

# Predicting NBA Regular Season Results Based on some popular statistics

Chenjie Li, Qiao Qiao, Giyoon Kwag

November 26, 2018

## 1. Problem Description

### 1.1 Background

NBA (National Basketball Association) is one of the most successful basketball league around the world. And since it is popular, there exists tons of websites recording the statistics related to the game. They not only provide simple metrics like Field Goal Percentage, Average Points, Average Assists, but also some advanced metrics like PER, True shooting percentage and much more.

Based on those statistics, people have been exploring what kinds of metrics can result in a team's success or failure as a season, or what metrics or combination of different metrics can result in team's success.

In this report, we investigated the relationship between the regular season results of the team and some popular metrics in two different perspectives:

- **Teams' Overall Effects: Team's Offensive Efficiency (TOE)**
- **Teams' Net Rating**

From those 3 metrics above, we proposed new TOE and NEW Teams' Net Rating and compared models by using "Cross Validation method", we found out that our new version of those two metrics are stronger correlated with teams' win ratios than the results for old versions. In this report, we collected data for the past 5 seasons (2014 – 2018). Based on

the results we get from 2014-2017 seasons, we made some predictions for the results of this 2018 season and compare the results between our predictions and the real results.

## 2.Data Source

There are many “data hubs” available, our data source is mainly from [insider.espn.com/nba/hollinger/statistics](http://insider.espn.com/nba/hollinger/statistics) ,<https://www.basketball-reference.com/> and <http://www.espn.com/nba/statistics>

To be able to work with data, we developed a “web crawler” using Python “Scrapy” module. By using the program we got the formatted CSV file so that we can easily input them in R.

Furthermore, since we are familiar with Database and SQL language, we also put the data in a database so we can manipulate the data easily by Queries.

## 3.Win Ratio Analysis

### 3.1 New TOE introduction and Experiments

Offensive Efficiency, as the name suggests, is a parameter measuring the efficiency of the offense. Here we explain this intuition by giving an example:

Player A took 5 shots and made 4 of them; Player B took 5 five shots, and made 2 of them. Therefore, we can conclude that player A is more efficient than Player B.

Of course, this is a oversimplified example.

There are a lot different kinds of definitions for Player and Team’s Efficiency: In Shea,Stephen M’s book<sup>[1]</sup>, They define the Offensive Efficiency for Team as:

$$TOE = \frac{FG}{FGA - ORB + TO}$$

Where FG is the field goals made, FGA is the field goals attempts, ORB is the number of offensive rebounds, TO represents the number of turnovers.

It is of course reasonable to guess that the more efficient a team’s offensive is, the more number of games a team is going to win.This idea was proved in their book: “The top 5 teams in TOE in the 2012-2013 season all won at least 56 games.”

This book was written and published five years ago. As time goes by, NBA statistics has been getting more and more comprehensive. Now we want to propose a new version of Team's offensive efficiency.

The formula we proposed is:

$$New\ TOE = \frac{2PTM + 1.5 * 3PTM}{FGA - ORB * ORBS + TO * OPPOTS}$$

Where  $2PTM$  means 2 points field goals made;  $3PTM$  means 3 points field goals made;  $ORBS$  means the offensive rebound scoring rate;  $OPPOTS$  means opponent scoring rate off the team's turnovers.

Now let's motivate our formula. An individual offensive possession could lead up to 4 possible results: taking a 2 point shot, taking a 3 point shot, missing a shot but getting an offensive rebound, and turning the ball over. We assume that a 3 point shot made is 1.5 times as great as a 2 point shot made. Offensive rebound could compensate the results of missing the shot. Since most (not all of them, of course) second chance points are 2 points, we ignored the effect of the 3 point second chance points. But here we times Offensive rebound by "offensive rebound scoring rate", because different teams might have different abilities to convert a offensive rebound to box score. Here we also ignored the effect of 3 points made by opponent off turnovers. And finally, we add a factor " $TO * OPPOTS\%$ " to take the effects of turnovers into consideration.

In this section all our data is from <https://stats.nba.com/><sup>[2]</sup>.

We calculated the results from Season 2014 - Season 2017, here we only show two groups of results. For comprehensive results, please check out the Appendix part.

As Figure 1 and Figure 2 show, our new TOE is much better (significantly higher  $R^2$  value) than evaluating the relationship between Win Ratio and Team Offensive Efficiency. We noticed the coefficient  $\beta_1$  values keep increasing in 3 consecutive years, whereas in 2017 the  $\beta_1$  drops again.

