

WAA Report

James Pinedo, II

3/2/2022

R Markdown

Overview

There is a trope in the United States that was once common and is now beginning to fade. That trope is the idea of a father and son playing baseball together. Cut off from this underlying source of interest and passion, baseball too is now fading. There are many reasons why what was once commonplace and expected is now rather rare. As Tyler Durden said in the hit film from the late 90's, *Fight Club*, "We're a generation of men raised by women." Perhaps the most popular singer of this period, Eminem, also rapped to the following thought, "My father? I never knew him. Never even seen a picture of him." I am sure it would be much more depressing study to back these specific experiences into data showing the reality for an entire generation, and I am sure that this has already been done quite well. That is not the point of this project, however.

Ironically, even while baseball is fading, the data analytics side of the sport is growing by leaps and bounds. As shown in the hit film, *Moneyball*, the analytics side of the sport has overtaken the very highest levels of each and every front office. At this point, every step that a professional baseball player takes on the field is now quantified and calculated for profit. The amount of analysis available is simply awe-inspiring. It would not be an overstatement to posit that baseball is blazing a trail not just with sports data science but also from the point of view of how man is working to understand how each and every action can be measured and tabled for future analysis. Needless to say, the amount of data is daunting.

Although this is well and good, all that data does not help anyone like me one bit. I am a father with a young son, and I would like to play baseball with my son. I also appreciate, thanks to Harvard's generous set of courses on the subject, how data analysis is an incredibly practical tool for understanding patterns in every day life and helping to visualize what set of actions are conducive to successful outcomes. I would like to pair this understanding in a practical way to playing baseball with my son, but a little league field is not wired with radar guns or cameras, nor is a little tee-baller capable hitting the ball over the fence (a "home-run" - the best thing that a offensive player can accomplish) or striking someone out (the best thing that a defensive player can accomplish), and these two outcomes are the ones that advanced baseball statistics such as "WAR" most value. In fact, baseball on the play ground has a lot more in common with the teams that played in 1875 than the teams that are currently refusing to play right now (in case the reader is unaware, baseball games are being cancelled now as professional baseball players and team owners are fighting over wages).

The point of this project is to meld these two things: a desire to play baseball with my son and my new-found data science skills. Through this combination, I would like to help my son improve at baseball, perhaps even helping other young players as well. I will look at very simple counting stats, something that a dad in the stands or coach with access to a piece of paper and a pencil, can keep track of and then compile into a single statistic that weights each stat based upon the historical value that these simple stats have contributed to generating winning teams.

I will look at every shred of available data, going back as far as 1871 and as recent as the past major league season, and I will look at as many different leagues as possible, including the famed "Negro Leagues"

of old and leagues that are long defunct. This compilation of data will never truly be completed, as any team data is useful, and baseball has been played anywhere and everywhere. However, since I need to present my findings at some point, I opted to use all available team hitting (offensive) data that could be scraped from Stathead Baseball's website (stathead.com). This data was then compiled into a spreadsheet and is available for use on my project's github page, you can also access it on the Stathead site here - <https://stathead.com/sharing/UqErJ>

Key steps that were performed on this data set were to first even out the data so that it could be compared on a average statistic per game basis. For example, modern Major League baseball plays a 162 game season. Most teams in this data set did not, so it was important to look at the averages of each stat on a per game basis. I first used linear regression to use these averaged-out statistics as predictors, and then used K Nearest Neighbors, and finally Random Forest.

Of these, the linear regression model was the most illuminating, for it showed how stats that are valued highly in the major leagues today - for example the walk (or "base-on-balls", which is when the defensive team allows a hitter on the offensive team to reach first base without that hitter team actually hitting the ball), are not effective predictors of winning teams from a historical perspective. I will speak more on that in the methods/analysis section.

In brief, this study was able to use simple offensive stats to predict how successful (or not) a team would be, by predicting the "Wins Above Average" (WAA). Average wins were determined by taking the winning percentages of all the teams in the data set and averaging it out to equal a 162 game schedule. Using only basic hitting statistics, I was able to predict the correct WAA within an average of 13.3 games given a 162 game season. This is an improvement of the standard error of 16 games should we guess the average at all times. I will speak more on how this could be improved in the conclusion.

