

# Using Logistic Regression to Predict Baseball Playoff Probabilities

Matthew Wankiewicz

26/11/2020

## Abstract

The book moneyball highlighted a statistics revolution that was occurring in baseball in the early 2000's. In this report, logistic regression to determine whether the advanced stats that stemmed from this revolution are useful for predicting a team's success. After running the model, we find that it successfully determines if a team will make the playoffs (or not) 88 percent of the time. This is a significant finding as it shows that these advanced metrics are extremely useful and teams should fully shift their focus on these metrics.

**Keywords:** Baseball, Moneyball, Logistic Regression, Analysis

Code supporting this analysis can be found at: <https://github.com/matthewwankiewicz/moneyball>

## 1 Introduction

In 2003, Michael Lewis wrote a book called Moneyball<sup>1</sup>. Moneyball is the story of the Oakland Athletics who, after losing players to the richer teams in the league, decided to focus on using the misfits of baseball to try to win a World Series. These "misfits" were players who never received support from teams because their traditional stats didn't look good but they excelled in the "advanced" stats that mattered.

In this report, I will use logistic regression to test if the basis of Moneyball was correct, can you throw payroll out the window and just focus on advanced statistics? Is it possible to build a cost-effective team that will make the playoffs in a baseball? I decided to put the benchmark of a successful season at making the playoffs because once a baseball team makes the playoffs, anything can happen. Often times, teams will just barely make the playoffs and still be able to win the championship, like the 2019 Washington Nationals.

In order to conduct this analysis using logistic regression, I have set the response variable to whether the team makes the playoffs and the predictors are various advanced metrics. The metrics I plan to use include weighted on base average (wOBA), fielding independent pitching (FIP), opponent batting average balls in play (oBABIP), walks plus hits divided by innings pitched (WHIP) and strikeout rate (SO%). Although these metrics do look confusing, I will explain what they mean and how to calculate them in the data section of the report.

After using the model, we find that the advanced metrics are fairly accurate in predicting a team's success, displaying an 88% accuracy. These results are significant because it shows that the statistical revolution displayed in Moneyball should be here to stay. Over the last few years, many teams have increased their budgets for analytics and those teams are appearing to be more successful than others. In a FiveThirtyEight article from 2016, it is shown that only 3 teams had analytics departments with more than eight people (Lindbergh and Arthur [2016]) and two of those teams were the best teams in baseball last year. This result helps confirm the purpose of the model, showing that advanced analytics can lead teams to more success.

The data collected for this analysis was from the Lahman package in R (Friendly et al. [2020]) and Baseball-Reference (LLC). Fangraphs (Fangraphs) was also used to help calculate one of the statistics.

---

<sup>1</sup>Moneyball was actually the main reason I decided to get into statistics. Applying baseball with numbers seemed like an absolute win to me

This data analysis was conducted using R (R Core Team [2019]), and in particular the packages Tidyverse (Wickham et al. [2019]) and Cowplot (Wilke [2019]) and was compiled using R markdown (Xie et al. [2018]). This paper includes 4 sections not including the introduction. The first section covers the data I used for my report, including a few plots to show how the distributions and relationships between various stats. The next section is about the model I used, this will include an in-depth breakdown of logistic regression. Next, I will display the results of my model, this will include the coefficients of the model and an interpretation of the model, along with an application of the model to see how accurate its predictions are. Lastly, there is a section discussing what we learn from this model and some next steps. Was Moneyball right? Can advanced metrics lead us to the playoffs?

## 2 Data

The data collected is from the `Lahman` package in R (Friendly et al. [2020]). The `Lahman` package contains yearly player statistics going all the way back to the late 1800s. For this project, I will only use the data from 1920 and onward because those years have the most observations for statistics like payroll data. When looking for payroll data, I found that some years did not have complete payroll data for all teams so these observations were removed<sup>2</sup>. After removing the years with missing entries, there were 2096 observations present for this analysis. In the `Lahman` package there are many different datasets present, for this project I will be using the `Teams` dataset. In order to collect the payroll data, I created a scraper that collected yearly Team stats from Baseball Reference (LLC). The scraper takes the data from each year and saved it into a larger dataset<sup>3</sup>. The payroll data will be crucial for my report because the basis of Moneyball is fielding the best possible team by spending the least amount of money.

As mentioned above, the `Lahman` package contains the `Teams` dataset which was used for the calculation of the basic statistics along with the more advanced ones. The statistics used for this analysis were Fielding Independent Pitching (FIP), Weighted on base average (wOBA), batting average balls in play of the opponents (oBABIP), walks plus hits divided by innings pitched (WHIP) and the strikeout rates (K%). The metrics used for this analysis will be explained below.

