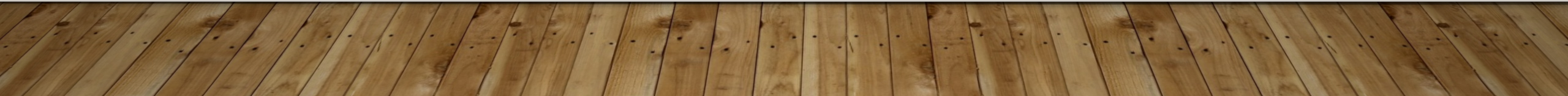


NCAA BASKETBALL

PREDICTING THE OUTCOMES OF BASKETBALL GAMES



CONTENTS

- Introduction
- The Data
- Cleaning the Data
- Analysis
- Machine Learning
- Conclusion

INTRODUCTION

- What variables are correlated with whether or not a team wins a game?
- Why is this important?
- How can we predict which team will win a game?

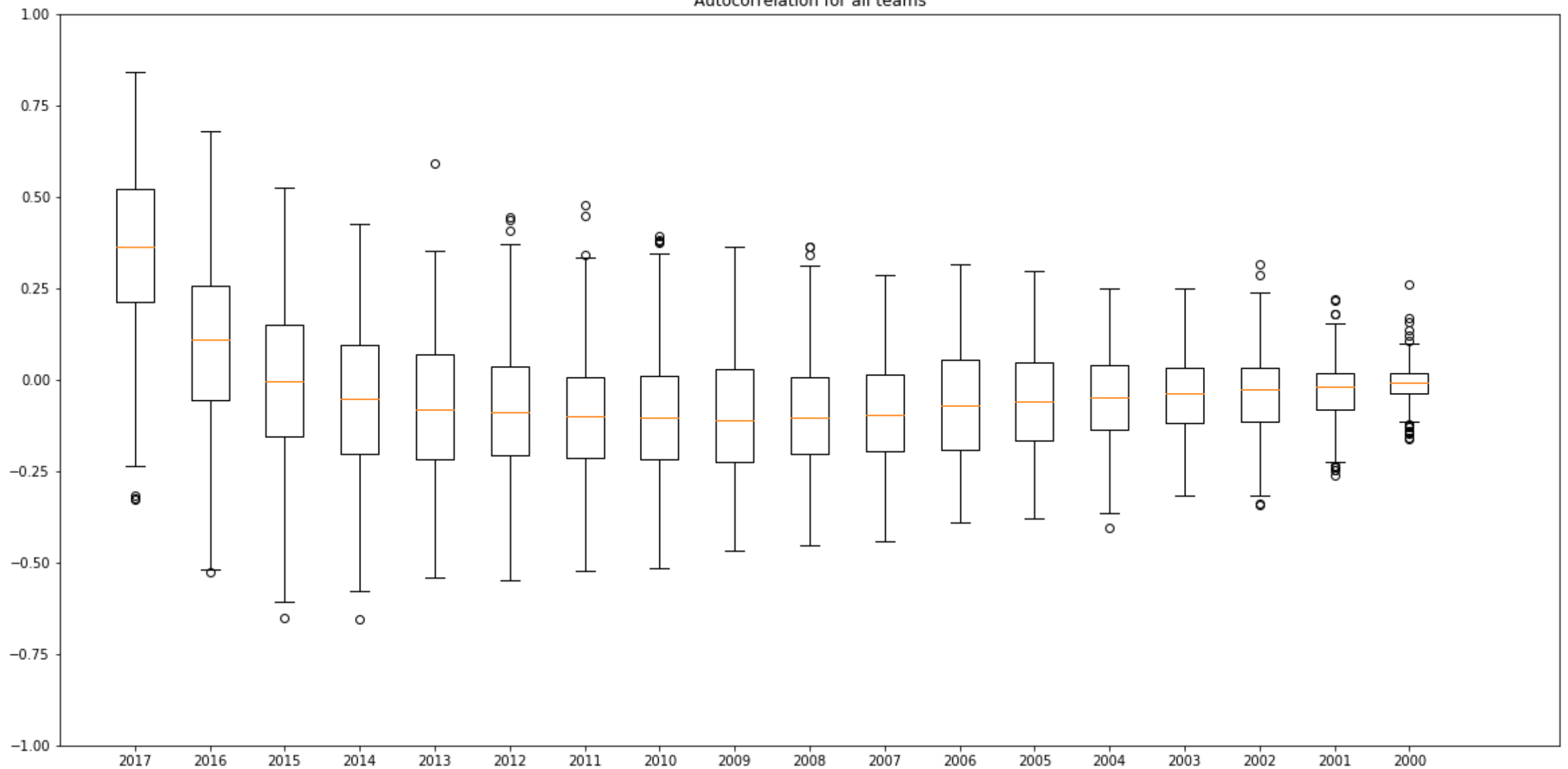
THE DATA

- Multiple csv files with team stats (from Kaggle)
- Overall ratings for teams ([Sports-reference.com](https://www.sports-reference.com))
- Merging multiple files, web scraping
- Data from past 18 years

CLEANING THE DATA

- Tools: pandas, numpy, matplotlib
- Extracting the data needed
- View autocorrelation to determine which data to use
- Preparing the data for merging

Autocorrelation for all teams



For the previous two years, there is a correlation, but after that the correlation is almost zero. This indicates that the data from the previous two years may be useful for predicting the outcomes for the current year.

CREATING FEATURES

- Raw data contains numbers, but we want rates
- Example: Field goal rate instead of number of field goals made
- Created function to calculate rates for each team for previous years
- Also calculated average number of steals, rebounds, blocks, etc.
- Average difference in scores for each game

