

# Logistic Regression shows that Moneyball was right

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, I will use logistic regression to determine whether the advanced stats that stemmed from this revolution are useful for predicting a team's success. According to the model (I still have to do this). The model correctly predicted that (insert) teams would make the playoffs using these metrics.

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

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 stats that mattered.

In order to test if Moneyball's basis was correct, I will be using logistic regression to predict whether or not a team will make the playoffs. The seasons that I will analyze go all the way back to 1920, with some missing years in the 70s because of missing values in the contract data. In order to predict if a team makes the playoffs or not, the variables I plan to use are the team's payroll, the team's batting average (BA), the team's earned run average (ERA), their weighted on base average (wOBA) and lastly their fielding independent pitching (FIP). These stats may sound scary but they are fairly simple and will be explained in the data section below.

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 stats like payroll data. 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>2</sup>. The payroll data will be crucial for my report because the basis of Moneyball is fielding the best team by using the least amount of money.

This data analysis was conducted using R (R Core Team [2019]), and in particular the packages Tidyverse (Wickham et al. [2019]) and was compiled using R markdown (Xie et al. [2018]). I have 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 a 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

---

<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

<sup>2</sup>This scraper can be found in the scripts folder of my GitHub repo.

the model to the results of the 2020 season 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? Does money impact a team’s ability to make the playoffs? Keep reading to find out!

## 2 Data

The data collected 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.

The **Lahman** package contains the **Teams** dataset which was used for the calculation of the basic statistics along with the more advanced ones. The basic statistics used for this analysis were Batting Average, Earned Run Average, Weighted On Base Average and Fielding Independent Pitching. Batting average is the rate at which a player will get a hit, it is calculated by  $Hits/Atbats$ . Earned Run Average is the amount of runs a player allows, per 9 innings (1 whole game), it is calculated by  $Runs/InningsPitched$ . Next, we get to the advanced statistics, 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>3</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>4</sup> Lastly, 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.

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.

Figure 2 shows the relationship between wOBA and the amount of teams a win gets, along with colours representing a team’s ability to make the playoffs. We see that as wOBA increases, the amount of wins increases and the number of teams making the playoffs increases as well. For the most part, it seems that if a team can get to about 87 wins in a season, they will likely make the playoffs. There were also some teams present which made the playoffs after winning less than 65 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>5</sup>.

Figure 3 shows that having a lower FIP means that a team is going to win more games. This makes sense because if a team gives up less runs, they are more likely to outscore their opponents, leading to a win.

## 3 Model

Table 1 shows us the coefficients from our logistic regression model. We see that teams with a higher wOBA have better chances of making the playoffs, along with team with higher payrolls. This makes sense because when teams will attempt to spend the most money to get the best players and when they score more runs they’ll win more games. We can also see that the coefficient for FIP is negative, this does not mean that FIP brings down the chances of making the playoffs. The best values for FIP are the ones closer to 0 and as the value begins to increase, it shows that the pitcher is not that good. What the coefficient is telling us is that

<sup>3</sup>The Book is seen as the best book for aspiring baseball analysts or “saberists” to get into the MLB analysis game.

<sup>4</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).