- **wOBA:** Weighted on base average is considered to be a “better” version of batting average. wOBA was created by Tom Tango and he writes more in depth about it in “The Book” Tango et al. [2014].<sup>4</sup> wOBA is a weighted average of each of the ways a player can get on base (walking, getting a single, double, triple and a home run) and is divided by the ways a player gets a chance to hit (at bats, walks and sacrifice flies). Usually the coefficients of wOBA change by year but for this report, I used a manipulation of other advanced statistics that were described in “The Book”,  $(2 * OBP + SLG)/3$ <sup>5</sup>
- **FIP:** Fielding independent pitching estimates a player’s run prevention if they did not have a defence, as opposed to ERA which includes defence. The equation for FIP is simple,  $(13 * HRA + 3 * BBA - 2 * SOA)/IP$  where HRA is home runs allowed, BBA is walks allowed and SOA is the number of strikeouts a player gets. As Fangraphs says, “FIP is a measurement of a pitcher’s performance that strips out the role of defense, luck, and sequencing, making it a more stable indicator of how a pitcher actually performed over a given period of time” (Slowinski).
- **oBABIP:** Batting average on balls in play against is pretty much the opposite of FIP. It looks at the batting average of only the balls that the defence has to field, as opposed to FIP which only looks at outcomes that the defence cannot control. The equation for oBABIP is:  $(HA - HRA)/(IPouts - SOA - HRA)$  where HA is hits allowed, HRA is home runs allowed, IPouts is the number of outs a team gets while pitching and SOA is the number of strikeouts a team gets while pitching.

---

<sup>2</sup>I did attempt to take the average payroll and apply it to the other teams, but some years only had 2 of 16 payrolls recorded

<sup>3</sup>This scraper can be found in the scripts folder of my GitHub repo.

<sup>4</sup>The Book is seen as the best book for aspiring baseball analysts or “saberists” to get into the MLB analysis game.

<sup>5</sup>OBP is on base percentage, the rate a player gets on base (includes walks). SLG is slugging percentage, the total bases (TB) a player gets divided by at-bats (a single = 1 TB, double = 2, triple = 3 and Home Run = 4).

- **WHIP:** Walks plus hits divided by innings pitched is kind of self explanatory. It takes the amount of walks a team gives up, plus the number of hits and divides it by the number of innings pitched:  $(BBA + HA)/IP$ . WHIP is useful because it tells us how good a team (or pitcher) is at preventing runners from getting on base. If a team allows more people to get on base, they are more likely to give up runs and losing a game.
- **SO%:** Strikeout rates tell us how often a team (or batter) strikes out when they are batting. Some players have a tendency to only hit the ball far or not hit it at all while other players hit the ball almost all of the time. SO% helps show what type of gameplay a team relies on and it will be interesting to see how it affects a team's playoff chances. To calculate SO% we take the number of strikeouts a team got while batting and divide that by the number of times they bat  $SO/AB$
- **ERA:** Earned Run Average represents the amount of runs a player gives up over a full game. Although I will not be including it in my model, I will be using it to show why FIP is a good metric so I thought it would be best to include it. ERA is the traditional measure of success for pitchers and if they have a low ERA, they're seen as great players.

Now that the metrics we're focusing on have been discussed, we can look at some plots to help visualize the data.

