# Predicting Win Shares in the NBA

By Anthony Le

# Purpose

-Predicting the win shares statistic based on player's stats

- Find which stats are important for win shares

-Big fan of the NBA

-WIll help appreciate the game more.

-Find which model will be the best to predict the win shares

# Win Shares

- player statistic which attempts to divvy up credit for team success to the individuals on the team
- An advanced statistic to see if the player is helping their team win.
- Add Offensive Win  Shares and Defensive Win Shares together to get the total WIn Shares.
- The formula for the WIn Shares are from this source. (https://www.basketball-reference.com/about/ws.html)
- The data collected for this analysis are not related to the formula.

# NBA Data

- Collected from the Basketball-Reference.com
  (https://www.basketball-reference.com/leagues/NBA_2019_advanced.html)
- Data collected website
- 2018-2019 NBA Season
- 708 players data collected ( couple players played for multiple seasons)
- Used the Advanced Stats

# Data Attributes

- 20 continuous variables
    - Age- Age; player age on February 1 of the given season.
    - G-Number of Games played
    - MP- Minutes Played
    - PER- Player Efficiency Rating
    - Ts%- True Shooting Percentage
    - 3PAr- 3-Point Field Goal Attempts Rate
    - FTr- Free Throw Attempts Rate
    - ORB%- Offensive Rebound Percentage
    - DRB%- Defensive Rebound Percentage
    - TRB%- Total Rebound Percentage
    - AST%- Assist Percentage
    - STL%- Steal Percentage

# Data Attributes(Cont.)

- 20 continuous variables
  - BLK%- Block Percentage
  - TOV%- Turnover Percentage
  - USG%- Usage Percentage
  - OBPM- Offensive Box Plus/Minus
  - DBPM- Defensive Box Plus/ Minus
  - BPM- Box Plus/ Minus
  - VORP- Value over Replacement Player

Source-(https://www.basketball-reference.com/about/glossary.html)

# Data Exploration

- Check to see any nulls in the dataset at all

```
nba_data.isnull().sum()*100/nba_data.isnull().count()
```

```
Player    0.000000
Age       0.000000
G         0.000000
MP        0.000000
PER       0.000000
TS%       0.847458
3PAr      0.847458
FTr       0.847458
ORB%      0.000000
DRB%      0.000000
TRB%      0.000000
AST%      0.000000
STL%      0.000000
BLK%      0.000000
TOV%      0.847458
USG%      0.000000
WS        0.000000
OBPM      0.000000
DBPM      0.000000
BPM       0.000000
VORP      0.000000
```
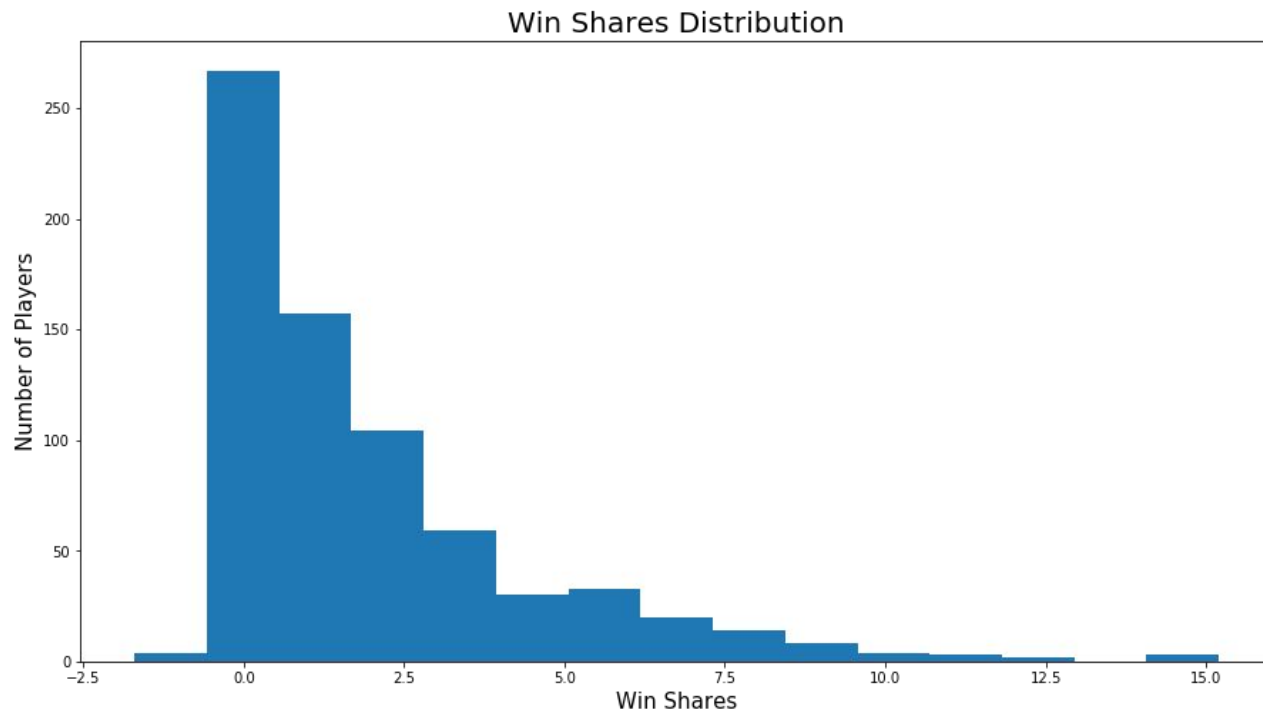
# Null