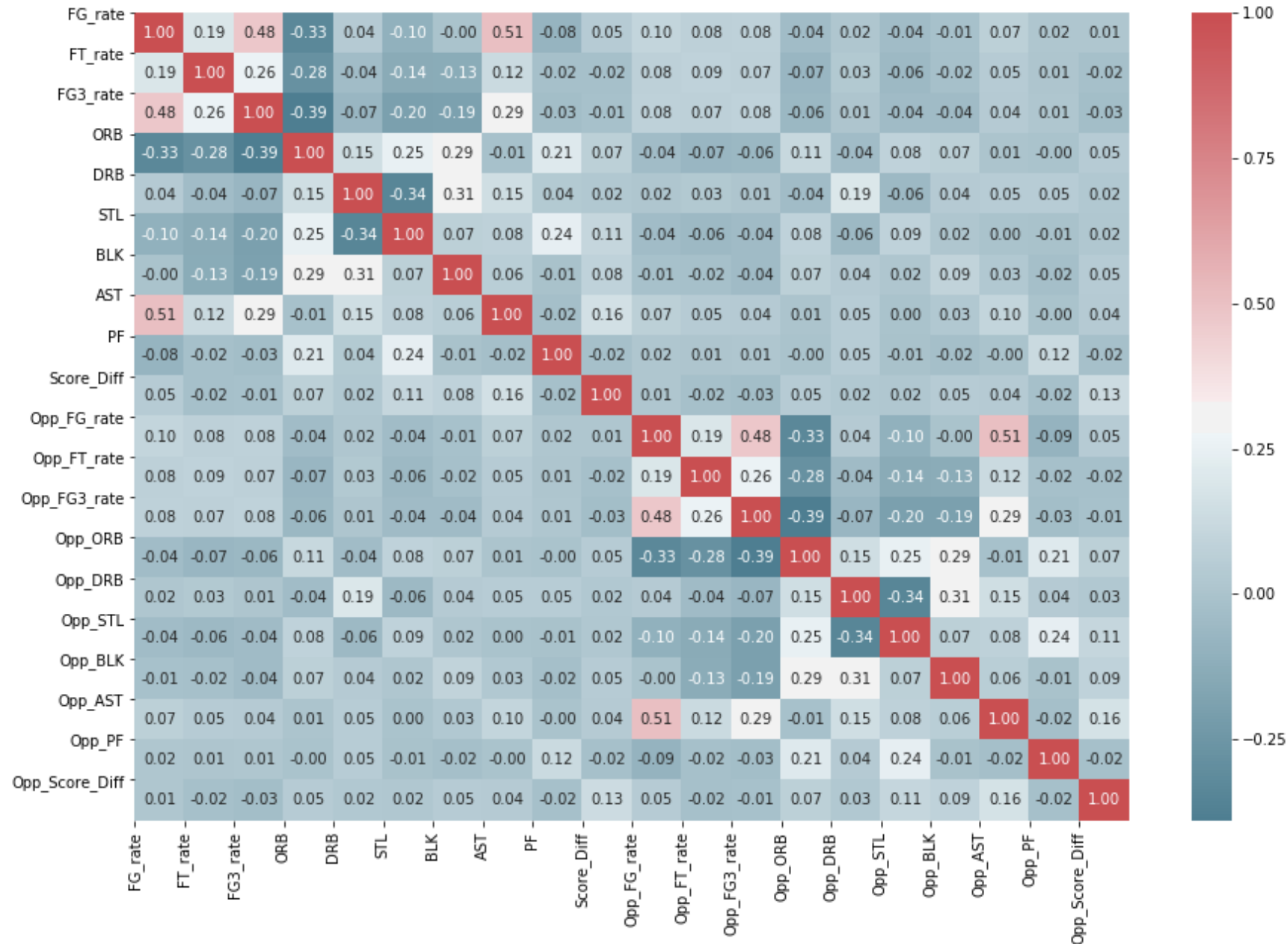
PREPARING THE DATA FOR ANALYSIS

- Merged the dataframes
- Created a column of ones and zeros, indicating whether the team won or lost
- Exported cleaned data as csv file

