

# Predicting Baseball Salaries

Greg Eastman

4/23/2021

Baseball franchises are both sports teams and businesses, as such they are trying to get the most value for their players. There is a lot that goes into player value, and since we do not have access to confidential sales data instead, we focus just on athlete performance. With data from the Sean Lahman Baseball Database<sup>1</sup>, we build 4 statistical models, a random forest, Ada boosted tree, gradient boosted tree, and a linear regression. The best of these models is then used for two tasks. First, it finds the most important factors in predicting player salaries. Second, it finds the players that most outperformed their salaries from 1996 to 2019. We then synthesize this information to gain insights into what aspects of a player's game can be leveraged to best increase their pay.

# Introduction

Baseball and statistics have a rich history, and since the Moneyball era started in the early 2000's that relationship has only deepened. Now analysts are a core part of any franchises with data driven insights factoring into player development, recruitment, sales, and so much more. Although there is no easy metric to evaluate athlete value, but since a franchise is a type of business, we can apply some economic logic to try and get some information.

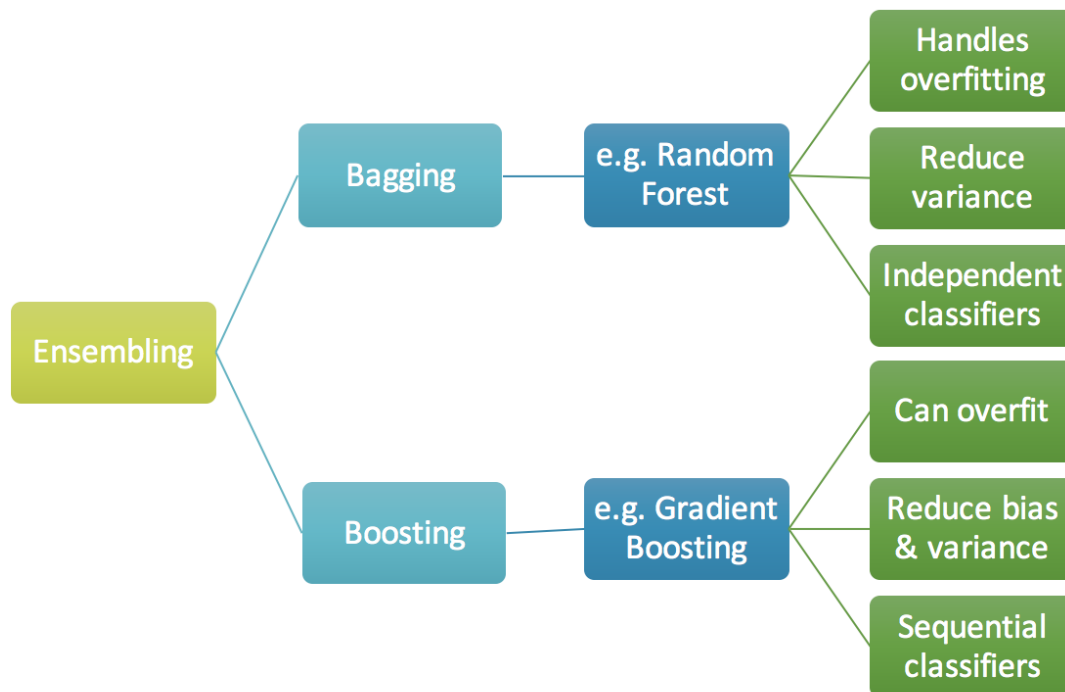
The MLB is an industry like any other, only their export is specifically baseball media, so we can assume then that they will financially follow the same patterns as any other corporation. This means that they will overall not pay a player more than that athlete bring them in revenue. At the very least, if an assumption of perfect competition is true, they will pay the exact amount that athlete brings them in value. This all happens because players can negotiate contracts, so we can assume that if a player could get a better offer, they would take it. This should mean that overall, the MLB pays players their worth. This is a rather simple assumption, which I will relax.

There are two major factors that create a less than perfectly competitive market, new franchises cannot just enter the market, and rookie contracts disallow negotiation for 3-4 years [add citation]. These factors imply that players have weaker bargaining power than the franchises. While the corporations compete for players, there are many players, but only a limited number of franchises. So, a franchise can choose to pass up one player for another, but a player may not have that option. And rookies will be forced into contracts that may severely undervalue them, because they cannot renegotiate when they perform well. Together these imply that salaries may underpay but should not overpay players. Mathematically we can denote this relationship as ***Player Salary***  $\leq$  ***Player Value***. This then begs the question, who were the most underpaid baseball players? We will compare multiple methodologies to get the best estimate of a player's salary and use then use these to estimate to find the 5 most underpaid players.

