# Predicting Team Success in the NBA using Oliver's Four Factors

Taaj Cheema

## 1 Introduction

How can we distill the success of a basketball team? Can the performance of a team be modeled by just a few key components? There is no shortage of box score and advanced statistics to characterize different aspects of a team's performance. Ideally, we would like to use a small selection of readily available and understandable statistics that can capture most of the variation in a given target metric of success. Furthermore, can we can use these features to predict a team's wins and average margin of victory using historical National Basketball Association data?

In this paper we will test if the four factors posited by Dean Oliver have statistically significant prediction power on the success of an NBA team. Additionally, we will use various machine learning models to predict the success of basketball teams in the NBA using the four factors as features. We define success with two output variables, number of wins in a season and average margin of victory. Lastly, we will calculate the respective weights of these four factors in predicting success and compare them to the weights proposed by Oliver and Ed Küpfer.

## 2 Background

The NBA has been undergoing an offensive renaissance over the past 15 years. We have seen teams prioritize 3-point attempts, improve their effective field goal percentage, and reduce turnovers, to an extent never seen before [5]. This focus on offense has led to record high points per game [17]. There are many factors that have contributed to this offensive boon, but the one that gets the most media attention is the 3-point shot [4] [7] [7] [18].

In his 2004 book, *Basketball on Paper: Rules and Tools for Performance Analysis*, Dean Oliver noted that for basketball, there are four key areas that winning teams excel at [13]. Oliver also notes that for basketball, points per possession, not points per game is the benchmark for offensive efficiency. Oliver asserts that the four factors encapsulate effective basketball possessions and can be used to measure the efficiency of a basketball team's offense or defense. The four offensive factors are effective field goal percentage, turnover percentage, offensive rebound percentage, and free throw factor. The corresponding four defensive factors are opponent effective field goal percentage, opponent turnover percentage, defensive rebound percentage, and opponent free throw factor [6].

$$eFG\% = \frac{FG + 0.5 * 3FG}{FGA} = \text{effective field goal percentage} \tag{1}$$

$$TOV\% = \frac{TOV}{FGA + 0.44 * FTA + TOV} = \text{turnover percentage} \tag{2}$$

$$OREB\% = \frac{OREB}{OREB + DREB_{OPP}} = \text{offensive rebound percentage} \tag{3}$$

$$FTF = \frac{FTM}{FGA} = \text{free throw factor} \qquad (4)$$

Each of the four factors has different weights representing its importance to a teams offensive or defensive efficiency. Oliver assigns a weight of 40% to eFG%, 25% to TOV%, 20% to OREB%, and 15% to FTF respectively [13]. In later analysis by Houston Rockets analyst Ed Küpfer, Küpfer assigns a weight of 10, 6, 3, and 3 to eFG%, TOV%, OREB%, and FTF respectively [11]. This works out to approximately 45.45%, 27.27%, 13.64% and 13.64%.

The four factors effectively evaluate how good a team is at shooting the ball, how good a team is at protecting the ball, how good a team is at winning rebounds, and how good a team is at getting to the free throw line and successfully converting these attempts. In theory, the four factors should be able to predict the success of a NBA team.

This research intends to:

1. Test if the four factors are statistically significant predictors of success in the NBA

2. Use the four factors to predict a team's number of wins and average margin of victory in the NBA using historical data and various machine learning models

3. Compare the observed weightings of the different factors in our models to those proposed by Dean Oliver and Ed Küpfer.

All of the code, the write up, as well as additional resources related to this paper are available at:
https://github.com/taajcheema/four_factors

## 3   Related Work

Similar research to model NBA games has been performed before.

In a 2003 paper, *A Multivariate Statistical Analysis of the NBA*, by Lori Hoffman and Maria Joseph, the researchers predicted if teams would qualify for the NBA playoffs using methods including principal components analysis and discriminant analysis. They identified 5 principal components that were used to accurately predict 26 out of 29 teams as either playoff or non-playoff teams for the 2002-2003 NBA season [8].

A 2008 paper by Beckler et al. titled *NBA Oracle* sought to predict the winner of individual NBA games from the 1991-1992 through 1996-1997 seasons using linear regression, logistic regression, support vector machines, and neural networks. They were able to correctly predict the winner 73% of the time using linear regression and a total of 60 features. This is compared to a baseline of 50% by choosing a winner randomly [3].

