

Predicting NBA Championship Winners Using Data from Past Championship Series

Ravi G and Rohit K

3/8/2022

Background:

Every year, the National Basketball Association (NBA) ends the season with a series of championship games between the two best teams in the league. The series is a best-of-7, where the first team to win 4 games is the winner of the series. The championship is a very important accolade, both teams and players are compared by the number of championships they have won. In these comparisons, statistics like field goal percentage, offensive rebounds, steals, and blocks can be used to determine a team's performance and be an indicator of how much better one team is than another.

Statistics in the NBA are highly analyzed, very often with the goal of assessing which areas a team can improve upon. Our group wants to create a model that will be able to tell us which key statistics from the NBA Championship series from 1980-2018 were important in deciding the championship. More specifically, we want to find weights for each statistic that can show us how important each statistic is to predicting the overall winner, which helps our group conduct a greater analysis of how different focuses and strategies can affect the outcome of NBA games.

We will be using the NBA Finals Team Stats Dataset, which has been analyzed and used to create models by several Kaggle Users. One project to note is a report written by Ziyu Liu called "Three pointers win championships", in which the author creates a model to see if the number of three point shots made by a team can predict whether or not the team wins the championship. In the study, the model achieves an accuracy of 59%. This tells us that, while three point shots are important, more statistics are required to be able to create a more accurate model.

Another example of analyzing NBA championship data comes from the article "Stat of the Week: How champions are built in the NBA" by Ryan Blackburn. In this article, Blackburn specifically analyzes the Denver Nuggets, and why the team has never won an NBA championship. In the paper, Blackburn ranks teams by offense, defense, and net-rating, and points out that only five out of past 20 championship winners have ranked outside of the top 3 teams in net-rating. This leads us to believe that a team needs both impressive defensive statistics and offensive statistics in order to win an NBA championship.

Given these studies and our group's own knowledge of the NBA, we believe that both offensive statistics like field goals, three pointers, offensive rebounds; and defensive statistics like steals and defensive rebounds will have a high impact on whether or not a team wins the championship.

Data:

The dataset with which we want to create our model is the NBA Finals Team Stats dataset on Kaggle uploaded by Dave Rosenman. The dataset contains final data from 1980 to 2018, and is divided into two tables. The first table contains the data of each winning team and the second contains the losing team.

Each observation includes data points like field goals made, field goals attempted, three point shots made, free throws made, total rebounds, assists, steals, turnovers, blocks, and many other statistics that will be covered in the data summary. The data takes averages from each game in the series, and its an average of the performance of the team in this category across all the games played in the series.

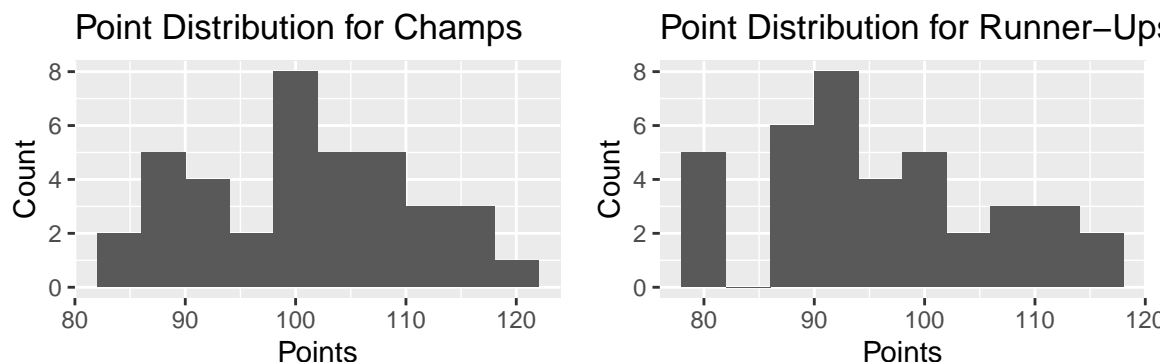
In order to create the dataset we are using in this study, we started with two separate datasets, one for all of the series winners (NBA Champions) and one of the runner-ups. We created a new column **win** with a 1 if the team won the series and a 0 if the team lost. This will be our predictor variable for the model. Next, we combined the two datasets and randomized the order of the entries. Our goal is to first analyze each of the variables to determine which will be the most useful in creating our model, then going through several iterations of models before choosing the most accurate one.

Exploratory Data Analysis

As mentioned before, we combined the datasets and added a **win** variable to tell us whether the team won a championship or not. Since there were over 20 variables, our goal was to choose key variables to analyze that we knew would have an impact on the game. To do this, we eliminated some variables that we did not wish to explore or view the effect they would have on the model. This includes statistics like **FTA** (Free Throw Attempts), **BLK** (Blocks). While in the game this statistics might be important, having a statistic like Free Throw Attempts or Blocks without any context about the other team's performance in the same category would most likely not be useful.

Now, we can start our EDA. First, we check the plots of each variable for both the losers and the winners to make sure each distribution is roughly normal.

