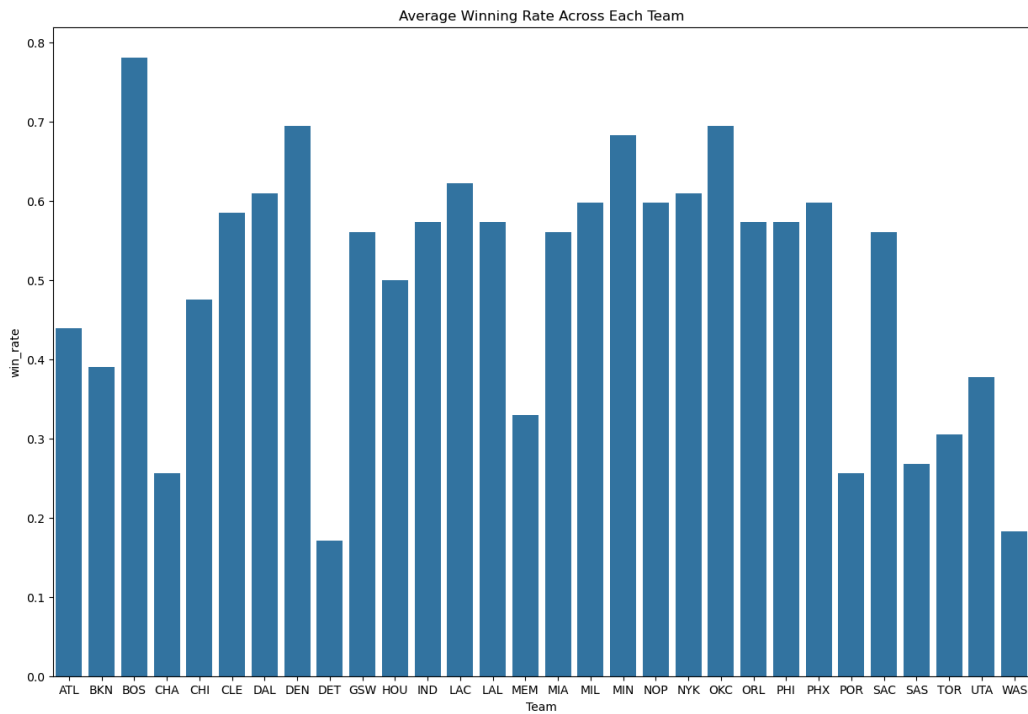
Figure 2: Winning Rate of Each Team

Taking a closer inspection, Team BOS has the highest winning rate of 0.7804878, while DET has the lowest winning rate of 0.1707317. Our task is to predict the win/loss outcome of each competition, given the historical data.

## 3      Feature Engineering

Considering that there exist variables that may affect the outcomes but weren't included in the dataset, we started data processing by adding the following variables.

### 3.1.    Home Advantage

As teams with home advantage are likely to perform better, we took this factor into consideration by creating a dummy variable that equals to 1 if the team is the home team in the competition, and equals to 0 otherwise. This variable will be denoted as "Home advantage."

### 3.2     Historical variance of the team's PTS

Historical variance of team's PTS captures the variability of team's scoring performance over time. Specifically, it calculates the variance from the first recorded game up to the current game in the dataset, and this feature is only available for games occurring on and after November 1st. Historical variance of team's points might have potential connection or relation to winning or losing the next game. This variable will be denoted as "variance."

### 3.3     Previous competition record between the two teams

Basically, it shows the historical outcomes of games between two teams before the current match. For example, if Team A and Team B have played 3 games before, and Team A won all of them, the previous record would indicate that Team A has a better chance of winning again. However, it's still possible for Team B to beat Team A, despite Team A's winning streak. To keep track of this, we start by setting the value to zero for both teams. Every time they play a game, we update the record: we add 1 to the winning team and leave the value as 0 for the losing team. This way, we get a simple and clear view of how well each team has performed directly against the other over time. It makes comparing the competition between the two teams more straightforward. This variable will be denoted as "previous record."

### 3.4     Historical average score

For each match between team A and team B, we computed the historical average PTS in every match between team A and team B prior to the current date. If the two teams never had a competition with each other before, we computed the mean PTS of all other matches as an approximation.

Building on the historical averages we calculated, we analyzed the PTS differences between teams across 2,460 matches. On average, the difference in scores between teams before each match was about 114.64, with most results clustering closely around this number. The range of score differences, from a low of 93.67 to a high of 143.00, shows some variability, but overall, the data is quite consistent. This consistency suggests that past performance data is a strong indicator of how teams might perform in future matches.