A 2013 paper by Amorim Torres, *Prediction of NBA games based on Machine Learning Methods*, also predicted the winner of individual NBA matchups using statistical methods. He was able to correctly predict the winner of 68.44% of NBA games in the 2006-2007 to 2011-2012 seasons using a multilayer neural network and 8 features [14].

A 2014 paper by Lin et al., *Predicting National Basketball Association Winners*, predicted the winner of NBA games using various machine learning methods. They used 16 independent variables to predict the winner. The researchers found that a basketball team's win record plays a central role in determining their likeliness of winning future games. When they removed the win record feature from their model, the classification accuracy decreased significantly. They concluded that traditional box score statistics fail to represent all of a team's success on the court [12].

A 2016 paper by Avalon et al., *Various Machine Learning Approaches to Predicting NBA Score Margins*, predicted the margin of victory of two NBA teams using linear regression, Gaussian discriminant analysis,

principal component analysis coupled with support vector machines, random forest and adaptive boosting. They started with 218 features that they reduced into the most important components using PCA. They trained their models on 1052 games and tested them on 264 games all from the 2013 - 2014 NBA seasons. With a PCA and SVM model, they were able to predict the margin of victory within a 2 point score margin of error for 90% of the teams in the test set [1].

Our paper differs from the related works in that we are attempting to predict the number of wins and average margin of victory for teams over a season, rather than predicting the winner in individual matchups. We are additionally testing the significance of the four factors of basketball success proposed by Dean Oliver, and using them as features in our models. Lastly, we are assigning a weighting to each of the factors using our models.

## 4   Data

All data was obtained from Basketball Reference [2]. The website has an incredible amount of NBA data available.

A dataframe is created including statistics for all NBA teams from 11 seasons spanning from the 2008-2009 season to the 2018-2019 season. The data is only for the regular season, it does not include information on postseason games. The features included in our dataframe are WINS, MOV (5), NRTG (8), eFG% (1), TOV% (2), OREB% (3), FTF (4), OPPeFG%, OPPTOV%, DREB% and OPPFTF. This data is used to create our predictive models and to calculate our weightings of the four factors.

$$MOV = \frac{TotPTS - TotPTS_{OPP}}{82} = \text{margin of victory} \tag{5}$$

$$ORTG = 100 * \frac{TotPTS}{TotPoss} = \text{offensive rating} \tag{6}$$

$$DRTG = 100 * \frac{TotPTS_{OPP}}{TotPoss_{OPP}} = \text{defensive rating} \tag{7}$$

$$NRTG = ORTG - DRTG = \text{net rating} \tag{8}$$

A dataframe is also created on NBA league averages per 100 possessions for 20 NBA seasons spanning the 2000-2001 to the 2019-2020 season. This data is used to visualize historical changes in the four factors.

## 5   Methods

The data needs to undergo certain preprocessing tasks before we can use it for modeling.

Certain features in the dataframe are scaled. The eFG%, FTF, OPPeFG% and OPPFTF features are multiplied by 100. This is done because in the original dataframe, these features are in decimal form, whereas the other features that will be used as predictors for the model are in percentage form. Furthermore, the 2011-2012 NBA season was a lockout season, so teams only played a total of 66 games. The wins feature for records in the 2011-2012 season is scaled by (82/66). Additionally, some of column names in the dataframe are changed to make them easier to work with.

Once our data is prepared, we analyze it graphically with ggplot2 by creating a variety of graphs. This is done to check the basic assumptions of mulitple linear regression. We create scatterplots of the predictor and response variables to see if there is a relationship between them. We also create boxplots to check for