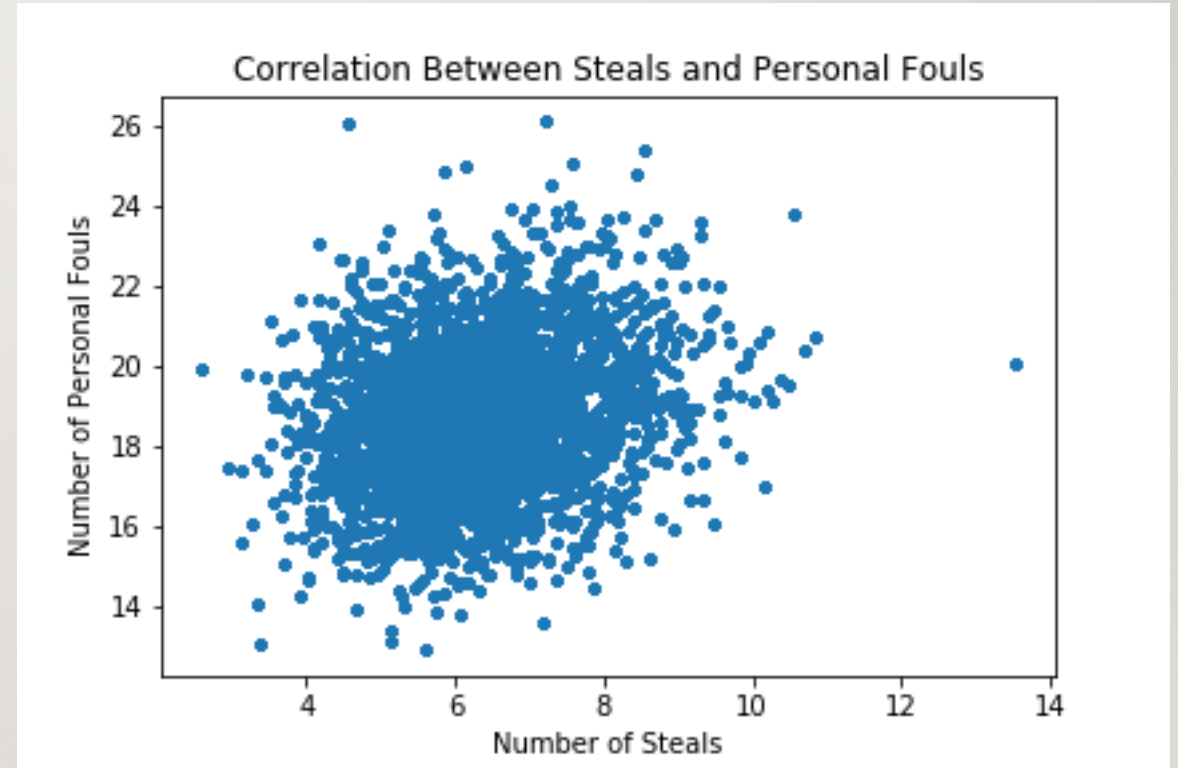
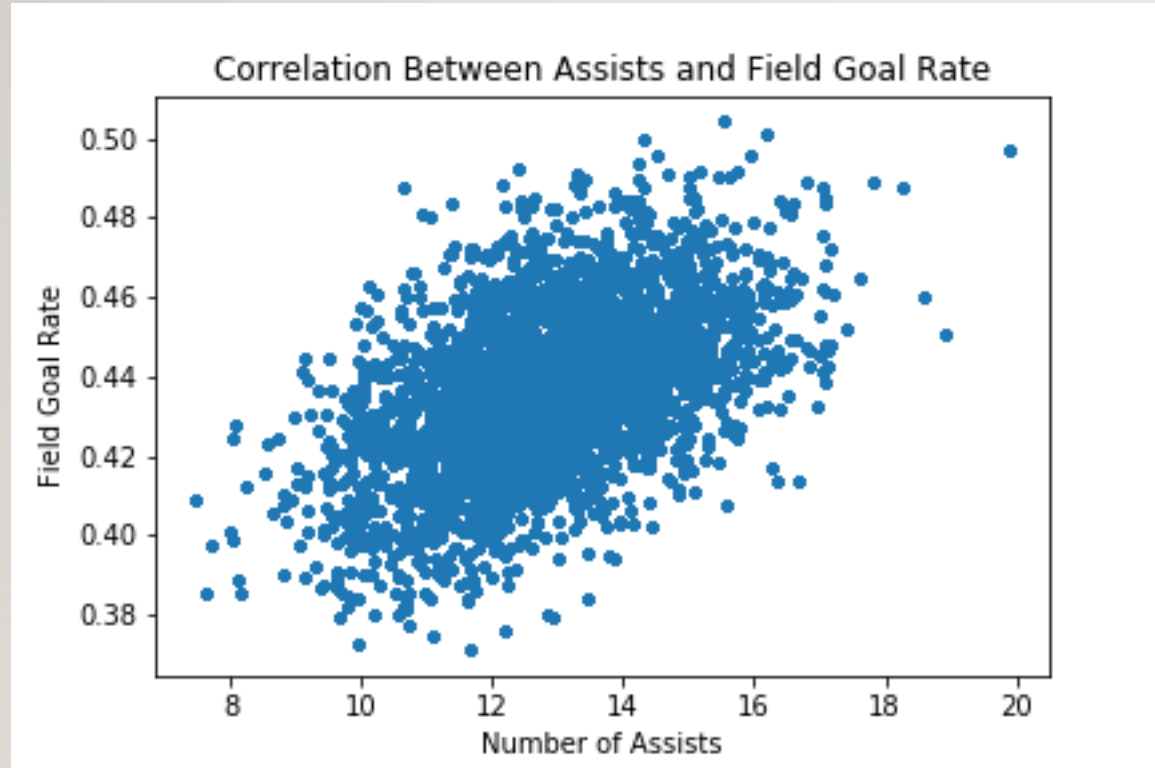
ANALYSIS

- Which features are correlated with each other?
- Which features are useful for predicting the outcome of a game?
- Testing the difference in a feature for winners versus losers.

