

Analyzing the Efficacy of the *RAPTOR* Statistic

Kevin Zhang
University of Pittsburgh
Information Science
Junior
kez27@pitt.edu

Chris Park
University of Pittsburgh
Information Science
Junior
chp117@pitt.edu

Abstract—The project name is “Analyzing the Efficacy of the *RAPTOR* Statistic” and the group members are Chris Park and Kevin Zhang. The research question we want to answer is “What effect do different categories of statistics have on NBA’s newly implemented *RAPTOR* score.” We decided to explore this research question regarding the NBA because we are both avid fans of the game of basketball and are interested in how data analytics works in the sports industry. Also, it is intriguing that sports analytics is seen everywhere these days in social media outlets such as *Instagram* and *Facebook*. Although basketball is for entertainment, it is ultimately a business and we thought it would be interesting to look at the NBA with a statistical business perspective.

In general, there are large sums of money involved when it comes down to predictions. Specifically, there are large monetary amounts involved in the sports world because of the long-term contracts that are being signed by professional athletes and the betting that is involved in each game that is played. This project is useful because it offers analysis on the newly implemented *RAPTOR* statistic that the NBA uses to determine a player’s value on a team. The model will reveal which category of statistics are most valued in predicting the *RAPTOR* statistic. This will then help people be more intelligent about certain decisions they make with their money.

The main beneficiaries of the model are sports managers and coaches. General team managers can use the model’s results to see which players to recruit or trade while coaches can see which game strategy to use. For example, team managers would want to sign certain players that are determined to bring the most value to the team. As for coaches, they can decide on a certain style of play or

determine which players to give more playing time. Other groups that would benefit include stakeholders and gamblers. Stakeholders and gamblers can use these models to see which team to invest their money in before the season even starts.

The first step of our plan for data analysis was to merge the NBA player statistics dataset found on *NBAStuffer* with the NBA player advanced metric dataset found on *FiveThirtyEight’s GitHub* to see both the box-score statistics and advanced-metric statistics. The next step was to clean the data such as dealing with missing values, changing data types, choosing our necessary features, dropping unnecessary columns and to understand the shape of the data through histograms and boxplots to scale or transform the data. Once the data was prepared, we built four linear regression models to fit the data by using Scikit-learn. After that, we compared the models based on the R^2 Score and the Mean Squared Error Score and decided which model (with its specific predictors) was the best model at predicting the *RAPTOR* statistic.

I. INTRODUCTION

For a while now, the NBA has been using the CARMELO (Career-Arc Regression Model Estimator with Local Optimization), which forecasts the careers of basketball players. It was used to project each player’s playing time and overall value on offense and defense, but not individual statistics¹. However, the NBA has moved onto a new predictor called RAPTOR.

The NBA statistic “RAPTOR” stands for **R**obust **A**lgorithm (using) **P**layer **T**racking (and) **O**n/**O**ff **R**atings². RAPTOR is a

- [1] ¹N. Silver, “We’re Predicting The Career Of Every NBA Player. Here’s How.” *FiveThirtyEight*, 09-Oct-2015. [Online]. Available: <https://fivethirtyeight.com/features/how-were-predicting-nba-player-career/>. [Accessed: 30-Nov-2020].
- [2] N. Silver, “How Our RAPTOR Metric Works,” 10-Oct-2019. [Online]. Available: https://fivethirtyeight.com/features/how-our-raptor-metric-works/?fbclid=IwAR0uC5Yipy-54_cPGg4STxqrPBNYj3UghRbj4fxgPrS2grcTlfl1iOxRrkL. [Accessed: 30-Nov-2020].

plus/minus statistic that measures a player's point contribution to his team based on offensive and defensive contributions per 100 possessions. There are two major components when calculating a player's RAPTOR score. The first component is the "box" component which uses a player's individual statistics. The second component is the "on-off" component which evaluates a team's performance when the player and combinations of his teammates are on or off the floor. As a descriptive statistic, RAPTOR is calculated solely based on a player's on-court performance and the performance of the player's teammates.