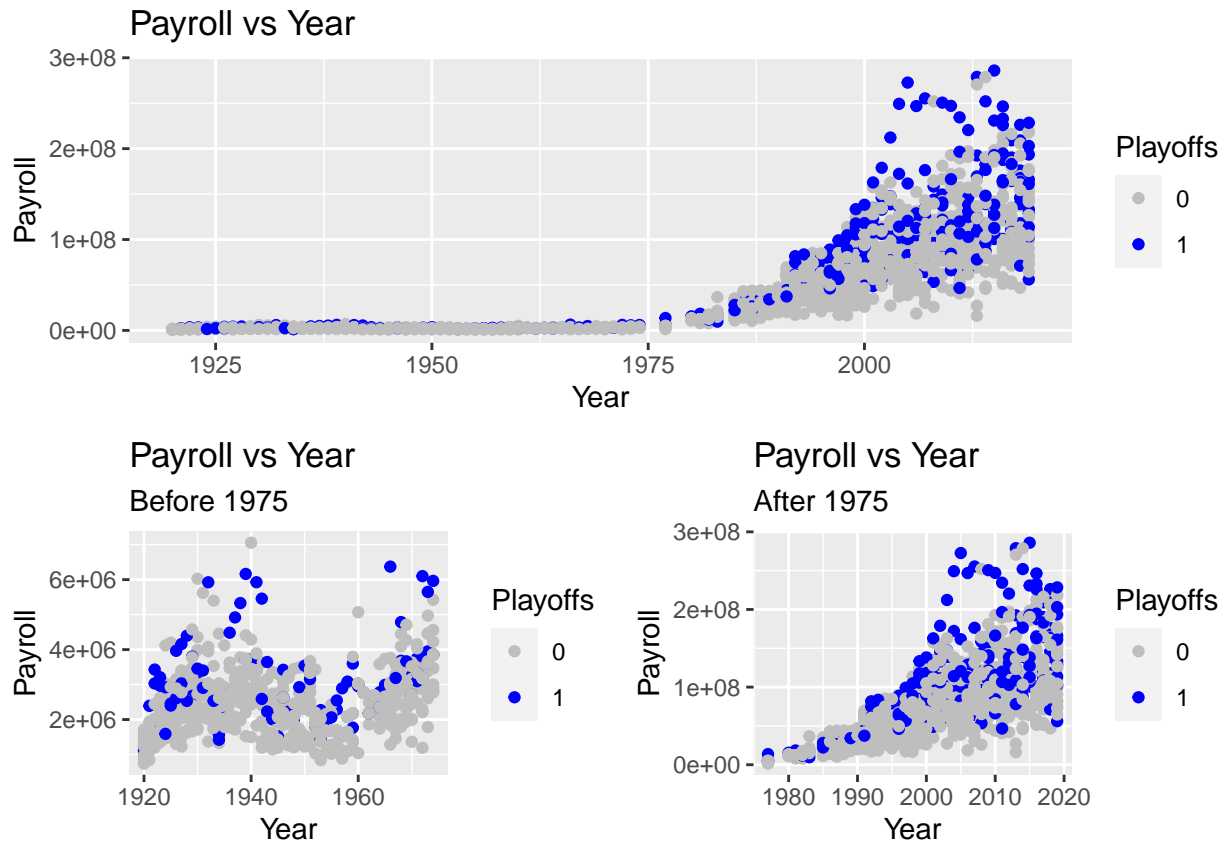


Figure 1: The increase in Payroll over the Years

Figure 1 shows us the increase in payroll over the years, starting at 1920 and ending at 2019. We see that up until 1960, payrolls were fairly constant but after the 60s we have seen a rapid increase in payroll. We can also see colours which represent if a team made the playoffs or not (light blue representing yes and dark blue meaning no). It is clear that in most seasons, teams that spend more make the playoffs. The salaries have been adjusted to represent the dollar in terms of what it would be in 2019 using inflation rates from the US Inflation Calculator (cit [2020]). Clearly, the adjustment does not do much as we see that teams were likely

just spending more compared to the past.