- Replace the null with zero since the null is very low

```
nba_data["TS%"].fillna(0, inplace =True)
nba_data["3PAr"].fillna(0, inplace =True)
nba_data["FTr"].fillna(0, inplace =True)
nba_data["TOV%"].fillna(0, inplace =True)
```

# Data Exploration

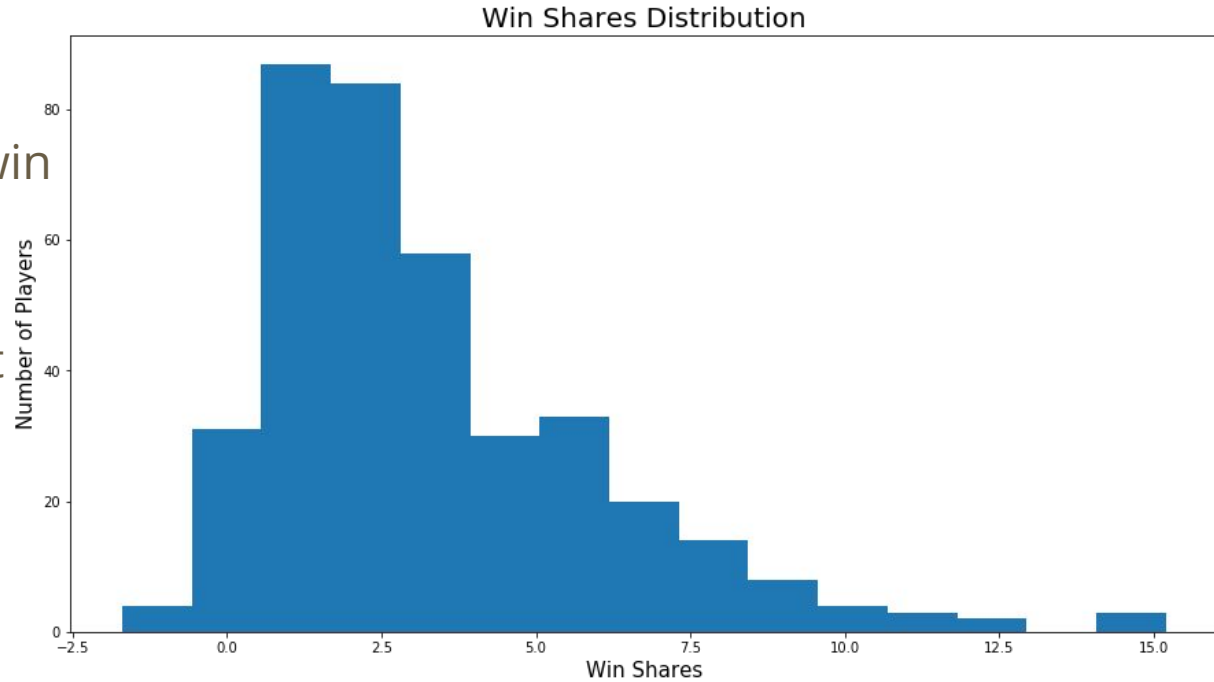- Need to see what the Win Shares Distribution looks like.