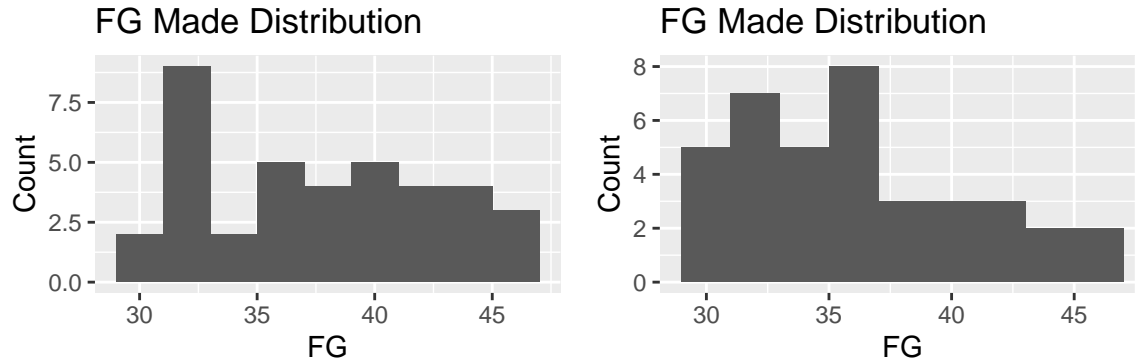
Distribution of Points:



Summary of Points:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	79.80	90.74	99.19	98.56	106.67	121.60

Distribution of Field Goals Made:



Summary of Field Goals:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	30.00	32.96	36.08	37.02	40.81	46.71

The distributions and summaries of the other variables can be seen in the appendix. All of the distributions appear to be somewhat normal with no outliers, We have a large sample size, so we can continue to make our model and check the residuals.

Creating the Model

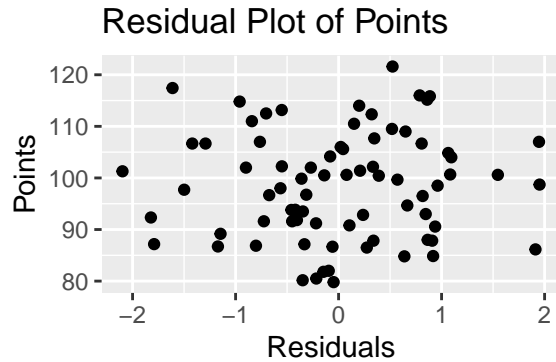
Model Refinement

Our model will be a binomial model (only options are 0 and 1). The first step will be to plot the residuals of each variable in the model to check the linearity assumption. Then, we will plot the Cook's distance and remove any high-leverage points. Finally, the VIF will be checked and any variables with a high VIF will be removed from the model.

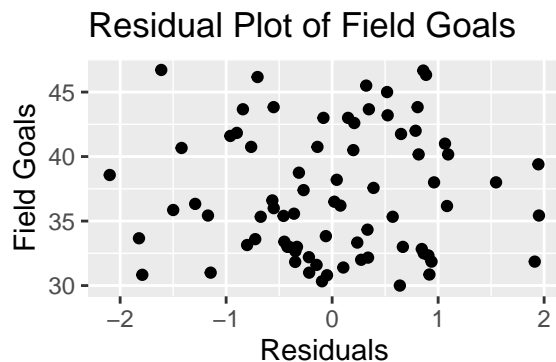
Summary of Model:

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	-11.372	91.053	-0.125	0.901	-195.869	168.984
PTS	0.058	0.113	0.515	0.607	-0.165	0.289
FG	0.365	2.311	0.158	0.875	-4.252	5.015
FGA	-0.637	1.052	-0.605	0.545	-2.785	1.428
FGP	0.337	1.961	0.172	0.864	-3.538	4.322
TP	0.150	0.739	0.202	0.840	-1.276	1.734
TPA	0.020	0.271	0.072	0.943	-0.550	0.551
TPP	0.042	0.084	0.496	0.620	-0.142	0.203
ORB	0.984	0.302	3.260	0.001	0.464	1.662
DRB	0.469	0.194	2.422	0.015	0.129	0.904
AST	-0.020	0.142	-0.143	0.886	-0.303	0.265
STL	0.712	0.298	2.385	0.017	0.166	1.359
BLK	0.011	0.333	0.035	0.972	-0.647	0.681
TOV	-0.256	0.192	-1.334	0.182	-0.658	0.110
PF	-0.051	0.153	-0.333	0.739	-0.360	0.255

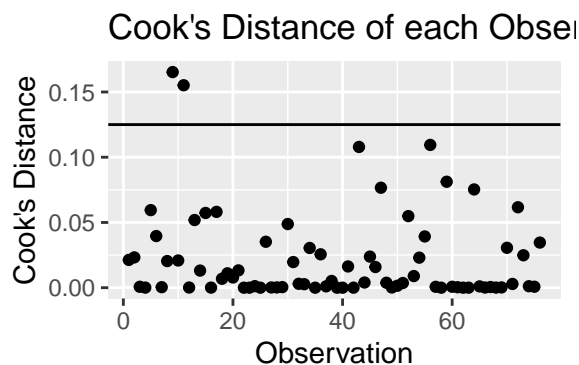
Residual Plot of Points:



Residual Plot of Field Goals:



The residual plots for each variable appear to be random and evenly dispersed (the rest can be seen in the appendix), which means that the linearity assumption is satisfied. Before we can test the accuracy of the model, we must also explore how these observations affect the model, and how the variables used in the model affect each other. First, we can plot the leverage (.cooks) of each observation to see if there are any high-leverage data points.



It is obvious that there are two high-leverage data points. If we use a threshold of 0.125, we can eliminate these two high-leverage points to make the model better at prediction. After this, the model must be trained on the newly-filtered data. Then, we re-train our model on the now filtered dataset.