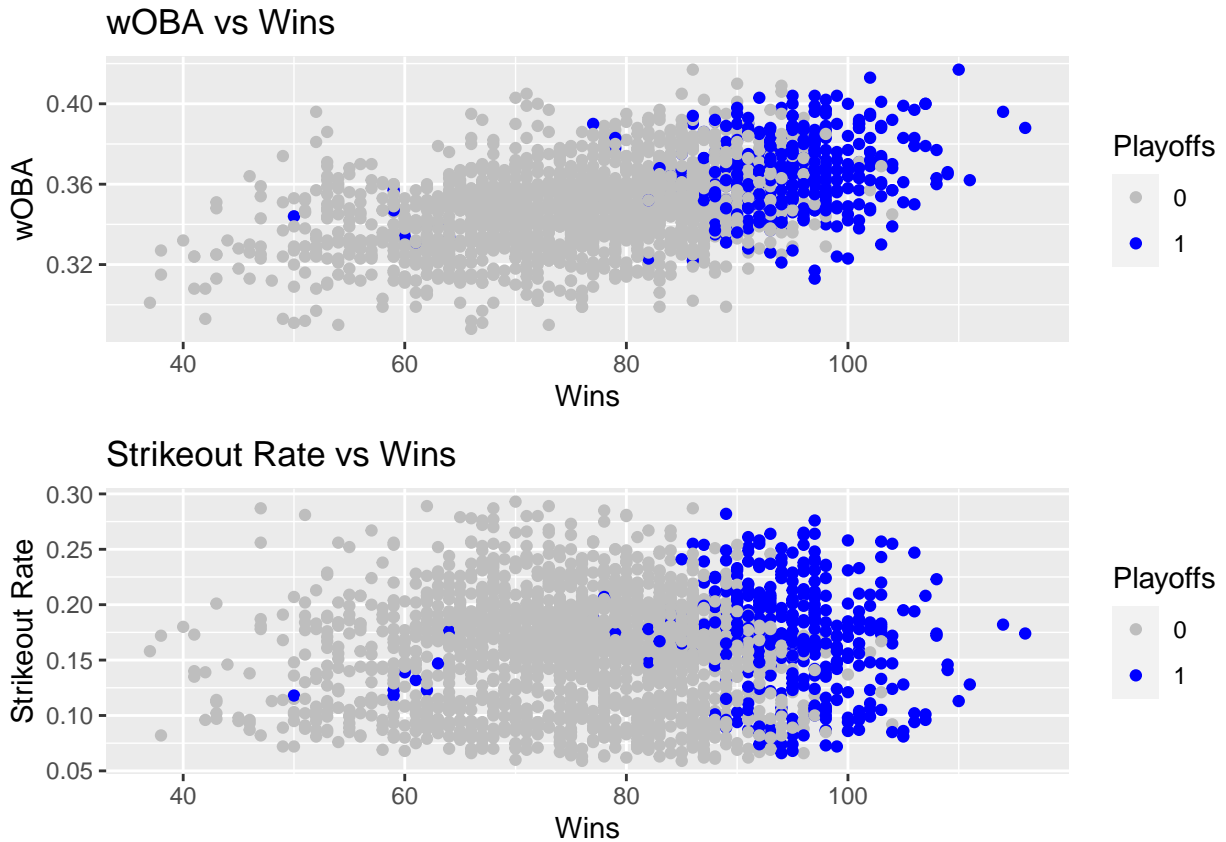


Figure 2: Relationship between Advanced Offense Metrics and Wins

Figure 2 shows the relationship between offensive advanced metrics and the amount of teams a win gets, along with colours representing a team's ability to make the playoffs. From the first graphs, we see that as wOBA increases, the amount of wins increases and the number of teams making the playoffs increases as well. We also see that there does not appear to be too much of a relationship between the number of wins a team gets and their strikeout rate. This makes sense because some teams are built on hitting home runs and strikeout a lot, versus teams who play more strategically and don't hit the ball as hard.

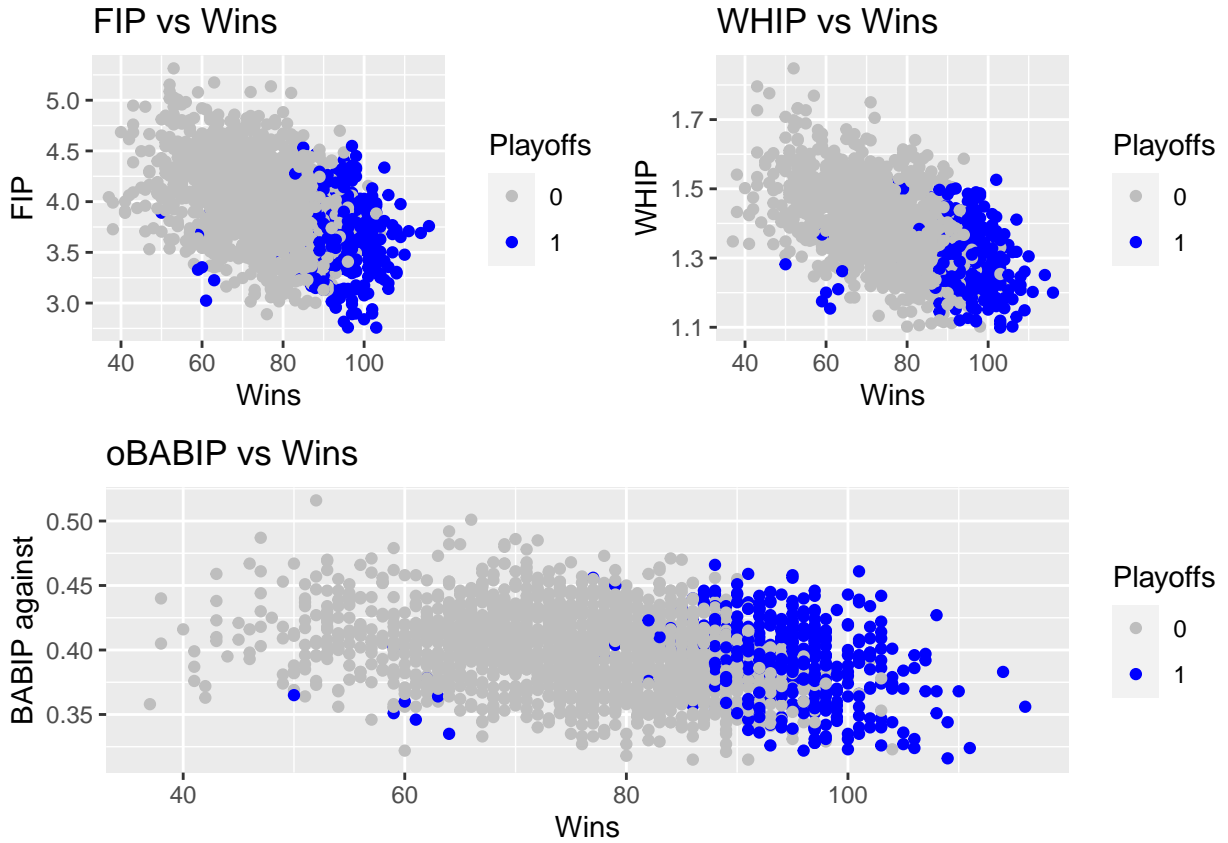


Figure 3: Relationship between Advanced Pitching Metrics and Wins

Figure 3 shows the relationship between some advanced pitching metrics and the amount of wins a team gets in a season. We see that WHIP appears to have a very strong impact on wins and if a team can keep their WHIP as low as possible, they're more likely to make the playoffs. It appears that FIP has the next strongest relationship with wins, better teams have lower FIPs. Lastly, it appears that BABIP against does have an impact on the amount of wins a team gets, just not as significant as FIP and WHIP.