The overall RAPTOR score is a combination of three different factors: Box RAPTOR Offense, Box RAPTOR Defense, and RAPTOR On-Off. Box RAPTOR Offense includes a few offensive statistics, but most of them fall into one of four categories: scoring and usage, passing, rebounding, and space creation. The measures of scoring and usage include points, usage rate, time of possession, and assisted field goals. The measures of passing include enhanced assists and net passes. The measures of rebounding include enhanced offensive rebounds, team offensive rebounds on missed shots, and positional opponents' defensive rebounds. The measures of space creation include defended 3-point attempts and isolation turnovers. In addition to these categorized offensive statistics, there are a few miscellaneous offensive metrics such as fast-break starts, non-shooting defensive fouls drawn, penalty fouls drawn, and opponents' defensive rating. Though the offensive statistics are intuitive, the metrics for Box RAPTOR Defense act as proxies for other unmeasured statistics. RAPTOR uses several variables in its defensive regression: steals, offensive fouls drawn, opponents' field goals made and attempted, enhanced defensive rebounds, positional opponents' points scored, positional opponents' offensive rebounds, distance traveled (perimeter defenders only), opponents' free throws made, fast break turnovers committed, penalty fouls committed, and opponents' offensive rating. RAPTOR On-Off uses a few metrics such as: the player's offensive and defensive rating, the player's court mates' weighted average offensive and defensive ratings, and the player's court mates' other court mates' weighted average offensive and defensive ratings. When combining the "Box" and "On-Off" components, the RAPTOR rating is equal to about 85 percent of "Box" RAPTOR plus 21 percent of "On-Off" RAPTOR. The combined values of the "Box" and "On-Off" components are greater than 100 percent because there is no redundant information. After the combination of the two components, the RAPTOR score is then adjusted in two ways: score effects adjustment and team effects adjustment. The score effects adjustment accounts of "junk time" statistics which occur when one team is way ahead of another. When a team is way ahead, the team tends to be less

efficient while the opposing team tends to be more efficient. If left unadjusted, players on good teams will become underrated while players on bad teams will become overrated. The team effects adjustment recalculates individual players' ratings so that they sum up to reflect the team's overall performance. The players who were most heavily involved with the offense or the defense receive more of the credit of blame than those who were not as heavily involved.

Ultimately, we decided to choose RAPTOR overall score or the RAPTOR+/- as our response variable because it is the best at showing which players are the most valuable and most efficient.

PLAYER	TEAM	POSITION(S)	MINUTES	BOX SCORE RAPTOR			ON/OFF RAPTOR			OVERALL RAPTOR		
				OFF.	DEF.	TOT.	OFF.	DEF.	TOT.	OFF.	DEF.	TOT.
1 James Harden	Rockets	PG, SG	2,931	+9.5	+1.2	+10.8	+1.6	+4.6	+6.2	+8.5	+2.0	+10.5
2 LeBron James	Lakers	PG, SF, PF	3,078	+6.6	+0.2	+6.9	+4.4	+5.7	+10.1	+6.5	+1.3	+7.8
3 Anthony Davis	Lakers	PF	2,900	+3.9	+4.7	+8.5	-0.8	+2.2	+1.4	+3.1	+4.3	+7.4
4 Kawhi Leonard	Clippers	SF	2,359	+6.3	+3.2	+9.5	+5.4	+2.7	+8.1	+6.4	+3.3	+9.8
5 Giannis Antetokounmpo	Bucks	SF, PF	2,194	+5.6	+2.6	+8.2	+2.9	+6.2	+9.1	+5.3	+3.4	+8.8
6 Jayson Tatum	Celtics	SF	2,955	+3.5	+1.3	+4.8	+2.1	+5.0	+7.1	+3.4	+2.1	+5.5
7 Damian Lillard	Trail Blazers	PG	2,617	+9.7	-2.2	+8.5	+3.7	+0.3	+3.9	+8.1	-1.8	+6.4
8 Nikola Jokic	Nuggets	C	3,030	+4.5	+0.9	+5.4	+3.8	-1.4	+2.5	+4.6	+0.5	+5.1
9 Jimmy Butler	Heat	SG, SF	2,785	+4.3	+1.8	+6.1	+2.0	+0.1	+2.1	+4.1	+1.5	+5.6
10 Rudy Gobert	Jazz	C	2,603	-0.7	+5.9	+5.3	+3.3	+2.2	+5.6	+0.1	+5.5	+5.6