Summary of Model on Filtered Dataset:

```
##
## Call:
## glm(formula = win ~ PTS + FG + FGA + FGP + TP + TPA + TPP + ORB +
```

```
##      DRB + AST + STL + BLK + TOV + PF, family = binomial, data = filter_data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.12838  -0.32438  -0.00962   0.66696   2.05077
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -13.34302   98.44346  -0.136  0.89218
## PTS          -0.03651    0.13731  -0.266  0.79033
## FG           1.03938    2.46225   0.422  0.67293
## FGA          -1.02089    1.12581  -0.907  0.36451
## FGP           0.30264    2.11119   0.143  0.88601
## TP           -0.32034    0.89001  -0.360  0.71890
## TPA           0.28922    0.33879   0.854  0.39328
## TPP           0.09593    0.10564   0.908  0.36384
## ORB           1.30905    0.41405   3.162  0.00157 **
## DRB           0.72101    0.28588   2.522  0.01166 *
## AST          -0.05620    0.18053  -0.311  0.75559
## STL           1.20375    0.42457   2.835  0.00458 **
## BLK           0.16852    0.38747   0.435  0.66362
## TOV          -0.29451    0.21214  -1.388  0.16506
## PF           -0.01274    0.18549  -0.069  0.94523
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 102.532  on 73  degrees of freedom
## Residual deviance:  47.107  on 59  degrees of freedom
## AIC: 77.107
##
## Number of Fisher Scoring iterations: 7
```

The next step is to check how the variables interact with each other. To measure this, we want to calculate the Variable Inflation Factor, or *VIF*.

```
##      PTS      FG      FGA      FGP      TP      TPA
## 14.554416 1116.928102 468.327785 307.483725 68.519839 65.480931
##      TPP      ORB      DRB      AST      STL      BLK
##  5.335238  8.903163  3.340391  4.573557  2.795145  1.839274
##      TOV      PF
##  1.638405  2.045241
```

Some of the variables have a very high *VIF*, so we can only include certain variables to keep the score lower. To see which variables are better included and not included, we must create multiple models and assess which one has the best accuracy. We can leave all of the variables who's *VIF* is lower than 10, but the others must be removed for higher accuracy. We will create three models, one that will include only the Field Goals made and the Three Point shots made; one that will only include the Field Goal percentage and Three Point percentage; and finally an attempts model that will include the number of Field Goal attempts and Three Point attempts. The, we check the *VIF* of each model to make sure that it has been reduced (that each variable has a *VIF* lower than 10)

Points Model:

```
##      PTS      FG      TP      ORB      DRB      AST      STL      BLK
## 10.603796 13.770485 2.563876 1.707288 1.471669 4.101363 1.333806 1.435059
##      TOV      PF
## 1.473768 1.403311
```

Percentage Model:

```
##      PTS      FGP      TPP      ORB      DRB      AST      STL      BLK
## 6.842594 5.870015 1.837503 3.514570 2.031535 3.559040 1.557436 1.702397
##      TOV      PF
## 1.451581 1.387433
```

Attempts Model:

```
##      PTS      FGA      TPA      ORB      DRB      AST      STL      BLK
## 7.070250 11.799342 2.191565 3.448663 1.576465 3.958743 1.645281 1.667579
##      TOV      PF
## 1.394461 1.461882
```

To evaluate the accuracy of each model, we will train each model to our dataset to make predictions and measure the accuracy of these predictions. We believe that the best way to evaluate the models is to calculate the *AIC*, or the Akaike information criterion.

Points Model:

```
## [1] 105.6141
```

Percentage Model:

```
## [1] 87.40903
```

Attempts Model:

```
## [1] 79.57669
```

We can see that the AIC of the attempts model is significantly lower than the AIC of the other two models. Since a lower AIC is the result of a better fitting model for the data, we will use only the attempts model in the process of refining our model.

Since the VIF of the model was close to 10, we want to check whether or not including the **PTS** variable in the model makes the model better or worse. To do this, we will create two models from our attempts model, the first including points and the second not including points. We will chose whichever model has the lower AIC to be our final model.

AIC of No Points Model:

```
## [1] 89.03652
```

The AIC of the points model has already been calculated (79.57669), so we know that the points model is a better fit so this will be our final model. So we consider our attempts model, with an AIC of 79.57669 to be our final model.