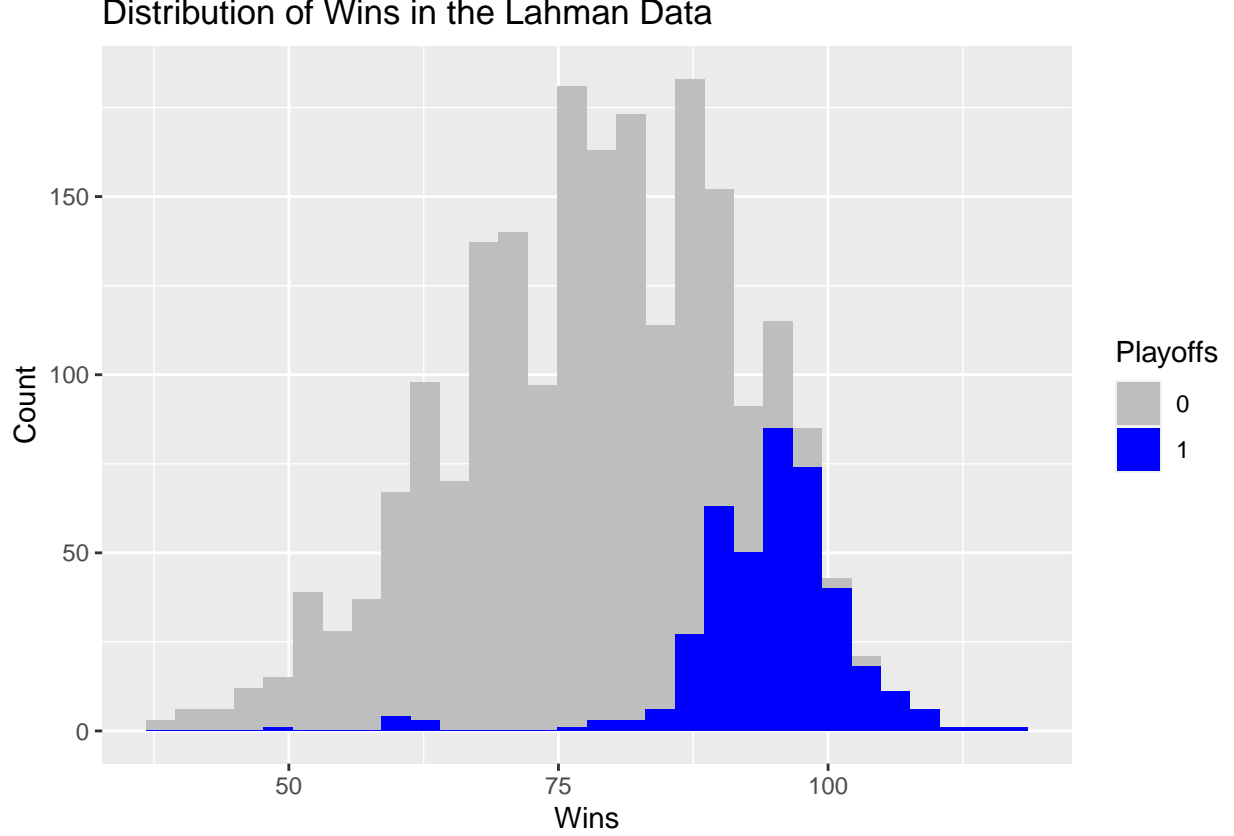


Figure 4: Distribution of Wins for Lahman data

Figure 4 shows the distribution of wins for teams present in the Lahman data. We see that majority of team make the playoffs when they win about 85 games although there is some spread in the distribution of wins for playoff teams. There are also some teams present which made the playoffs after winning less than 75 games. This occurred because of dispute between the players and the league which led to a 50 day strike and led to teams winning less games in total (Bumbaca [2020])<sup>6</sup>.

### 3 Model

In order to get playoff probabilities, I will use logistic regression. Logistic regression is a type of statistical analysis that instead of giving a prediction as to what a value will be, it gives the probability that an event will happen. In this case, logistic regression is used to determine the probability that a team will be making the playoffs. Logistic regression takes a sum of various factors and then after manipulating the equation, we will end up with a probability. The logistic regression will be in the form:

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \beta_0 + \beta_1 x_{FIP} + \beta_2 x_{wOBA} + \beta_3 x_{oBABIP} + \beta_4 x_{SOrate} + \beta_5 x_{before75} + \beta_6 x_{payroll} \quad (1)$$

where  $x_{FIP}$  represents the team's FIP,  $x_{wOBA}$  represents the team's wOBA,  $x_{oBABIP}$  represents the team's BABIP against,  $x_{SOrate}$  represents the team's strikeout rate,  $x_{before75}$  represent whether the team in question was playing before 1975 or not and  $x_{payroll}$  represents the team's payroll.

