

Predicting NBA All-Stars Using Linear Regression Models

Executive Summary:

Each season, 24 of the NBA's elite are selected to the All-Star Team: a commemoration of the most popular and talented players in the league. We analyzed the performance of players during the 2020-2021 season in order to find out which attributes had the greatest impact on a player's selection to the All-Star squad. After devising a suitable model, we used players' performance during the second half of the prior season in order to predict All-Stars in the current season, and conducted a corresponding analysis of our results.

Introduction:

Every season, the NBA hosts the All-Star Weekend, a three-day event where a collection of the league's stars and rising talents showcase their skills throughout a series of exhibition games and events. The centerpiece of the All-Star Weekend is the All-Star Game: a full-length exhibition match played between the best players in the Eastern and Western Conferences.

Selection to this event is a testament to both a player's dominance and longevity within the league, as evidenced by Kareem Abdul-Jabbar, Kobe Bryant, and LeBron James topping the record books with 19, 18, and 17 selections respectively. However, the award holds importance beyond merely a player's legacy, as oftentimes salary incentives are structured around performance criteria such as All-Star selection. Because of the nontrivial implications of being an All-Star, the NBA revised the selection process in 2017 to reduce fan impact on voting - from 100% fan-based to now include player and media opinion.

The current weighted ranked vote consists of 50% popular fan votes, 25% current players, and 25% media representatives. This weighted average is computed by ranking each player within their conference and position (backcourt vs frontcourt) across each of the voting channels (fan, player, media). For example, budding superstar Trae Young ranked 6th among fans, 11th among players, and 6th among media within Eastern Conference guards and did not make the All-Star team. The top two guards and top three frontcourt players among each conference are selected as starters, with the remaining seven players on each team being chosen by head coaches within the conference. If a player selected as an All-Star is injured prior to the event, a replacement is selected by the All-Star coach or the NBA Commissioner. For the purposes of this project, we ignored players selected to the team in this fashion.

Problem Statement:

Considering the importance and year-to-year variability of being selected as an All-Star, either as a starter or reserve, predicting a player's likelihood of selection has far-ranging implications. These include future salary and potential legacy, their team's success and popularity, as well the ability to attract free agents to the franchise. As such, forecasting the chance of a player's selection to the All-Star roster is a compelling proposition.

To do so, we conducted linear regression using box score data on all NBA players from the first half of the 2020-2021 NBA season, with the key variable of interest being a player's weighted All-Star Voting rank (the average of the fan, player, and media vote). After assessing the relative

value of a range of attributes such as Win/Loss, Points per Game, Age, and so forth, we devised a model to predict whether or not a player will make the All-Star Game by comparing their forecasted All-Star Voting Rank to the necessary score to make the game in their position and conference. By computing a new All-Star Rank using the latter half of the season's data, we predict whether or not each player would similarly make the All-Star Game during the 2021-2022 season. Of course, we would not yet be able to verify this assumption as the game occurs in February of 2022 and is based upon player performance during the current season rather than the former.

Data Description:

Our core dataset on player performance comes from *NBA.com* - the official source of NBA-related statistics. We concatenated the primary data with All-Star Voting dating sourced from *Basketball-reference.com*, the premier source for basketball stats across leagues. The initial data set displays seasonal performance on a per-player basis up until the 2021 All-Star game. We then take the second half of this same dataset (up until the end of the regular season) to predict All-Star Rank. Examples of the variables included in this dataset are age, wins, minutes per game, points per game, and so forth.

One issue that arose while joining the player performance data with All-Star Voting data related to the players with exceptionally forgettable seasons - those most likely to not receive any votes. Such players who did not receive All-Star votes did not appear in the All-Star voting dataset, and as such could not be immediately mapped to seasonal performance. While we had the option of simply assigning them the lowest All-Star Voting rank in the dataset (157.5), we elected not to do so as the presence of such outliers could have skewed our regression models. In any case, the number of players not included in the All-Star data was negligible: only 20 of the bottom feeders in the league were excluded relative to the 500 players considered in the initial dataset.

Following is a description of the variables included in our analysis on the 2020-2021 NBA season; not all of them are incorporated into the final model. Certain variables become redundant as they are a composite of other attributes. e.g. $\text{Rebounds} = \text{Offensive Rebounds} + \text{Defensive Rebounds}$.