<sup>5</sup>The players eventually returned and the playoffs continued

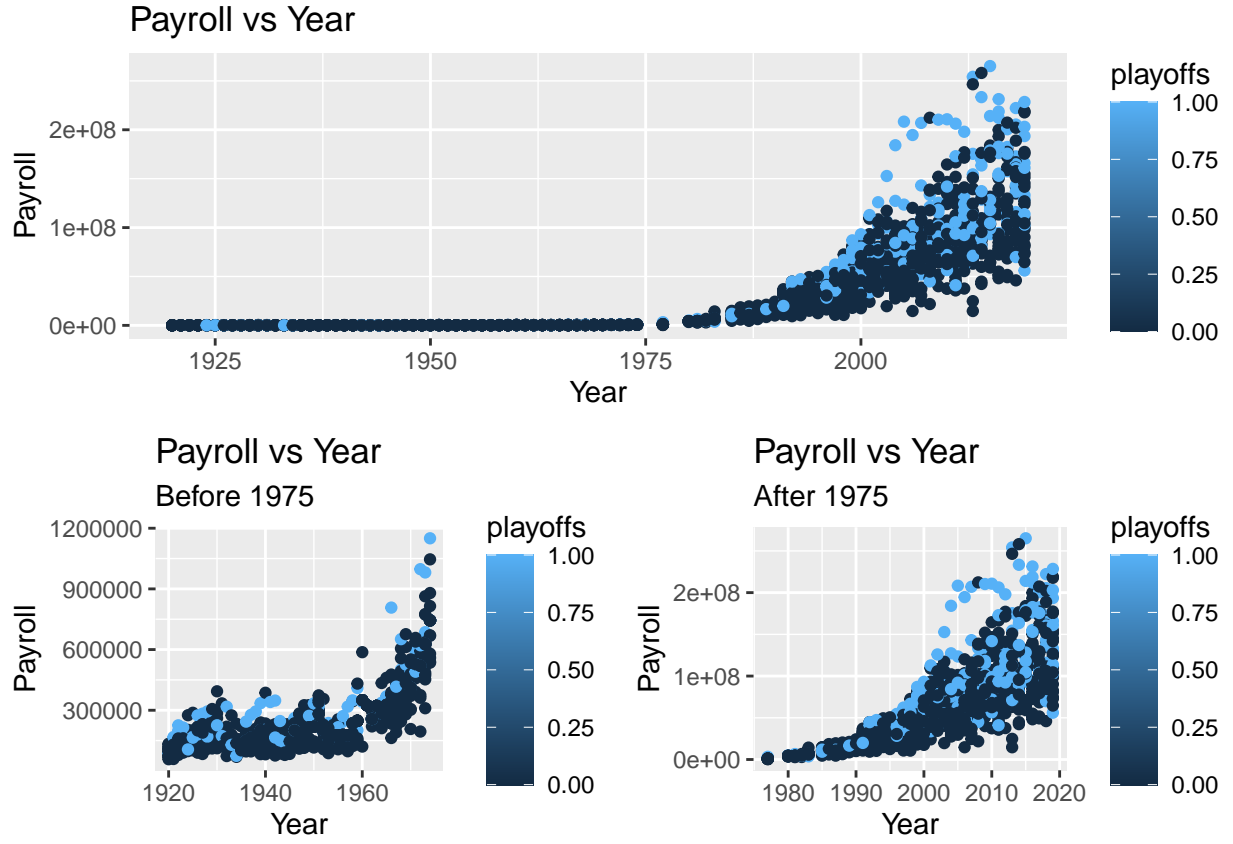


Figure 1: The increase in Payroll over the Years

Table 1: Regression Model Coefficients

term	estimate	std.error	statistic	p.value
(Intercept)	-21.4913	1.8914	-11.3628	0.0000
wOBA	140.3117	9.8677	14.2193	0.0000
FIP	-0.3222	0.3760	-0.8570	0.3914
payroll	0.0000	0.0000	4.5658	0.0000
BA	-39.4284	10.2847	-3.8337	0.0001
ERA	-4.7689	0.3403	-14.0123	0.0000
before_75	-1.2911	0.2352	-5.4902	0.0000