outliers, along with density plots to make sure the features follow a normal distribution. We additionally use the Durbin-Watson test to check for potential autocorrelation in the residuals of our regression models. Finally, we use pairplots and correlation plots to check for multicollinearity.
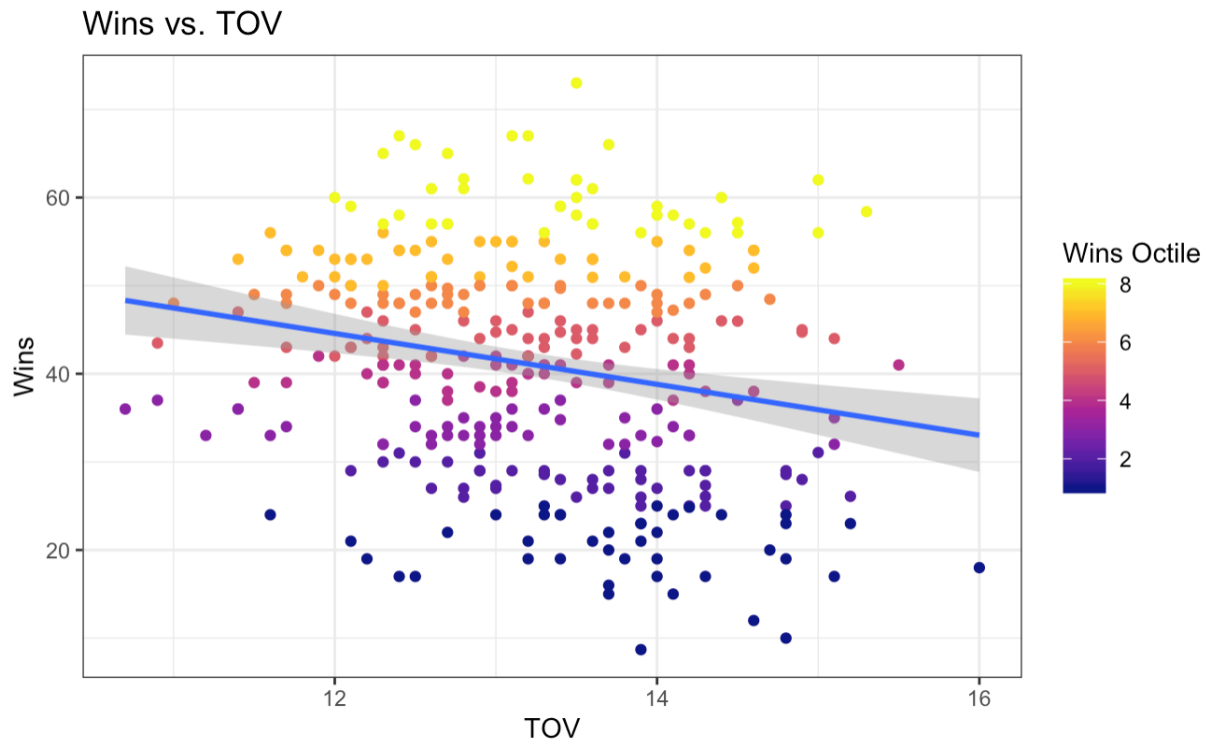


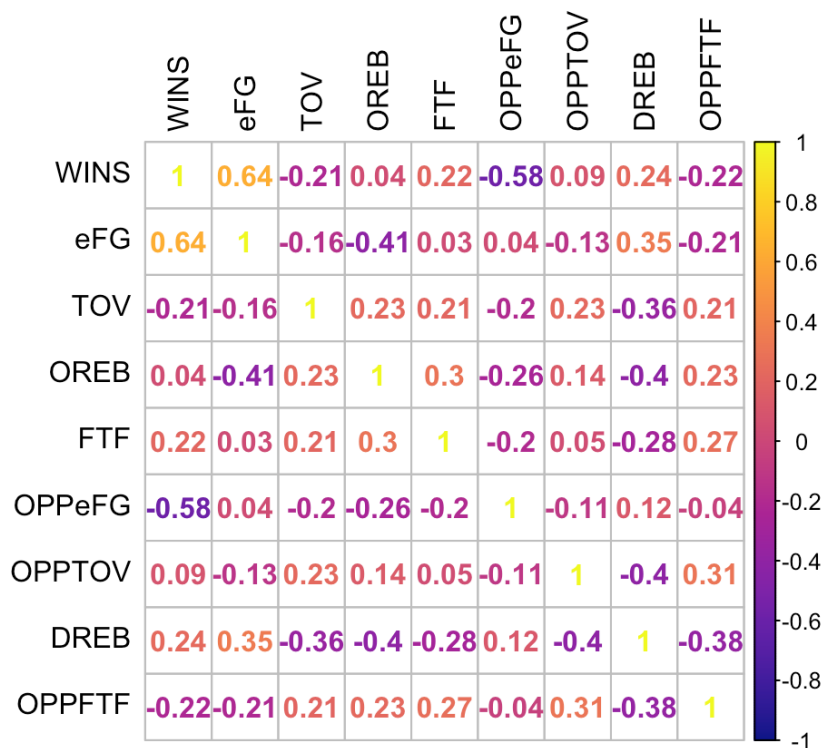**Figure 1:** Scatter plot of WINS vs TOV%



**Figure 2:** Correlation matrix

4

This data is then split into random train and test subsets. 70% of the original dataset is randomly partitioned to ff_train while the other 30% is partitioned to ff_test. This is done using the sample function in R. The set.seed() function is used to create a reproducible output.