Looking deeper, we found that half of the matches had score differences between 111.81 and 117.29, which shows that most games fall within a predictable range. This tight clustering around the average means the model we built based on historical data does a solid job of capturing team performance patterns.

What this tells us is that historical averages are not just numbers, and they provide real insight into team dynamics. These patterns could be incredibly useful for coaches, analysts, or even fans looking to understand and predict game outcomes. As we move forward, adding more layers to the analysis, like player stats or game-day conditions, could make the predictions even more accurate and insightful. Therefore, we included the historical average score into our dataset. This variable will be denoted as "average score."

## 4        Data Pre-processing

In the prediction task, we are only given the historical record of the two teams. Therefore, we constructed a new data set by computing the average difference between the two teams in each variable before each competition, and made variable selection and predictions based on the new dataset. As early matches have limited previous records, we included only the matches on or after November 1st in our model, which gave us a dataset of 2352 observations.
Besides, we cleaned the data by removing a row with missing value, converting the outcome variable into binary values (1 for win, 0 for loss), and eliminating irrelevant categorical columns like identifiers and dates.

## 5        Variable Selection

In variable selection, we used two methods: lasso classification and random forest variable importance ranking.

### 5.1     Lasso classification

We took the competition result as the response variable, and the other numerical variables as the predictors. From the result of lasso, we concluded that MIN (The total time of the game, measured in minutes), PTS (Points scored), FGM (Field Goals Made), FGA (Field Goals Attempted), FG% (Field Goal Percentage), 3PA (Three-Point Attempts), 3P% (Three-Point Percentage), FTM (Free Throws Made), FT% (Free Throw Percentage), OREB (Offensive Rebounds), REB (Rebounds), AST (Assists), STL (Steals), BLK (Blocks), TOV (Turnovers), PF (Personal Fouls), +/- (Plus/Minus), variance (historical variance of the team's PTS), average score (historical average score) and previous record (previous competition records) should be included in the model. With these variables, the testing accuracy of lasso is 0.658, which serves as a baseline of our project.

### 5.2     Random forest variable importance ranking

We used random forest to rank the importance of each variable, and the result presented in table 1 shows that the importance of each variable are roughly equal, except for the previous competition records and +/-. This is aligned with the lasso result, as this algorithm also indicated that they are among the most important variables. Therefore, we decided to use the variables selected in the lasso classification.

## Table 1: Ranking of the Variable Importance

| Feature | Score | Feature | Score |
|---------|-------|---------|-------|
| Cumulative Wins | 0.07368 | Average Score | 0.03984 |
| +/- | 0.07580 | FT% | 0.03583 |
| 3P% | 0.05106 | REB | 0.03646 |
| TOV | 0.04395 | Home Advantage | 0.00902 |
| DREB | 0.04424 | BLK | 0.03914 |
| FGA | 0.04176 | PF | 0.04249 |
| Variance | 0.04363 | STL | 0.03700 |
| FG% | 0.04594 | 3PM | 0.03887 |
| PTS | 0.03949 | FGM | 0.03660 |
| AST | 0.03863 | OREB | 0.03410 |
| FTA | 0.04158 | MIN | 0.03528 |
| FTM | 0.03793 | | |
| 3PA | 0.03768 | | |

## 6    Model
To predict the outcome of each competition, we used logistic regression, QDA, decision tree, ridge, and Support Vector Machine (SVM).

### 6.1    logistic regression
Logistic Regression is a statistical method used for binary classification problems. Here, a logistic regression model is employed to predict whether games resulted in a "Win" or a "Loss" based on the selected features.

The logistic regression model was trained using Leave-One-Out Cross-Validation (LOOCV) to ensure robust evaluation. The model achieved an accuracy of 65.31%, with a Kappa value of 0.306, indicating moderate agreement between the predicted and actual outcomes.