<sup>6</sup>The players eventually returned and the playoffs continued

In Equation (1), the  $\beta$  values represent a coefficient determined by the `glm` function, they will either be positive or negative, depending on if a higher  $x$  value increases or decreases a team's chances of making the playoffs. Each  $\beta$  value will be multiplied by its corresponding  $x$ , for example, as we saw in the data section, teams with higher  $wOBA$ 's tend to make the playoffs more often, so the  $\beta$  value for  $x_{wOBA}$  will likely be positive. Once all of the coefficients are multiplied we take the sum of the equation and insert it into the equation below:

$$\frac{e^{sum}}{1 + e^{sum}} \quad (2)$$

Equation (2) is just a manipulation of equation (1), where  $e$  is the exponential equation and  $sum$  is the sum of the right side of equation (1). We can see that as the sum of the equation increases, the probability of a team making the playoffs increases as well.

The logistic regression model is run using the `glm()` function in R (R Core Team [2019]). The decision to run this model over other models like linear regression was made by the fact that we are predicting a binary variable of a team's success in the season. Since there are only two possible options our data will likely follow an S shape and a straight line equation will not be helpful to model this relationship. As a straight line regression will likely miss many data points.

One weakness present with a logistic regression model is that since the response variable must be binary (either make the playoffs or don't), we cannot expand it to determine if a team wins their division or just scrapes into the playoffs.

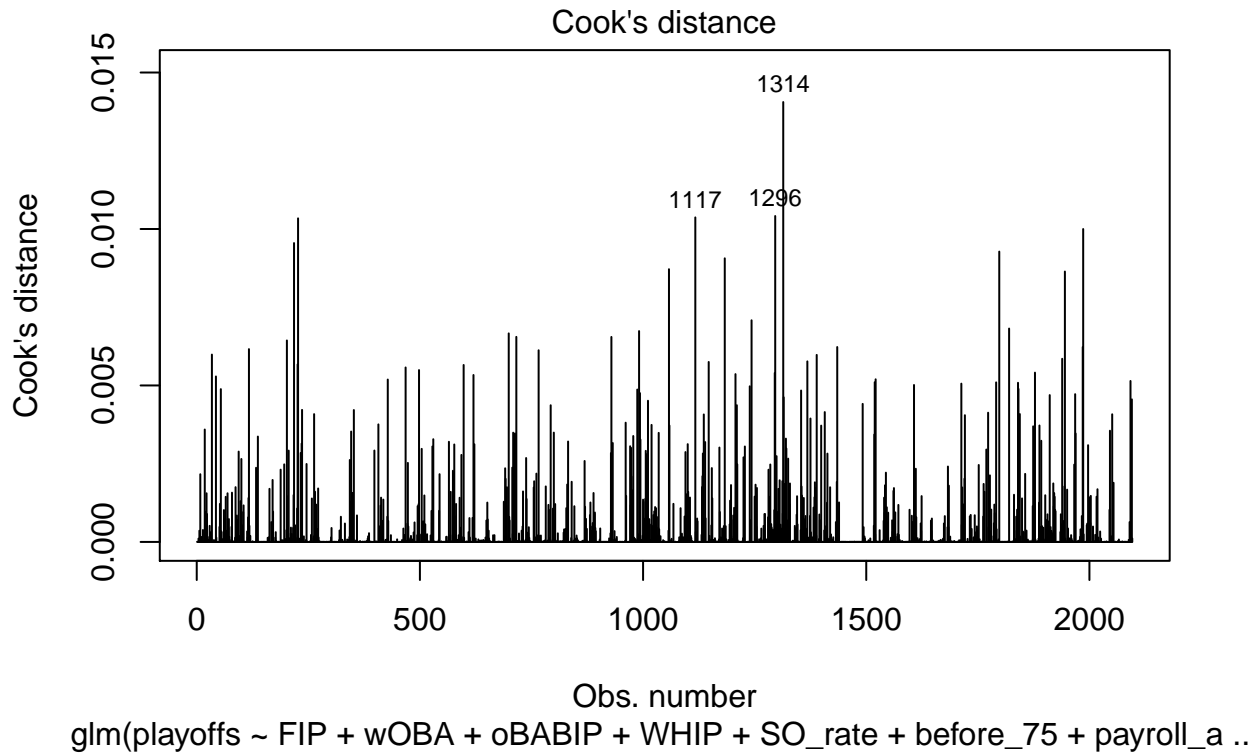


Figure 5: Cook's Distance for the Model

Figure 5 shows us the observed cook's distances for the model. Cook's distance tells us how far a point is from the predicted value of it, telling us which points negatively impact our model. We see that these cook's

Table 1: Regression Model Coefficients

term	estimate	std.error	statistic	p.value
(Intercept)	-6.69671	1.88881	-3.54547	0.00039
FIP	-2.86375	0.43896	-6.52392	0.00000
wOBA	113.29413	6.73092	16.83189	0.00000
oBABIP	-38.70566	6.28537	-6.15805	0.00000
WHIP	-8.18975	2.36440	-3.46377	0.00053
SO_rate	10.30124	3.02687	3.40327	0.00067
before_75	-0.40806	0.24141	-1.69034	0.09096
payroll_adj	0.00000	0.00000	2.67947	0.00737

Table 2: Accuracy of the Model

Data Type	result	Accuracy (%)
Data	18.989	-
Model	15.076	88.263358778626

distances are very low (most are under 0.01), since our distances are fairly low, we can conclude that the model is fitting the data well.