Two linear regression models are fit using the lm() function in R using the ordinary least squares method. The OLS method fits a linear model and minimizes the residual sum of squares between the observed and predicted targets in the dataset [9].

The first model fits the four factors as features to predict wins.

$$\begin{aligned} WINS = \beta_0 &+ \beta_1 * eFG\% + \beta_2 * TOV\% + \beta_3 * OREB\% \\ &+ \beta_4 * FTF + \beta_5 * OPPeFG\% + \beta_6 * OPPTOV\% \\ &+ \beta_7 * DREB\% + \beta_8 * OPPFTF + \epsilon \end{aligned} \tag{9}$$

The second model fits the four factors as features to predict margin of victory.

$$\begin{aligned} MOV = \beta_0 &+ \beta_1 * eFG\% + \beta_2 * TOV\% + \beta_3 * OREB\% \\ &+ \beta_4 * FTF + \beta_5 * OPPeFG\% + \beta_6 * OPPTOV\% \\ &+ \beta_7 * DREB\% + \beta_8 * OPPFTF + \epsilon \end{aligned} \tag{10}$$

After training our model we get our fitted equations

$$\begin{aligned} WINS = &-53.185 + 3.899 * eFG\% - 3.392 * TOV\% + 1.111 * OREB\% \\ &+ 0.708 * FTF - 3.765 * OPPeFG\% + 2.902 * OPPTOV\% \\ &+ 0.887 * DREB\% - 0.729 * OPPFTF + \epsilon \end{aligned} \tag{11}$$

$$\begin{aligned} MOV = &-38.851 + 1.436 * eFG\% - 1.323 * TOV\% + 0.438 * OREB\% \\ &+ 0.296 * FTF - 1.420 * OPPeFG\% + 1.170 * OPPTOV\% \\ &+ 0.386 * DREB\% - 0.288 * OPPFTF + \epsilon \end{aligned} \tag{12}$$

Summary statistics for the lmWINS and lmMOV models can be seen below in figures 3 and 4.

```
Call:
lm(formula = WINS ~ eFG + TOV + OREB + FTF + OPPeFG + OPPTOV +
    DREB + OPPFTF, data = ff_train)

Residuals:
     Min      1Q   Median      3Q      Max
-10.6982  -2.1067  -0.1126   2.0078   9.0728

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -53.18468   14.54796  -3.656  0.00032 ***
eFG           3.89934    0.11042  35.314  < 2e-16 ***
TOV          -3.39168    0.24879 -13.633  < 2e-16 ***
OREB          1.11070    0.09013  12.323  < 2e-16 ***
FTF           0.70791    0.09304   7.609 7.89e-13 ***
OPPeFG       -3.76483    0.11439 -32.913  < 2e-16 ***
OPPTOV        2.90168    0.21229  13.668  < 2e-16 ***
DREB          0.88666    0.11764   7.537 1.22e-12 ***
OPPFTF       -0.72897    0.10377  -7.025 2.62e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.158 on 221 degrees of freedom
Multiple R-squared:  0.9435,    Adjusted R-squared:  0.9414
F-statistic: 461.2 on 8 and 221 DF,  p-value: < 2.2e-16
```

**Figure 3:** Summary statistics of the lmWINS model

```
Call:
lm(formula = MOV ~ eFG + TOV + OREB + FTF + OPPeFG + OPPTOV +
    DREB + OPPFTF, data = ff_train)

Residuals:
     Min       1Q   Median       3Q      Max
-2.87201 -0.41263 -0.00396  0.46326  1.89270

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -38.85131    3.16528  -12.27   <2e-16 ***
eFG           1.43555    0.02402   59.75   <2e-16 ***
TOV          -1.32334    0.05413  -24.45   <2e-16 ***
OREB          0.43792    0.01961   22.33   <2e-16 ***
FTF           0.29640    0.02024   14.64   <2e-16 ***
OPPeFG       -1.42002    0.02489  -57.06   <2e-16 ***
OPPTOV        1.17049    0.04619   25.34   <2e-16 ***
DREB          0.38560    0.02559   15.07   <2e-16 ***
OPPFTF       -0.28811    0.02258  -12.76   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.6872 on 221 degrees of freedom
Multiple R-squared:  0.9806,    Adjusted R-squared:  0.9799
F-statistic:  1398 on 8 and 221 DF,  p-value: < 2.2e-16
```