# Data Transformation

- Decide to filter out players that barely played
  - 41 games or more
- There's ton of player with 0 win shares.
- Only 15 players per team
  - Injuries, trades, etc. affects the data set

# Data Transformation

- 381 players
- Less players with 0 win shares
- More normalized compared to the last distribution
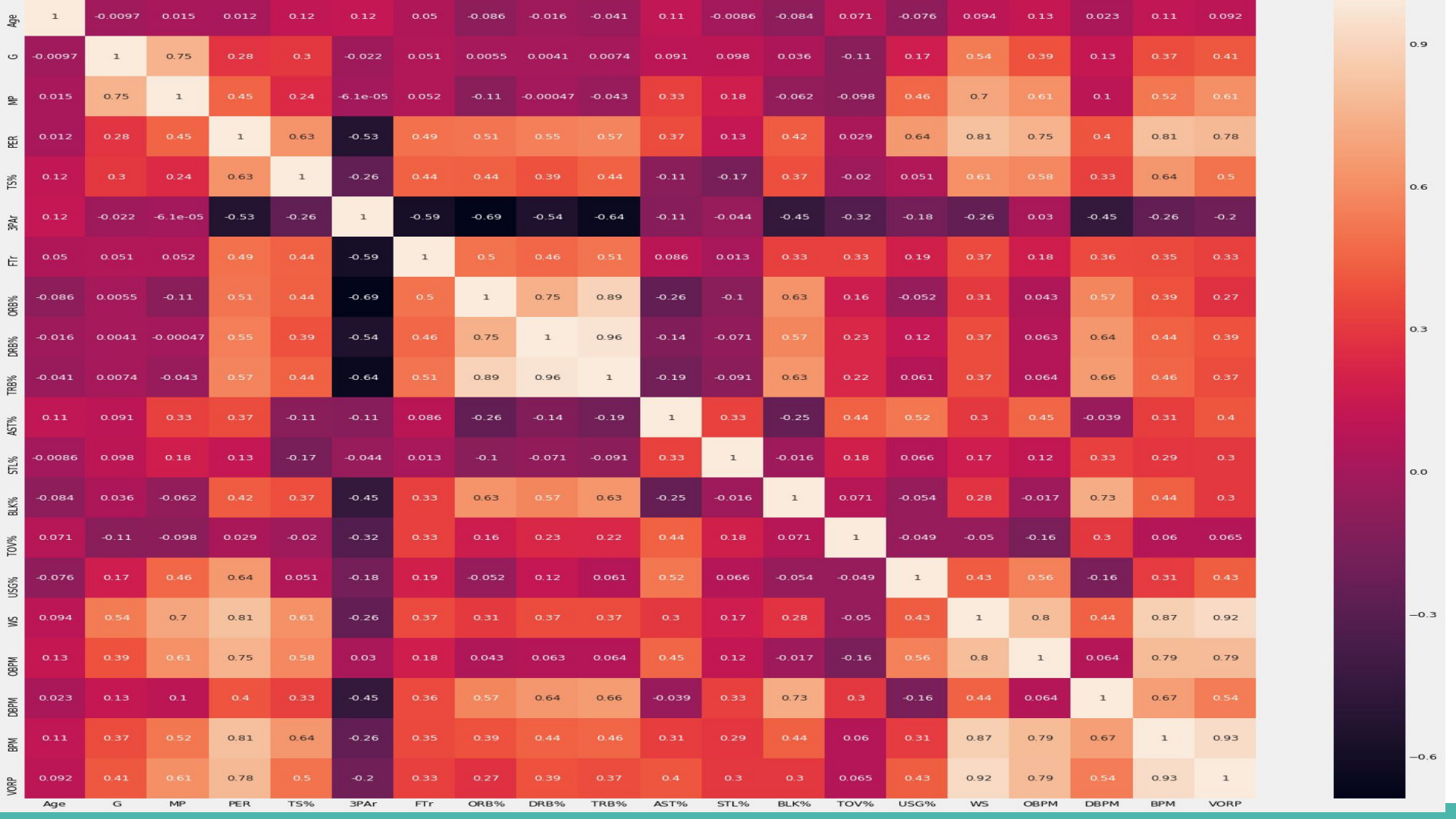


Win Shares Distribution

# Correlation

- Finding which variables correlated well with WIn Shares.
- Use correlation function and heatmaps to determine correlation.