## 4 Results

Table 1 shows us the coefficients from our logistic regression model<sup>7</sup>. We can see that all variables have statistically significant p-values, meaning that we know for sure that the variables have an effect on a team's playoff chances. An important point to note from Table 1 is that the p-value for payroll is less significant than the advanced metrics. This means that although it is significant to predicting success, we are less confident in its abilities compared to the advanced metrics.

We see that using the logistic regression equation, the lowest possible probability of making the playoffs is  $4.23 \times 10^{-5}$  percent. This was calculated using the worst possible stats present in the data used and it is extremely unlikely a team would be this bad! The highest possible probability of making the playoffs is 99.99 percent. This was taken from the maximum stats present in the data, and although it would be tough for a team like this to exist, there have been some teams that were close.

Table 2 displays the proportion of teams which made the playoffs from the observed data vs the results of our model's predictions. We can see that the model is very strict with its results, only predicted about 15 percent of teams to make the playoffs, less than the observed 18.9 percent.

We can also look at the overall accuracy of the model. In order to find the accuracy of the model, I compared the predicted results to the actual results and we find that the model successfully predicted if a team would make the playoffs about 88 percent of the time. These results are very significant because it shows that the model can fairly successfully determine what a playoff team is. This also bodes well for our attempt to see if the results for the 2020 predictions will be correct.

## 5 Discussion

Now that the model has been created (and tested) and we have our results, what can we make of them? Well, the first major point we can see is that when using advanced statistics, we can accurately predict if a team

<sup>7</sup>Created using the `tidy` function from the `broom` package (Robinson et al. [2020])



will make the playoffs 88 percent of the time. This result is truly significant because it means that if a team put their focus on acquiring players who excel in advanced statistics but are average in the ‘regular’ ones, they can drastically improve their chances of making the playoffs.

When looking at our batting stats, we see that wOBA is very influential in determining a team’s playoff status. Taking a quick look at Fangraphs’ wOBA leaderboard from 2019<sup>8</sup> we see that 6 of the top ten players in wOBA were not in the top 10 for batting average. If teams are still evaluating players using traditional statistics, they have a much higher chance of paying for a player who may not have a significant impact on their playoffs chances. If teams make the shift to stats like wOBA, they certainly won’t be hurting their chances at making the playoffs.

We can also take a look at the strikeout rate of teams. As we saw in Figure 2, it didn’t appear as if a team’s strikeout rate had any impact on the amount of games a team won. This makes sense due to the various different play styles that have been popularized over the years. Some teams have adopted the style of go big or go home when they go up to bat, either hit a home run or strikeout. This type of play has been fairly popular over the past few years and was especially popular during the “Steroid Era” in the early 2000s when the league saw an rise in offense due to steroids (Woltring et al. [2018]). Table 1 shows that having a high strikeout rate actually improves a teams chances of making the playoffs, so according to the model, the go big or go home strategy appears to lead to success.

Next, we can also take a look at FIP, now that we know it is significant. FIP is much more difficult to evaluate compared to the traditional stats because of how it is calculated. As we saw in the data section, FIP is calculated only by true outcomes: home runs, strikeouts and walks. Since FIP only deals with these “true outcomes”, it tends to overlook players who don’t rely on strikeouts but instead rely on their teammates to help them out. Regardless, finding players with an above average FIP<sup>9</sup> and below average traditional metrics is fairly simple. Once again, from fangraphs we find that there are many pitchers with high ERAs and above average FIPs. This is significant because FIP predicts success better than other metrics. In a Sporting News article from 2017, John Edwards talks about a player named Michael Pineda’s 2016 season. He started off the season with a high ERA but a low FIP and as the season went on, his ERA began to decrease (Edwards [2017]). This is one of many examples where FIP is more valuable than the traditional statistics and should be emphasized when evaluating teams and players.

We can talk about oBABIP and WHIP together because they both yield similar results. As shown in Figure 3, having lower oBABIPs and WHIPs both lead to teams winning more games and thus increasing their chances of making the playoffs. This result is confirmed by Table 1 because it of their negative term. This means that as the pitching numbers increase, the teams will have less of a chance of making the playoffs and both these metrics are backed up by very significant p-values. These two metrics have a similar argument to that of FIP where if a team allows less runs, they’ll win more games and in turn, will make the playoffs.

Lastly, we can look at a Team’s payroll and how that impacts their playoff chances. Figure 1 shows us that along the increasing payrolls, teams that had higher payrolls appeared to make the playoffs more often than teams who didn’t. We see that in Table 1, the coefficient for payroll is 0, but upon further inspection it is actually greater than 0 but is very small. This is likely due to the fact that the lowest payroll after making the adjustments was about 780 thousand. This means that for the variable to not completely shift the model’s results, it should be very small. Regardless, since the coefficient is small, we can conclude that according to the model having a higher payroll increases a team’s chances. This makes sense because one of the only reasons why a team would spend more money would be to try and compete and make the playoffs. We also see in Table 1 that other than the variable `before_75`, payroll has the highest p-value. This is interesting because it shows that we are more confident that the advanced metrics we have been looking at impact a team’s playoff chances more than payroll does.