**Figure 4:** Summary statistics of the lmMOV model

We also performed diagnostic tests to make sure that we do not violate any assumptions of multiple linear regression. Using the plot() function in R, we created plots for residuals vs. fitted, normal Q-Q, scale-location, and Cook's distance. We also checked the variance inflation factor of our models to make sure there is no multicollinearity in our features. The results of the diagnostic tests can be seen in figures 5 and 6.
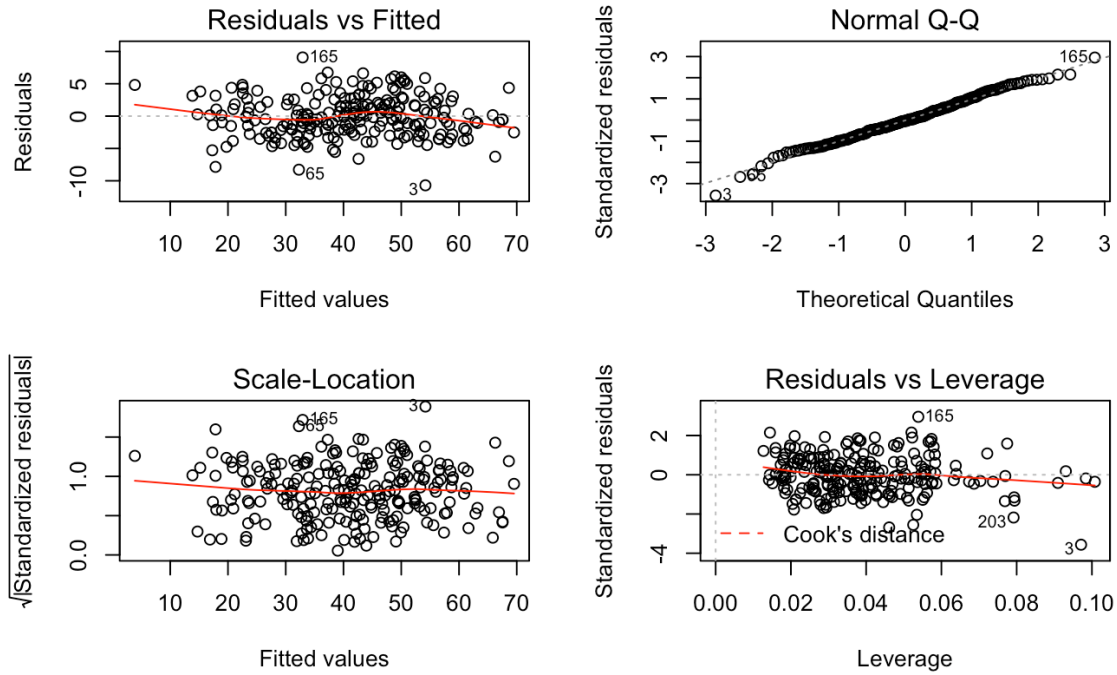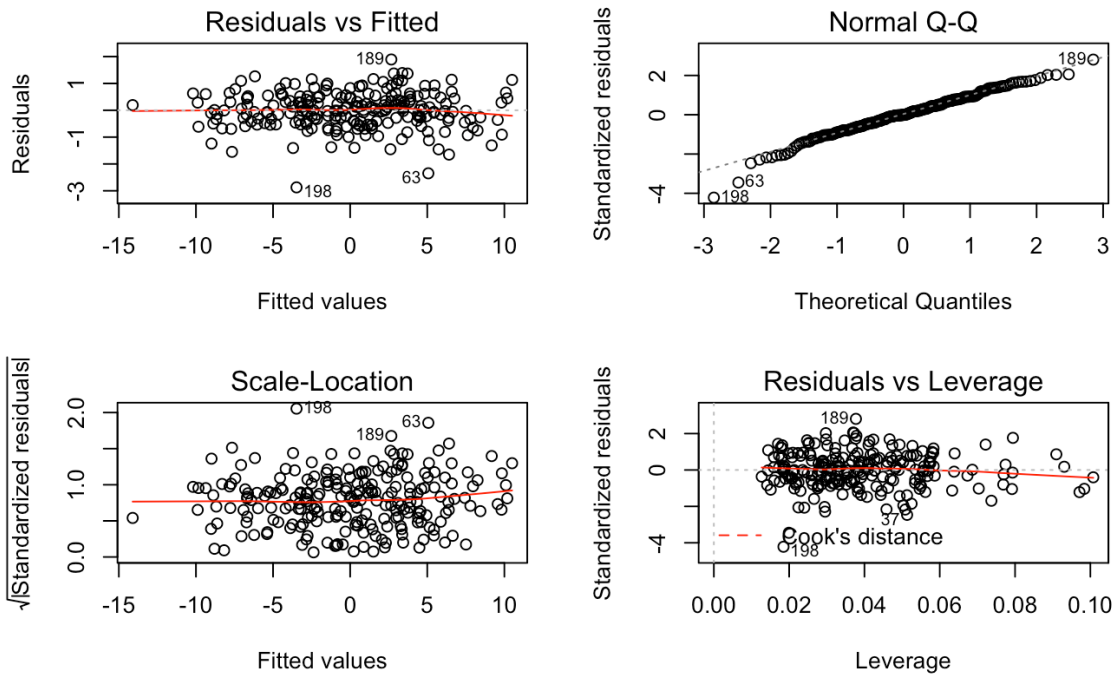