Methods/Analysis

Before starting out on the methods used to start predicting, I will comment on the methods used for obtaining the data and then cleaning it for use in analysis. I relied on the data found on the Stathead Baseball site for this project. Unfortunately, the data there is split up into many pages, so I downloaded the data from each page and combined them into a spreadsheet. This spreadsheet was then uploaded to Github and then imported into my R project using the RCurl package. The data was a little messy upon import, with the "Year" column being imported with an improper name, but a simple mutate function fixed that.

Now that the data was in the project, I found the average win-loss percentage ("WL_perc", the number of wins divided by number of games, an near universally tracked average that was already in the data set) for all teams, which was 0.4964048 (teams on average won about half their games, which makes perfect sense - every time one team wins, another has to lose). Then a column was added to show the wins above average percentage (WAAP) for each team. For example, if a team played 10 games and won seven of them, their WL_perc would be .7, thus their WAAP would be 0.2035952. This average is a little difficult to visualize, so I opted to add another column to the data, one that took WAAP and multiplied it by 162, this showing what each team's WAA would be should they have played 162 games or "WAA_162". This is the column that we will be trying to predict moving forward, with using simple counting stats as predictors.

A quick note on the predictors (you may skip this paragraph if you are already familiar with baseball), when playing baseball the way that a team wins is by scoring more "runs" than the other team. Each team has nine tries to score runs and they switch off between being the offensive and defensive side throughout these nine tries (each set of tries by both teams is called an "inning"). While one team is on the offensive, they are the "hitting" team, as this is the time that they are attempting to hit the baseball from the starting position called "home plate". The defensive team is called the "fielding" team - this is when their player, called the "pitcher", is throwing the baseball - trying not to let the hitting team hit it but also needing to keep the ball within a narrowly defined area (called the "strike zone") over home plate. A baseball field is shaped in a diamond with four safe areas (or "bases") that a player on the hitting team may touch safely, and if the fielding team catches them off of this base they are called "out" and must leave the field. A run is scored every time the hitting team manages to safely run around the baseball diamond without being called out three times in one inning. The best chance that the hitting team has at running around these bases (starting from the home base [called "home plate"], to first base, second base, third base and then back to

home) is by hitting the baseball where the fielders are not. This is called a “hit”. Each hitters attempt to do this is called a “plate appearance” (PA). Should the hitter hit the ball where the fielders are not and make it all the way to second base, that is called a double. Should they make it to third base that is called a triple. Should they manage to run around all the bases in single hit that would be called a home run. There are other ways that a hitter may reach base, mostly due to the fielding team making a mistake, but for the purposes of this study, I will only explain only one of those - the hitter may reach base without hitting the ball when the opposing team’s pitcher fails four times to throw the baseball within the strike zone. Should this happen, the batter is awarded with first base and this is counted as a “walk” or “base-on-balls”. A hitter is credited for a “run-batted-in” if they hit the ball and allow one of their team mates that were already stationed on one of the bases to run home. All of these actions are counted and have been counted reliably since the 1870’s and are of great use to this project.

Of course, I needed to average out these stats to reflect that most teams in the data set did not play a 162 game season. To that end, I divided the hits (H), doubles (X2B), triples (X3B), home runs (HR), runs batted in (RBI), and walks (BB) by the number of games that each team played. Easily tracked by any parent sitting the stands watching their child play, these averages will be used as predictors, as well as the following averages: the average amount each hitter reached base per plate appearance (on-base percentage or OBP); the batting average (BA) or average hit per “at-bat” (at bats are plate appearances minus sacrifices [an action that was not tracked reliably in the early days of baseball, so we did not use it here] or walks); and the average amount of bases that the player manages to run around per attempt (called the slugging percentage or SLG). Again, these averages are easy to determine with simple division, so these are the nine predictors that we will use.

A note on why I didn’t choose other numbers in the data set for predictors - the point of this project is to help us understand better through data analysis the value of each action that any baseball player (whether playing in the major leagues on on a local little league) is able to accomplish. This rules out the use of several predictors that would definitely make our algorithm more accurate, but would defeat the purpose of the project. The easiest example of this is the year. At the end of the analysis, I did use the year as a predictor, and was able to be more accurate, but this would be of no use to a parent in the stands, so I do not count it in the larger project. That guiding principal - hoping to be of some practical use, also rules out the use of games, plate appearances, and league - none of these would help us be able to determine the value of an action that my young son playing baseball is taking on his baseball team, for example.