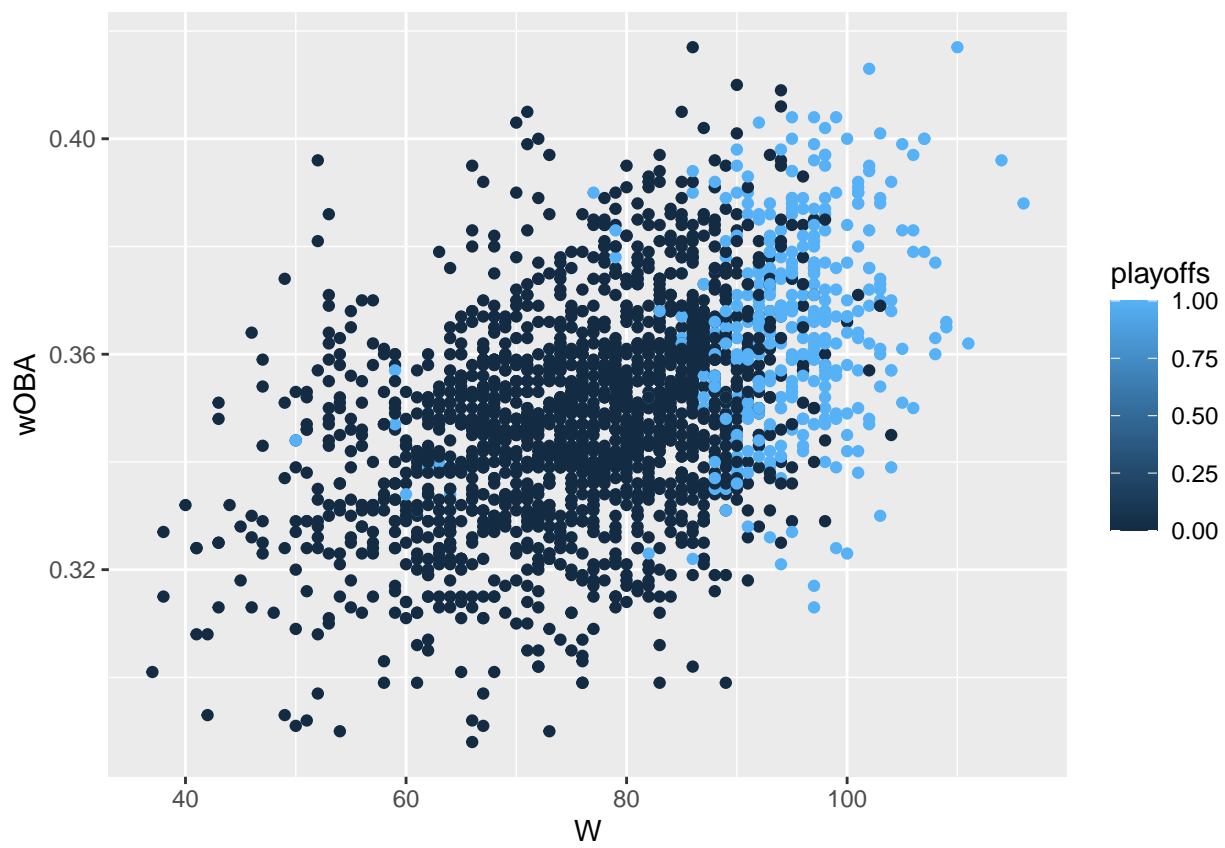


Figure 2: Relationship between wOBA and Wins

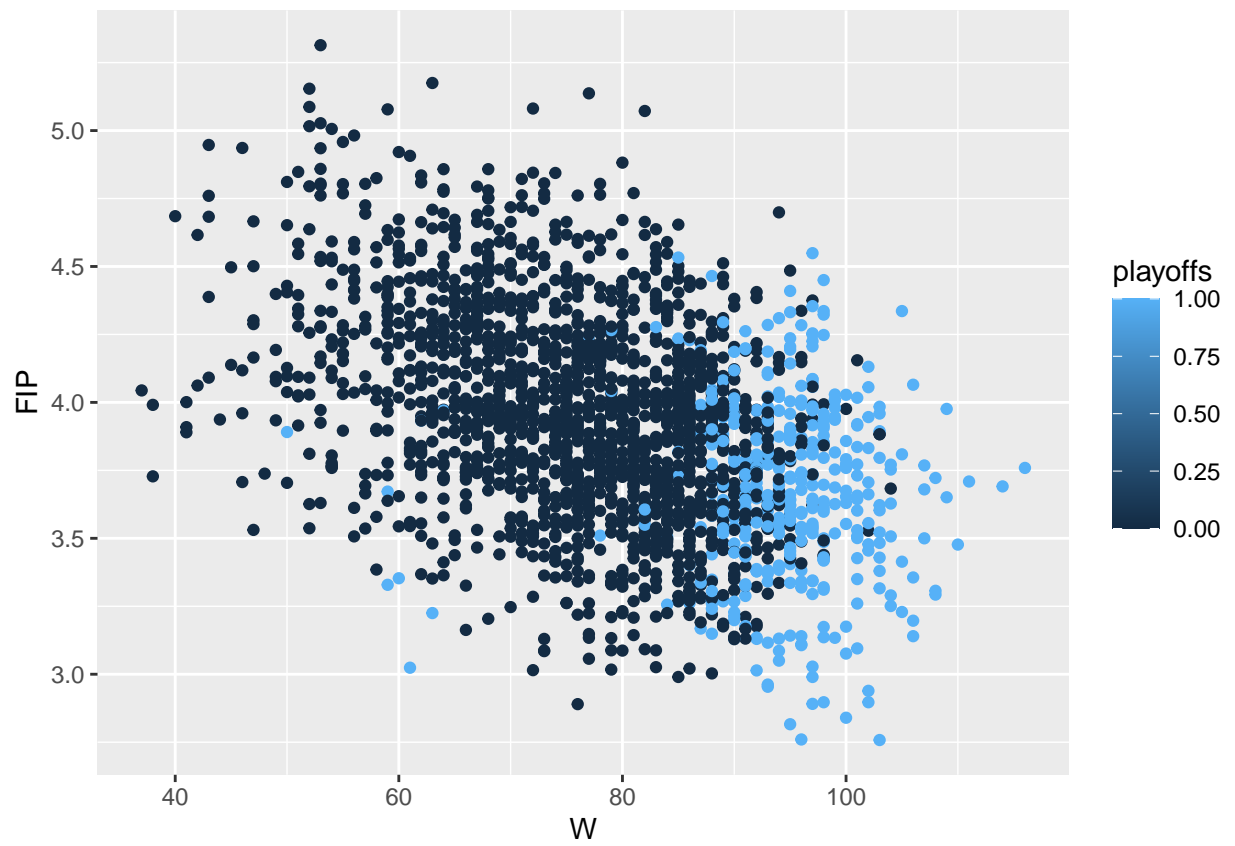


Figure 3: Relationship between FIP and Wins

lower FIP's give higher chances of making the playoffs while higher FIP's bring down the chances of making the playoffs.

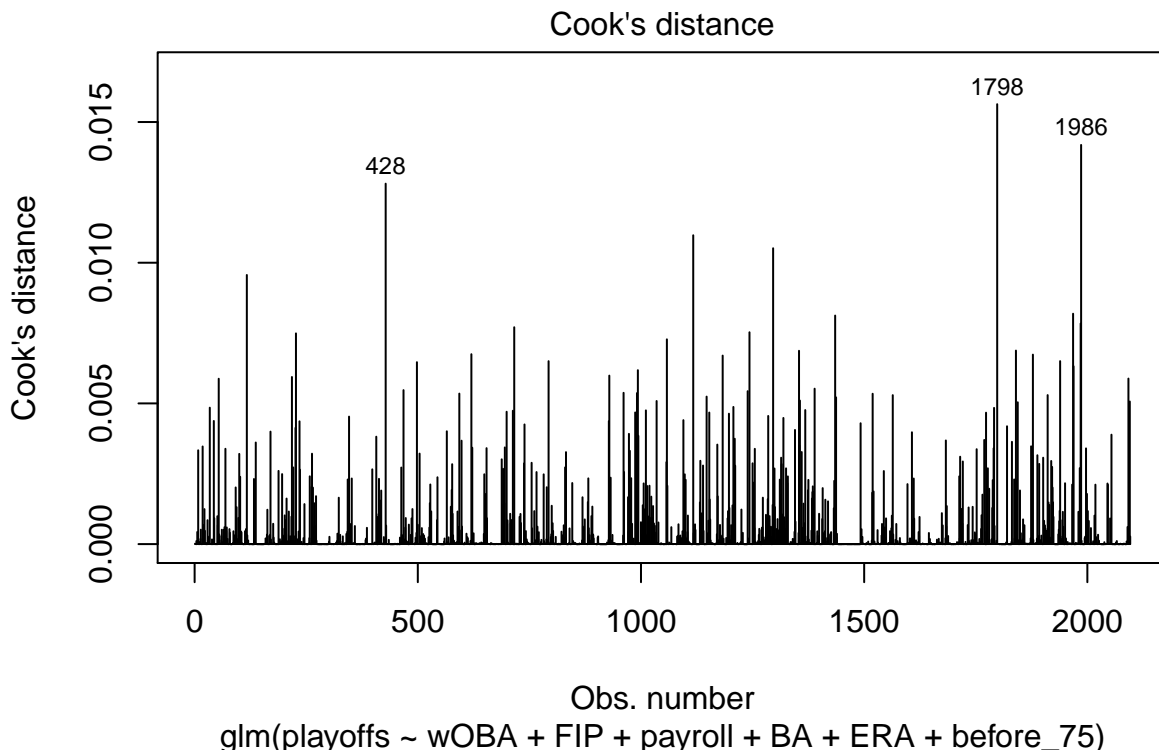


Figure 4: Cook's Distance for the Model

Figure 4 shows us the observed cook's distances for the model. Cook's distance gives us the difference between the ... . We see that these cook's 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 fairly well.

## 4 Results

Data Type	Accuracy (%)	result
Data	-	18.989
Model	88.8358778625954	15.649

Table ?? 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 15.6 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 88.8 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 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>6</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.

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>7</sup> and below average traditional 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 ERA. 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 outperforms ERA and this shows...

## 6 Weaknesses and Next Steps

---

<sup>6</sup>We’re going to focus on 2019 because the 2020 season was just a really small sample size

<sup>7</sup>under about 4.2

## References

- 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.
- 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/>.
- 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 Grolmund, 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, Kohske 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.
- Yihui Xie, J.J. Allaire, and Garrett Grolmund. *R Markdown: The Definitive Guide*. Chapman and Hall/CRC, Boca Raton, Florida, 2018. URL <https://bookdown.org/yihui/rmarkdown>. ISBN 9781138359338.