**Figure 5:** Diagnostic plots for the lmWINS model



**Figure 6:** Diagnostic plots for the lmMOV model

We additionally made a random forest model and a gradient boosting machine model to predict Wins and margin of victory.

# 6    Results

Wins and margin of victory were predicted using three different models

1. Multiple Linear Regression

2. Random Forest

3. Gradient Boosting Machine

For the final model, we decided to use multiple linear regression. The linear regression models we fit performed better than the random forest and gradient boosting machine models. Had we further tweaked and tested different hyperparameter configurations for the other models, they may have been able to perform better than the linear regression model.

Another reason we picked multiple linear regression as the final model type, is because it is the simplest model. Simple models are often easier to explain. In this case, the other models are black-box models, which can potentially provide more accuracy, but at the expense of explainability compared to a model like multiple linear regression.

## 6.1   Are the four factors statistically significant predictors of success in the NBA?

The p-values for both the lmWINS and lmMOV models are practically zero. This is strong evidence that we can reject the null hypothesis that there is no relationship between our response and predictor variables. Additionally, both the lmWINS and lmMOV models produce statistically significant p-values for all included features at a 0.001 significance level. All of the included predictors have statistically significant value in predicting wins and margin of victory respectively.

From the plots in figures 5 and 6 we can see that we reasonably meet the assumptions of linear regression. The residual plot is mostly flat for both the lmWINS and lmMOV models. There are a few points skew the line, but overall the relationship is linear. In the scale-location plot, points are normally spread. There is no evidence of potential heteroskedasticity. In the Q-Q plot, points are very close to the diagonal line. We can assume that values are normally distributed. From the Cook's Distance plot, we can see there are a few outliers. We choose to leave these points in because we know that they were real data, and not errors. We also checked the variance inflation factor to make sure there is no multicollinearity in our features. All of the VIF's for our features were below 2 and we can conclude that there is no multicollinearity in our data.

Effective Field Goal percentage has been trending upwards in the NBA over the past twenty years. Turnover percentage, offensive rebound percentage and free throw factor have all been trending downward.

## 6.2   Model Observed vs. Predicted

The lmWINS model was successful in predicting NBA wins and achieves an adjusted R-squared value of 0.9435. The can be interpreted as our model is able to explain 94.35% of the variation in wins. When we fit the lmWINS model on the test data, we achieve a an adjusted R-squared value of 0.9300. This shows that our lmWINS model generalizes well. The linear model fit well and could potentially be used to predict future wins.

The lmMOV model achieved a R-squared value of 0.9785. The can be interpreted as our model is able to explain 97.85% of the variation in margin of victory. When we fit the lmMOV model to the test data, we achieve a R-squared value of 0.9707. The model again generalizes well on the test data.

