
National Basketball Association salary prediction: a data-driven Linear Regression analysis

Jiahe Zhang
freddie_1534@163.com

Abstract.

National Basketball Association (also known as NBA which will be used as an abbreviation throughout this report) as one of the few most successful professional sports leagues in the world, it is well known by the fierce physical competition among the most talented and competitive players in the realm of basketball. However, there is something else to discover behind the dazzling crossovers and sensational clutches, hidden in the reflection on the O'Brien Cup.

This report will be discussing the prediction on NBA players' salary using multiple supervised machine learning algorithms based on a dataset of players' on-court data and achievements, aiming for an objective result that can be used for both prediction and evaluation of players' contracts. The result of this report reveals which specific variables are useful for predicting players' salary using Multiple Linear Regression

1. Introduction

In NBA, the best way to make profit has been to win the championship. Naturally, by recruiting or trading valuable players to build a championship-winning team is the goal for every club. While, to make the competition fair, there is a limit on the amount of money that a team can spend on players' salaries, which is so called 'salary cap', if such limit was reached, luxury tax is required to pay towards the association. Therefore, making the budget manageable and spending the money efficiently is crucial. This is where a prediction algorithm comes in. Not only to manage the money flow, but also to evaluate the values of current players' contracts in the team so that the decision-makers can construct the team better.

The structure of this report is as follows: Literature review, Methodology, Conclusion, References

2. Literature review

Data and statistics have been heavily influencing every sport in this era, which can objectively enhance decision makers' work for a grand amount. Sports analytics is an emerging field that grabs data to optimize the decision-making process. With it, teams can have better on-court winning strategies, off-court training plans and other approaches to improve athletes' performance. By the year of 2028, it is expected that the sports analytics industry will profit \$3.4 billion globally [1]. There have already been some great papers researching on the factors influencing NBA players' salary. A paper in 2018 claimed that the determining factors of NBA players' pay are experience, points, rebounds, assists and fouls, whereas 3-point shots made and Hollinger's player efficiency rating (PER) are insignificant [2]. Also, [3] aims to explore the best player selection strategies by studying the player statistics, team performance, and the salary cap. It explains basic NBA statistics concepts and how the players' efficiency is measured and their relationship with the team performance. However, there are only 450 players in NBA in 2021, it is relatively low (compared with 1,696 players in The National Football League, 780 players in Major League Baseball), it is difficult to build a good model. Besides, data from video games of basketball (2K series) is as useful as that from real basketball world. As confirmed in [4], the variables in the video game 2K20 profoundly contributed to predict NBA players in the real world. As a matter of fact, the predictions were very close to the salaries in season 2021-2022.

3. Methodology

48

49
50
51
52

53
54
55
56
57

58
59
60
61
62

63

65

66
67
68
69



71

72

There were 2 regression models used in this report: Simple Linear Regression and Multiple Linear Regression

3.2.1 Simple Linear Regression

Simple Linear Regression is a Linear Regression model with only one independent variable as the predictor to predict the dependent variable which always is Average Salary in this report. The first model starts with Average Salary versus Total Score:

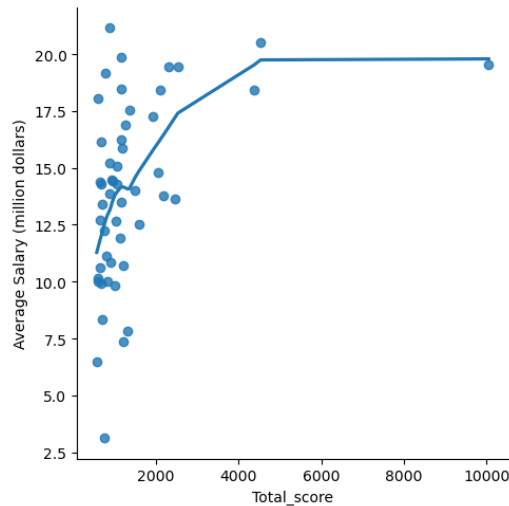


Figure 3: Simple Linear Regression results

Figure 3 contains the scattered data and the curve describing it, which indicates that there is not a strong linear relationship between Average Salary and Total Score.

Size of train set = **40** (80% of dataset)
Size of test set = **10** (20% of dataset)
Mean Squared Error \approx **11.97**

Size of train set = **45** (90% of dataset)
Size of test set = **5** (10% of dataset)
Mean Squared Error \approx **19.62**

Figure 4: two results from Simple Linear Regression

Figure 4, they are two results of different ways splitting the dataset. They have Mean Squared Error (an indicator on the deviation of data, as model error increases, MSE increases) of 11.97 and 19.62, respectively. However, both results are not satisfying since the dependent variable, Total Score, is the sum of all the 17 variables from [6], and some of the variables help to predict Average Salary, but some do not. Hence, Simple Linear Regression does not suffice to be a good model in this case.