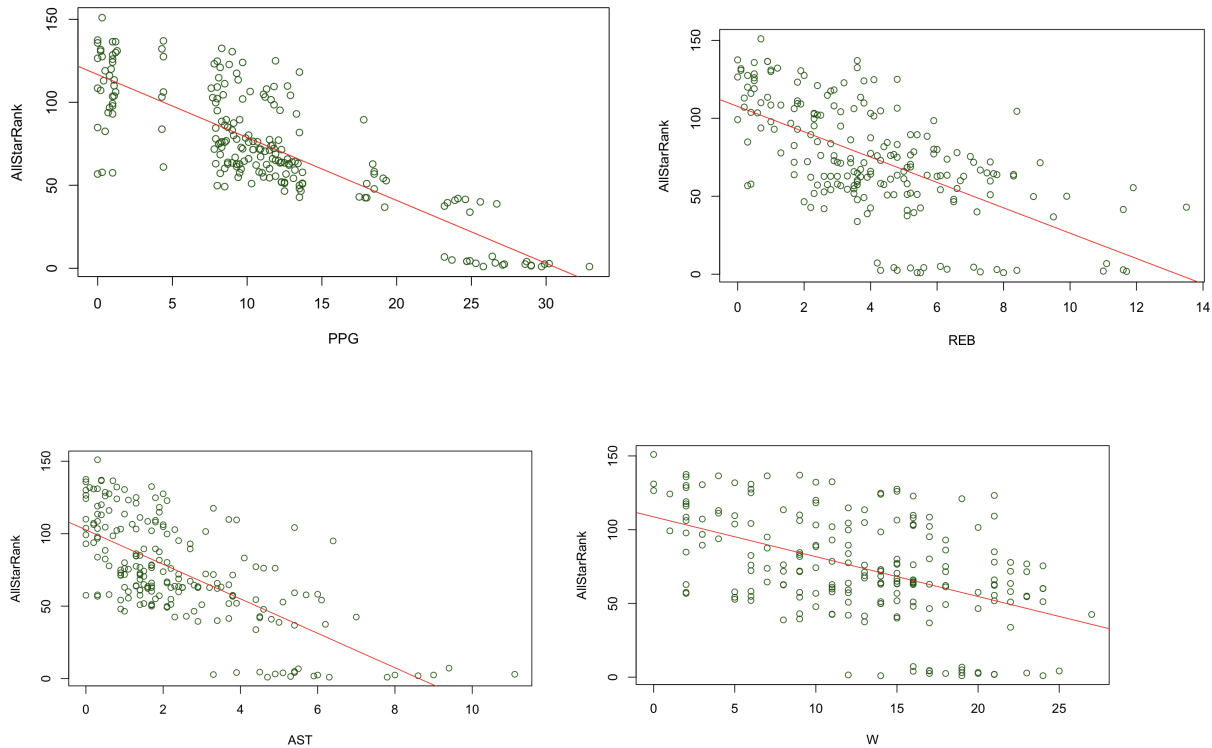
ID (player ID), PLAYER (player name), TEAM (team), AGE (age), GP (games played), W (wins), L (losses), MIN (minutes), PPG (points per game), FGM (field goals made), FGA (field goals attempted), FG% (field goal %), 3PM (three pointers made), 3PA (three pointers attempted), 3P% (three point %), FTM (free throws made), FTA (free throws attempted), FT% (free throw %), OREB (offensive rebounds), DREB (defensive rebounds), REB (rebounds), AST (assists), TOV (turnovers), STL (steals), BLK (blocks), PF (personal fouls), FP (fantasy points), DD2 (double-doubles), TD3 (triple-doubles), Plus/Minus (box plus/minus), AllStarRank (weighted average All Star rank), AllStar (binary All-Star selection), Conference (East/West), Position (Front/Backcourt).

All variables are quantitative (integer or numeric) except for PLAYER, TEAM, Conference, and Position (qualitative).

