



NBA Game Match-up Prediction Tool

Yotan Demi-Ejegi, Diego Escalera, Ethan Garbow, Dan Goldin, Sebastián de la Hoz and Michael Lamontagne



INTRODUCTION

According to ESPN, NBA teams and the sports betting industry spend well in excess of \$3.1 billion in researching surrounding player stats and predictions in order to gain a competitive advantage in the respective fields.

We will be developing a visual tool that evaluates the head-to-head NBA match-up of a game to predict a winner.

NBA teams and coaches can use this to identify match-up benefits or issues. Also, for sports betting seeking to find betting advantages it can be used to improve sports betting lines or predict winners which has a financial impact

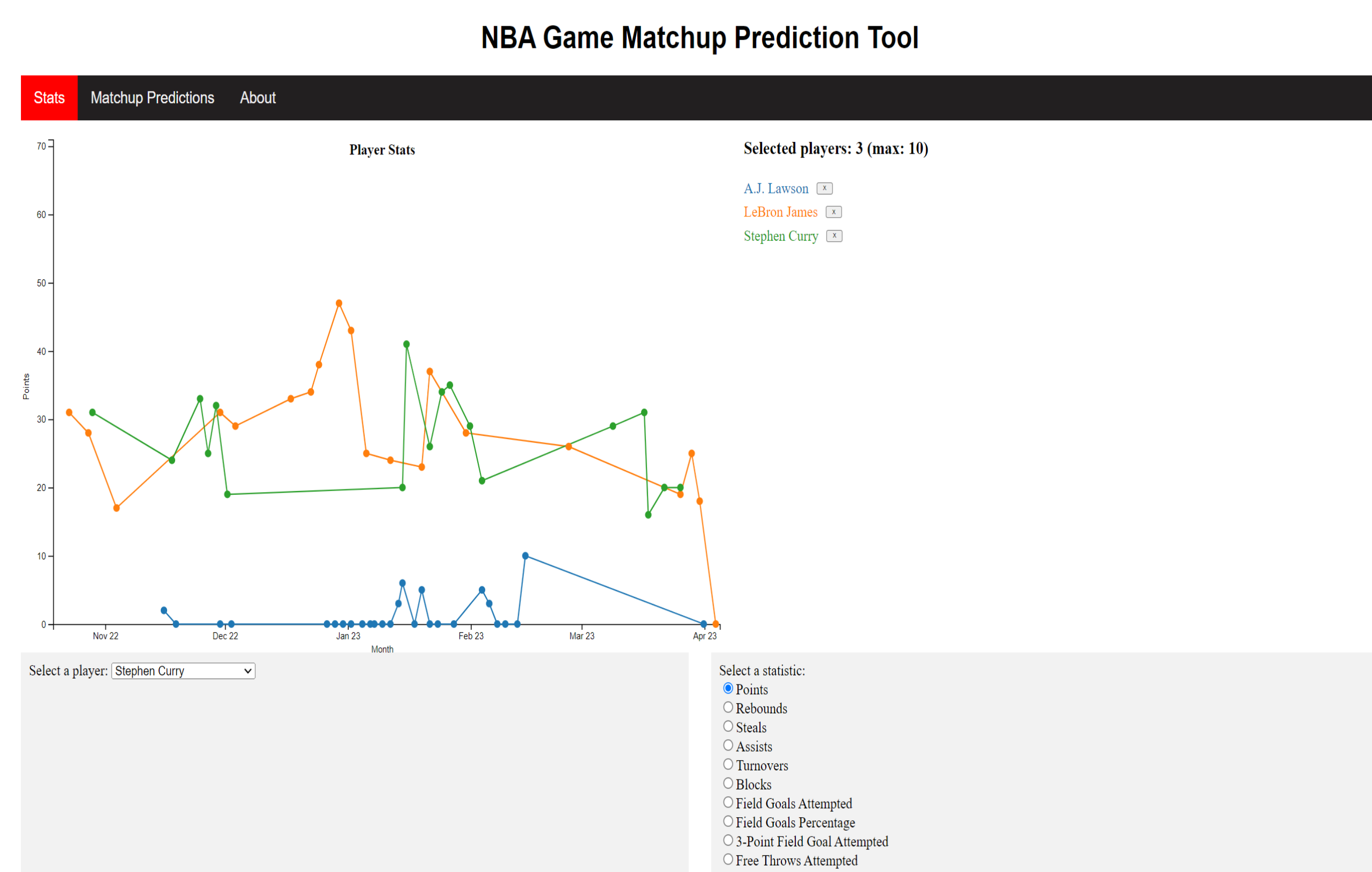
THE APPROACH

The algorithm involves running a multiple linear regression model looking at historical game performance to determine baseline projections of points per game, assists, rebounds, etc. We then took past 7-days of social sentiment via Twitter and created a sentiment score which was incorporated into the model to estimate potential changes to impact per game. And then run a multiple linear regression with the players' points per game as dependent variable and all teams and all players as categorical variables.