Fig. 1 The top NBA Players are revealed through the RAPTOR Statistic
(Taken from FiveThirtyEight)

II. METHODOLOGY

A. Reading & Merging the Datasets

The first step to creating the model was to read in the datasets and to merge them. The first dataset was found on *NBASTuffer* (<https://www.nbastuffer.com/2019-2020-nba-player-stats/>) and the second dataset, with the NBA player advanced metric dataset, was found on *FiveThirtyEight's GitHub* (<https://github.com/fivethirtyeight/nba-player-advanced-metrics>) which provided both the box-score statistics and advanced-metric statistics. *NBASTuffer* provides sports datasets in excel for NBA basketball analytics and *FiveThirtyEight* is a website that focuses on opinion poll analysis by providing data from *GitHub*. The first dataset had unclear column names, so every column had to be renamed. The first column within the first dataset was also deleted because it provided no information. The second dataset included statistics from 1977 to 2020. To remediate this, the dataset was queried to only include statistics from 2020 which allowed it to match up with the first dataset. The two datasets were merged using an inner

statistics because we wanted actual player statistics instead of pace adjusted statistics. After dropping the columns that contained statistics that would not be used, columns that had ambiguous names were renamed.

Fig. 2 2019-2020 Season NBA Player Dataset from *NBAstuffer*

Fig. 3 1977-2020 NBA Player Advanced-Metric Dataset from *FiveThirtyEight*

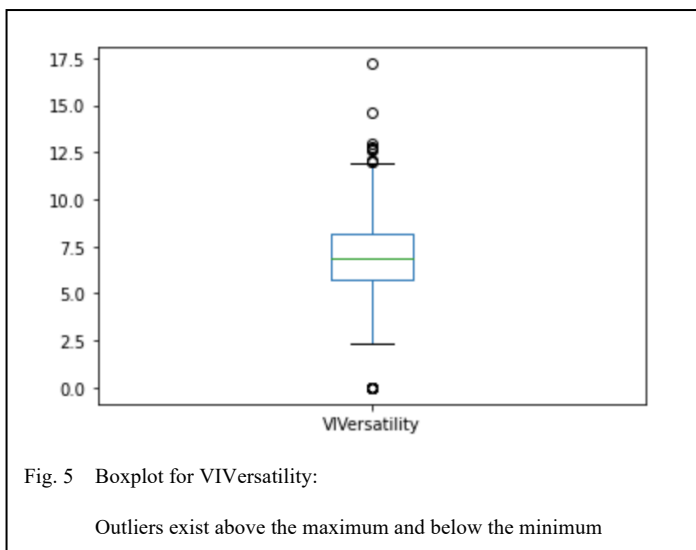
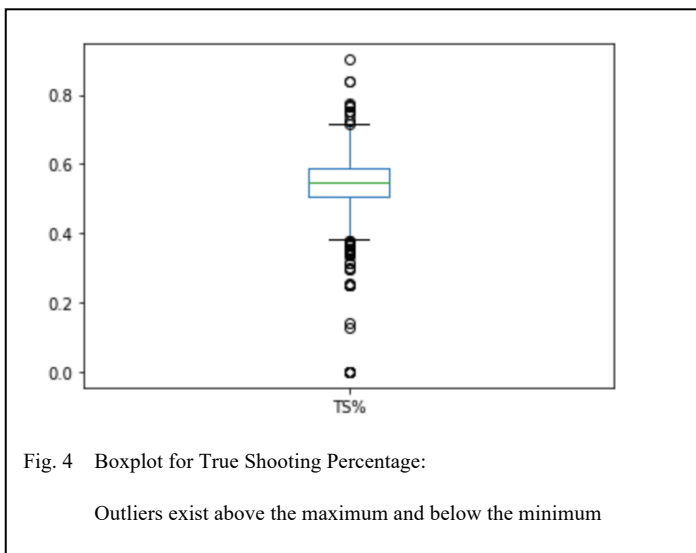
C. Check For Missing Values