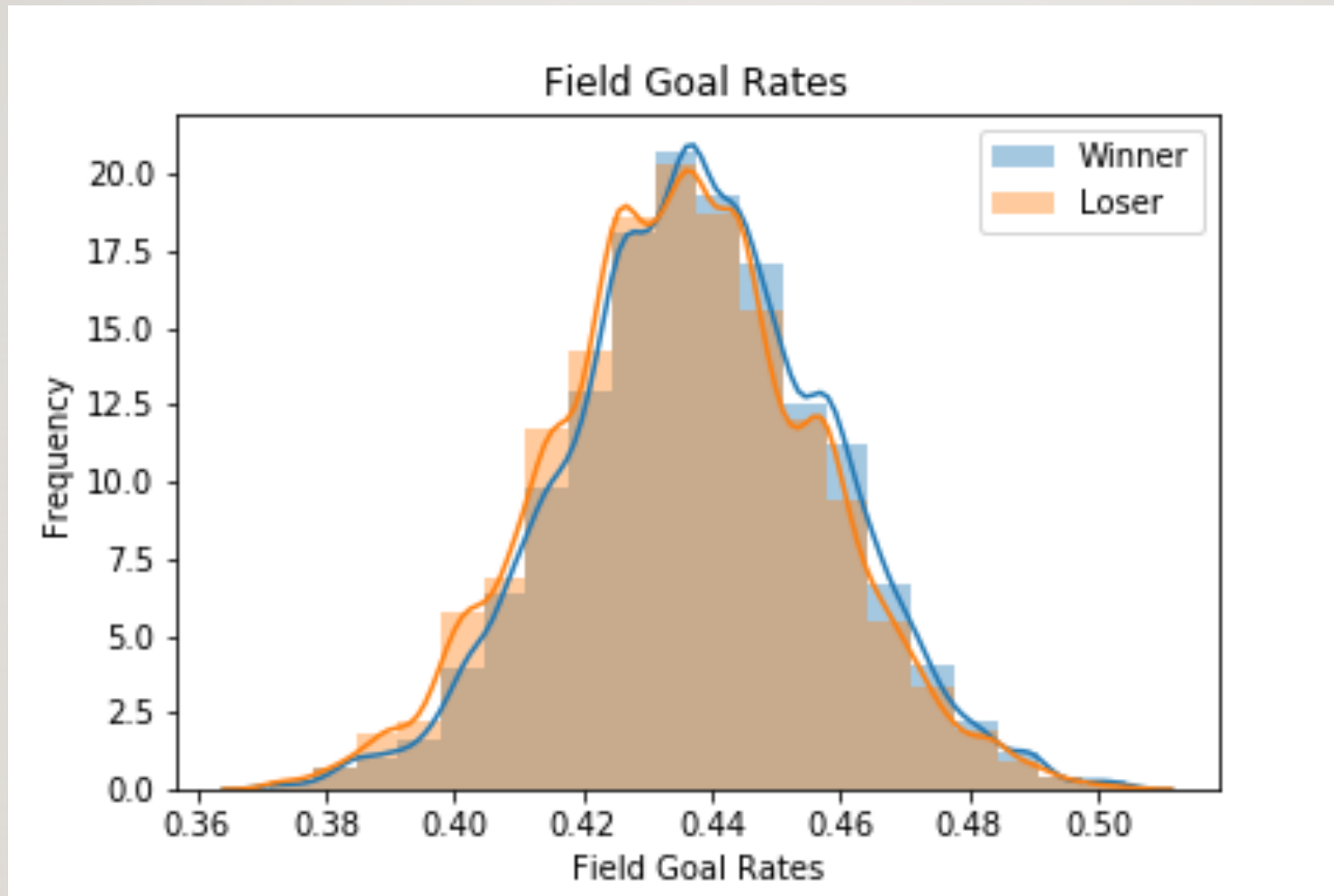
Heat map of the correlation between features



Correlation between features



Histogram of field goal rates for winning team versus losing team.