Using combination of baseline and team/player models, we estimated points per game for each player. For an upcoming game, we will aggregate the players points for each team to determine the winner

We are placing value on the individual team member and seeking to understand how individual performance is impacted based on match-ups. Each final predicted score is the aggregate of many models each customized to the individual player.

Current prediction tools only show which team will play against which other team and show numbers to predict the most likely. We also have gone extra mile of creating an interactive display using D3 that allows users to explore and compare different players, teams, and attributes.



THE DATA

Data used to train the model was gotten from

- Social discussion data was captured from Twitter firehouse via the Tweepy Python Library
- NBA datasets from KAGGLE and NBA api. All games from 2004 season to end of 2023 season (April 2023) from api was pulled using python code against exposed end points

Data after collection was stored in GitHub and contains below characteristics

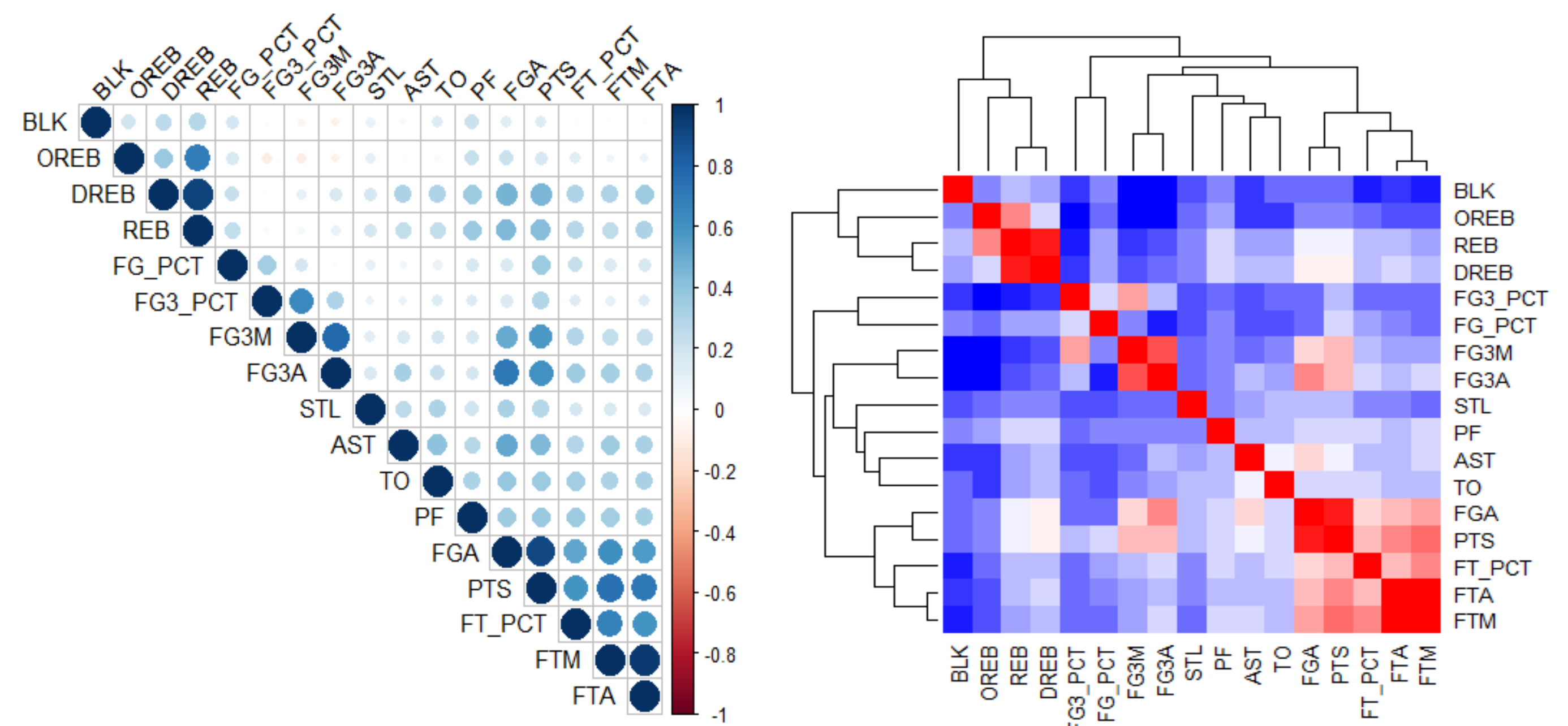
- Games Data containing total number of 21 variables of both continuous and categorical including game date, game id, home and away team ids, aggregate game stats for each team.
- Game Details data with 29 variables including Game ID, player ID, Team Abbreviation, Player Name, Game Statistics (Minutes played, field goal attempted and made, points, rebounds, assists etc.)
- For each player we took a sample of fifty tweets that had the following breakdown of fields like Tweet Text, Username, user_followers, attribute_score etc.