To demonstrate the model's performance, a decision boundary was plotted using the features Average Score Before Match (Points) and Positive/Negative with other features taking their average values. The upper region is where the model predicts a "Win" (blue) and the bottom region predicts a "Loss" (red). The decision boundary illustrates how the model separates the two outcomes.
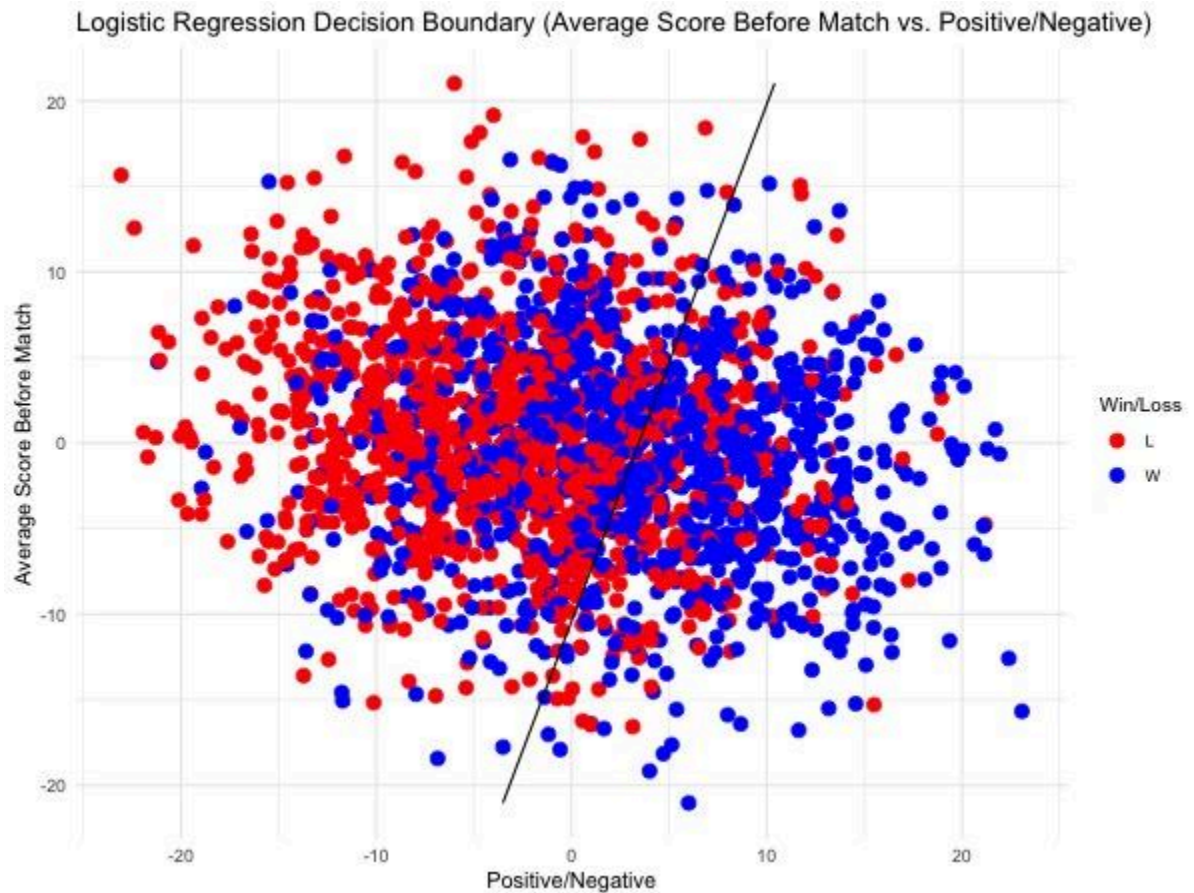


Figure 1: Illustration of Logistic Regression Decision Boundary

6.2    QDA

Quadratic Discriminant Analysis (QDA) is a classification method that assumes each class in the data follows its own Gaussian distribution. Unlike simpler models, QDA allows each class to have unique covariance structures, which helps it handle non-linear relationships more effectively. By creating quadratic decision boundaries, QDA can classify data points based on the probabilities of belonging to different classes.

For this project, a QDA model was built to predict whether basketball games resulted in a "Win" or a "Loss." The data was split into training and testing sets, with 70% used for training the model and 30% for testing its performance. After training, the QDA model was evaluated on the