# Methodology and Models

This work used four models and compared them against each other so that the best could be used to find the most underpaid players and analyze the most essential statistics for a player's pay. We will go through the way that each of the models works. This next section assumes that the reader has familiarity with decision trees. For more resources outside of this work this [link](#)<sup>2</sup> may be useful. Additionally, we will not walk through the math as doing so for each model would be superfluous to the purpose of this work. Before going into a random forest, the diagram below will provide some important information.

3



This diagram has a lot going, but we will break the essential elements. The overall methods this work uses is called Ensembling. The basic idea is ideologically aligned with the saying “many heads are better than one”. These techniques work by creating many weak decision trees and conglomerating their predictions. This helps reduce variance because having more models making guesses allows any individual overcompensation to be corrected. Although the two methods do this differently. The biggest difference between bagging and boosting is that in bagging the models are grown independent from each other. In boosting the models are built sequentially with one picking up where the previous left off. Both approaches have benefits and deficits we will discuss further.

## Random Forest

A random forest is a bagging technique that grows decision trees from different pieces of data to avoid over fitting. The idea of bagging has been discussed, and it is a powerful predictive technique with proven performance<sup>4</sup>. It works by first randomly sampling the training data and building a model on each of the created subsets. This helps prevent it from overtraining by running a different one on the same data. The next thing it does is that when a tree considers splitting a node, it randomly subsets a feature. This random element helps it to try new splits to check performance and avoid all the trees from trying to converge on just one feature. Since each tree is made independently from each other on different subsets of data, and with different feature splits the model is robust.

The benefits of this technique are two-fold, as shown in the diagram. First, its classifiers are independent, so the models do not affect one another, which helps it to be more robust. Leading into number two, handling over fitting. By grouping a set of different models, made on different data, the model can avoid becoming too dependent on trends in the training data. The models have all generated with random splits so when it groups them together, even if one model over fit in a certain way, it is compensated for by the large number of other independently grown trees. This allows researchers to feel more comfortable when training a model that they do not have to worry a large performance loss on the testing data. Next, we will go into boosting models.

## Ada Boost

Ada boost is the original boosting technique<sup>5</sup>, and it works by sequentially growing trees one on top of another. This is a sequential technique where each tree works off how the previous missed. They all link together to form the final full model. To do this we make two sets of weights, the first is for the data and the others are for the trees. The weighting of the data is the key to this process. After each link in the chain of trees the data is given a new distribution. Values that were misclassified in previous steps get a higher weight. This makes them more likely to appear later and to become more influential in the loss minimization of future trees. The weights for each tree are calculated with a function that creates a weighted average. The more useful trees get a higher weight, and the total sum of the weights is always equal to one. This approach has several advantages and disadvantages.

The advantage of this technique is that it generally converges faster than a forest and that it reduces bias, but it is inflexible and prone to overfitting. The sequential nature of the model means that there is less randomization that needs to be done. The guesses are more intentional, this means that the model should get to the desired value quicker than a random forest will.

Additionally, by constantly improving on missed guesses instead of just the overall error, the model helps to reduce bias in the error term. This makes its importance weights very useful because it is more robust to incorrectly signing and effect. Although, the models have an issue with overfitting. Since they are all built on the same data, and more so they are built on each other, they can quickly start learning trends that are specific only to the sample and not the population. This is usually addressed by tuning the learning rate and number of samples. Although, the model still struggles with the fact that it has a specific loss function that it can use. This means that it may struggle with modeling different types of data because it is more rigid in its construction. The next technique addresses this issue.

## Gradient Boost

Gradient boosting is like ada boost, the only major difference is that instead of using the ada boost function to minimize error it uses gradient descent. Through this method the weights on the data and the trees are trained. After each iteration, the model will reduce the loss by walking down its gradient towards the minima. This method is more flexible for several reasons. First, in ada boost the weights are additively bounded, which is not the case here, as the weights are chosen iteratively without a built-in constraint. Additionally, you can perform this methodology over any loss function instead of just the original one. This helps to predict different types of data, and update with newer techniques. The last model is one of the most classic statistical models.