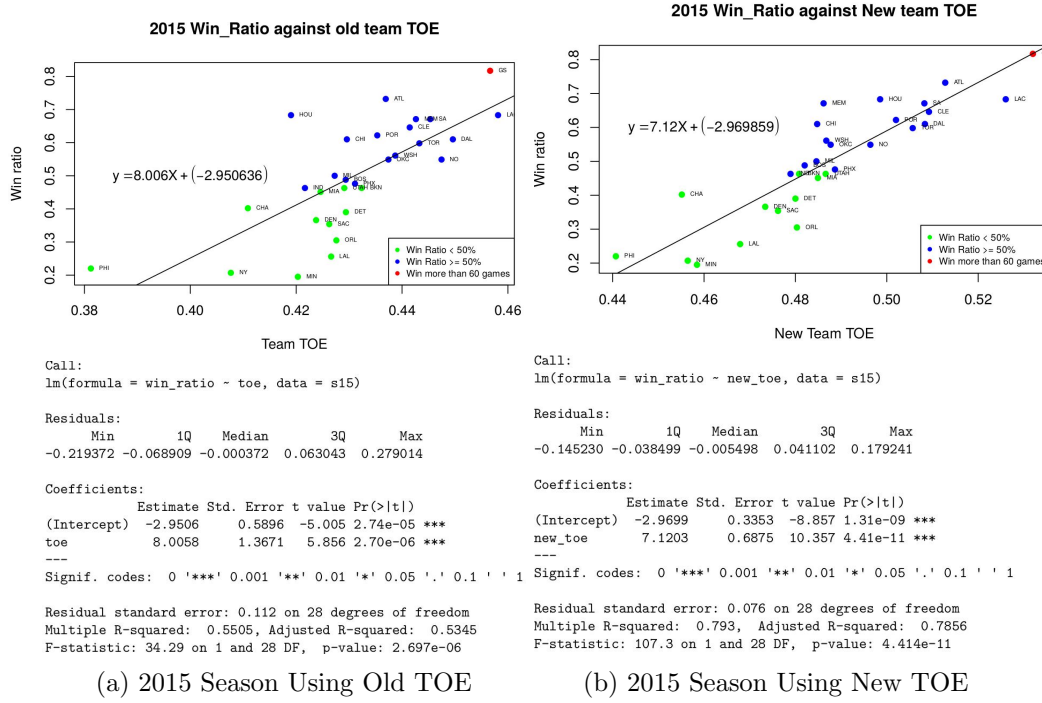


Figure 1: Comparision Between TOE and New TOE in 2015 Season

Season	Old TOE	New TOE	Season	Old TOE	New TOE
2014	0.438	0.602	2014	7.558	6.623
2015	0.551	0.793	2015	8	7.12
2016	0.608	0.776	2016	8.874	7.811
2017	0.45	0.667	2017	6.645	6.399

Figure 2: Comparison of  $R^2$ (left) and  $\beta_1$ (right)

### 3.2 Module comparison and validation

#### Train Test and $AIC_p$

Now we pay our attention to the evaluation of the model. In this section we use “Correlation Accuracy”. A simple correlation between the actuals and predicted values can be used as a form of accuracy measure. A higher correlation accuracy implies that the actuals and predicted values have similar directional movement, i.e. when the actuals values increase the predicted values also increase and viceversa<sup>[3]</sup>.

Also, we calculated  $AIC_p$ , Akaike information criterion<sup>[4]</sup> to compare between old model and the new model.

Statistic	Train:Test Ratio	Old TOE	New TOE
Train,Test Score	5:1	0.757	0.810
Train,Test Score	4:1	0.757	0.812
Train,Test Score	3:1	0.754	0.825
Train,Test Score	Avgerage	0.757	0.816
$AIC_P$	5:1	-152.2682	-186.1824
$AIC_P$	4:1	-145.4626	-178.7008
$AIC_P$	3:1	-138.7349	-166.051
$AIC_P$	Average	-136.489	-176.978

Table 1: Tran Test Accuracy And  $AIC_p$

From Table 1 we could clearly see that our new TOE has better correlation accuracy and smaller  $AIC_p$  values.

### K-Fold Cross Validation

In this section, we compare the module performances by using Cross Validation, which, in essence, is just a systematic evaluation by implementing “train and test” procedure. We splitted our data into ‘k’ mutually exclusive random sample portions. Each time, we keep one of the portions as test data, we build the model on the remaining (k-1 portion) data and calculate the mean squared error of the predictions. This is done for each k random samples. In this report, we chose  $K = 5$  and  $K = 4$  as our “fold number”, which means we divided our data(120 points) into 5 and 4 groups, respectively.

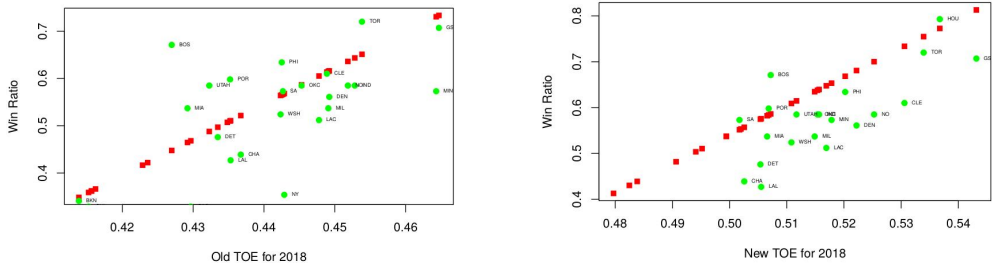
The results are shown below:

K value	Old Toe	New Toe
4	$0.45 \pm 0.18$	$0.50 \pm 0.23$
5	$0.47 \pm 0.12$	$0.51 \pm 0.27$

Table 2:Correlation Accuracy(  $\pm 95\%$ CI) for Old TOE and New TOE

After model validation procedures, we believe that our proposed new model has stronger correlations with win ratios. Now it's time to see the prediction results for 2018 seasons.