Regression Analysis:

i) Plots of Variables:

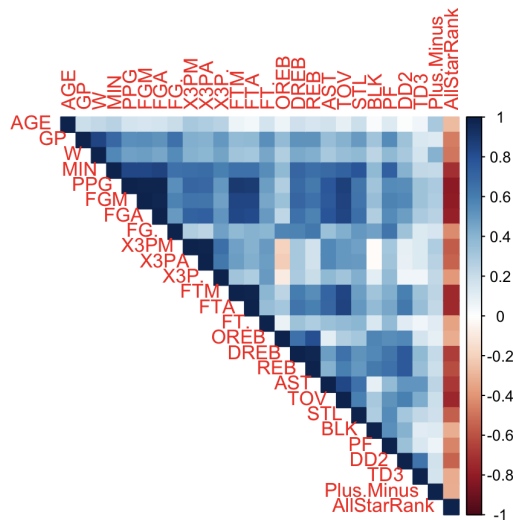
To start off our analysis, we decided to plot the major components of the box score (Points Per Game, Rebounds, Assists) as well as Wins. The primary stats within the box score typically constitute the “eye test” for a player’s skill level, while wins serve to differentiate “skilled” players from “valuable” players. We hypothesize that these variables will be important and influential for the model and as such we want to look at their distribution.



From the above plots, we can see that there is a strong negative correlation between the various stats and All-Star Rank.

ii) Correlation:

We then plotted the correlation between the variables used in the model to check for multicollinearity.



As you can see from the above figure, there is a high degree of positive correlation between some of the variables such as field goals made and field goals attempted. This makes sense as many of these variables have an inherently causal relationship. For example, more field goals attempted should lead to more field goals made. Because of this, we will likely not need two such variables which tell us similar information for the final model.

IX) Model Selection Process:

		r ²	r ² _adj	S	CP	AGE	GP	W	MIN	PPG	FGM	FGA	FG	X3PM	X3PA	X3P	FTM	FTA	FT	OREB	DREB	REB	AST	TOV	STL	BLK	PF	DD2	TD3	Plus.Minus
1	(1)	"66.4"	"66.2"	"20.41"	"56.4"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
1	(2)	"64.5"	"64.3"	"20.99"	"71.3"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
2	(1)	"68.7"	"68.4"	"19.74"	"40.4"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
2	(2)	"68.6"	"68.3"	"19.79"	"41.6"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
3	(1)	"70.4"	"69.9"	"19.26"	"29.7"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
3	(2)	"70.3"	"69.8"	"19.29"	"30.2"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
4	(1)	"71.3"	"70.8"	"18.99"	"24.2"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
4	(2)	"71.2"	"70.7"	"19.03"	"25"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
5	(1)	"72.5"	"71.9"	"18.63"	"16.8"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
5	(2)	"72.4"	"71.7"	"18.67"	"17.7"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
6	(1)	"73.4"	"72.6"	"18.38"	"12.1"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
6	(2)	"73.4"	"72.6"	"18.4"	"12.5"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
7	(1)	"74.1"	"73.2"	"18.17"	"8.4"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
7	(2)	"74.1"	"73.2"	"18.19"	"8.8"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
8	(1)	"74.6"	"73.6"	"18.05"	"6.7"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"
8	(2)	"74.5"	"73.5"	"18.1"	"7.8"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"	"

We choose model 8(1). This has the highest adjusted R^2 value and lowest C_p . The 8 variables are Age, Games Played, Wins, Points per Game, Field Goals Attempted, Field Goal Percentage, Rebounds, Assists.

Multicollinearity:

After selecting these variables, we wanted to check multicollinearity as not to have multiple variables giving us the same information.

AGE	GP	W	PPG	FGA	FG.	REB	AST
1.075780	4.523346	3.579074	29.854102	27.831255	2.322817	2.033955	2.667912

As is shown by the VIF scores above, PPG and FGA both have a high degree of multicollinearity ($VIF > 10$). Keeping this in mind, we run backward stepwise variable selection using AIC as the metric. The final step of this process is shown below:

Step: AIC=1209.51
 AllStarRank ~ AGE + GP + W + PPG + FGA + FG. + REB + AST + DD2 +
 Plus.Minus