## Linear Regression

Linear regression is a staple of the statistical world. It works finding the line through the data that minimizes the sum of squared errors. The math for this can be found [here](#)<sup>6</sup>, and is one of the most classic statistical minimization problems. The model has a strict set of assumptions that are usually not fully met, but its easy interpretability makes it a powerful tool. This latter fact has made it an essential tool for causal research in disciplines ranging from chemistry to biology to sociology to economics. We will not go into how to interpret this model because this work is purely predictive. For the sake of completeness, we will discuss the assumptions.

Linear regression has four major assumptions, linear relationship, no multicollinearity, homoskedasticity, and normality. The linear relationship requirement is often the laxest of all the requirements. It is often met by forming scatter plots and transforming the data so that the relationship becomes linear. The tradeoff is complex transformations to get linearity may affect

interpretability. The next assumption is that the model must not be multicollinear. In other words, no two predictors should be highly correlated. This is usually tested for with variance inflation factors, a VIF above 3 is usually cause for concern, above 5 is a problem, and above 10 is a total violation. The next assumption is homoskedasticity, which is tested by looking at a residual vs. fit plot. If the data has a pattern, such as a cone around the  $y=0$  line, then there is an issue, if it is randomly distributed then the assumption is met. Finally, the normality assumption is tested with a QQ-plot. The pattern of data should follow the normal print on the plot, if it has a large deviation from the line its assumption is not met. With all the models covered, we will move onto the experiment.

# Experiment

We took a 3-step approach to answer the research questions by finding the best model, looking at the important variables, and then finding the most profitable players. Although, before going into the specifics of the experiment we will go over the data.

## Data

This work uses data from the Sean Lahman database. This source is from a sports reporter who has been documenting baseball statistics since 1996. We used only data from that year and beyond because the earlier years are filled with missing salary and other values. Additionally, for more robust analysis we got inflation data from the Federal Reserve<sup>7</sup> and used that to adjust the salary data to the monetary value of 2020. The salaries were then logged to account for the natural skew in the data. Next, we removed all players that had less than 50 games and all pitchers. We wanted to keep only the starters and not substitutes, and pitchers are not paid for fielding and batting ability. The final data cleaning step was to bin all continuous variables so that the tree can learn on them properly. With the data covered we will go over the results of the models.

## Model Results

The four models were measured on mean squared error and  $R^2$ , which in the end both told the same story, the gradient boost was the best model.

Model	rMSE	$R^2$
Random Forest	.81	.64
Ada Boost	.85	.6
Gradient Boost	.79	.65
Linear Regression	.81	.63

We will break down the table one model at a time by comparing them to the basic multiple linear regression, beginning with the random forest. The forest outperformed the control model of a linear regression by a noticeable margin of half a percentage point in both rMSE and  $R^2$ . This improvement over a regression model shows that a random forest is better at handling predictions in a space that may not be linear. Additionally, this was the second-best model overall, being second only to the gradient boost.