The Results are shown below:



because of today's NBA trend: Most Teams in the league are paying more attention on their offense, whereas they are ignoring the importance of the defense. First, let's define Net Rating:

In this Net Rating definition, we added what is missing in TOE part: The opponents' effect.

Net Rating:

$$100 * (\frac{FGM}{FGA - ORB * + TO} - \frac{OPFGM}{OPFGA - OPORB + OPTO})$$

Where:

$$OPTOE = \frac{OPFGM}{OPFGA - OPORB + OPTO}$$

New Net Rating:

$$100 * (\frac{2FGM + 1.5 * 3FGM}{FGA - ORB * ORBS + TO * OPPOTS} - \frac{OP2FGM + 1.5 * OP3FGM}{OPFGA - OPORB * OPORBS + OPTO * POTS})$$

Where:

$$New OPTOE = \frac{OP2FGM + 1.5 * OP3FGM}{OPFGA - OPORB * OPORBS + OPTO * POTS}$$

Just as what we've defined in TOE and New TOE, we added opponents' performance. OP2FGM is a team's opponents' average 2 point field goals made per gam. OP3PM is a team's opponents' average 3 point field goals made per game; OPFGA is a team's opponents' average field goal attempts; OPORB is a team's opponents' average offensive rebound; OPORBS is a team's opponents' "offensive rebound scoring rate" against this team being evaluated; OPTO is a team's opponents' average turnovers; and at last, POTS the team's scoring rate off the opponents' average turnovers.

After subtraction operation, we times the result by 100, so this means **this team's net win margin in 100 possessions**

## 4.2 Experiments

As Section 3 for TOEs, we also tested our Net Ratings on 2014-2017 seasons. The results for 2015 season are shown below, for comprehensive results, please check out the Appendix

part.

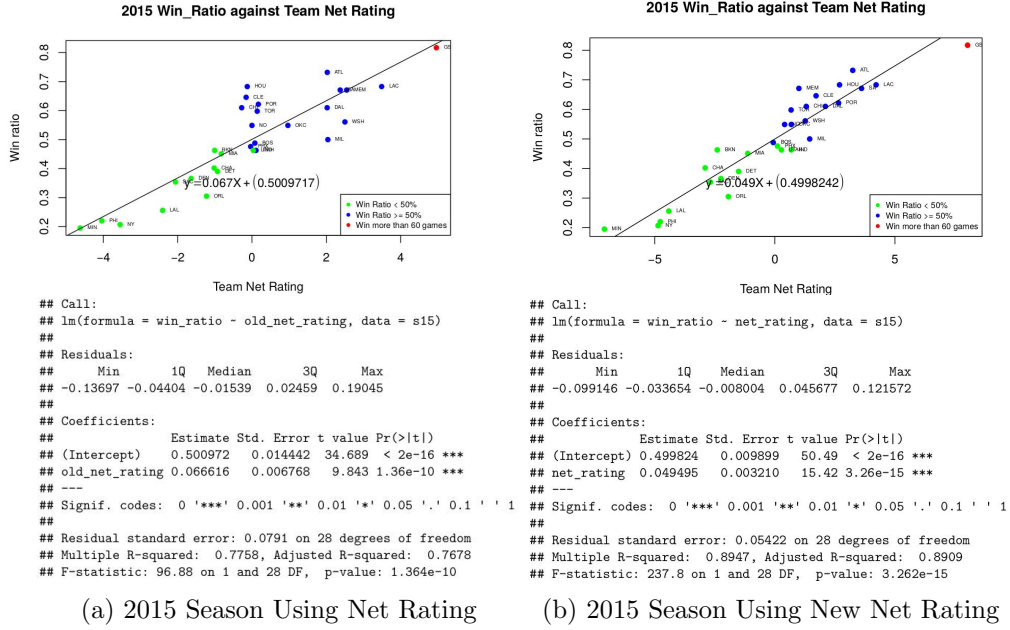


Figure 4: Comparison Between Net Rating and New Net Rating in 2015 Season

Season	Net Rating $R^2$	New Net Rating $R^2$	Season	Net Rating $\beta_1$	New Net Rating $\beta_1$
2014	0.784	0.895	2014	0.065	0.052
2015	0.776	0.895	2015	0.067	0.049
2016	0.822	0.882	2016	0.064	0.048
2017	0.699	0.838	2017	0.061	0.048

Table 3: Comparison of  $R^2$ (left) and  $\beta_1$ (right) for Net Ratings

As Figure 4 and Table 3 show, our New Net Rating is better (higher  $R^2$  value) than evaluating the relationship between Win Ratio and Team Offensive Efficiency. One thing we noticed is that  $\beta_1$  and  $R^2$  are all very steady across those 4 years.