The results were very good. The four factors proposed by Oliver can be highly effective in predicting wins and margin of victory in the NBA.

| Team | Season | Observed Wins | Predicted Wins | Observed Margin of Victory | Predicted Margin of Victory |
|---|---|---|---|---|---|
| Milwaukee Bucks | 2018 - 2019 | 60 | 61.01426 | 8.87 | 7.5019316 |
| Denver Nuggets | 2018 - 2019 | 54 | 48.44225 | 3.95 | 2.7918936 |
| Boston Celtics | 2018 - 2019 | 49 | 51.48235 | 4.44 | 3.8579222 |
| Brooklyn Nets | 2018 - 2019 | 42 | 41.65049 | 0.71 | 0.1291505 |
| Charlotte Hornets | 2018 - 2019 | 39 | 38.43315 | -0.23 | -0.8615088 |
| Washington Wizards | 2018 - 2019 | 32 | 35.05902 | -2.90 | -2.4180218 |
| Cleveland Cavaliers | 2018 - 2019 | 19 | 17.24568 | -9.34 | -9.0182489 |
| Phoenix Suns | 2018 - 2019 | 19 | 18.57312 | -9.61 | -8.7940497 |
| Houston Rockets* | 2017 - 2018 | 65 | 59.85960 | 8.48 | 7.2595134 |
| Toronto Raptors* | 2017 - 2018 | 59 | 58.05547 | 7.78 | 6.3801864 |
| Golden State Warriors* | 2017 - 2018 | 58 | 58.90024 | 5.98 | 6.3904133 |
| Portland Trail Blazers* | 2017 - 2018 | 49 | 44.46868 | 2.60 | 1.3311391 |
| Indiana Pacers* | 2017 - 2018 | 48 | 44.94865 | 1.38 | 1.4623198 |
| San Antonio Spurs* | 2017 - 2018 | 47 | 47.55992 | 2.89 | 2.6330804 |
| Milwaukee Bucks* | 2017 - 2018 | 44 | 42.03005 | -0.30 | 0.3325235 |
| Chicago Bulls | 2017 - 2018 | 27 | 22.34545 | -7.04 | -6.9689690 |
| Boston Celtics* | 2016 - 2017 | 53 | 46.97189 | 2.63 | 2.1073394 |
| Memphis Grizzlies* | 2016 - 2017 | 43 | 39.29831 | 0.49 | -0.4522467 |
| Miami Heat | 2016 - 2017 | 41 | 45.34191 | 1.06 | 1.5203129 |
| Portland Trail Blazers* | 2016 - 2017 | 41 | 40.81297 | -0.52 | -0.2177054 |
| Charlotte Hornets | 2016 - 2017 | 36 | 40.88568 | 0.20 | 0.2178854 |
| Minnesota Timberwolves | 2016 - 2017 | 31 | 37.46738 | -1.11 | -1.2279486 |
| Philadelphia 76ers | 2016 - 2017 | 28 | 26.72349 | -5.70 | -5.5494226 |
| Phoenix Suns | 2016 - 2017 | 24 | 27.68881 | -5.63 | -4.9422247 |
| Brooklyn Nets | 2016 - 2017 | 20 | 25.10779 | -6.73 | -6.2136536 |

**Figure 7:** Observed vs. predicted values for wins and margin of victory for the first 25 observations in our test data

## 6.3   Weightings

We obtained weights for the different features of the lmWINS and lmMOV models by averaging the absolute value of each of the coefficients. An example where the relative weight of the eFG% factor for the lmWINS model is derived can be seen in equations 13, 14 and 15.

$$
\begin{aligned}
sum\_all\_lmWINS = |3.89934| + |-3.39168| + |1.11070| + \\
|0.70791| + |-3.76483| + |2.90168| + \\
|0.88666| + |-0.72897|
\end{aligned}
\tag{13}
$$

$$
sum\_lmWINS = |3.89934| + |-3.76483|
\tag{14}
$$

$$
weight\_eFG\_lmWINS = 100 * (sum\_eFG\_lmWINS / sum\_all\_lmWINS)
\tag{15}
$$

The order of importance of the weights that we obtained with our model agree with both Oliver and Küpfer. We found shooting or EFG% to be the most important of the four factors to a team's success. Turnovers, rebounding and free throw factor followed respectively. However, the weights we obtained for the four factors are different. We created our model on regular season NBA games from the 2008 - 2009 season to the 2018 -

2019 season. Oliver originally created his model in 2002 and Küpfer in 2005, so it is possible the importance of the weights has shifted as the NBA has changed.

| Factor | Oliver | Küpfer | lmWINS |
|---|---|---|---|
| Shooting | 40.00% | 45.45% | 44.07% |
| Turnovers | 25.00% | 27.27% | 36.19% |
| Rebounding | 20.00% | 13.64% | 11.48% |
| Free Throw Factor | 15.00% | 13.64% | 8.26% |

**Figure 8:** Comparison of weights for the four factors

# 7    Conclusion

A multiple linear regression model to forecast NBA teams' wins and margins of victory using Oliver's Four Factors was successfully created. The Four Factors provide us an effective way to represent a team's overall performance in a relatively simple and easy to understand model.