# Correlation-2

- Use variables that correlation that is .5 or higher

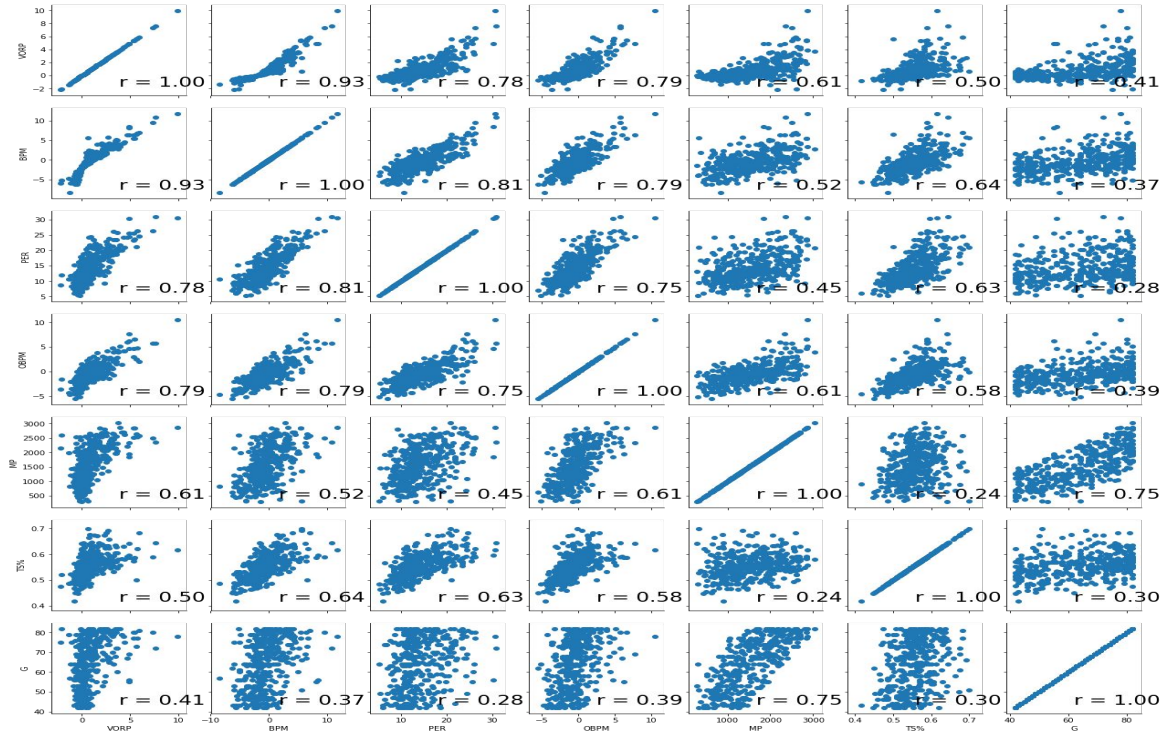| | index | WS |
|---|---|---|
| 15 | WS | 1.000000 |
| 19 | VORP | 0.924396 |
| 18 | BPM | 0.868633 |
| 3 | PER | 0.813769 |
| 16 | OBPM | 0.795985 |
| 2 | MP | 0.698498 |
| 4 | TS% | 0.612444 |
| 1 | G | 0.537833 |
| 17 | DBPM | 0.441935 |
| 14 | USG% | 0.425792 |
| 9 | TRB% | 0.373042 |
| 8 | DRB% | 0.367111 |
| 6 | FTr | 0.365449 |
| 7 | ORB% | 0.305793 |
| 10 | AST% | 0.298471 |
| 12 | BLK% | 0.276973 |
| 11 | STL% | 0.172317 |
| 0 | Age | 0.094430 |
| 13 | TOV% | -0.049696 |
| 5 | 3PAr | -0.255229 |

# Multicollinearity

- Need to check that other variables correlated with each too closely.
- Won't skew the model
- Use graphs and table to check

# Multicollinearity

- VORP and BPM have a correlation of .93
- BPM would be dropped because VORP used BPM in their formula