### 4.3 Module comparison and validation

#### Train Test and $AIC_p$

Just as what we did for TOEs, we also evaluated thos two different Net Ratings using  $AIC_p$ , Train Test and K-Fold Cross validation.

Statistic	Train:Test Ratio	Net Rating	New Net Rating
Train,Test Score	5:1	0.885	0.965
Train,Test Score	4:1	0.888	0.962
Train,Test Score	3:1	0.891	0.956
Train,Test Score	Avgerage	0.888	0.961
$AIC_P$	5:1	-235.3074	-289.7162
$AIC_P$	4:1	-225.8254	-279.0621
$AIC_P$	3:1	-212.9546	-261.1197
$AIC_P$	Average	-224.6958	-276.6327

Table 4: Tran Test Accuracy And  $AIC_p$

From Table 1 we could clearly see that our new Net Rating has better correlation accuracy and smaller  $AIC_p$  values than those of Net Rating.

#### K-Fold Cross Validation

In this section, we compare the module performances by using

The results are shown below:

K value	Net Rating	New Net Rating
4	$0.74 \pm 0.12$	$0.86 \pm 0.06$
5	$0.75 \pm 0.12$	$0.86 \pm 0.07$

Table 5:Correlation Accuracy(  $\pm 95\%$ CI) for Net Rating and New Net Rating

### 4.4 Prediction And Evaluation

Now it's time to see the prediction results for 2018 seasons.

The Results are shown below:

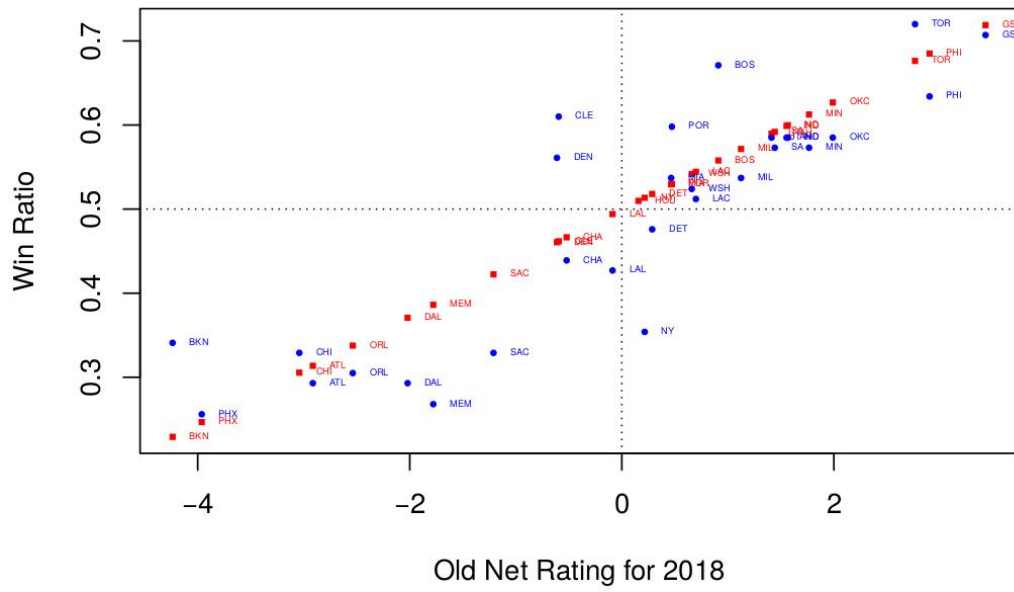


Figure 5: Predictions for 2018 season using Net Rating

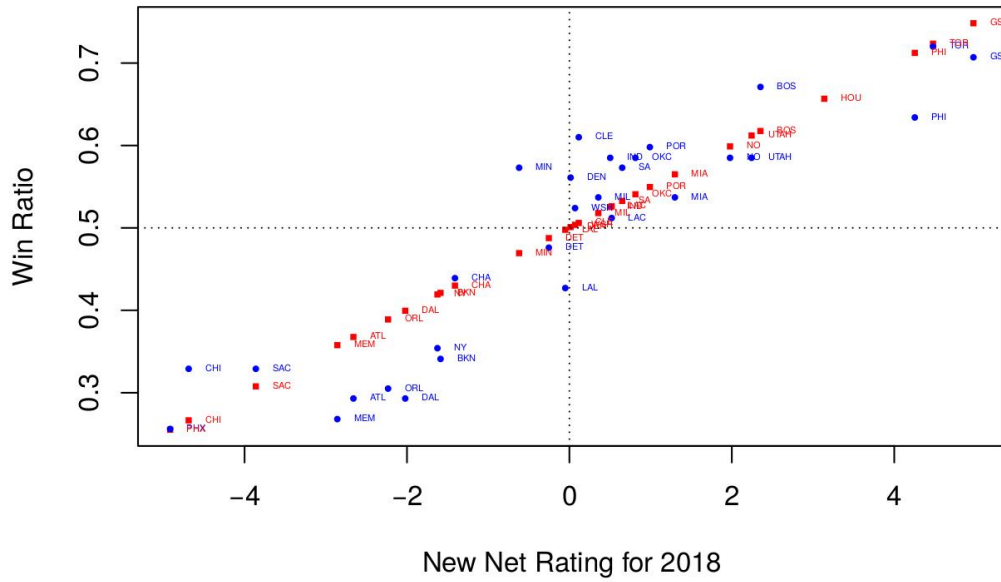


Figure 6: Predictions for 2018 season using Net Rating

From the results above, we calculated  $R^2$  for both of the model.  $R^2_{old} = 0.816$ ,  $R^2_{new} = 0.663$ .

Here are the comparison between two versions of Net Ratings and the Real Results:

New Net Rating				Net Rating			
Western Conference		Eastern Conference		Western Conference		Eastern Conference	
Team	Win Ratio	Team	Win Ratio	Team	Win Ratio	Team	Win Ratio
GS	0.748	TOR	0.724	GS	0.719	PHI	0.685
HOU	0.657	PHI	0.712	OKC	0.627	TOR	0.676
UTAH	0.612	BOS	0.618	MIN	0.612	IND	0.599
NO	0.599	MIA	0.565	NO	0.600	MIL	0.572
POR	0.550	IND	0.525	SA	0.592	BOS	0.558
OKC	0.541	MIL	0.518	UTAH	0.590	WSH	0.542
SA	0.533	CLE	0.506	LAC	0.544	MIA	0.529
LAC	0.526	WSH	0.504	POR	0.530	DET	0.518

Table 6: 2018 End Of Season Playoffs Prediction

(Yellow means correct prediction if this team is in the first 8 seeds in its conference)

2018 Playoffs Standings(Actual Results)			
Western Conference		Eastern Conference	
Team	Win Ratio	Team	Win Ratio
HOU	0.793	TOR	0.720
GS	0.707	BOS	0.671
POR	0.598	PHI	0.634
OKC	0.585	CLE	0.610
UTAH	0.585	IND	0.585
NO	0.585	MIA	0.537
SA	0.573	MIL	0.537
MIN	0.573	WSH	0.524

Table 7: The Actual Playoffs Rankings

## 4.5 Analysis And Conclusion

From the results shown above, we can clearly see that the prediction is very accurate in terms of “who make the playoffs”, except some exceptions:

Houston Rockets, who was actually top seed in Western Conference, in “Net Rating Mdel”, are surprisingly left out from the Playoffs, whereas in “New Net Rating” model, Houston ranked in the second seed.

Another surprise is Lebron James’ Cavs team, which ranked 4th at the end of last regular season, are both ranked very low in our two models(7th in New Net Rating Model,11th in the Net Rating Model). I think Lebron is really someone special who can turn things around.

In terms of win ratio predictions, New Net Rating Model has

<div>Category Win Ratio</div>	New Net Rating	Net Rating	Actual Results
Team	Win Ratio	Team	Win Ratio
$\geq 0.7$	3	1	3
$\geq 0.6$	3	5	3
$\geq 0.5$	10	10	10

Table 8: Distribution of Teams in terms of Win Ratio

From the table shown above, we can clearly see that our New Net Rating model has exactly same number of teams in different win ratio ranges.

## 5. Conclusions

In this report, we proposed modified models, New TOE and New Net Rating, and for each new model we made some predictions using past 4 years data and did some model validations and model comparisons.

The results show when evaluating the models to predict NBA regular season results, both teams' offense and defense should be considered so that we can make more accurate predictions.

## References

1. Shea, S. M., Baker, C. E. (2013). Basketball analytics: Objective and efficient strategies for understanding how teams win. CreateSpace Independent Pub. Platform.
2. Teams General Misc Statistics. Retrieved from *stats.nba.com*
3. Simple linear regression: A complete introduction with numeric example, Retrieved from <http://rstatistics.net/linear-regression-with-r-a-numeric-example/>
4. Neter, J., Kutner, M. H., Nachtsheim, C. J., Wasserman, W. (1996). Applied linear statistical models (Vol. 4, p. 359). Chicago: Irwin.