test set, achieving an accuracy of 0.6671. This accuracy was calculated by comparing the model's predictions with the actual game outcomes.
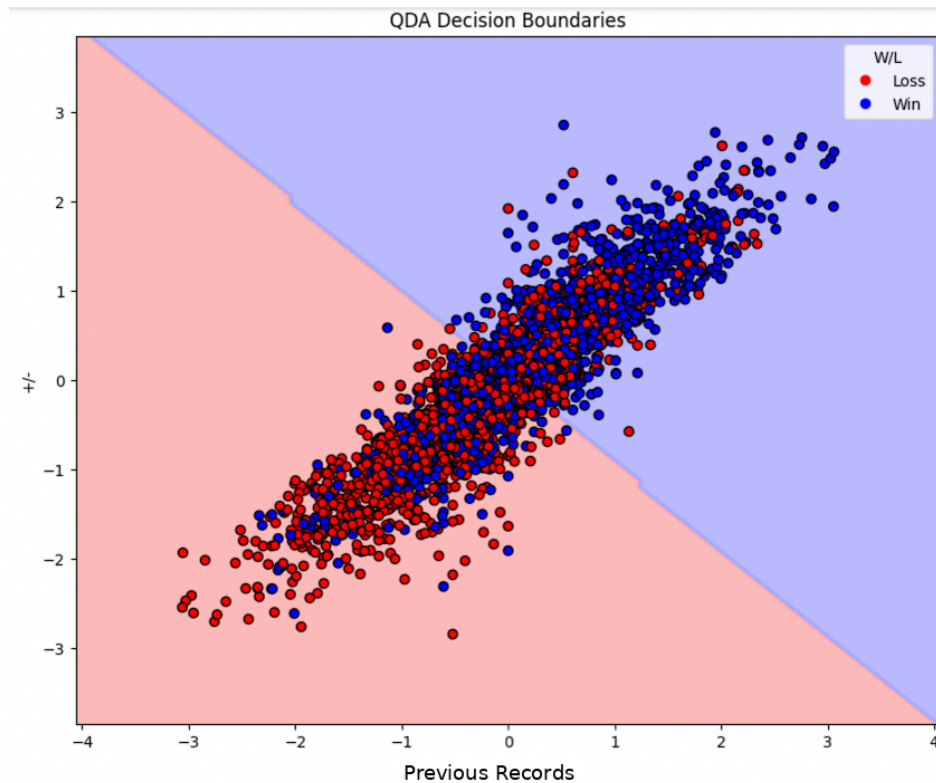


Figure 2: Illustration of QDA Decision Boundary

To better understand how the QDA model works, a decision boundary plot was created using the features +/- and previous record as an example. The plot visually shows the regions where the model predicts a "Win" (blue) or a "Loss" (red). Data points were overlaid on this background, with blue and red points representing actual "Win" and "Loss" outcomes, respectively. The decision boundary highlights how the model separates the two outcomes.

## 6.3     Decision Tree

A decision tree model with XGBoost was used to predict the outcomes of games, specifically focusing on whether a team would "win" or "lose" based on some selected features.

The decision tree classifier was trained and tested using a split of the dataset (70% for training and 30% for testing) to effectively measure the model's performance. The model reached an accuracy of 65.16%, which means it could reasonably distinguish between winning and losing outcomes based on past performance and current metrics. Interestingly, this accuracy was the highest we could achieve; adding any other features actually made the accuracy drop.

To better understand how the model was making its predictions, we generated a decision tree plot to show the sequence of splits and the importance of each feature. The decision nodes clearly

show how crucial "previous record" and "+/-" are for predicting whether a team would win or lose. By simplifying the decision tree to focus only on these key features, the model became more interpretable, making it easier to see the impact of previous games and score differences on game outcomes.
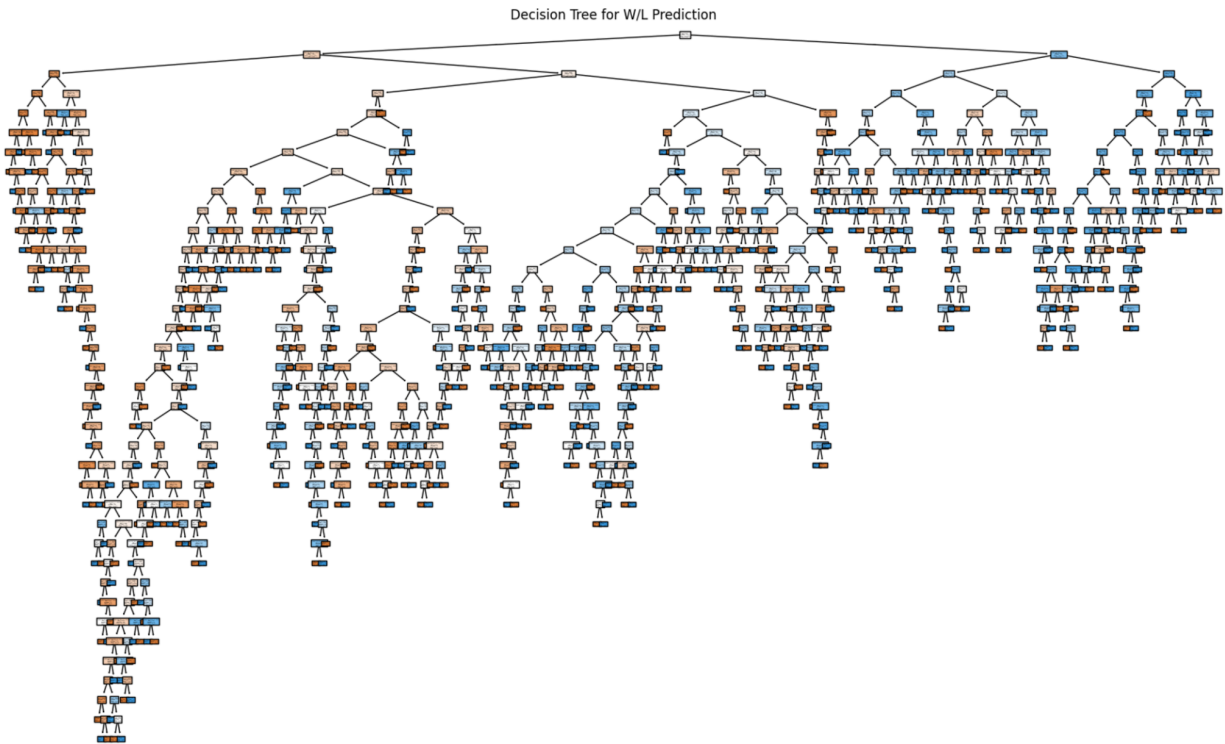