Conclusion

```
final_model <- attempts_model
tidy(final_model, conf.int = TRUE, exponentiate = FALSE) %>%
  kable(digits = 3, format = "markdown")
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	-0.025	6.400	-0.004	0.997	-13.094	12.528
PTS	0.250	0.089	2.821	0.005	0.095	0.446
FGA	-0.626	0.167	-3.742	0.000	-1.009	-0.341
TPA	0.025	0.055	0.448	0.654	-0.081	0.139
ORB	0.744	0.252	2.951	0.003	0.311	1.332
DRB	0.442	0.169	2.621	0.009	0.138	0.811
AST	0.160	0.147	1.091	0.275	-0.119	0.468
STL	0.805	0.297	2.715	0.007	0.266	1.452
BLK	0.343	0.326	1.052	0.293	-0.281	1.011
TOV	-0.328	0.181	-1.806	0.071	-0.709	0.015
PF	-0.132	0.142	-0.928	0.353	-0.424	0.144

Our final model had an AIC of 79.57669. Using the coefficients of the final model, we can see that Offensive Rebounds and Steals are the two most important factors in whether or not a team from 1980-2018 won the NBA championship.

Observations

We found it interesting that the attempts of field goals and three pointers provided the most accurate model for predicting who won the series. It intuitively makes sense that the field goals scored matter more than attempts because basketball games are decided by the score, not the attempts; however, attempts might signal how ineffective a team's offense really can be.