## 6. Appendix

### § Glossary

**TOE:** Teams Offensive Efficiency

**FG:** Field Goals Made

**FGA:** Field Goal Attempts

**ORB:** Offensive Rebounds

**TO:** Turnovers

**2PTM:** 2 Points Total Made

**3PTM:** 3 Points Total Made

**ORBS:** Offensive Rebound Scoring Rate

**OPPOTS:** Opponent Points Scored Off Team's Turnovers Scoring Rate

**OPFGA:** Opponent Field Goal Attempt

**OPORB:** Opponent Offensive Rebound

**OPTO:** Opponent Turnovers

**OPTOE:** Opponent Team Offensive Efficiency

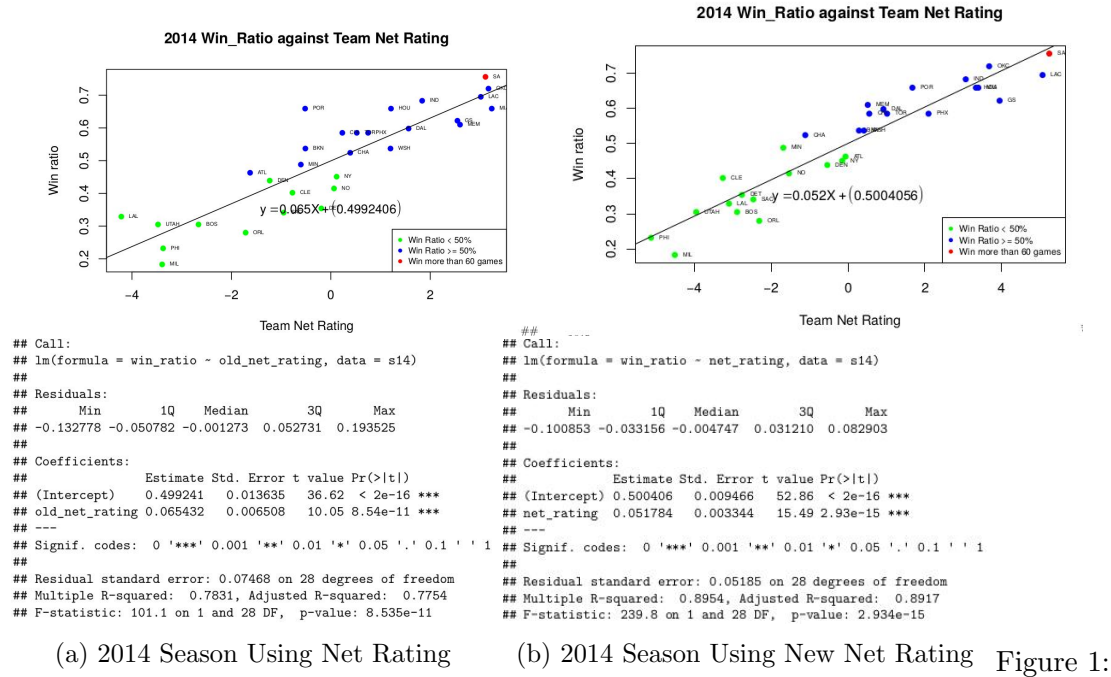
**OP2FGM:** Opponent 2 Points Field Goals Made

**OPORBS:** Opponent Offensive Rebound Scoring Rate

**POTS:** Points Off Scored Off Opponent Team's Turnovers Scoring Rate

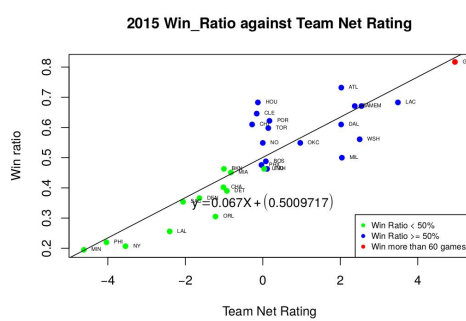
## §Results of The Other Years

### Net Ratings from 2014-2017 Seasons



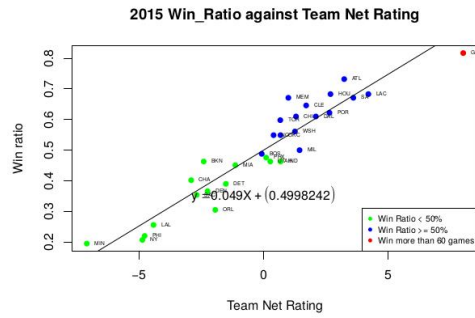
### Comparison Between Net Rating and New Net Rating in 2014 Season

### TOEs from 2014-2017 Seasons



```
## Call:
## lm(formula = win_ratio ~ old_net_rating, data = s15)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.13697 -0.04404 -0.01539  0.02459  0.19045
##
## Coefficients:
##      (Intercept)      0.500972      0.014442      34.689 < 2e-16 ***
##      old_net_rating  0.066616      0.006768      9.843 1.36e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0791 on 28 degrees of freedom
## Multiple R-squared:  0.7758, Adjusted R-squared:  0.7678
## F-statistic: 96.88 on 1 and 28 DF,  p-value: 1.364e-10
```