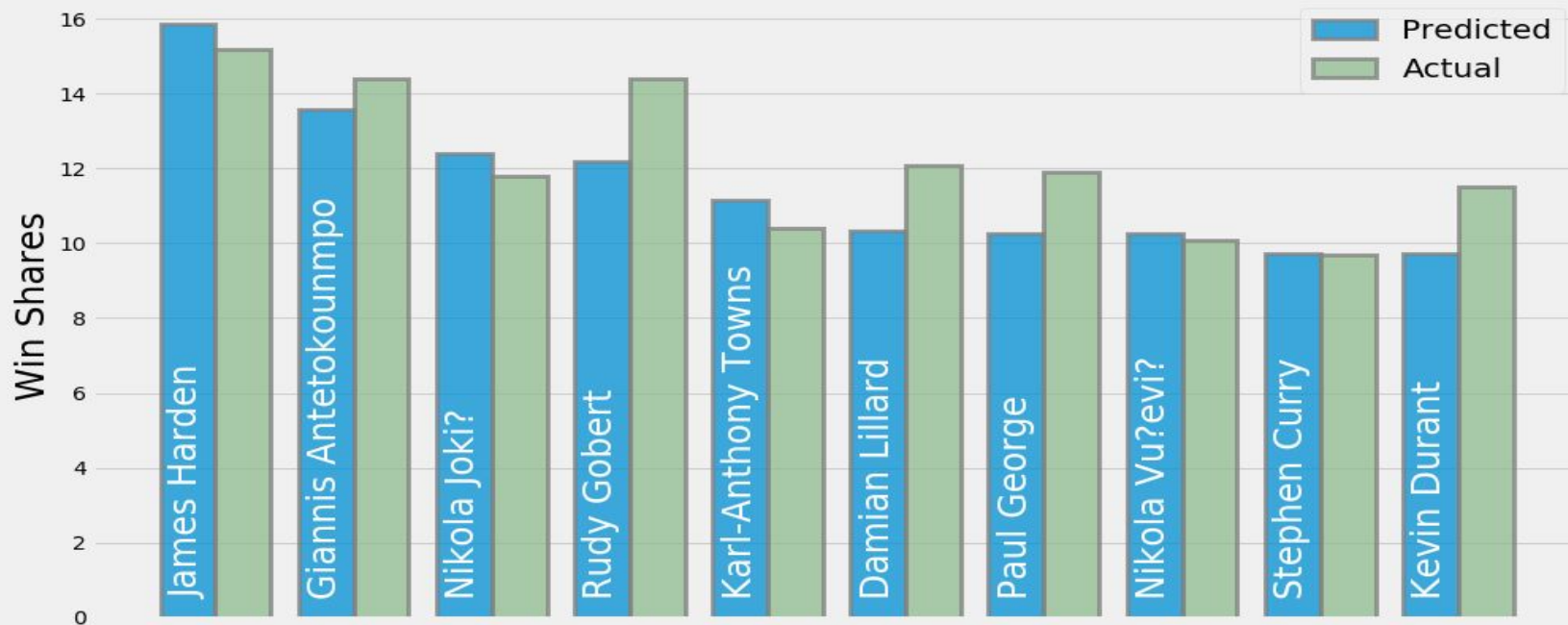|      | VORP     | BPM      | PER      | OBPM     | MP       | TS%      | G        |
|------|----------|----------|----------|----------|----------|----------|----------|
| VORP | 1.000000 | 0.929908 | 0.779450 | 0.794319 | 0.606268 | 0.504991 | 0.408920 |
| BPM  | 0.929908 | 1.000000 | 0.806153 | 0.785744 | 0.521236 | 0.637544 | 0.371151 |
| PER  | 0.779450 | 0.806153 | 1.000000 | 0.749606 | 0.448797 | 0.625866 | 0.281484 |
| OBPM | 0.794319 | 0.785744 | 0.749606 | 1.000000 | 0.612428 | 0.581164 | 0.387401 |
| MP   | 0.606268 | 0.521236 | 0.448797 | 0.612428 | 1.000000 | 0.235015 | 0.749390 |
| TS%  | 0.504991 | 0.637544 | 0.625866 | 0.581164 | 0.235015 | 1.000000 | 0.296773 |
| G    | 0.408920 | 0.371151 | 0.281484 | 0.387401 | 0.749390 | 0.296773 | 1.000000 |