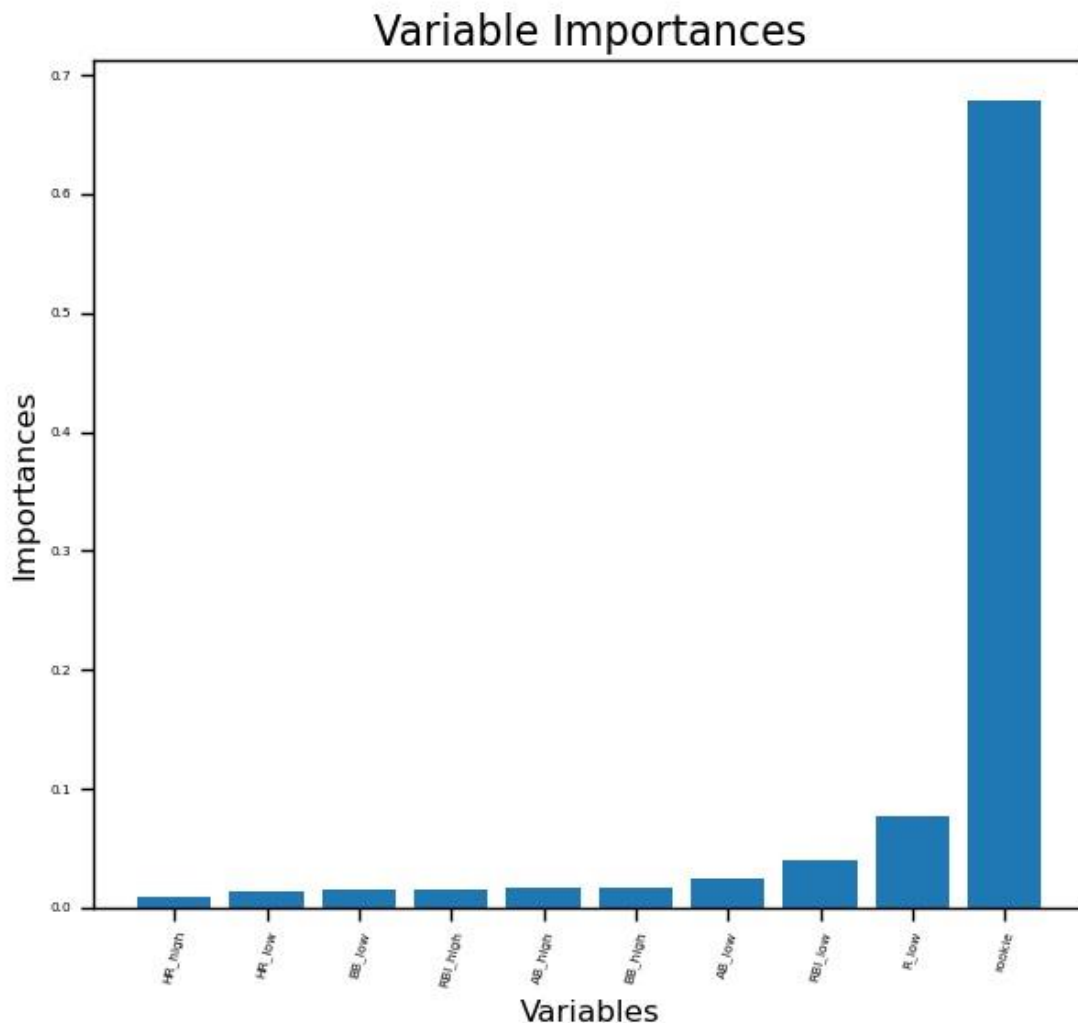
The gradient model had by far the best performance by explaining over 65% of the variation and having the lowest rMSE. The variance explanation, or  $R^2$  is a common metric for model performance as explaining more of the variation in the data is good. Although, it is a measure relative to a null model, and not in absolute terms. This is where rMSE comes in, as it measures the standard deviation of the error term in the units of log salary. So having an rMSE of less .8 is strong and seeing as it outperformed the regression by over 2 percentage points in both categories, it will be our best predictive model to use. Although, there is one other tree-based method to still discuss.

The ada boost model had the worst performance overall. This model did worse than the linear regression. This is surprising because it is more complex, and since it is not making a simple linearity assumption should better predict the data. Although, we see in the table that this is not the case because the regression beat it by 3 percentage points in both rMSE and  $R^2$ . This could be just that the salaries are complex, and the simple function to refine the tree is just not quite up to that task. Meanwhile a simpler linear model may be more successful due to the large number of inputs and data causing it to split the difference better but not guess on specific points as well. Since the regression was just a baseline model, we will not go over it much and will instead look at what 10 variables were the most important in predicting a player's salary.

## Importance Weights

The importance weights, or the weights given to how useful each variable is in predicting the log of salary, were those generated by the gradient boost model. This model was chosen due to its superior performance. The chart on the next page will help to draw some meaning on what variables are the best determinants of salary.





The first thing to note is that the rookie variable is by far the most important as it holds more weight than the next 9 combined. This makes sense as rookies are signed with contracts before they get into the pros for well under the million-dollar contracts they may grow into. Since these players cannot renegotiate or go to another team there is no bargaining power. This goes back to the idea that the league holds more power in the negotiations than the players. Regardless, it is not surprising that being a rookie is the largest determiner of salary, but the magnitude is still notable. The next finding involves multiple stats.

The most important determiners of player worth are runs, at bats, walks, RBI's, and homeruns. The first stat to talk about is runs, which scoring a lot does not seem important, but players that do not score many runs seems to matter. This likely means that teams do not pay well for those who do not deliver. This is shown by the "R\_low" bar. Next, having low RBIs seems

to do the same, but unlike runs, having a lot of RBI's does help to increase pay. Then at bats seems to determine salary. Although this is an important variable to have, its importance presents an endogeneity scenario because it is highly likely that teams want to play those that they pay more and not just the other way around. Then we are left with walks and homeruns. Both are extremely important and having lots of homers and base on balls drastically changes pay. With this all covered, we will move onto what is not in the chart.