It is worth noting that, while the coefficients for most variables make sense (ie. turnovers having a negative coefficient and points having a positive coefficient), others have a surprising effect on the model.

For example, the coefficients for field goals attempted (FGA) is negative, implying that more field goal attempts indicate a lesser likelihood that a team won the game. This can be attributed to the fact that basketball is about making baskets rather than taking shots. If you aren't making these shots you are taking, you most likely won't win the game.

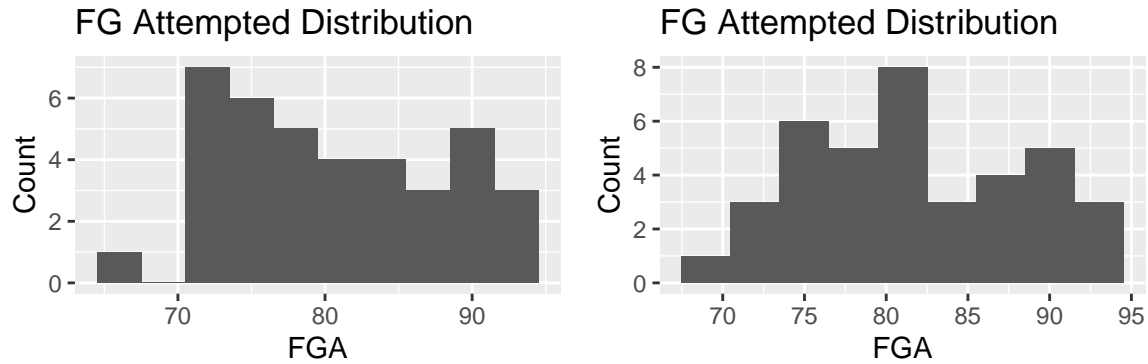
Some variables with high positive coefficients include: Points (ORB), Offensive Rebounds (ORB), Defensive Rebounds (DRB), and Steals (STL). These all seem to be good at indicating whether or not teams won the series, which is fascinating because they cover different aspects of the game. A team cannot solely rely on shooting, defense, or size (in the case of offensive and defensive rebounds) to win a championship; they must have all aspects of the game.

The two variables with low negative coefficients include: Field Goals Attempted (FGA) and Turnovers (TOV). These also seem good in indicating whether or not teams won the series. Having low turnovers is good as