Now that I have predictors that are evened out for all the teams in the data set - I can go to work!

In order to avoid over training a linear regression model, I split up the data into a test set and training set. The test set was made from a random sample of 10% of the overall data.

First, I set a baseline. I found that if I guessed the average WAA_162 (which on our train set is -0.001531599) for every team in the test set, I would be on average about 16 games off (I use the root mean squared error or “RMSE” formula for this). An RMSE of 16 is what we will need to get better.

Results (and Methods continued)

Now, I started with all nine predictors to train a linear model on the train set, and then used the predict function with the fitted model to make predictions on the test set. I then checked the RMSE of the those predictions on the actual WAA_162 of the test set and found that this gave me a RMSE of 13.40354.

Thinking that I perhaps had some noisy predictors, I then started cycling them out one by one to see if I could improve my RMSE. I did this several times with different slices of the data, for the slice that I presented in the accompanying code, I found that only by removing the RBI/G (rbi_g) predictor did the overall RMSE improve to 13.29213. Removing any of the other predictors did not make for a better score. There were a few notes of interest at this point, however. By far the predictor that was the least helpful was walks, which only improved our predicted WAA_162’s RMSE by 0.00492. In fact, in other slices of the data this predictor was not helpful at all. Also of noted interest is that in triples/game only improved our RMSE in this slice of the data by 0.02256. In other slices of the data it was not helpful either.

Given than our linear model only resulted in an improvement to an RMSE of 13.29213, meaning that we are about 13 games off on average, I went beyond linear regression in hopes of finding some better results.

First, I tried K-Nearest Neighbors. Using our top 8 predictors for this data slice, I fitted a model using the `knn3` function on the train set. I then used this fitted model to make predictions on the test set and found that these predictions resulted in a RMSE of 16.70832. Not what I was hoping for, but perhaps I could choose a better tuning parameter for this function? Understanding that this tuning parameter will randomly sample segments of the data many times, I opted to train it on the entire data set. I used the `train` function to with the tuning parameter set to `k = seq(5, 99, 2)`. I then plotted the results to find that the K value with the lowest RMSE was `K=61`. This tune resulted in a RMSE of 14.60331 - not an improvement on our linear model.

So I turned to the Random Forest package. Again understanding that it would take multiple random samples of the data, I used the entire data set for it. Using the top 8 features, I made predictions on the data set and found that this resulted in an RMSE of 13.98472 - better than KNN but not than the linear model. Thinking that the nature of Random Forest might help with noisier predictors, I tried a random forest fit with all nine predictors. This resulted in a better RMSE of 13.66735.

Finally, to see the outcome of adding the year and number of games to the algorithm, I added those two features to the random forest fit and found that the predictions that results improved. The resulting RMSE was 12.84503. This makes sense. Random Forest was tracking that there is a correlation in accurately predicting `WAA_162` based upon the year and the number of games played. The algorithm was now tracking that the game of baseball was changing with the years, but given that a young boy playing baseball with his local team has statistically more in common with a team of 1875 than a major league team of 2021 (he will hit an average home run total closer to 1875 than now, for example) this improvement only shows that I am using the data analytics tools at my disposal properly, but not much more.

Conclusion

In summary, the predictions resulting from an algorithm built on a linear model of eight predictors (Home Runs per game, batting average, on-base percentage, slugging percentage, hits per game, walks per game, triples per game, and doubles per game) gave me the best results. However, there is much more work that can be done. Due to the constraints in my data set, defensive data (pitching and fielding) was not used. This algorithm could easily be plugged into other data sets (the Lahman data set for example) that have pitching and fielding data. Armed with this additional data, I believe it would be possible to get the RMSE below 8 games, if not, more. In addition, this algorithm is simple enough for others to use to train on their local league's stat sheet, meaning that it is easy to check to see if it applies for your individual situation.

Once you have a predicted `WAA_162` that can accurately reflect the reality of your situation, I would suggest that we take the amount that each stat was able to accurately lower the resulting RMSE and weight that stat accordingly. Then that weighted score could be compiled into a single individual player stat that is able to simply and reliably show how much that individual player has contributed to wins above average - that is, the success of the team as a whole. For example, if a young player had 5 walks, we would multiply 5×0.00492 to their over all WAA score. This will allow us to compile and weight statistics in way that has more direct bearing than other modern advanced statistics.