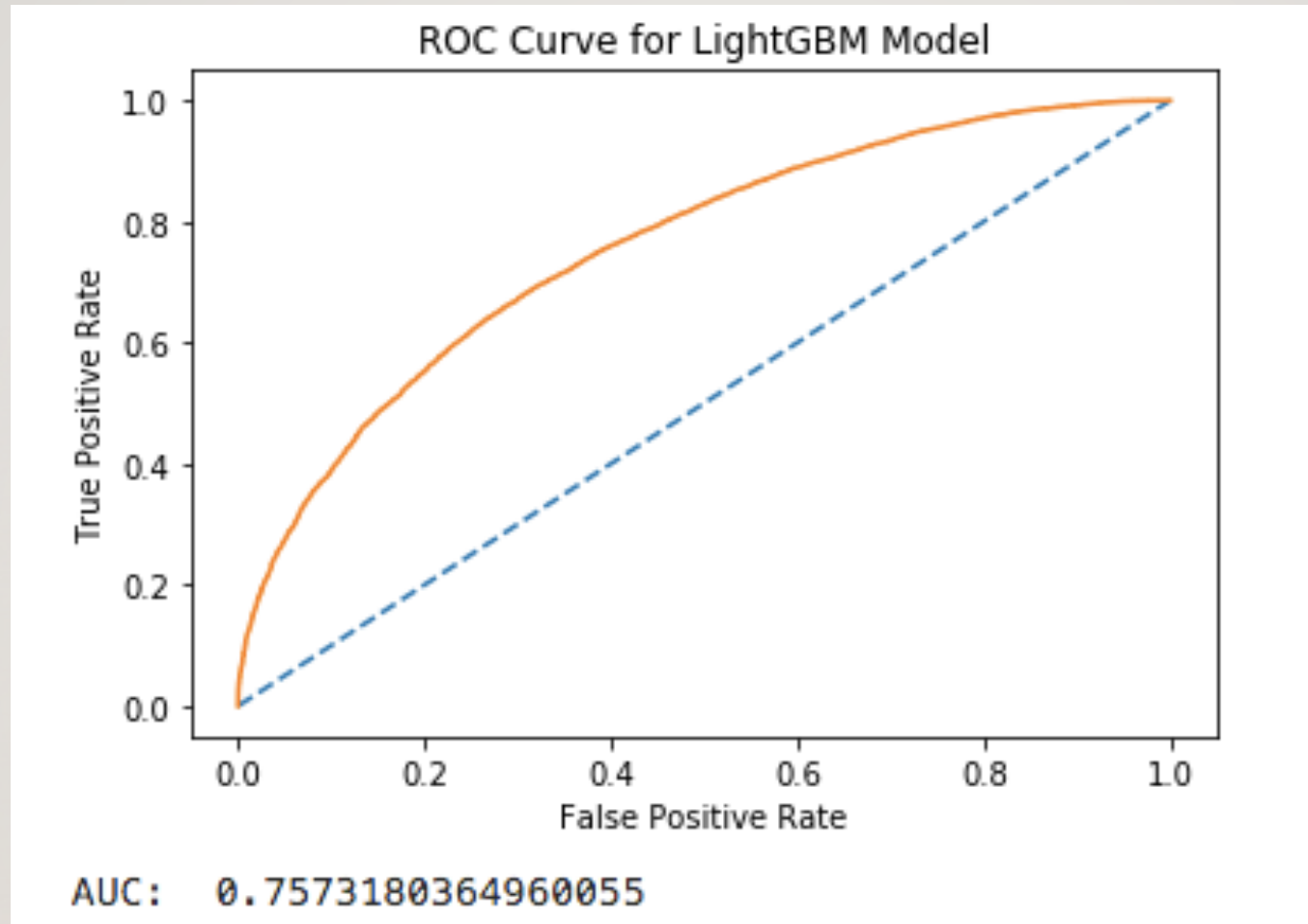
MACHINE LEARNING

- LightGBM (can handle missing data)
- Gradient boosting decision tree
- Binary classification
- Categorical features: team id, season, played in previous tournament (binary)
- Scikit-learn (to split into train and test sets)

FITTING THE MODEL

- Scikit-learn Randomized Search CV
- Randomized Search to find best values for the hyperparameters
- Used these parameters to train the model
- Converted probabilities to binary values (to compute accuracy)

Measuring the performance of the model



The area under the curve is about 0.757, so this is a good model.

CONCLUSION

- We have discovered which features are useful in predicting the outcome of a game
- We can use data from the previous two years to predict whether a team will win or lose.
- We have a model that is about 68.6% accurate
- Knowing which features are useful for predicting wins is useful for those attempting to predict the outcome of the NCAA tournament

FUTURE WORK

- Create additional features
- Use individual player data (may be difficult to access)
- Model could be used for similar problems: predicting wins for other sports