it will give you more opportunities to shoot the ball and prevent free baskets. The other variables don't have that big of an impact on the final result, but that does not mean those variables do not impact the final result of a championship.

Appendix

Distributions of variables for Winners and Losers

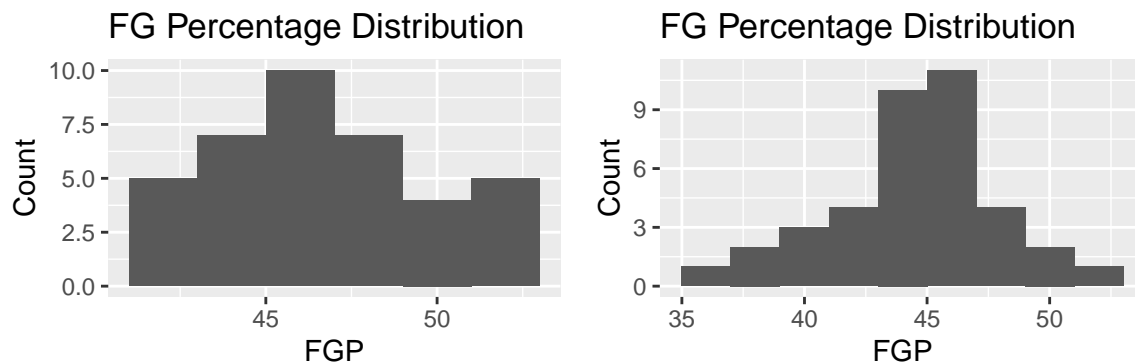
Distribution of Field Goals Attempted:



Summary of Field Goal Attempts:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	67.40	75.42	80.60	81.18	87.00	92.86

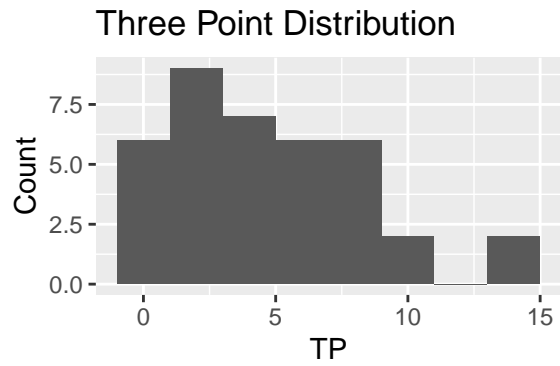
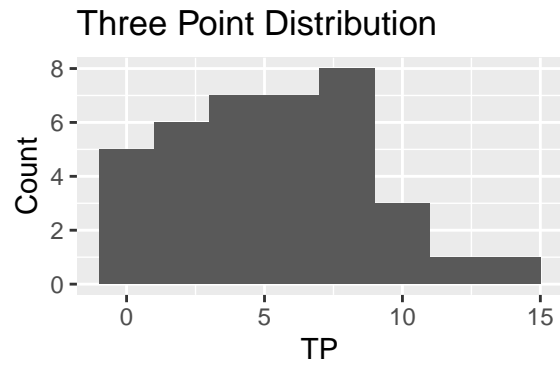
Distribution of Field Goal Percentage:



Summary of Field Goal Percentage:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	36.99	43.20	45.40	45.52	47.51	52.76

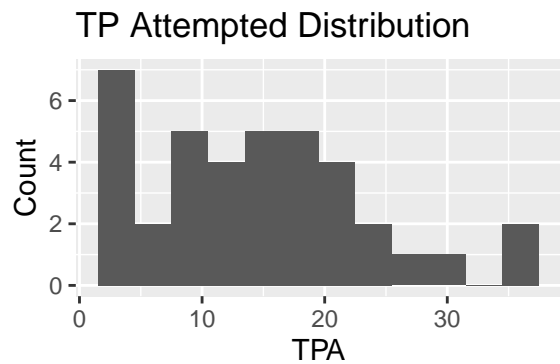
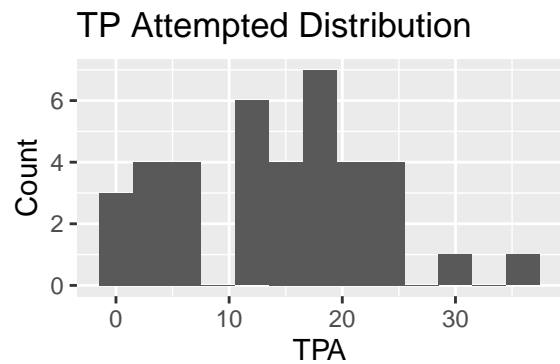
Distribution of Three Point Distribution:



Summary of Three Point:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.000	2.167	4.857	4.987	7.259	14.200

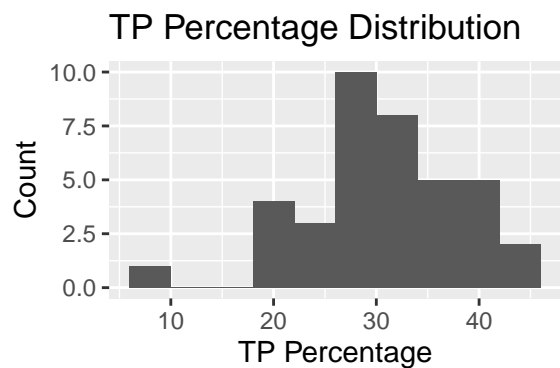
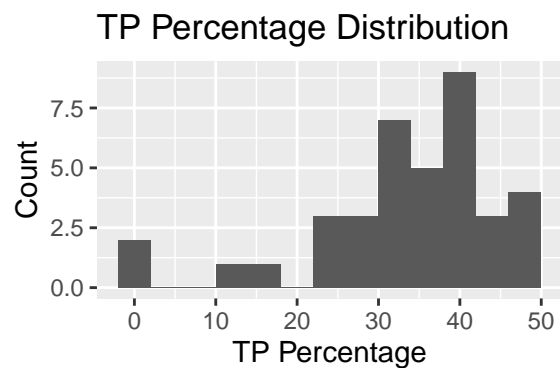
Distribution of Three Point Attempted:



Summary of Three Point Attempts:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.6667	6.9583	14.7083	14.3523	20.2500	37.2000

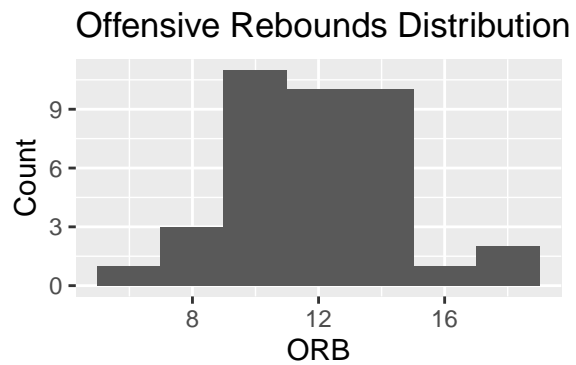