# Multicollinearity

# Chosen Features

- VORP
- PER
- OBPM
- MP
- TS%
- G
- 20% tested and 80% trained

# Model 1- Linear Regression



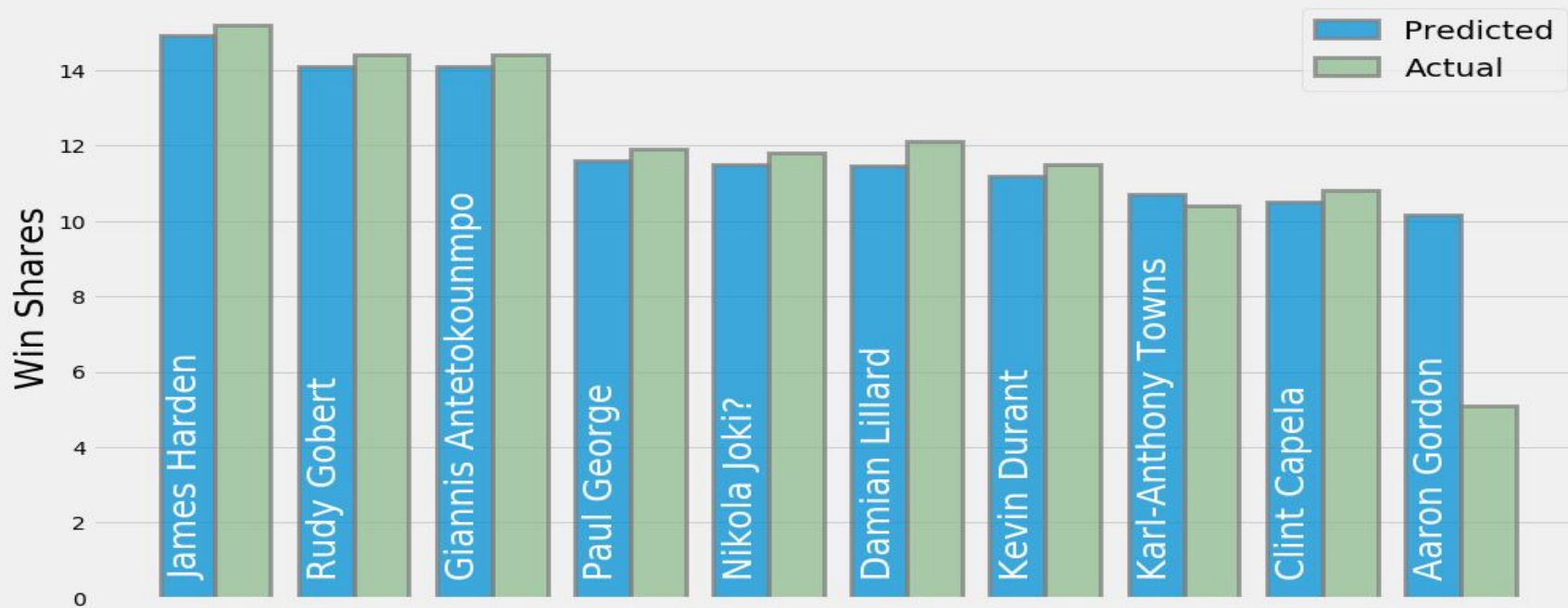## 2018 NBA Predicted vs Actual Win Shares - Top 10 Players
Wins shares are predicted with Linear Regression model

**Win Shares**

Legend:
- Predicted
- Actual

Players (x-axis):
- James Harden
- Giannis Antetokounmpo
- Nikola Joki?
- Rudy Gobert
- Karl-Anthony Towns
- Damian Lillard
- Paul George
- Nikola Vu?evi?
- Stephen Curry
- Kevin Durant

National Basketball Association

Source: Basketball-Reference.com

# Model 2- Support Vector Regression



## 2018 NBA Predicted vs Actual Win Shares - Top 10 Players
Wins shares are predicted with Support Vector Regression model

Legend: Predicted, Actual

Win Shares (y-axis): 0, 2, 4, 6, 8, 10, 12, 14

Players: James Harden, Rudy Gobert, Giannis Antetokounmpo, Paul George, Nikola Joki?, Damian Lillard, Kevin Durant, Karl-Anthony Towns, Clint Capela, Aaron Gordon

# Model 3- K-Nearest Neighbors Regression



## 2018 NBA Predicted vs Actual Win Shares - Top 10 Players

Wins shares are predicted with K-nearest Neighbors Regression model

Legend:
- Predicted
- Actual

Players (x-axis): Paul George, Damian Lillard, P.J. Tucker, Tobias Harris, Tobias Harris, Kemba Walker, Bradley Beal, James Harden, Ben Simmons, DeMar DeRozan, Kevin Durant