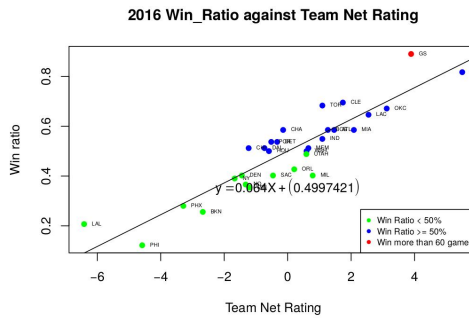
(a) 2015 Season Using Net Rating



```
## Call:
## lm(formula = win_ratio ~ net_rating, data = s15)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.099146 -0.033654 -0.008004  0.045677  0.121572
##
## Coefficients:
##      (Intercept)      0.499824      0.009899      50.49 < 2e-16 ***
##      net_rating    0.049495      0.003210      15.42 3.26e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.05422 on 28 degrees of freedom
## Multiple R-squared:  0.8947, Adjusted R-squared:  0.8909
## F-statistic: 237.8 on 1 and 28 DF,  p-value: 3.262e-15
```

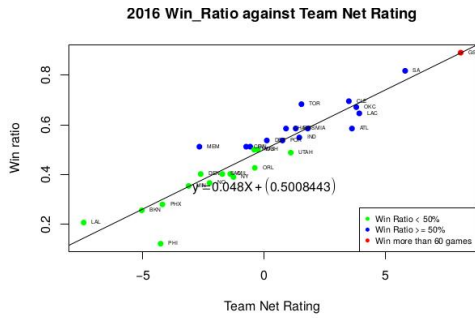
(b) 2015 Season Using New Net Rating

Figure 2: Comparison Between Net Rating and New Net Rating in 2015 Season



```
## Call:
## lm(formula = win_ratio ~ old_net_rating, data = s16)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.14813 -0.04903 -0.01327  0.05864  0.14214
##
## Coefficients:
##      (Intercept)      0.49974      0.01326      37.69 < 2e-16 ***
##      old_net_rating  0.06375      0.00560      11.38 5.1e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07263 on 28 degrees of freedom
## Multiple R-squared:  0.8224, Adjusted R-squared:  0.816
## F-statistic: 129.6 on 1 and 28 DF,  p-value: 5.102e-12
```

(a) 2016 Season Using Net Rating

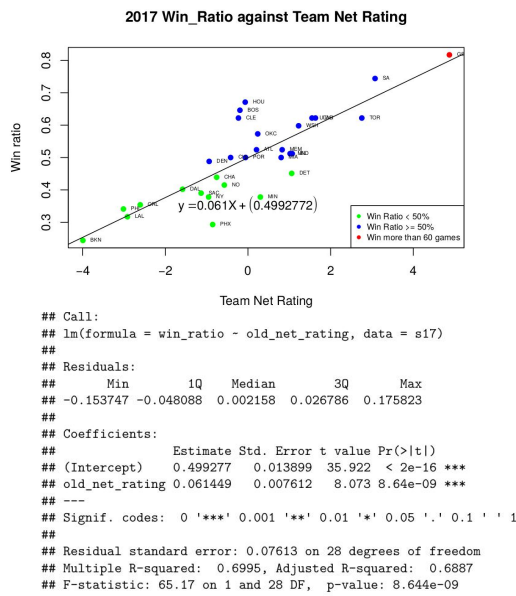


```
## Call:
## lm(formula = win_ratio ~ net_rating, data = s16)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.174523 -0.025737  0.000448  0.029789  0.138892
##
## Coefficients:
##      (Intercept)      0.500844      0.010815      46.31 < 2e-16 ***
##      net_rating    0.047943      0.003317      14.46 1.64e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.05924 on 28 degrees of freedom
## Multiple R-squared:  0.8818, Adjusted R-squared:  0.8776
## F-statistic: 209 on 1 and 28 DF,  p-value: 1.639e-14
```