After dropping and renaming columns, the dataset was checked for any missing values. The dataset had two missing values for TS% and three missing values for ORtg. Before addressing these missing values, the columns needed to be converted to the correct data type. Many columns that had numeric values were of an object type. This was fixed by converting the columns to their appropriate numeric data type. After ensuring that the columns had the correct data type, we needed to decide whether to impute based on the mean or median. To make this decision, a boxplot of TS% and ORtg were created and it was used to see if there were any outliers that would skew the data. After analyzing the boxplots, we decided to impute based on the mean for both columns because there are outliers both below the minimum and above the maximum. After imputation based on mean, there were no missing values within our dataset.

The next step was to check for outliers. An outlier of every numeric column/statistic was created (AGE, GP, MPG, USG%, TS%, PPG, APG, RPG, SPG, BPG, VIVersatility, tmRtg, Raptor+/-, PIE%, ORtg, DRtg, Raptor O, Raptor D). We chose not to drop any outliers because it makes sense for there to be outliers below the minimum and beyond the maximum. There is a large discrepancy between players in the NBA. Those outliers that are below the minimum are likely the players who are not that good and do not contribute to their teams which results in a low accumulation of stats. The outliers that are beyond the maximum are likely star players who generate most of the statistics for their team. In addition, dropping outliers would heavily change the distribution. For example, in Figure 4, the boxplot for "TS%" percentage shows outliers both above the maximum and below the minimum. "TS%" stands for true shooting percentage which is a measure of shooting efficiency that takes into account field goals, 3-point field goals, and free throws. There are outliers above the maximum because there are players who take few shots and make those shots which results in a high true shooting percentage. There are outliers below the minimum because there are players who either take many shots or a few shots and miss a lot which results in a low true shooting percentage.