There is a surprising non-finding to report, that the player's team was not a significant determinant of pay. It is shocking that teams are not a big predictor of pay, because baseball is well known for having some extremely wealthy and "poor" teams. The Yankees have historically had a mountain of money that they spend on their roster, and the Oakland A's have had the opposite. Yet neither are meaningful determinants of pay. This implies that the league is at a strong equilibrium in player evaluation, and that salary caps may be working to keep things more even across franchises. With all this information covered we will take a quick look at the most underpaid players since 1996.

## The Five Most Underpaid Players

The five most underpaid players were all undervalued by roughly the same amount, these athletes in order of profit for their franchise were, Jorge Pasada (2000), Tino Martinez (1996), Scott Brosius (1998), Bernie Williams (1996), Ryan Ludwig (2008). All these players were all stars in these years, and most of them for the first time. Additionally, three of them (Martinez, Williams, and Ludwig) were MVPs in those years and had not been previously. This implies that they all performed unexpectedly well in those years, and when looking at the stats they were all the best years that they had played in those times. Furthermore, it is interesting to note that even though the data goes to 2019 all the misvaluations were around the early 2000s with only one in the latter half of the decade. The early 2000s is significant because that is when Moneyball started its rise and statistics began choosing players instead of just scouts. Since the contracts have gotten closer to the value a player brings as time has passed, we can see an implication that statistics has made players and teams into more informed and accurate actors when signing contracts.

# Conclusion

Baseball is one of the most statistically driven sports, and there is much more to learn about the sport with each study. In this work we looked at several different modeling strategies to predict the value of a player. We found that gradient boosting performed the best. Then using this finding, we investigated the most important determinants of a player's success. We found that being a rookie has a massive impact, and that underperforming is more impactful than an overachieving player. The next thing of note was that teams did not play a huge role in a player's pay. Finally, we found that as time went on player pay got closer to the value their performance brought their teams each year.

These findings could be important for players trying to take that next step in their baseball careers. Teams seem to pay most for players that hit on par with other pros rather than for defensive excellence, or other offensive traits. The teams should not be a factor in determining where a player wants to go, because financially they are not impactfully different for each player. Outside of this advice the paper just found that statistics appears to have made baseball more economically efficient for teams and players by lowering the number of outliers in players being underpaid.

Going forward other work could build on this by taking a different modeling approach. Networks may be able to lend a performance boost, or clustering methods could lend insights on what traits the top paid players exhibit. Also, Bayesian statistics may be able to utilize priors to better leverage the information in this data. Finally, combining this data with sales and other metrics for player valuation could help to increase model accuracy and refine the findings of this paper.

## References

- 
- <sup>1</sup> "Download Lahman's Baseball Database." SeanLahman.com. February 16, 2021. Accessed April 25, 2021. <http://www.seanlahman.com/baseball-archive/statistics/>.
- <sup>2</sup> "Decision Tree Algorithm, Explained." KDnuggets. Accessed April 25, 2021. <https://www.kdnuggets.com/2020/01/decision-tree-algorithm-explained.html>.
- <sup>3</sup> Grover, Prince. "Gradient Boosting from Scratch." Medium. August 01, 2019. Accessed April 25, 2021. <https://medium.com/mlreview/gradient-boosting-from-scratch-1e317ae4587d>.
- <sup>4</sup> Fawagreh, Khaled, Mohamed Medhat Gaber, and Eyad Elyan. "Random forests: from early developments to recent advancements." *Systems Science & Control Engineering: An Open Access Journal* 2, no. 1 (2014): 602-609.
- <sup>5</sup> Freund, Yoav, and Robert E. Schapire. "A decision-theoretic generalization of on-line learning and an application to boosting." *Journal of computer and system sciences* 55, no. 1 (1997): 119-139.
- <sup>6</sup> "5.4 - A Matrix Formulation of the Multiple Regression Model." 5.4 - A Matrix Formulation of the Multiple Regression Model | STAT 462. Accessed April 25, 2021. <https://online.stat.psu.edu/stat462/node/132/>.
- <sup>7</sup> "Inflation, Consumer Prices for the United States." FRED. March 03, 2020. Accessed April 25, 2021. <https://fred.stlouisfed.org/series/FPCPITOTLZGUSA>.