EXPERIMENT/RESULTS

As part of the data exploratory, we performed a correlation matrix for each basketball stat in order to identify and visualize any pattern in the data.

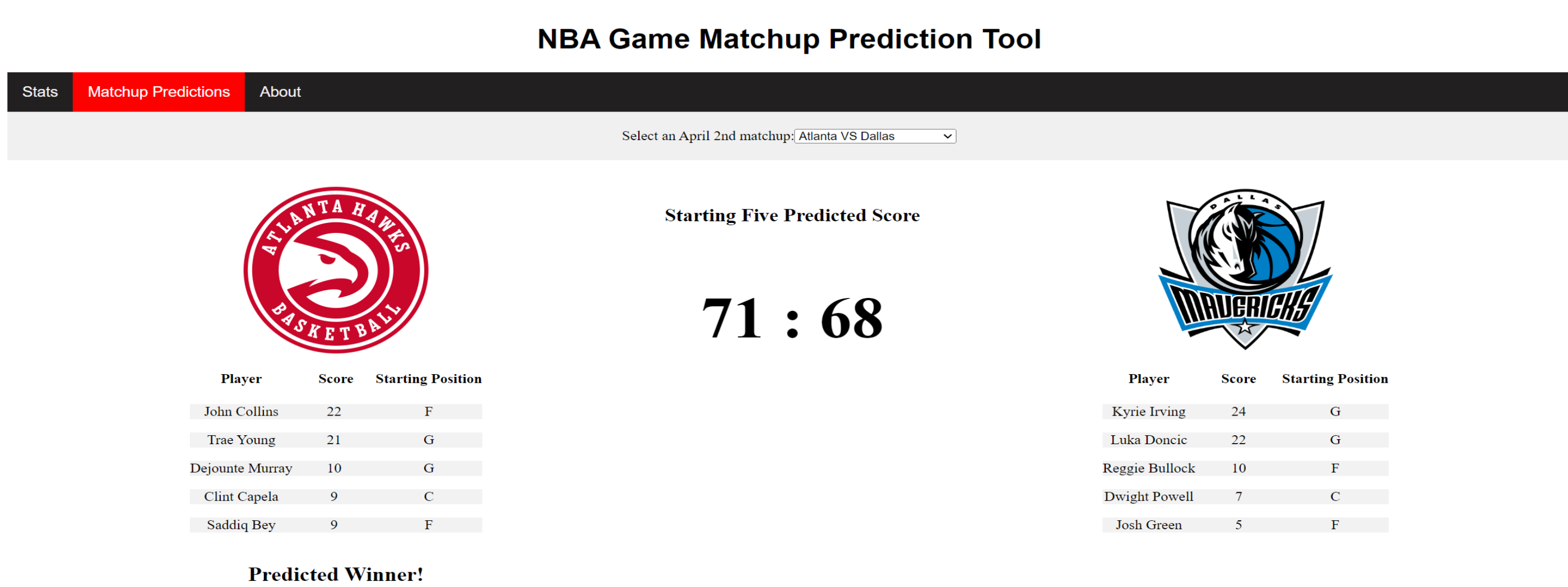
Figure below presents the correlation matrix plot. Most variables related to the act of shooting such as: FGA, FG3A, FT, among others, are highly related to the variable of interest, points per game. Similarly, the heat map presents clearly those findings as well as the different clusters based on the correlations results. We can see that there are about 4 clusters, confirming the results of the player cluster analysis on the progress report. The 4 clusters consist of: defensive stats such as blocks and rebounds, 3-pointers, possession and shooting stats.

Since teams combine player clusters to build balanced lineups, players' performance could vary throughout season matchups. These findings helped to support using non-traditional metrics to predict output such as the team matchup.



Our modeling process was used to predict the winner of games on specific date (April 2nd). This does represent a small sample but however gives a promising outcome. We focused on the starting five players of each team and aggregated the predicted points. The teams with the highest points were predicted to be the winner. The model prediction accuracy was 83%

GAME_ID	GAME_DATE	MATCHUP	ACTUAL WINNER	PREDICTED WINNER
22201167	2023-04-02	DAL @ ATL	ATL	ATL
22201165	2023-04-02	MEM @ CHI	CHI	MEM
22201166	2023-04-02	POR @ MIN	POR	POR
22201163	2023-04-02	CHA @ TOR	TOR	TOR
22201164	2023-04-02	UTA @ BKN	BKN	BKN
22201168	2023-04-02	WAS @ NYK	NYK	NYK



We tested 10 players to predict their points per game using the model for the games occurring after January 1, 2023. Our per game points prediction for each of the 10 players tested yields an Mean Squared Error (Prediction - Actual)² in the range of 33-300 (average - 97.62).

Many NBA game prediction models have been created but focus on team results and statistics. Team-based outcome models largely ignore player performance and contributions. Player-based models largely ignore the impact of a specific matchups.

Our model successful by focuses on player performance accounting for factors that are critical in today's NBA.