Figure 3: Illustration of Decision Tree

## 6.4    Ridge Classifier

This model aimed to predict game outcomes (win or loss) using Ridge Regression. Categorical features were one-hot encoded, and all predictors were standardized to ensure consistent scaling. The dataset was then split into 70% training and 30% testing, and Ridge regression was applied with cross-validation over a range of alpha values (0.1, 1.0, 10.0, 100.0, 1000.0, 10000, 100000).

Table 2: Table of Coefficients in Ridge Classifier

| Feature | Estimated Coefficient |
|---|---|
| Previous Record | 0.063 |
| Average Score | -0.057 |
| AST | - 0.053 |
| FG% | 0.057 |
| +/- | 0.044 |
| Variance | 0.037 |
| STL | 0.023 |
| FGA | -0.016 |
| 3PA | 0.029 |
| PTS | 0.026 |
| FT% | -0.014 |
| TOV | -0.013 |
| BLK | 0.011 |
| PF | -0.009 |
| MIN | -0.007 |
| OREB | 0.004 |

The model achieved 65.58% accuracy on the test set, a moderate result suggesting some patterns were captured but with room for improvement. The estimated coefficients are presented in table 2. Key insights from the feature coefficients revealed that the previous records and field goal percentage (FG%) had the strongest positive impact on predicting a win. Interestingly, some features, like the average score before a match and assists, had negative coefficients, possibly pointing to unexpected influences. Other predictors, such as offensive rebounds (OREB), home advantage, and three-point percentage (3P%), showed little to no impact on game outcomes.

While Ridge Regression provided clear insights and helped avoid overfitting, the model's accuracy suggests it may miss deeper patterns or nonlinear relationships in the data. The unexpected behavior of some features highlights the need for better  key factors.

6.5     Support Vector Machine

With the cleaned data, we trained the Support Vector Machine, an algorithm that makes classification decisions based on the separating hyperplane. We tested several choices of C (0.1, 0.5, 1.0, 5, 10.0, 100.0, 1000.0) and selected the model with the best performance on the testing dataset. Overall, the model performed moderately with a testing accuracy of 67.14%.

## 7     Conclusion

Above all, we converted the original data into a dataset containing the historical records before each match, and trained 5 models on the new dataset. Since the dataset is balanced, we use testing accuracy as the metric to evaluate the performance of each model, which is presented in table 3.

Table 3: Summary of Model Accuracy

| Model | Testing Accuracy |
|---|---|
| Logistic Regression | 66.67% |
| QDA | 66.71% |
| Decision Tree | 65.16% |
| Ridge | 65.58% |
| SVM | 67.14% |

As shown in the table, the model with best performance on the testing dataset is SVM, with an accuracy of 67.14%, outperforming our baseline while still having room for improvement.

To achieve better model performance, we may enlarge our dataset. Currently, we are training on a dataset with only 2352 observations, which limit our model choice and provide insufficient training to each model. Furthermore, we can consider more features, such as collecting data on the performance of each player.