	Df	Sum of Sq	RSS	AIC
<none>			62737	1209.5
- DD2	1	971.0	63708	1210.7
- Plus.Minus	1	1377.2	64114	1212.0
- FGA	1	2369.8	65107	1215.2
- GP	1	2898.4	65636	1216.9
- AST	1	3363.2	66100	1218.4
- AGE	1	3439.9	66177	1218.6
- FG.	1	3743.5	66481	1219.6
- W	1	6444.0	69181	1227.8
- REB	1	6796.0	69533	1228.9
- PPG	1	8985.7	71723	1235.3

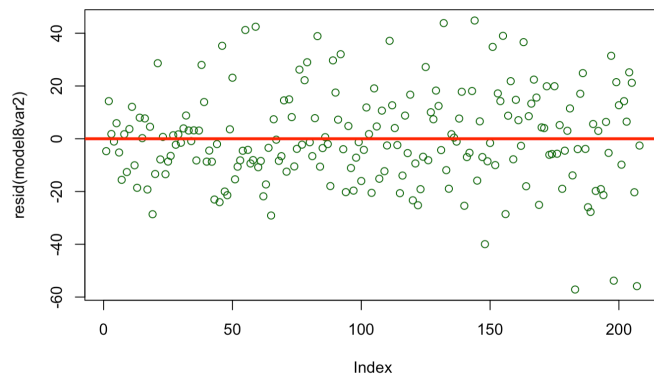
Given that PPG was in this list of variables and FGA was not, we decided to remove FGA to try to alleviate the issue of multicollinearity. We also added DD2(double-doubles) and Plus.Minus(Plus-Minus) as the backward stepwise regression showed these variables could be important. After this change in variables, we again run the VIF scores.

AGE	GP	W	PPG	FG.	REB	AST	DD2	Plus.Minus
1.198464	6.010411	5.535492	3.865659	2.442023	4.437507	2.977787	3.280499	1.829877

There is no apparent multicollinearity within this new model as all of the VIF scores are less than the threshold of 10. This new model (model8var2) also retains a high adjusted R^2 value of 0.7335.

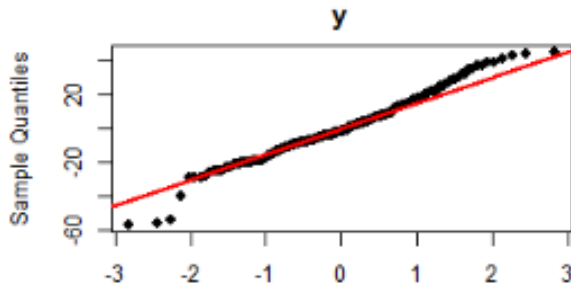
iii)Residual Analysis:

One of the assumptions in the use of a linear model is that the residuals show constant variance and do not grow or decrease in a noticeable pattern. We plot the residuals below measured against the player index(which is sorted by PPG).



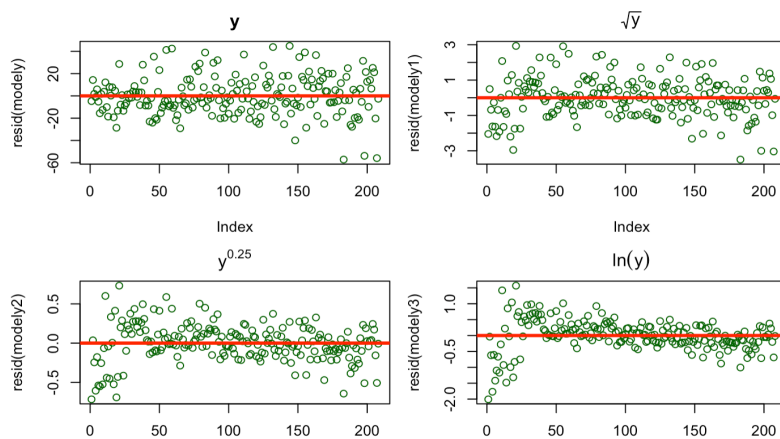
We can see that the constant variance assumption is violated as there is a growth in the spread of the residuals as the index increases.

Another assumption that we make in using the linear model is that the residuals are normally distributed. We also graph the normal probability plot which in theory can tell us if there are places in the data in which the residuals significantly veer off of the Normal Distribution. The variable y represents AllStar Rank.