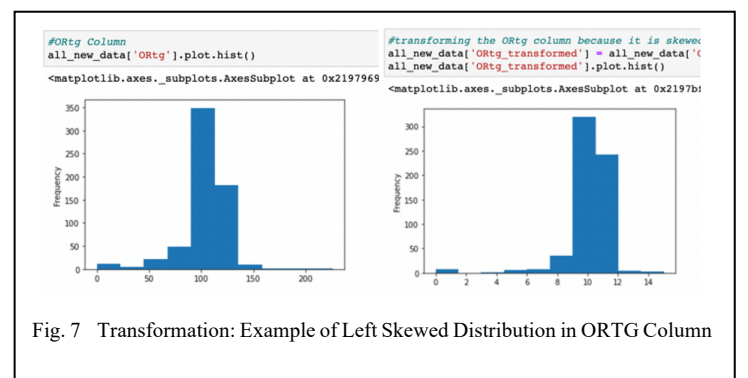
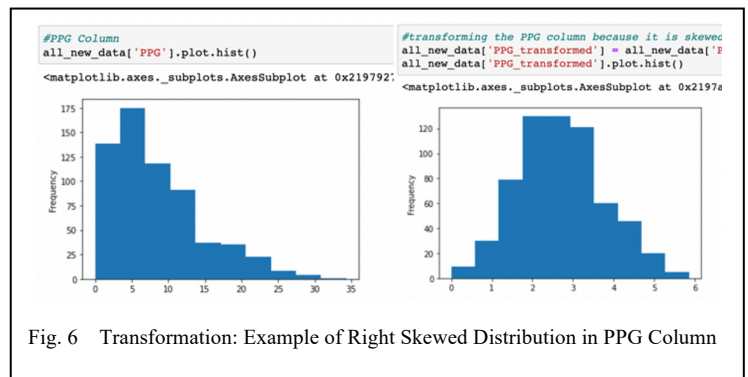
The next step was to drop the columns that were not going to be used and to rename the columns. We chose to drop columns for two reasons: the statistic would not provide any value to the model or the statistic was repeated and needed to be dropped. Of the columns that were not repeated, we dropped MIN%(percentage of team's minutes used by player), TO%(number of turnovers a player commits per 100 possessions), FTA(free throw attempts), 2PA(two-point attempts), 3PA(three-point attempts), eFG%(effective shooting percentage), TOPG(turnovers per game), P/36(points per 36), A/36(assists per 36), R/36(rebounds per 36), SB/36(steals+blocks per 36), and TO/36(turnovers per 36). We chose to drop these columns because the statistics provided would have already been included in other statistics such as PPG. There was no need for columns that calculated "per 36"

There were some boxplots that had many outliers above the maximum and had one outlier at zero. For example, VIVersatility is an index that measures a player's ability to produce points, assists, and rebounds. The average player will score around five on the index while a star player will score above ten. There is an outlier at zero because this represents that population of players who do not receive any playing minutes and therefore have no contribution to the team. There are many outliers above the maximum because star players accumulate much more stats than the typical NBA player. These outliers are not strange to see because there is such a discrepancy between players in the NBA. Statistics that require individual contributions will likely be filled up by star players which results in many outliers.



E. Visualize, Transform, and Scale the Data

The next step is to visualize the data and to transform/scale appropriately. We chose not to scale anything because it makes sense for there to be a large discrepancy between the maximum and minimum. The minimum was so low because of players who do not contribute to their teams while the maximum was so high due to star players who have the highest usage rates and accumulate the most statistics for their teams. We did transform several columns. We transformed the PPG, APG, RPG, SPG, and BPG columns because they were skewed to the right. The PIE%, ORtg, and DRTg columns were transformed because they were skewed to the left. The columns that were not transformed were AGE, GP, MPG, USG%, TS%, VIVersatility, tmRtg, Raptor O, Raptor D, Raptor+/- because they all had fairly normal distributions. After transforming the data, we checked for missing values again. The only column that had missing values was PIE%_transformed and we imputed by the mean because there were outliers below the minimum and beyond the maximum.