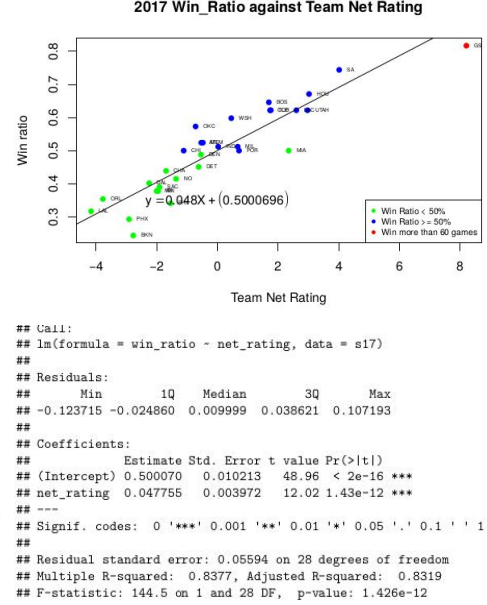
(b) 2016 Season Using New Net Rating

Figure 3: Comparison Between Net Rating and New Net Rating in 2016 Season



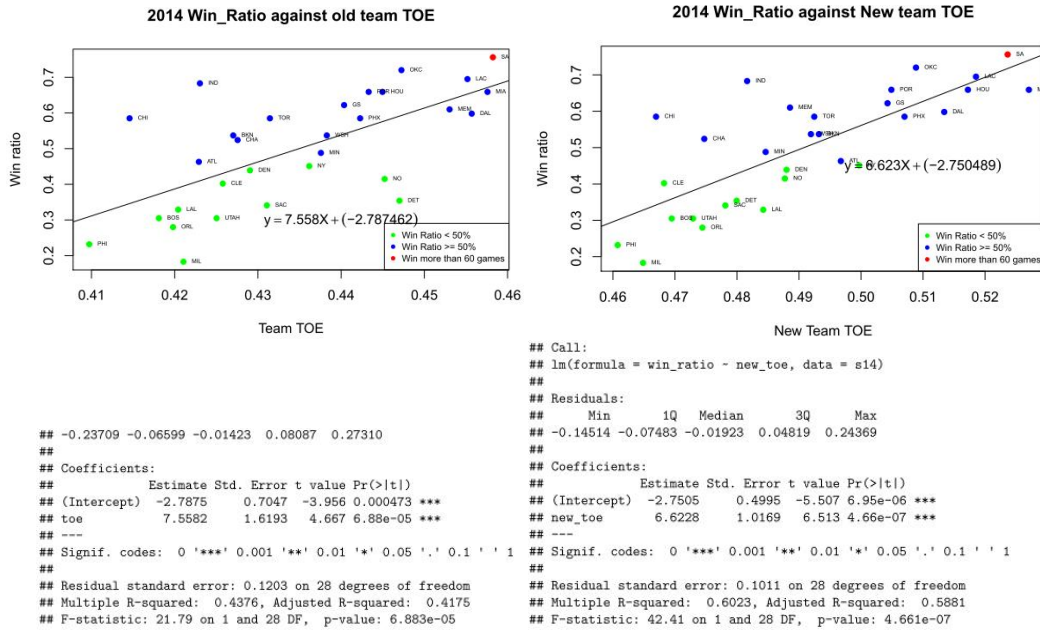


(a) 2017 Season Using Net Rating



(b) 2017 Season Using New Net Rating

Figure 4: Comparison Between Net Rating and New Net Rating in 2017 Season



(a) 2015 Season Using Old TOE

(b) 2014 Season Using New TOE

Figure 5: Comparison Between TOE and New TOE in 2014 Season

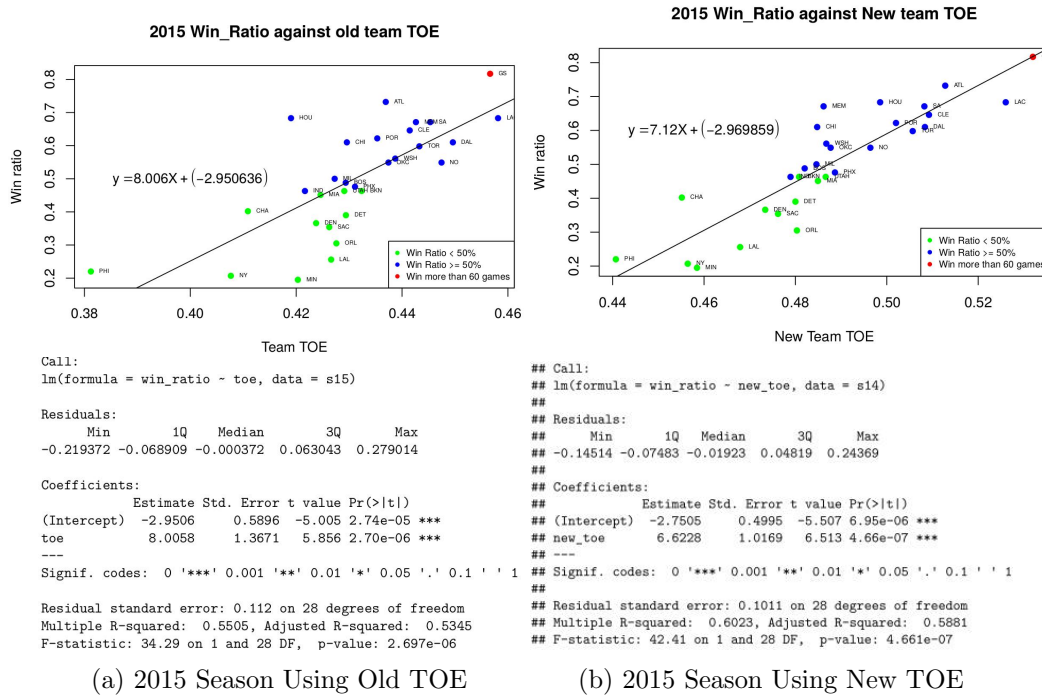


Figure 6: Comparison Between TOE and New TOE in 2014 Season