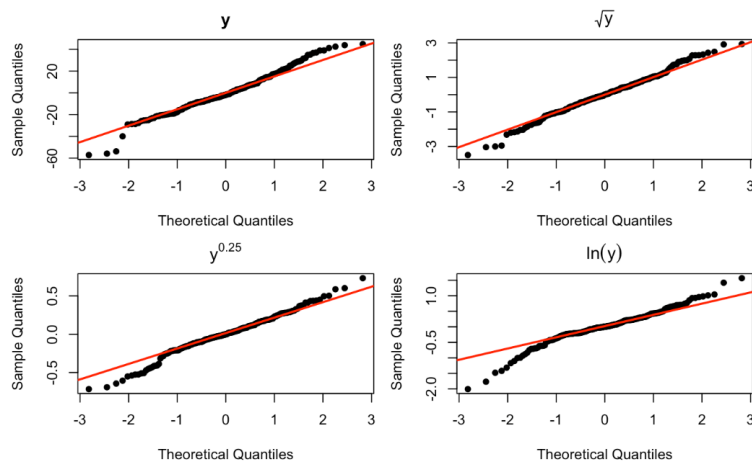
As can be seen from the above plot, the residuals have a significant difference from that of the normal distribution along with both extremes.

To fix these issues, we perform several transformations on the AllStar Rank. These transformations are square root, quartic root, and natural log. We show the residual plots under the transformations below.



While the natural log and quartic root transformation seem to make the problem of

differing variance worse, the square root transformation leads to a higher degree of constant variance. We then plot the normal probability plots under the transformations and show those results below.



Again, the quartic root and natural log transformation seem to maintain or exaggerate the problem of non-normally distributed residuals. The square root transformation however seems to generally alleviate the issue.

V) Analysis of Outliers:

Using the leverage values and their respective cutoff, we find the below players to be outliers with respect to the independent variables.

1	2	5	10	13	16	21	30	39	73	118	157	174
0.10036467	0.13564050	0.10669667	0.19285459	0.12873161	0.18174767	0.15400890	0.09671702	0.12036058	0.10411135	0.10769084	0.12994419	0.10764578
188												
0.11944982												

Utilizing studentized residuals and their respective cutoff, we can see the below players are outliers with respect to the square root of the AllStarRank.

19	21	55	183	198	207							
-2.630673	2.799746	2.659600	-3.142106	-2.675935	-2.747921							
19	21	38	46	55	59	79	83	132	148	183	198	207
-2.630673	2.799746	2.032860	2.005310	2.659600	2.187781	2.079804	2.139746	2.030545	-2.033182	-3.142106	-2.675935	-2.747921

VI) Influential Points:

We computed the quantile over the F-distribution to determine the possible influential points based on the Cook's Distance. Our results found no influential points ($\text{CooksD} > F_{0.5}$) or points with the potential to be influenced ($\text{CooksD} < F_{0.5}$ and $\text{CooksD} > F_{0.8}$)

VII) Performance Measures:

Our model had an MSE (Mean Squared Error) of 1.276448, MAPE (Mean Absolute Percentage Error) of 16.03014, and a MAD (Mean Absolute Deviation) of 0.8710102. Given the high degree of volatility in predicting all-stars, we believe that these metrics show that our model is fairly robust and has utility.

Conclusion/Summary/Recommendations:

After passing in the second half of the 2020-2021 season's player data to the optimized regression model, we observed the new AllStarRank for each player. Using the second half of the season as a proxy to determine All-Star selections for the 2021-2022 season, we identified the five starters on each team. We sorted the new AllStarRank in ascending order, accounting for position and conference to capture the following 15 highest ranked players (shown in appendix).

As seen in the table, the top two backcourt and top three frontcourt players in each conference make the starting roster. Considering the changes in the starting roster vs. the actual All-Stars we see that among Eastern frontcourt players Jayson Tatum (4th in the original standings) took the spot from Kevin Durant. Among Western frontcourt players, Zion Williamson (8th) and Karl Anthony-Towns (19th) took spots from LeBron James and Kawhi Leonard. Among Eastern backcourt players, James Harden (originally 3rd) received a starting spot over Kyrie Irving. Furthermore, Donovan Mitchell (4th) and Shai Gilgeous-Alexander (8th) received spots over Stephen Curry and Damian Lillard.