The weightings we obtained for the four factors were similar to those proposed by Dean, with our model prioritizing turnovers and deprioritizing rebounding and free throw factor.

A future work will focus on predicting winners and margin of victory of individual NBA matchups, rather than overall season predictions, using a similar four factors model. This would entail training a model on 1230 games per NBA season, rather than the 30 data points of summary statistics per NBA season we trained the models for this paper on. Another future approach will involve looking at only NBA playoff games. It would be interesting to see how well the four factors can predict success in the NBA postseason as it is often much more competitive and high stakes than the regular season. One final future approach will use a binary logistic regression model to predict probabilities of two opposing teams winning in matchups, using the four factors as features.

# References

[1] Avalon, Grant, et al. "Various Machine Learning Approaches to Predicting NBA Score Margins.".
    http://cs229.stanford.edu/proj2016/report/Avalon_balci_guzman_various_ml_approaches_NBA_Scores_report.p

[2] "Basketball Reference". *Basketball Reference*.
    www.basketball-reference.com.

[3] Beckler, Matthew, et al. "NBA Oracle".
    https://www.mbeckler.org/coursework/2008-2009/10701_report.pdf.

[4] Cohen, Ben. "The First Shots of the NBA's 3-Point Revolution". *The Wall Street Journal*, 12 Mar. 2018,
    www.wsj.com/articles/the-first-shots-of-the-nbas-3-point-revolution-1523542076.

[5] Dubin, Jared. "The NBA's Other Offensive Revolution: Never Turning The Ball Over". *FiveThirtyEight*,
    14 Mar. 2019,
    www.fivethirtyeight.com/features/the-nbas-other-offensive-revolution-never-turning-the-ball-over/.

[6] "Four Factors". *Basketball Reference*.
    www.basketball-reference.com/about/factors.html.

[7] "Gametime: 3-Point Revolution Part One". *NBA*.

[8]  Hoffman, Lori, et al. "A Multivariate Statistical Analysis of the NBA ".
     http://www.units.miamioh.edu/sumsri/sumj/2003/NBAstats.pdf.

[9]  James, Gareth, et al. *An Introduction to Statistical Learning: With Applications in R*. Springer, 2017.

[10] Kram, Zach. "The 3-Point Boom Is Far From Over". *The Ringer*, 27 Feb. 2019,
     www.theringer.com/nba/2019/2/27/18240583/3-point-boom-nba-daryl-morey.

[11] Kubatko, Justin; Oliver, Dean; Pelton, Kevin; and Rosenbaum, Dan. (2007) "A Starting Point for
     Analyzing Basketball Statistics". *Journal of Quantitative Analysis in Sports*. Volume 3: Issue 3, Article
     1.

[12] Lin, Jasper, et al. "Predicting National Basketball Association Winners".

[13] Oliver, Dean. *Basketball on Paper: Rules and Tools for Performance Analysis*. Potomac Books, 2004.

[14] Torres, Renato. "Prediction of NBA games based on Machine Learning Methods".
     https://homepages.cae.wisc.edu/ ece539/fall13/project/AmorimTorres_rpt.pdf.

[15] Uudmae, Jaak. "CS229 Final Project: Predicting NBA Game Outcomes".
     http://cs229.stanford.edu/proj2017/final-reports/5231214.pdf.

[16] Winston, Wayne. *Mathletics: How Gamblers, Managers, and Sports Enthusiasts Use Mathematics in
     Baseball, Basketball, and Football*. Princeton University Press, 2012.

[17] Woike, Dan. "107 Is the New 100.' How Scoring Trends in the NBA Are Setting a New Bar for What a
     Good Defense Is.". *Los Angeles Times*, 29 Oct. 2018,
     www.https://www.latimes.com./sports/nba/la-sp-nba-defenses-20181029-story.html.

[18] Young, Shane. "The NBA's 3-Point Revolution Continues To Take Over.". *Forbes*, 1 Dec. 2019,