3.2.2 Multiple Linear Regression with Backward Selection

Multiple Linear Regression is a prediction model that takes multiple variables as predictors to regress with the dependent variable, which shows the content each predictor helps explain the dependent variable. Figure 5 was attained after transforming and fitting Multiple Linear Regression of all 17 variables against Average Salary:

	coef	std err	t	P> t
intercept	-13.51880	10.04800	-1.34500	0.18800
Accom_awards	0.00810	0.00900	0.91200	0.36900
Accom_teams	0.01290	0.01500	0.88500	0.38300
Accom_ranks	0.01050	0.01300	0.81800	0.41900
Accom_mvp_shares	-0.05950	0.03300	-1.78600	0.08400
Accom_champ_factors	13.38820	9.13500	1.46600	0.15300
Accom_playoffs	-0.09750	0.05400	-1.81900	0.07800
Commit_GMSC_regular	-0.11150	0.08900	-1.25600	0.21800
Commit_GMSC_playoff	0.38750	0.13700	2.82100	0.00800
Commit_WS_regular	0.00880	0.10400	0.08400	0.93300
Commit_WS_playoff	-0.31850	0.13200	-2.40600	0.02200
Prime_GMSC_regular	0.06820	0.03900	1.76100	0.08800
Prime_GMSC_playoff	0.03180	0.03500	0.89600	0.37700
Prime_WS_regular	0.02700	0.05400	0.50400	0.61800
Prime_WS_playoff	0.03900	0.04900	0.79600	0.43200
Legacy_regular	0.03600	0.07600	0.47400	0.63900
Legacy_playoff	-0.03470	0.05700	-0.61000	0.54600
Legacy_final	0.16320	0.14700	1.11000	0.27500

Figure 5: Multiple Linear Regression predictors information

In Figure 5, $\text{coef}(\beta)$ refers to the amount Average Salary changes for one unit increase in predictor X_i ; std err refers to Standard Error; t-statistic refers to indicates the linearity of each predictor with Average Salary; p-value indicates the correlation of each predictor with Average Salary. With Figure 5 and the formula, $\text{Average Salary} = \beta_0 + \beta_1 * \text{'Accom_awards'} + \beta_2 * \text{'Accom_teams'} + \beta_3 * \dots + \beta_4 * \text{'Legacy_final'} + \epsilon$ (random error), a player's predicted salary can be calculated if the values of all the variables and coefficients are known. Nevertheless, as mentioned before, some predictors are not helpful as much as others. For instance, 'Commit_WS_regular' has high p-value of 0.93, which indicates that it does not have a linear relationship with the dependent variable. Variables like this need to be dropped one by one and fitting the model again so that the model can be optimised, such method is so called Backward Selection. Average p-value was 0.3238 with all 17 variables used. After 9 times of Backward Selection, the variables left are (5)championship factors, (6)playoffs, (7)Game Score in regular seasons, (8)Game Score in playoffs, (10)Win Share in playoffs, (11)Game Score in regular seasons at prime, (14)Win Share in playoffs, (17)all-time leading records in finals. The lowest average p-value was obtained of value approximately 0.0747, which decreased 76.93% from 0.3238. While the Mean Squared Error was 7.9, which declined 34% from 11.97 and 59.7% from 19.62.

4. Conclusion

The performance of Multiple Linear Regression with Backward Selection was much better than that of Simple Linear Regression for the reason that Backward Selection eliminated the variables that were valuable for predicting salaries. Multiple Linear Regression had a result of average p-value 0.0747 and Mean Squared Error 7.9. And the most useful 8 variables, among the original 17 variables, predicting salary were: 1. championship factors, 2. playoffs, 3. Game Score in regular seasons, 4. Game Score in playoffs, 5. Win Share in playoffs, 6. Game Score in regular seasons at prime, 7. Win Share in playoffs, 8. all-time leading records in finals.

However, there are some other factors affecting the accuracy of the models. 1. Rookie player: Some of the players in the dataset had very short experience in NBA, they surely could not have signed big contracts. 2. Voluntary Pay cut: There were few some players who volunteer to have pay cut on themselves so that the team can sign new players. 3. Subjective factor: Decision makers were not always right. It has been highly common that players signing contracts with astronomical figure but have ordinary performance on court.

144 **References:**

- 145 1. Alyssa Schroer, 'How Sports Analytics Are Used Today, by Teams and Fans', 2018
146 (updated with additional reporting by Brian Nordli in 2021.)
- 147 2. Kevin Sigler, William Compton "NBA Players' Pay and Performance: What Counts?",
148 United States Sports Academy, 2018
- 149 3. R. Nagarajan and L. Li, "Optimizing NBA Player Selection Strategies Based on Salary
150 and Statistics Analysis," 2017 IEEE 15th Intl Conf on Dependable, Autonomic and
151 Secure Computing, 15th Intl Conf on Pervasive Intelligence and Computing, 3rd Intl
152 Conf on Big Data Intelligence and Computing and Cyber Science and Technology
153 Congress(DASC/PiCom/DataCom/CyberSciTech), Orlando, FL, USA, 2017, pp. 1076-
154 1083, doi: 10.1109/DASC-PiCom-DataCom-CyberSciTec.2017.175.
- 155 4. Intelligent and Fuzzy Techniques for Emerging Conditions and Digital Transformation:
156 Proceedings of the INFUS 2021 Conference, Held August 24-26, 2021. Volume
157 2. Switzerland: Springer International Publishing, 2021.
- 158 5. <https://www.basketball-reference.com/players/>
- 159 6. 博闻爱打球, 'NBA 史上球员排行', <https://bbs.hupu.com/39307300.html>, 2020