Our findings show that the basis of Moneyball was in fact, correct. This is displayed from the suprisingly strong accuracy displayed in Table 2, to the very significant variables from Table 1. The point is also validated by the fact that payroll’s p-value is the highest of the variables, implying that the model believes more in using advanced metrics to predict playoff probabilities as opposed to payroll.

---

<sup>8</sup>We’re going to focus on 2019 because the 2020 season was just a really small sample size

<sup>9</sup>under about 4.20

## 6 Weaknesses and Next Steps

This analysis has had some great results, but there are some weaknesses we must address which may have impacted the analysis. Firstly, it would have been more preferable to have a larger data set. As mentioned in the data section, I had to remove entries because of the missing payroll data and even before that, there was only about 2200 observations. Ideally, we would want more observations but since there are only 30 teams playing each year, we have to make it work. Another weakness I encountered was that the payroll data did not fully translate from 1920 to 2019. Even though I did adjust the payrolls to account for inflation of the dollar, there was still a large difference in spending now compared to 100 years ago. This just shows that teams decided to spend more money than they did in the past but this was accounted for in the model by adding a term to tell us if the team played before 1975 (just after the spike in pay began).

When looking at potential future steps for analyses after this one, there are many different routes we can take. One potential route we could take is looking at the statistical value of a single player's contract and also look at their statistical impact on the team. This analysis looked at the impact of a whole teams payroll and performance but going in depth on a players impact would certainly be an interesting topic to research. Another topic we could look at is using this model to predict the next season's results. Right now the model is constructed to use current year metrics to predict playoff probabilities for that same season, but if it was possible to predict expand that to the next season, it would once again be an amazing application of the model.

## References

- Us inflation calculator, Dec 2020. URL <https://www.usinflationcalculator.com/inflation/consumer-price-index-and-annual-percent-changes-from-1913-to-2008/>.
- Chris Bumbaca. Explaining the 1981 mlb season: How baseball survived shortened year, Mar 2020. URL <https://www.usatoday.com/story/sports/mlb/2020/03/15/1981-mlb-season-coronavirus-delay-baseball/5054780002/>.
- John Edwards. Stat to the future: Move over, era, it's time for fip, Sep 2017. URL <https://www.sportingnews.com/us/mlb/news/what-is-fip-era-pitching-baseball-mlb-stats-statistics-advanced-sabermetrics/d5p48p6z6us51lqbo0wvvgigto>.
- Fangraphs. Guts!: Fangraphs baseball. URL <https://www.fangraphs.com/guts.aspx?type=cn>.
- Michael Friendly, Chris Dalzell, Martin Monkman, and Dennis Murphy. *Lahman: Sean 'Lahman' Baseball Database*, 2020. URL <https://CRAN.R-project.org/package=Lahman>. R package version 8.0-0.
- Ben Lindbergh and Rob Arthur. Statheads are the best free agent bargains in baseball, Apr 2016. URL <https://fivethirtyeight.com/features/statheads-are-the-best-free-agent-bargains-in-baseball/>.
- Sports Reference LLC. *Baseball-Reference.com*. Major League Statistics and Information. <https://www.baseball-reference.com/>.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2019. URL <https://www.R-project.org/>.
- David Robinson, Alex Hayes, and Simon Couch. *broom: Convert Statistical Objects into Tidy Tibbles*, 2020. URL <https://CRAN.R-project.org/package=broom>. R package version 0.7.1.
- Steve Slowinski. Fip. URL <https://library.fangraphs.com/pitching/fip/>.
- Tom M. Tango, Mitchel G. Lichtman, and Andrew E. Dolphin. *The book: playing the percentages in baseball*. TMA Press, 2014.
- Hadley Wickham, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Golemund, Alex Hayes, Lionel Henry, Jim Hester, Max Kuhn, Thomas Lin Pedersen, Evan Miller, Stephan Milton Bache, Kirill Müller, Jeroen Ooms, David Robinson, Dana Paige Seidel, Vitalie Spinu, Kokske Takahashi, Davis Vaughan, Claus Wilke, Kara Woo, and Hiroaki Yutani. Welcome to the tidyverse. *Journal of Open Source Software*, 4(43):1686, 2019. doi: 10.21105/joss.01686.
- Claus O. Wilke. *cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'*, 2019. URL <https://CRAN.R-project.org/package=cowplot>. R package version 1.0.0.
- Mitchell T Woltring, Jim K Rost, and Colby B Jubenville. Examining perceptions of baseball's eras: A statistical comparison, Oct 2018. URL <https://thesportjournal.org/article/examining-perceptions-of-baseballs-eras/>.
- Yihui Xie, J.J. Allaire, and Garrett Golemund. *R Markdown: The Definitive Guide*. Chapman and Hall/CRC, Boca Raton, Florida, 2018. URL <https://bookdown.org/yihui/rmarkdown>. ISBN 9781138359338.