III. RESULTS

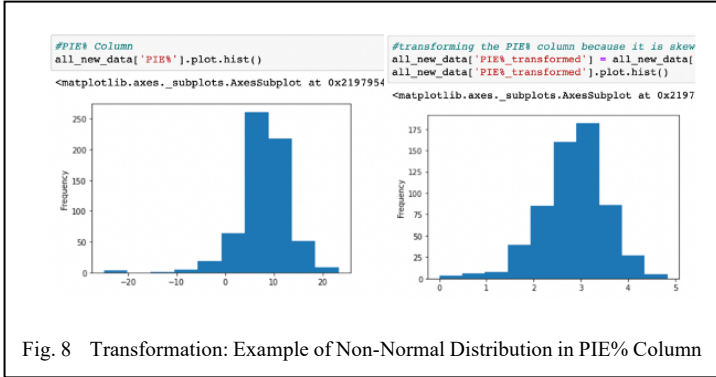


Fig. 8 Transformation: Example of Non-Normal Distribution in PIE% Column

F. Create Linear Regression Models

The last step is to create a linear regression model. We chose to do a linear regression model because it reveals how sets of variables are related to each other. We created four linear regression models that had different predictor variables. The first regression model included predictor variables that were focused on the individual player (PPG, USG%, TS%, PIE%, VIVersatility, ORtg, Raptor O). This was done to see how accurate the RAPTOR score would be if it only focused on individual statistics. The second regression model included all team-driven non-scoring statistics as the predictors (APG, RPG, BPG, SPG, DRtg, Raptor D). This focused more on the team component of the RAPTOR score. The third regression model only included box-score statistics. Box-score statistics are the basic individual accomplishments (PPG, APG, RPG, BPG, SPG, MPG). The last model included all advanced statistics as the predictor variables (USG%, TS%, VIVersatility, tmRtg, PIE%, ORtg, DRtg, Raptor O, Raptor D).

Linear Regression with Individual-Driven Statistics

```
lm = LinearRegression()
x = all_new_data[['PPG_transformed', 'USG%', 'TS%', 'PIE%_transformed',
                  'VIVersatility', 'ORtg_transformed', 'Raptor O']]
y = all_new_data['Raptor+/-']
lm.fit(x, y)

LinearRegression()

r_squared = lm.score(x, y)
print('Coefficient of determination: ', r_squared)

Coefficient of determination: 0.5732435446441481

import math
from sklearn.metrics import mean_squared_error
y_predict = lm.predict(x)
regression_model_mse = mean_squared_error(y_predict, y)
math.sqrt(regression_model_mse)

4.5941426347018375
```

Fig. 8 First Linear Regression Model: Calculating R^2 and RMSE Scores

Linear Regression with Team-Driven Non-Scoring Statistics

```
lm = LinearRegression()
x = all_new_data[['APG_transformed', 'RPG_transformed', 'BPG_transformed',
                  'SPG_transformed', 'DRtg_transformed', 'Raptor D']]
y = all_new_data['Raptor+/-']
lm.fit(x, y)

LinearRegression()

r_squared = lm.score(x, y)
print('Coefficient of determination: ', r_squared)

Coefficient of determination: 0.3697621064951383

import math
from sklearn.metrics import mean_squared_error
y_predict = lm.predict(x)
regression_model_mse = mean_squared_error(y_predict, y)
math.sqrt(regression_model_mse)

5.582987827669979
```

Fig. 9 Second Linear Regression Model: Calculating R^2 and RMSE Scores

Linear Regression with Box-Score Statistics

```
lm = LinearRegression()
x = all_new_data[['PPG_transformed', 'APG_transformed',
                  'RPG_transformed', 'BPG_transformed',
                  'SPG_transformed', 'MPG']]
y = all_new_data['Raptor+/-']
lm.fit(x, y)

LinearRegression()

r_squared = lm.score(x, y)
print('Coefficient of determination: ', r_squared)

Coefficient of determination: 0.12825762057613177

import math
from sklearn.metrics import mean_squared_error
y_predict = lm.predict(x)
regression_model_mse = mean_squared_error(y_predict, y)
math.sqrt(regression_model_mse)

6.566115303144572
```

Fig. 10 Third Linear Regression Model: Calculating R^2 and RMSE Scores

Linear Regression with Advanced Statistics

```
lm = LinearRegression()
x = all_new_data[['USG%', 'TS%', 'VIVersatility', 'tmRtg',
                  'PIE%_transformed', 'ORtg_transformed',
                  'DRtg_transformed', 'Raptor O',
                  'Raptor D']]
y = all_new_data['Raptor+/-']
lm.fit(x, y)

LinearRegression()

r_squared = lm.score(x, y)
print('Coefficient of determination: ', r_squared)

Coefficient of determination: 0.9999502357202907

import math
from sklearn.metrics import mean_squared_error
y_predict = lm.predict(x)
regression_model_mse = mean_squared_error(y_predict, y)
math.sqrt(regression_model_mse)

0.04961046369229128
```