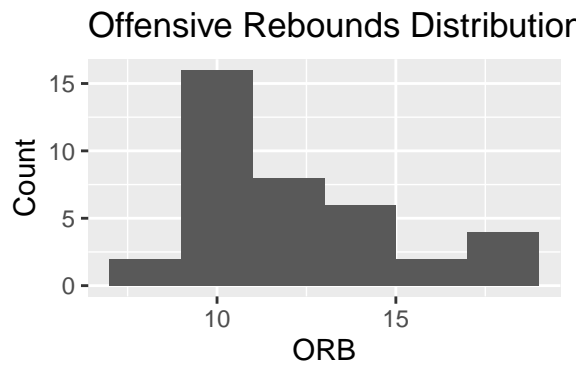
Distribution of Three Point Percentage:



Summary of Three Point Percentage:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.00	28.20	32.47	31.89	38.29	48.00

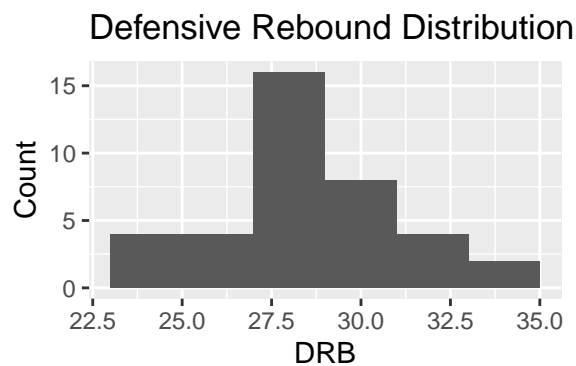
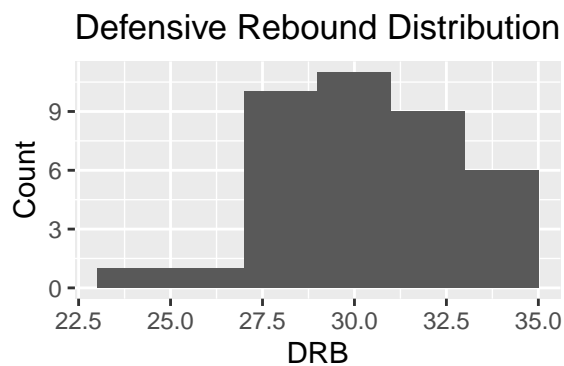
Distribution of Offensive Rebounds:



Summary of Offensive Rebounds:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	5.60	10.37	11.73	12.16	13.71	18.67

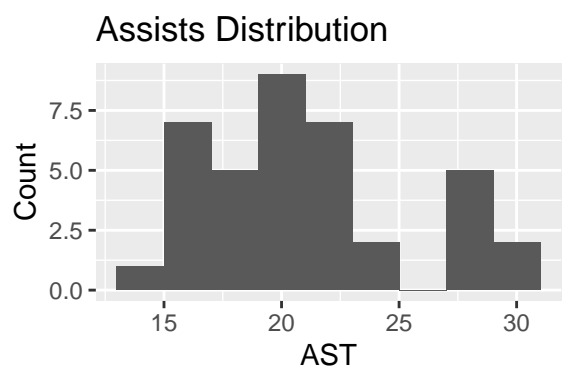
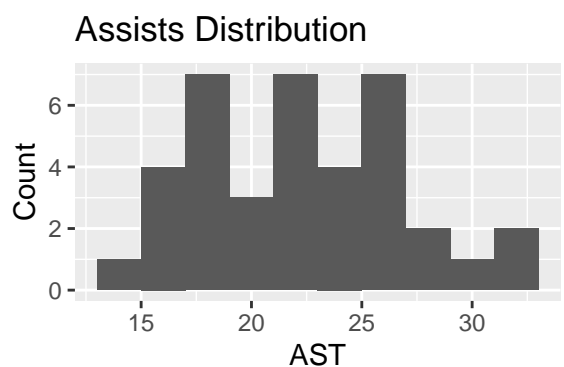
Distribution of Defensive Rebounds:



Summary of Defensive Rebounds:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	23.80	27.83	29.31	29.45	31.30	35.00

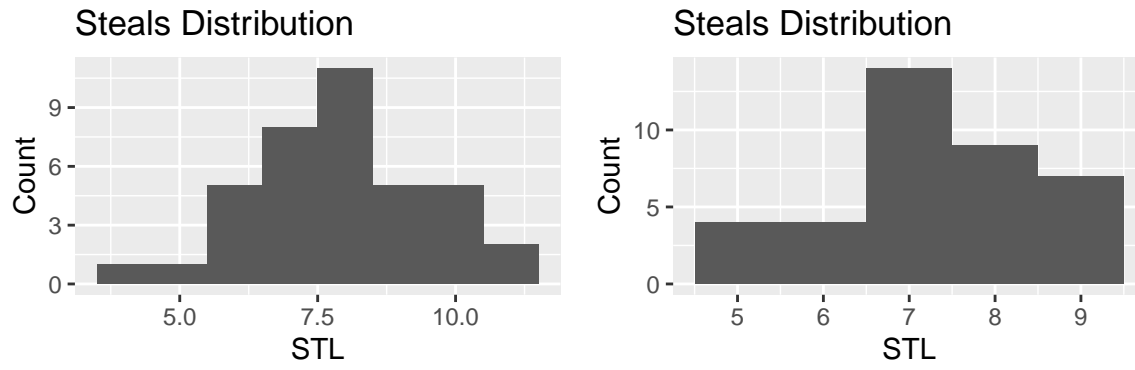
Distribution of Assists:



Summary of Assists:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	14.00	18.33	21.33	21.81	24.73	32.00

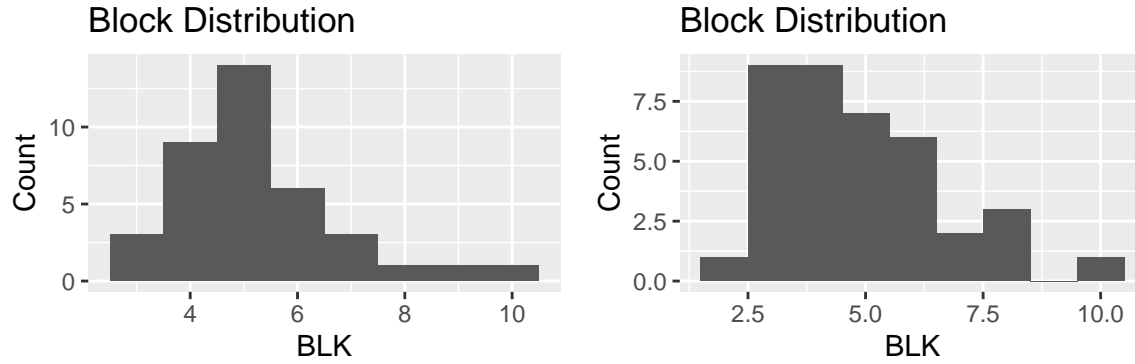
Distribution of Steals:



Summary of Steals:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	4.000	6.833	7.633	7.618	8.488	11.000

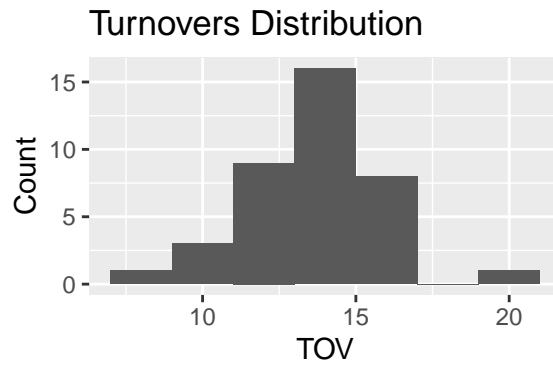
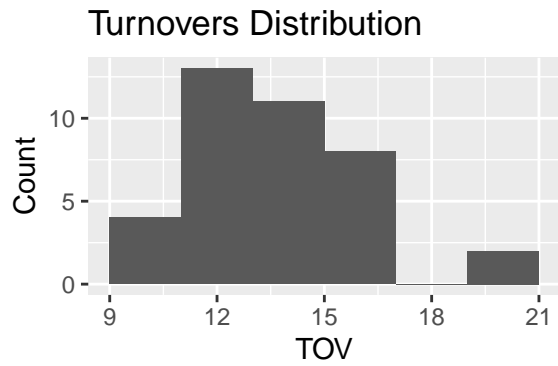
Distribution of Blocks:



Summary of Blocks:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	2.200	4.000	5.000	5.100	5.808	10.000