Y-axis: Win Shares

National Basketball Association

Source: Basketball-Reference.com

# Model 4- Random Forest Regression



## 2018 NBA Predicted vs Actual Win Shares - Top 10 Players
Wins shares are predicted with Random Forest Regression model

Legend: Predicted, Actual

Y-axis: Win Shares

Players: James Harden, Giannis Antetokounmpo, Rudy Gobert, Nikola Joki?, Karl-Anthony Towns, Damian Lillard, Paul George, Kevin Durant, Nikola Vu?evi?, Stephen Curry

National Basketball Association

Source: Basketball-Reference.com

# Performance Evaluation

| Model | Mean Squared Error | Mean Absolute Error | Variance Score |
|---|---|---|---|
| Linear | 0.360 | 0.448 | 0.926 |
| Support Vector | 2.878 | 1.270 | 0.406 |
| k-Nearest Neighbors | 3.560 | 1.332 | 0.265 |
| Random Forest | 0.331 | 0.451 | 0.932 |

# Best Model

- Random Forest  proven to be the best model
- Predict the top ten players accurately
- Has the lowest value of Mean Squared Error and Mean Absolute Error
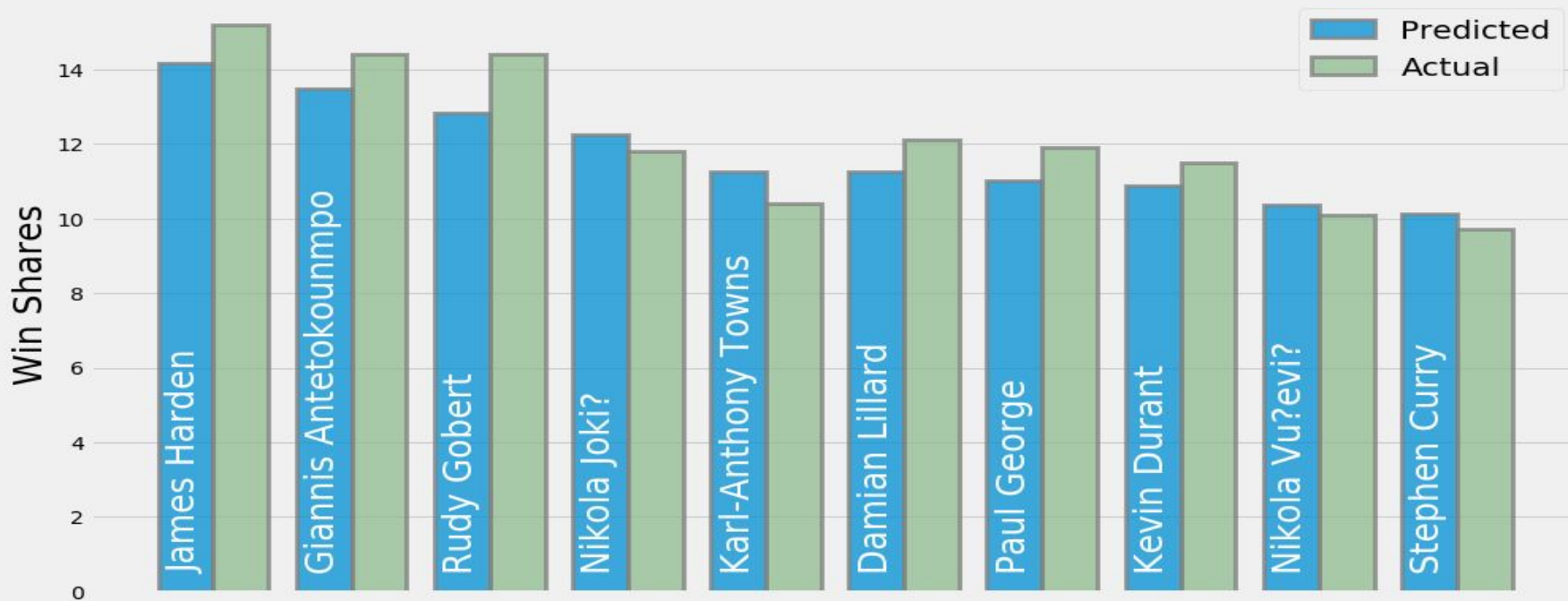- Has the highest value of Variance score.