Noting that all “replacement players” were typically within the original top 10 in their respective conference and position type, we conclude that the model is not obviously misrepresentative. However, we recognize that this mechanism is likely underestimating the effect of Wins and superstardom. For example, Zion, KAT, SGA, and Fox were all players on heavily losing teams, which heavily decreases their chances of being selected as starters. Furthermore, we think the effect of raw popularity is being undervalued; when LeBron and Steph are playing at an elite level, there is a negligible chance that Zion or Donovan Mitchell would displace them as starters. SGA also figures as an outlier as he was injured post-All-Star Break and therefore deactivated as a player - meaning his averages came from limited sample size and should preclude him from selection.

Of course, the latter half of the prior season's performance might not be a perfect indicator for a player's current season, so we would recommend rerunning the numbers with updated data. One option would be to merge the second half of last season's data with data from the 2021-2022 season thus far, but that creates uncertainty as all players' performance has not been consistent across seasons: regression and breakouts as misleading trends and retirements and rookies posing outliers (as their data would be unmergeable).

We also elected not to predict reserves because of the wider variability in their selection, considering that coaches can arbitrarily choose them. In order to conduct further analysis, we could identify a grouping of players that fall into a range which makes them likely candidates for selection as a reserve. By evaluating the All Star Rank of historically selected reserves as well as potentially assigning greater weight to certain variables (e.g. age, as all-time players playing into the twilight of their careers might get a selection, as well as veterans previously snubbed, etc.) we could develop a separate model to similarly predict All-Star reserves.

Considering the reasonability of our model, we recommend that players aspiring to receive All-Star recognition consider the variables we highlighted in order to cement their status as legends and maximize earnings potential.

Appendix:

```
anova(modely1)
...

Analysis of Variance Table

Response: y1
      Df Sum Sq Mean Sq F value    Pr(>F)
AGE      1  72.23    72.23   53.8664 5.424e-12 ***
GP       1 187.49   187.49  139.8194 < 2.2e-16 ***
W        1  55.60    55.60   41.4618 8.902e-10 ***
PPG      1 705.79   705.79  526.3494 < 2.2e-16 ***
FG.      1  10.02    10.02    7.4734 0.006829 **
REB      1  30.85    30.85   23.0070 3.169e-06 ***
AST      1  14.00    14.00   10.4431 0.001442 **
DD2      1   4.61     4.61    3.4355 0.065296 .
Plus.Minus 1   0.41     0.41    0.3062 0.580644
Residuals 198 265.50     1.34
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

MSE is 1.34
      FUA      2.77017  2.22933  1.243  0.21502
FG.      0.45259  0.18547  2.440  0.01564 *
X3PM     -4.18538  20.24226 -0.207  0.83643
X3PA      2.07783  4.06922  0.511  0.61023
X3P.      0.01376  0.14946  0.092  0.92676
FTM      1.22565  18.56767  0.066  0.94744
FTA      0.49676  5.20176  0.095  0.92402
FT.      0.08997  0.07190  1.251  0.21245
OREB     1.50240  27.90161  0.054  0.95712
DREB     4.48605  27.90651  0.161  0.87247
REB     -8.29419  27.95926 -0.297  0.76707
AST     -2.49225  1.67652 -1.487  0.13886
TOV     -3.91012  3.99588 -0.979  0.32911
STL     -1.81771  5.44379 -0.334  0.73884
BLK     -3.08453  4.22817 -0.730  0.46662
PF       4.77374  2.65694  1.797  0.07404 .
DD2      0.62583  0.60964  1.027  0.30599
TD3      0.86632  1.92107  0.451  0.65256
Plus.Minus 1.08048  0.52849  2.044  0.04235 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.15 on 182 degrees of freedom
Multiple R-squared:  0.7651,    Adjusted R-squared:  0.7329
F-statistic: 23.72 on 25 and 182 DF,  p-value: < 2.2e-16
```

```
summary(model8var2)

Call:
lm(formula = AllStarRank ~ AGE + GP + W + PPG + FG. + REB + AST +
    DD2 + Plus.Minus)

Residuals:
    Min       1Q   Median       3Q      Max
-57.149 -10.173  -1.816   10.193   44.815

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 142.4542     8.9336   15.946 < 2e-16 ***
AGE          -1.0971     0.3268   -3.357 0.000944 ***
GP             0.9021     0.2923    3.086 0.002321 **
W            -1.9941     0.4664   -4.276 2.96e-05 ***
PPG          -2.7343     0.3283   -8.329 1.34e-14 ***
FG.           0.3815     0.1452    2.628 0.009269 **
REB          -3.8755     0.9994   -3.878 0.000143 ***
AST          -2.7045     1.0787   -2.507 0.012975 *
DD2           0.2329     0.4253    0.548 0.584598
Plus.Minus    0.7955     0.4965    1.602 0.110688
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.13 on 198 degrees of freedom
Multiple R-squared:  0.745,    Adjusted R-squared:  0.7335
F-statistic: 64.29 on 9 and 198 DF,  p-value: < 2.2e-16
```