This may not seem like an earth shattering study or find, but I believe that it can be a practical help to me and my son as we play baseball together. I hope that you, the reader, may find it of some use to you as well.

Including Plots

Below is a plot of the K-Nearest Neighbors tune. In future studies, features pitching and fielding statistics would be sure to help.

```
## Loading required package: tidyverse
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```

## v ggplot2 3.3.5      v purrr  0.3.4
## v tibble  3.1.3      v dplyr  1.0.7
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   2.0.0      v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

## Loading required package: caret

## Loading required package: lattice

##
## Attaching package: 'caret'

## The following object is masked from 'package:purrr':
##
##      lift

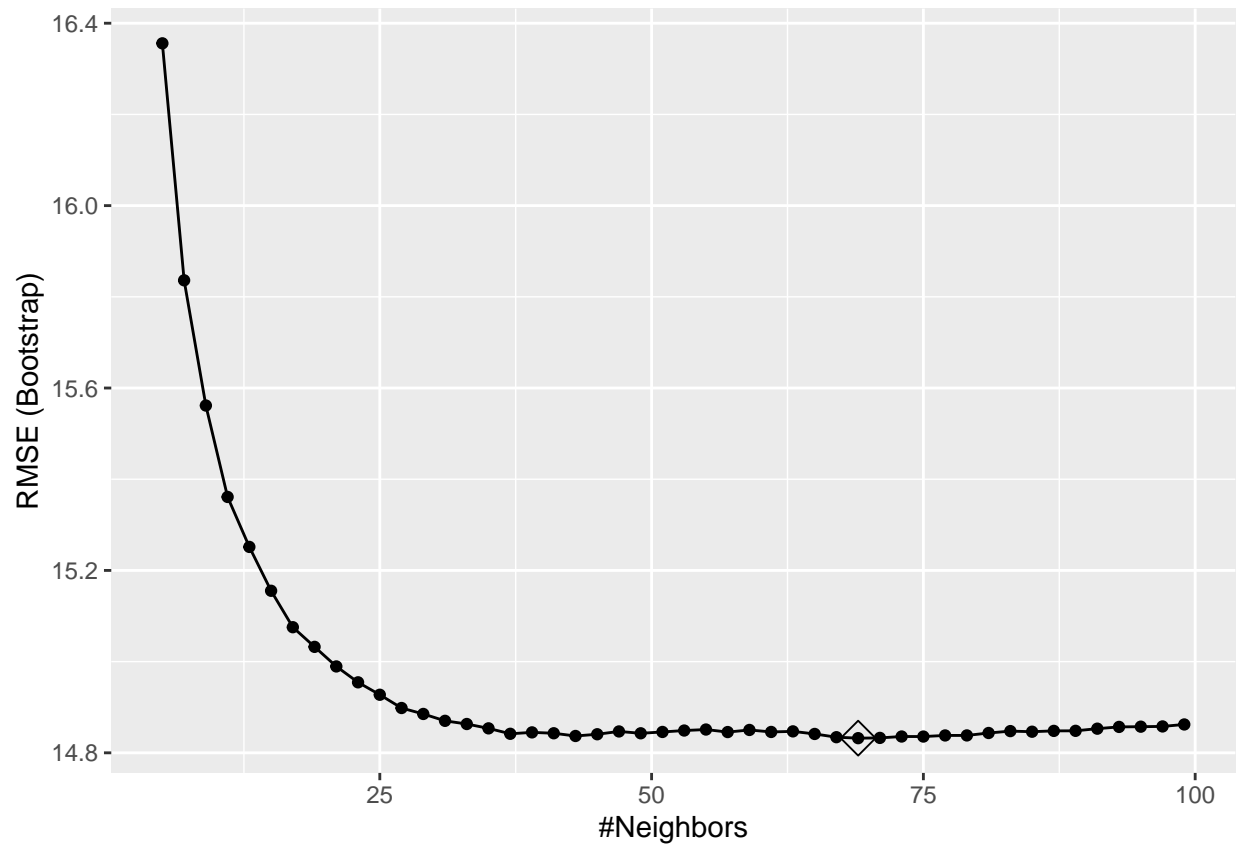
## Loading required package: RCurl

## Warning: package 'RCurl' was built under R version 4.1.2

##
## Attaching package: 'RCurl'

## The following object is masked from 'package:tidyr':
##
##      complete

```



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.