# Random Forest

- Need to see the Random Forest was overfitted or not
- To see there was a good generalization or not.

# 50%



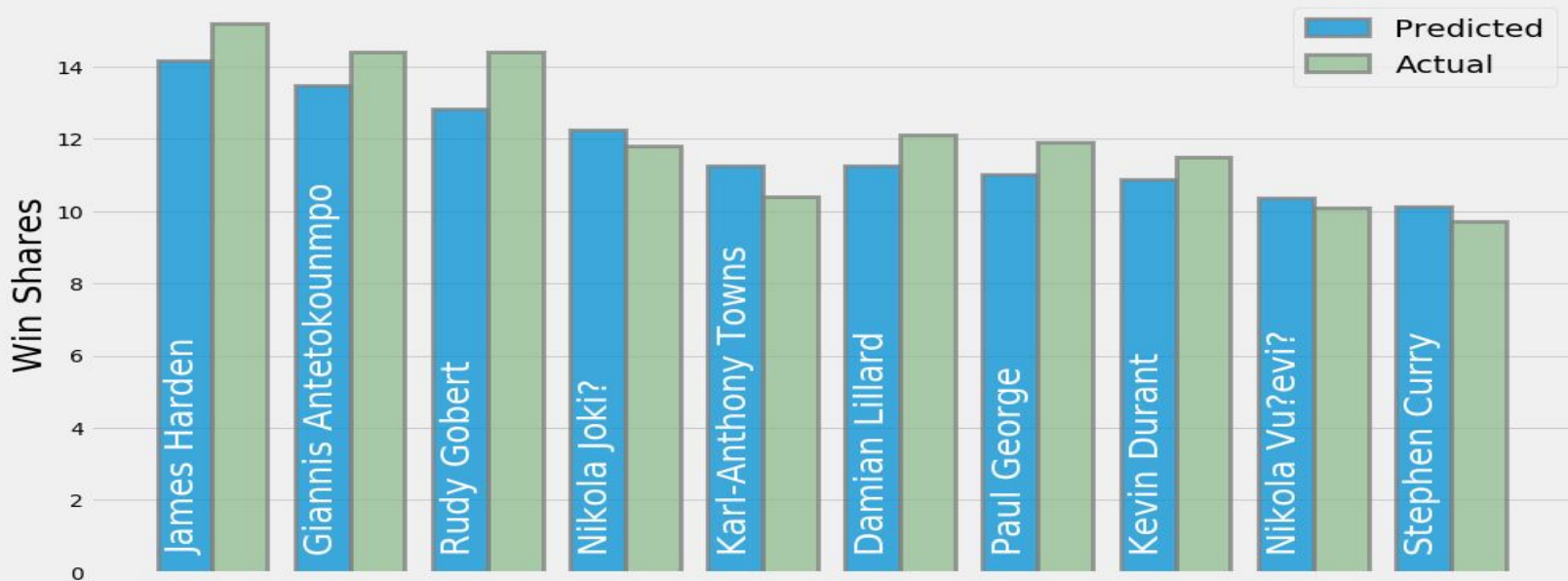**2018 NBA Predicted vs Actual Win Shares - Top 10 Players**
Wins shares are predicted with Random Forest Regression model

National Basketball Association    Source: Basketball-Reference.com

# 90%



**2018 NBA Predicted vs Actual Win Shares - Top 10 Players**

Wins shares are predicted with Random Forrest Regression model

# Random Forest

| Test | Mean Squared Error | Mean Absolute Error | Variance Score |
|------|--------------------|--------------------|----------------|
| 10%  | 0.319              | 0.471              | 0.940          |
| 50%  | 0.689              | 0.567              | .0902          |
| 90%  | 1.020              | 0.713              | 0.863          |

# Random Forest

- At 10% seem to be the closest for performance for 20%
- 10% and 50% have  similar results for top ten players
- 50% has higher error though
- This shows that the random forest wasn't overfitted for 20% because 90% would have similar results.

# Conclusion

- Random Forest best model for this dataset based on performance
- VORP has the most correlation
- Age didn't have an impact
- Not the best stat to judge individuals' performance
  - Team fit and player personnel
- Need more datasets
  - Used past seasons
  - Used more features