Optimal Linear Regression Model & ANOVA Table

ID	AllStarRank	Name	Conference	Position	All-Star Starter?
2	1.180892	Donovan Mitchell	West	Backcourt	Yes
10	1.41396	Shai Gilgeous-Alexander	West	Backcourt	Yes
16	1.530637	Nikola Jokic	West	Frontcourt	Yes
5	1.60616	De'Aaron Fox	West	Backcourt	No
14	1.916234	Joel Embiid	East	Frontcourt	Yes
3	2.262009	Bradley Beal	East	Backcourt	Yes
8	2.271601	Giannis Antetokounmpo	East	Frontcourt	Yes
4	2.621146	Zion Williamson	West	Frontcourt	Yes
6	2.862073	Jayson Tatum	East	Frontcourt	Yes
1	3.046502	Stephen Curry	West	Backcourt	No
26	3.223893	James Harden	East	Backcourt	Yes
12	3.439643	Devin Booker	West	Backcourt	No
13	3.44279	Karl-Anthony Towns	West	Frontcourt	Yes
9	3.502565	Luka Doncic	West	Backcourt	No
21	3.510287	Trae Young	East	Backcourt	No

Predicted All-Stars

<i>Name</i>	ID	PLAYER	TEAM	AGE	GP	W
<i>Description</i>	Player ID ordered by PPG	Player name	Player team	Player age	Games played	Wins
<i>Type</i>	Integer	Character	Character	Integer	Integer	Integer

<i>Name</i>	L	MIN	PPG	FGM	FGA	FG%
<i>Description</i>	Losses	Minutes per game	Points per game	Field goals made	Field goals attempted	Field goal percentage
<i>Type</i>	Integer	Numeric	Numeric	Numeric	Numeric	Numeric

<i>Name</i>	3PM	3PA	3P%	FTM	FTA	FT%
<i>Description</i>	Three pointers made	Three pointers attempted	Three point percentage	Free throws made	Free throws attempted	Free throw percentage
<i>Type</i>	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric

<i>Name</i>	OREB	DREB	REB	AST	TOV	STL
<i>Description</i>	Offensive rebounds	Defensive rebounds	Rebounds	Assists	Turnovers	Steals
<i>Type</i>	Numeric	Numeric	Numeric	Numeric	Numeric	Numeric

<i>Name</i>	BLK	PF	FP	DD2	TD3	
<i>Description</i>	Blocks	Personal fouls	Fantasy points	Double doubles	Triple doubles	
<i>Type</i>	Numeric	Numeric	Numeric	Integer	Integer	

<i>Name</i>	AllStarRank	AllStar	Conference	Position	Plus/Minus
<i>Description</i>	Weighted average All Star Ranking	All-star selection (binary)	Player conference	Player position	Box plus/minus (point differential when player in on court)
<i>Type</i>	Numeric	Integer	Character	Character	Numeric

Variable descriptions

References:

2021 NBA All-Star Game Voting Scores. Basketball. (n.d.). Retrieved November 20, 2021, from https://www.basketball-reference.com/allstar/NBA_2021_voting.html

Players advanced. NBA Stats. (n.d.). Retrieved November 20, 2021, from <https://www.nba.com/stats/players/advanced/?sort=GP&dir=-1>.