Fig. 11 Fourth Linear Regression Model: Calculating R^2 and RMSE Scores

IV. DISCUSSION

Looking at the models that we created, the best model was the Advanced Statistics Linear Regression Model with an R^2 score of 0.9999 and an RSME score of 0.04961. Coming in second was the Individual-Driven Statistics Linear Regression with an R^2 score of 0.5732 and an RSME score of 4.5941. Next was the Team-Driven Non-Scoring Statistics Linear Regression with a R^2 score of 0.36976 and RSME score of 5.5829. The worst model was the Box-Score Statistics Linear Regression with a R^2 score of 0.1282 and a RSME score of 6.5661. Given what we know about the Overall Raptor Score (+/-), it made sense that the Advanced Statistics model was the best at predicting the RAPTOR +/- . This was because it consisted of all predictors that measure in-game productivity and efficiency of a player. Not only that but these predictors are evaluated by the level of possessions which is crucial in the game of basketball.

According to our models, another finding was that individual-driven statistics are valued more than team-driven non-scoring statistics. Initially we made a hypothesis that the team-driven statistics would overshadow the individual-driven statistics because we were attempting to predict the efficacy of the RAPTOR +/- which is how much value a player contributes to the team. It made sense that a player who generates positive outcomes through other teammates would be more beneficial to the team. However, our hypothesis was proven to be wrong.

Lastly, seeing that the Box-Score Statistics Model was the worst, it solidified the fact that not everything seen on paper (NBA stat sheets) is directly correlated to an athlete's value on his/her team. This is seen in many instances in the NBA where some mediocre teams who are on the verge of missing the playoffs have great players with deceiving box-score statistics also known as stat padding. Stat padding is when a player's statistics are magnified but of little benefit to the team's chance of winning. For example, Russell Westbrook looked like the best player on paper in 2018. However, he was not the best team player which was proven by their team's ranking (6), barely leading his team to the playoffs.

Ultimately, people should look at advanced statistics more than any other statistic that is given. Secondly, individual-driven statistics proved to show some importance and should be valued more than team-driven statistics and box-score statistics.

V. ATTRIBUTION

The milestones of the project were done collaboratively through Google Docs such as the abstract assignment and data story. The data analysis on *Jupyter Notebook* was done through Zoom and in-person meetings. For the final paper, the abstract, figures, and discussion were done by Chris Park and the introduction, methodology, and references were done by Kevin Zhang. Lastly, the discussion portion was done together.

VI. REFERENCES

- [1] "2019-2020 NBA Player Stats," *13th Anniversary*. [Online]. Available: <https://www.nbastuffer.com/2019-2020-nba-player-stats/>. [Accessed: 30-Nov-2020].
- [2] "The Best NBA Players This Season, According To RAPTOR," *FiveThirtyEight*, 12-Oct-2020. [Online]. Available: <https://projects.fivethirtyeight.com/2020-nba-player-ratings/>. [Accessed: 30-Nov-2020].
- [3] Fivethirtyeight, "fivethirtyeight/nba-player-advanced-metrics," *GitHub*. [Online]. Available: <https://github.com/fivethirtyeight/nba-player-advanced-metrics>. [Accessed: 30-Nov-2020].
- [4] N. Silver, "How Our RAPTOR Metric Works," 10-Oct-2019. [Online]. Available: https://fivethirtyeight.com/features/how-our-raptor-metric-works/?fbclid=IwAR0uC5Yipiy-54_cPGg4STxqrPBNYj3UghRbj4fxgPrS2greTlflOxRrkI. [Accessed: 30-Nov-2020].
- [5] N. Silver, "We're Predicting The Career Of Every NBA Player. Here's How.," *FiveThirtyEight*, 09-Oct-2015. [Online]. Available: <https://fivethirtyeight.com/features/how-were-predicting-nba-player-career/>. [Accessed: 30-Nov-2020].