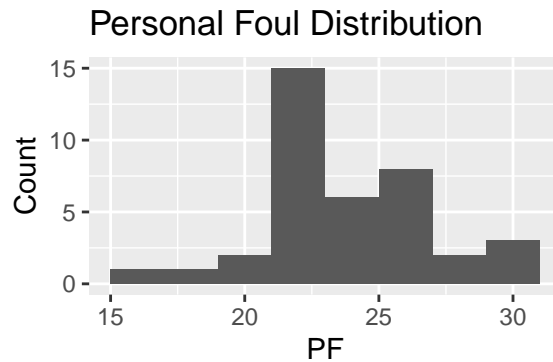
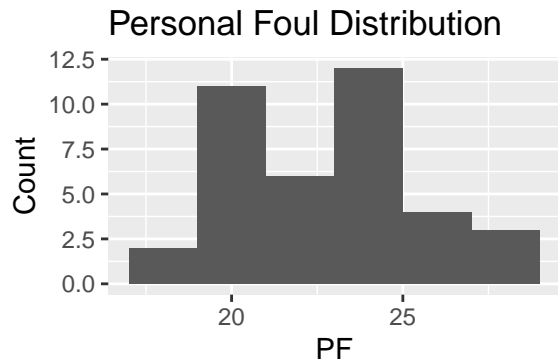
Distribution of Turnovers:



Summary of Turnovers:

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	8.714	12.421	13.310	13.728	15.042	20.000

Distribution of Personal Fouls:

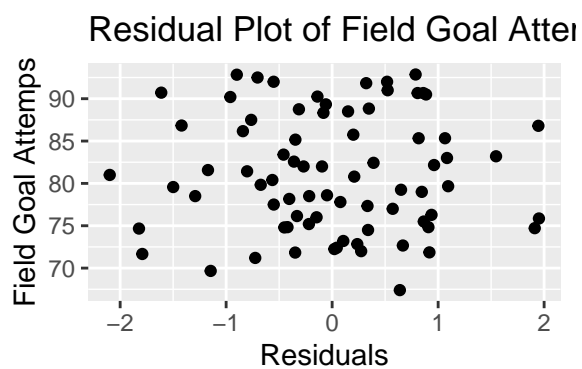


Summary of Personal Fouls:

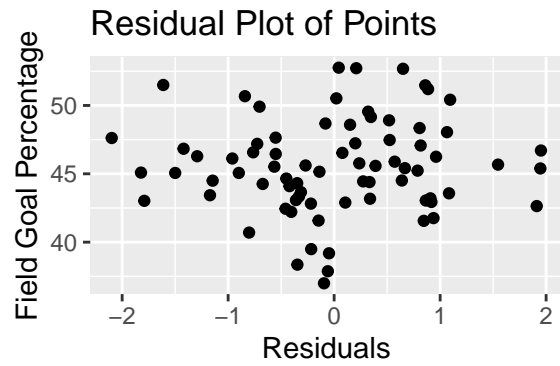
##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	16.86	21.32	23.08	23.39	25.38	30.00

Residuals of Each Variable

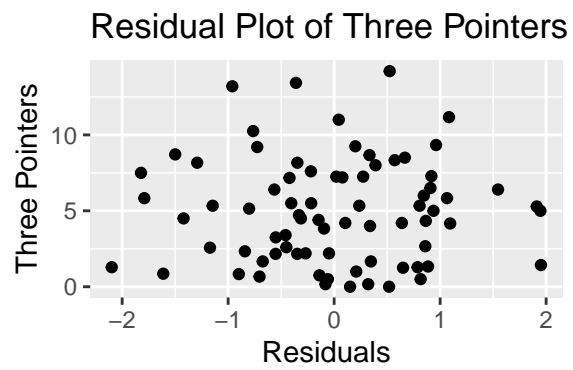
Residual Plot of Field Goal Attempts:



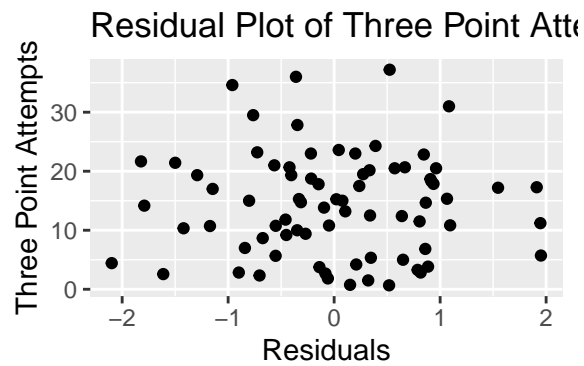
Residual Plot of Points:



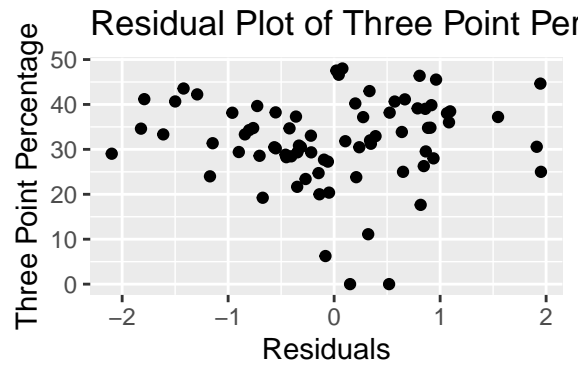
Residual Plot of Three Pointers:



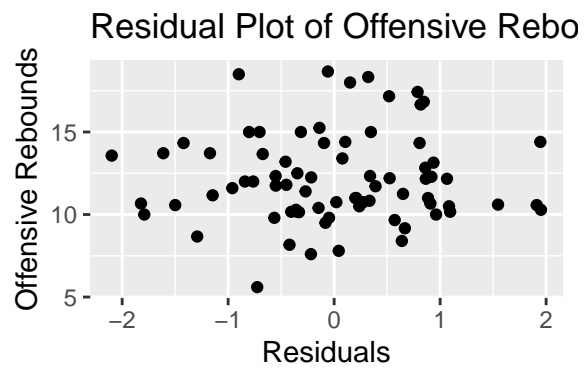
Residual Plot of Three Point Attempts:



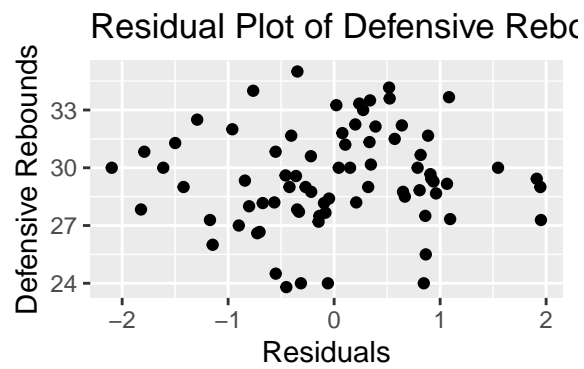
Residual Plot of Three Point Percentage:



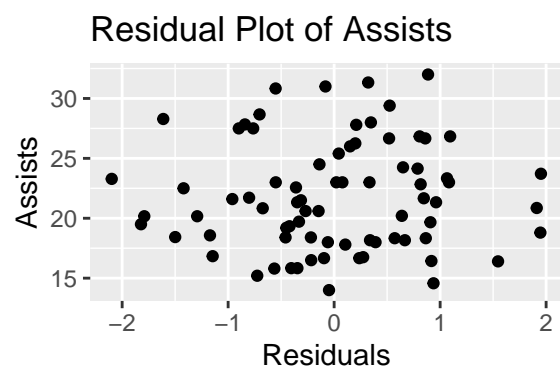
Residual Plot of Offensive Rebounds:



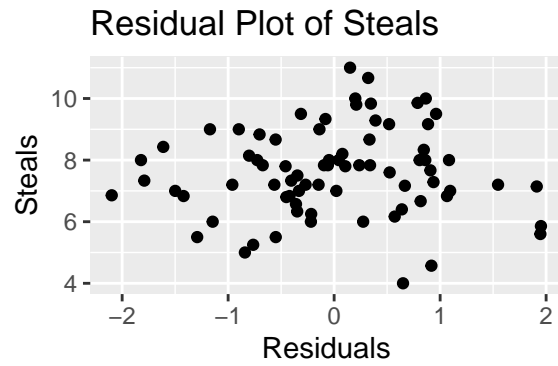
Residual Plot of Defensive Rebounds:



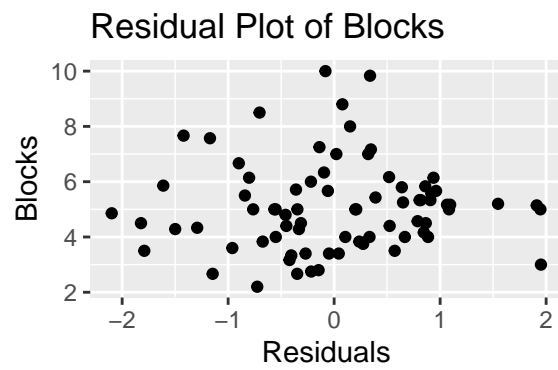
Residual Plot of Assists:



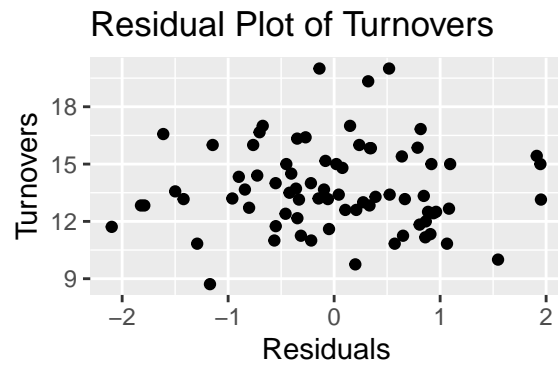
Residual Plot of Steals:



Residual Plot of Blocks:



Residual Plot of Turnovers:



Residual Plot of PF:

