

DSC 630: Predictive Analytics  
Milestone 5 - Final Project Paper

08/05/2020

By: Thanh Nguyen-Duong

Tai Ngo

## EXECUTIVE SUMMARY

The National Basketball Association (NBA), which is an American men's professional basketball league, is composed of 30 teams. It is one of the four major sports league in the United States, and it is considered to be the premier men's basketball league worldwide. Each NBA team is given a budget worth millions of dollars to build a talented team that aims to win the NBA championship. However, there are disparities in the budgets among different teams, and the teams that spend the most do not always end up winning the championship. Many NBA teams decide to form a strategy to spend their budget optimally with the goal of getting the best possible rosters with their given budget. Teams need to figure out a way to operate efficiently and recruit intelligently within their budget constraints. The problems that most teams have to encounter are the way that the player salaries are determined. Keeping underperformed, overpaid players is a detrimental, costly mistake for any teams. This project is created with the purpose of using existing data to predict whether players are overpaid or underpaid. This project will be considered a success if it can predict player salaries based on their performance accurately. The statistics for all players will be analyzed, the predictive models will be built to predict the players' salary based on their statistic indicators. At the end of the project, a graph of overpaid and underpaid players will be shown, the users can see how a player is paid according to their performance. If NBA teams decide to use predictive analytics as a primary tool to determine player salaries, their results within the leagues will be improved, and their budgets will be managed more efficiently.

## INTRODUCTION

Many NBA teams are interested in using their player performance indicators and the predictive analytics to predict whether a player is overpaid or underpaid. In order to optimally utilize their budgets, operate and recruit players efficiently, teams want to possess as many underpaid players as possible. A lot of players, who are exceptionally good, may be underpaid due to them being new, being in a mediocre team or being underachieved. Predictive analytics is a powerful tool in aiding the NBA teams for their player assessment and recruitment process. Having players who are overpaid and underperform is counter-productive and hurts the chances of winning a NBA championship. Teams that focus on recruiting players who are underpaid will result in better budget management and improve the team's ranking. On top of that, predictive analytics can bring many more benefits aside from monetary purposes and achievements such as better understanding of players' performance indicator, insights into training program to help existing players, and an overview of how player salaries are mismanaged in general.

In this project, the dataset is NBA player statistics dataset, which is obtained from Kaggle. It contains information about players such as their game plays, points and various performance indicators. The player statistics originating from 2017-2018 team rosters, but the rosters are more aligned with 2016-2017 data. In this dataset, there are 52 variables. Irrelevant variables will be discarded, and relevant variables will be used to predict player salaries. A detailed description for the variables will be also included.

## METHODS

To build predictive models, the approach is broken down into three phases. Each phase contains methods and tasks that need to be performed before moves on to the next phase.

- Phase 1: an exploratory data analysis is performed to better understand the dataset and its variables. Since the dataset contains 52 variables, we have to evaluate the variables to see which variables are relevant for the predictive models. In this phase, we will check if there are missing values and outliers. We also create a histogram for each variable to understand their data distribution. Predictive models tend to assume the model variables to be normally distributed. In addition, a histogram can help us examine the outliers in the data.
- Phase 2: this phase is dedicated to feature/variable selection. Meaningful features from 52 variables will be picked to build predictive models. We use various methods to select features such as a random forest classifier, a recursive feature elimination, an extra tree classifier, a chi squared analysis, and a Lasso regression. These methods create a ranking system for each feature. When some features end up with the same score, we use confusion matrix to reduce redundant features.
- Phase 3: after final features are selected, they will be used to build the predictive models in this phase. The model will be run and trained by using existing data. The results of the predictions will be summarized.

## RESULTS

### Phase 1: Exploratory data analysis (EDA)

A portion of the statistical summary of the dataset is shown in Figure 1. This summary lets us see the missing values among the variables; variables with less than 486 values are likely to have missing values. At the first glance, we can identify a variable with missing value for the variable True\_Shooting\_Pct. The range of the variables is shown in the min/max values.

	Number	Season_Start	Salary	Age	Games_Played	Games_Started	Minutes_Played	Player_Efficiency_Rating	True_Shooting_Pct
count	486.000000	486.0	4.840000e+02	486.000000	486.000000	486.000000	486.000000	486.000000	485.000000
mean	24392.170782	2017.0	6.717244e+06	26.405350	53.783951	25.308642	1223.051440	13.020782	0.526944
std	171.748286	0.0	7.376188e+06	4.345194	24.835638	28.715875	842.438143	5.762420	0.089771
min	24096.000000	2017.0	1.722400e+04	19.000000	1.000000	0.000000	1.000000	-17.600000	0.000000
25%	24242.250000	2017.0	1.471382e+06	23.000000	35.250000	1.000000	449.500000	9.800000	0.502000
50%	24392.500000	2017.0	3.343830e+06	26.000000	62.500000	11.000000	1197.500000	12.800000	0.537000
75%	24541.750000	2017.0	1.004073e+07	29.000000	75.000000	49.750000	1942.250000	15.800000	0.576000
max	24690.000000	2017.0	3.468255e+07	40.000000	82.000000	82.000000	3048.000000	31.500000	0.799000

Figure 1. Statistical summary of the NBA dataset.

Creating a histogram for each variable allows us to see their data distribution. The histograms for the variables are shown in the Figure 2.

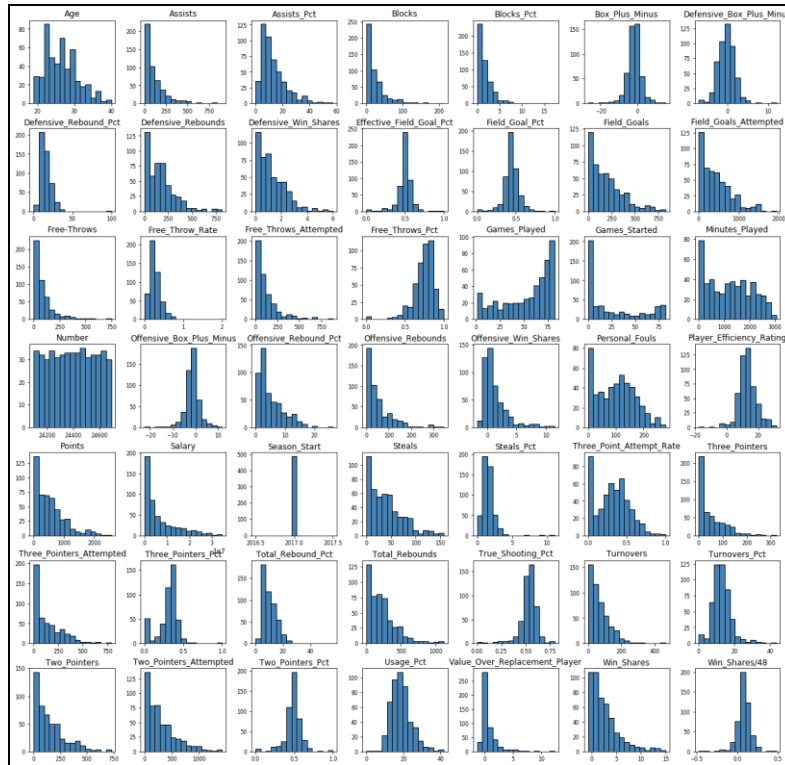


Figure 2. Histograms of all variables in the NBA dataset.

Many histograms display a normal distribution for variables such as Field\_Goal\_Pct and

Player\_Efficiency\_Rating. The histograms do not show any outliers in the data for all variables.

Examining the histograms gives us more insights into the variables/features, but we still need to conduct more steps to choose relevant features.

## Phase 2: Feature Selection

- **Random Forest Classifier**—this method utilizes the RandomForestClassifier from the sklearn library. Random forest is a very popular machine learning algorithm. They have good predictive capability with low overfitting, and easy interpretability due to boxes of information. This interpretability is given by the fact that it is straightforward to derive the importance of each variable on the tree decision. It is easy to compute how much

each variable is contributing to the decision. Each of our variables was subjected to this algorithm, and ranked on importance to the model. The output of this process is shown in Figure 3. The first column is the data frame column number, the index column is the name of the dataset variable, and the RF column is the Random Forest value. Due to the length of the output, only the first 10 values are shown.

	index	RF
47	Usage_Pct	0.033442
0	Age	0.031442
15	Free_Throw_Rate	0.031111
40	Total_Rebounds	0.025462
6	Defensive_Box_Plus_Minus	0.025320
19	Games_Started	0.025170
43	Turnovers_Pct	0.025037
16	Free_Throws_Attempted	0.024108
4	Blocks_Pct	0.023944
34	Team	0.023466

Figure 1. Output from Random Forest variable importance calculation.

- Recursive Feature Elimination (RFE)**—this technique is imported from the sklearn library. RFE is a backward selection of the predictors. This technique begins by building a model on the entire set of predictors and computing an importance score for each predictor. The insignificant predictors are removed, the model is re-built, and importance scores are computed again, hence the recursive nature of the process. The output from this process is shown in Figure 4. The first column is the data frame column number, the index column is the name of the dataset variable, and the RFE column is the Recursive Feature Elimination result. The values of the RFE column are shown as True because they have not been eliminated in the RFE selection process. In other words, these remaining variables are considered important to the overall

predictive model. The output below (in alphabetical order) shows all of the variables that returned a True value.

	index	RFE
0	Age	True
2	Assists_Pct	True
3	Blocks	True
4	Blocks_Pct	True
5	Box_Plus_Minus	True
6	Defensive_Box_Plus_Minus	True
7	Defensive_Rebound_Pct	True
14	Free_Throws	True
16	Free_Throws_Attempted	True
18	Games_Played	True
19	Games_Started	True
22	Offensive_Box_Plus_Minus	True
23	Offensive_Rebound_Pct	True
27	Player_Efficiency_Rating	True
30	Position	True
32	Steals	True
33	Steals_Pct	True
34	Team	True
47	Usage_Pct	True
49	Win_Shares	True

Figure 2. Output from RFE calculation.

- Extra trees Classifier**—this technique is imported as ExtraTreesClassifier from the sklearn library. The Extra Trees Classifier is similar to a Random Forest Classifier. The only difference is the way that the decision trees are constructed in the forest. Each decision tree in the extra trees forest is constructed from the original training sample. After that, at each test node, each tree is provided with a random sample of k features from the feature-set. From this, each decision tree must select the best feature to split the data. This random sample of features leads to the creation of multiple de-correlated decision trees. The output from this process is shown in Figure 5. The first column is the Data frame column number, the index column is the name of the dataset



variable, and the Extratrees column is the algorithm result. Due to the length of the output, only the first 10 values are shown.

	index	Extratrees
0	Age	0.037527
43	Turnovers_Pct	0.025375
28	Player_Name	0.024393
27	Player_Efficiency_Rating	0.023118
20	Minutes_Played	0.023043
14	Free-Throws	0.022644
34	Team	0.022525
16	Free_Throws_Attempted	0.022187
29	Points	0.022117
47	Usage_Pct	0.021970

Figure 3. Output from Extra Trees Classifier.

- **Chi Square**—this technique is imported from the sklearn library. The Chi Square Test is used in statistics to test the independence of two events. In feature selection, the two events are occurrence of the feature and occurrence of the class. In other words, we want to test whether the occurrence of a specific feature and the occurrence of a specific class are independent. When the two events are independent, the observed count is close to the expected count, thus a small chi square score. So a high chi square value indicates that the hypothesis of independence is incorrect. In other words, the higher value of the chi square score, the more likelihood the feature is correlated with the class, thus it should be selected for the model. The output of the chi square test is shown in Figure 6. The first column is the Data frame column number, the index column is the name of the dataset variable, and the Chi\_Square column is the calculated chi square value. Due to the length of the output, only the first 10 values are shown.

	index	Chi_Square
20	Minutes_Played	31888.49
13	Field_Goals_Attempted	30812.33
29	Points	30580.51
45	Two_Pointers_Attempted	27591.87
12	Field_Goals	27179.60
21	Number	26747.82
28	Player_Name	26366.52
37	Three_Pointers_Attempted	26266.40
40	Total_Rebounds	25068.39
44	Two_Pointers	23515.41

Figure 4. Output from the Chi Square calculation.

- **Lasso Regression (L1)**—this technique is imported from the sklearn library.

Regularization consists in adding a penalty to the different parameters of the machine learning model to reduce the freedom of the model and in other words to avoid overfitting. In linear model Regularization, the penalty is applied over the coefficients that multiply each of the predictors. From the different types of Regularization, Lasso or L1 has the property that is able to shrink some of the coefficients to zero. Therefore, that feature can be removed from the model. The output from our Lasso Regression is shown in Figure 7. The first column is the Data Frame column number, the index column is the name of the dataset variable, and the L1 column is the Lasso Regression result. Like with the RFE output, the values of the L1 column are shown as True because they have not been eliminated in the Lasso Regression selection process. In other words, these remaining variables are considered important to the overall predictive model. The output below (in alphabetical order) shows all of the variables that returned a True value.

	index	L1
0	Age	True
1	Assists	True
2	Assists_Pct	True
3	Blocks	True
7	Defensive_Rebound_Pct	True
8	Defensive_Rebounds	True
12	Field_Goals	True
13	Field_Goals_Attempted	True
14	Free_Throws	True
16	Free_Throws_Attempted	True
18	Games_Played	True
19	Games_Started	True
20	Minutes_Played	True
21	Number	True
24	Offensive_Rebounds	True
26	Personal_Fouls	True
27	Player_Efficiency_Rating	True
28	Player_Name	True
29	Points	True
31	Season_Start	True
32	Steals	True
34	Team	True
36	Three_Pointers	True
37	Three_Pointers_Attempted	True
40	Total_Rebounds	True
42	Turnovers	True
43	Turnovers_Pct	True
44	Two_Pointers	True
45	Two_Pointers_Attempted	True

Figure 5. Output from Lasso Regression (L1) calculation.

- **Scoring Table of Final Results**—a scoring table was constructed that tallied all of the different feature selection methods. From this table, we were able to determine which variables should be included in our model. The scoring table, listing only values with a final\_score above 2, is shown in Figure 8.

	index	RF	Extratrees	Chi_Square	RFE	L1	final_score
0	Age	1	1	0	1	1	4
27	Player_Efficiency_Rating	0	1	0	1	1	3
20	Minutes_Played	0	1	1	0	1	3
3	Blocks	0	0	0	1	1	2
16	Free_Throws_Attempted	0	0	0	1	1	2
43	Turnovers_Pct	0	1	0	0	1	2
45	Two_Pointers_Attempted	0	0	1	0	1	2
34	Team	0	0	0	1	1	2
32	Steals	0	0	0	1	1	2
29	Points	0	0	1	0	1	2
28	Player_Name	0	1	0	0	1	2
2	Assists_Pct	0	0	0	1	1	2
19	Games_Started	0	0	0	1	1	2
18	Games_Played	0	0	0	1	1	2
40	Total_Rebounds	1	0	0	0	1	2
13	Field_Goals_Attempted	0	0	1	0	1	2
12	Field_Goals	0	0	1	0	1	2
6	Defensive_Box_Plus_Minus	1	0	0	1	0	2
7	Defensive_Rebound_Pct	0	0	0	1	1	2
47	Usage_Pct	1	0	0	1	0	2
14	Free-Throws	0	0	0	1	1	2

Figure 6. Scoring table of all feature selection results.

We are obviously going to include Age and Player\_Efficiency\_Rating since they scored the highest. After that, however, we have many other variables that all returned a sum score of 2. We eliminated some of these features by looking at the correlation matrix.

- **Confusion (Correlation) Matrix**— once the features were scored a correlation matrix was used to test the features correlation. The table in Figure 9 is a heatmap showing the correlation of the features scoring at least a 2.

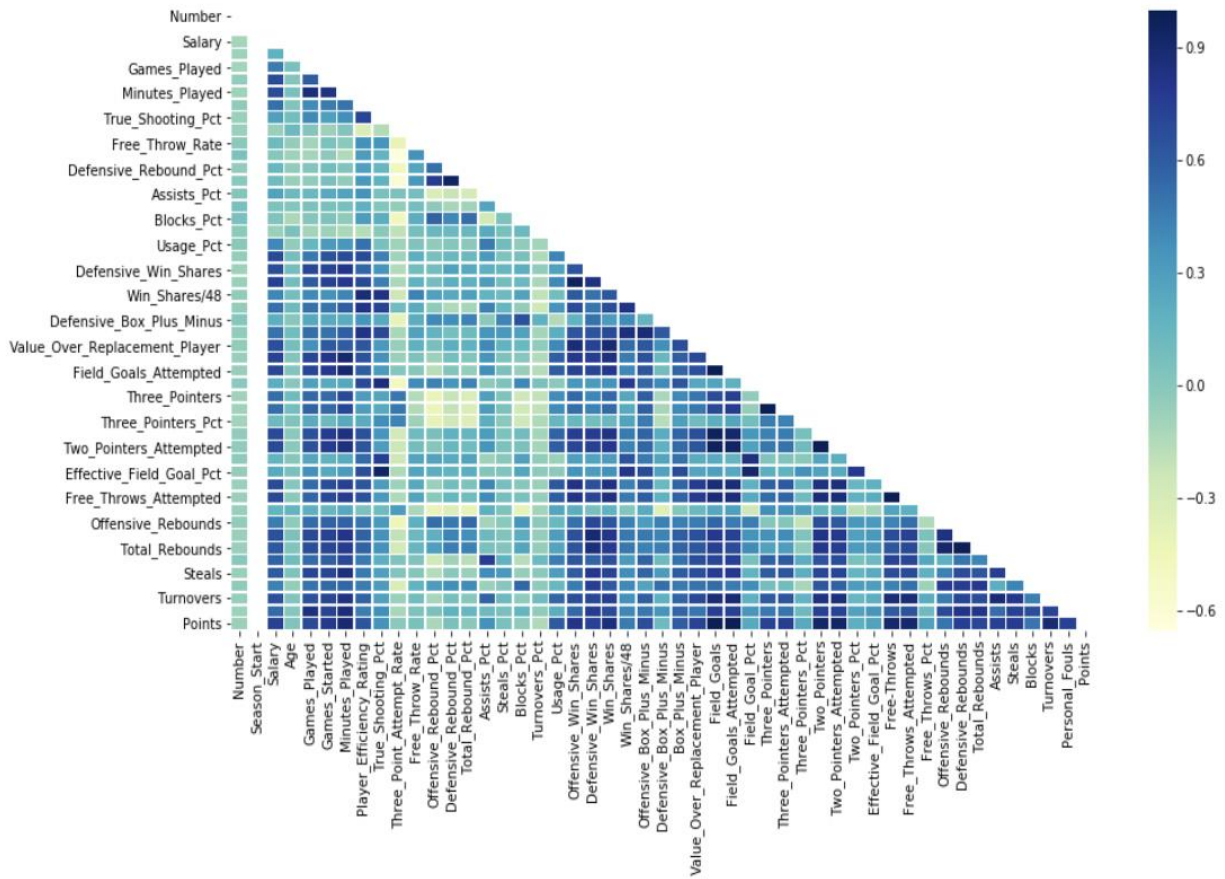


Figure 9. The Output from Initial Correlation Heatmap

It is clear in this heatmap that there is quite a bit of correlation between features. This is not surprising because a good basketball player is good at everything and a bad one is bad at everything. These are all important parts of the game so if someone is great at assists they are also likely great at field goals. That being said, we still want to remove as much unnecessary correlation as we can. We first removed all features marking “attempts”. They were too highly correlated with successes, so they did not add much value. The other feature we removed was Minutes Played. This scored highly in our L1 test, but it is very highly correlated with most features. That makes sense because the more time you play means the more chances you have at blocking, rebounding, taking

shots, etc. Removing these features gave us a heatmap with the 14 features in Figure 10.

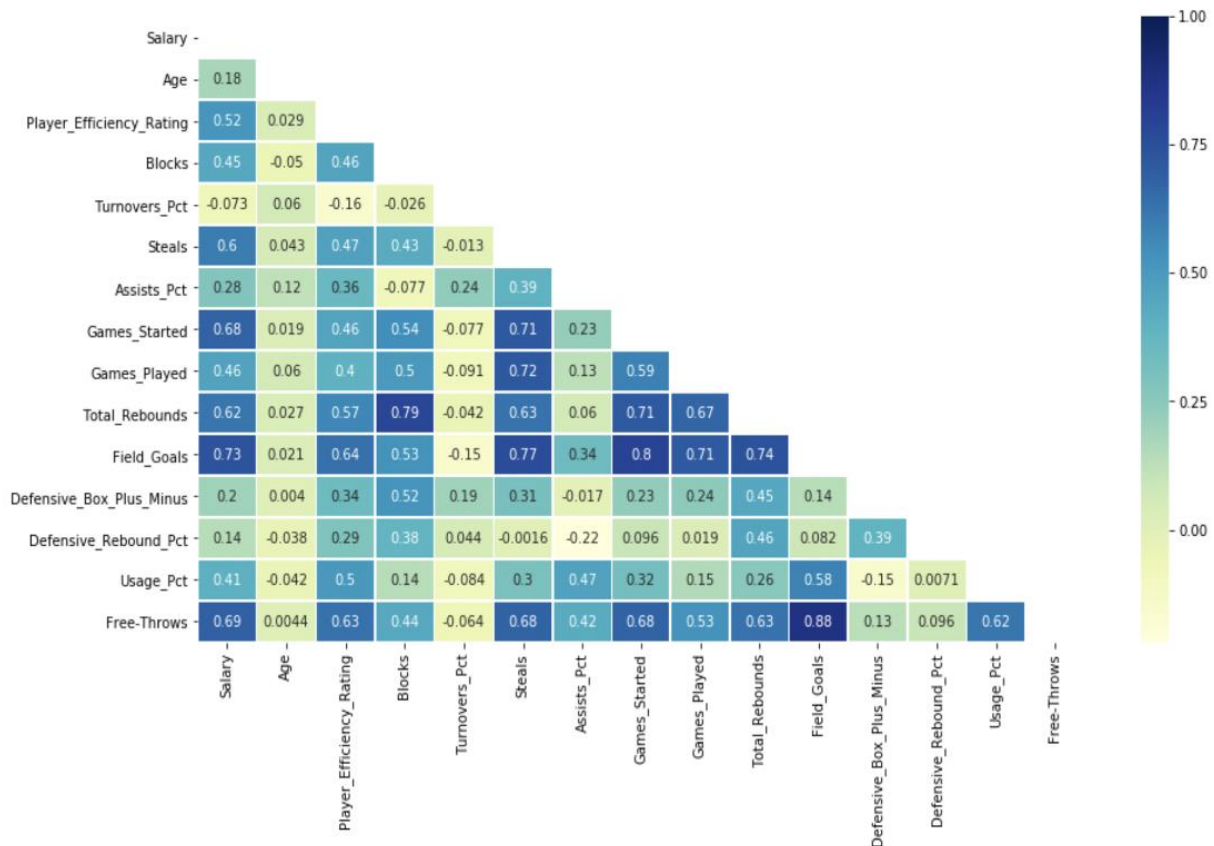


Figure 10. The Output from Final Correlation Heatmap

After highly correlated features are removed, the number of remaining features has been reduced. There is still some high correlation that we have left in for now. For example, Field Goals and Free Throws are highly correlated but field goals do not include free throws so there is not that direct relationship like we had in other features.

### Phase 3: Feature Finalization and Model Build

**Feature Finalization**—a features summary is shown below in Table 1. This table summarizes each feature selection method and was used to determine the final features that were used for the preliminary model run in Phase III.

Feature	Summary Statistics			Histograms	Correlation Matrix	Feature Selection
	Complete Data?	Reasonable Max Value?	Reasonable Min Value?	Normal Distribution?	Uncorrelated?	Final Score
Age	Yes	Yes	Yes	No	Yes	4
Player_Efficiency_Rating	Yes	Yes	Yes	Yes	Yes	3
Minutes_Played	Yes	Yes	Yes	No	No	3
Blocks	Yes	Yes	Yes	No	Yes	2
Free_Throws_Attempted	Yes	Yes	Yes	No	No	2
Turnovers_Pct	No	Yes	Yes	Yes	Yes	2
Two_Pointers_Attempted	Yes	Yes	Yes	No	No	2
Team	Ignored. Not a measurable player metric.					2
Steals	Yes	Yes	Yes	No	Yes	2
Points	Yes	Yes	Yes	No	No	2
Player_Name	Ignored. Not a measurable player metric.					2
Assists_Pct	Yes	Yes	Yes	No	Yes	2
Games_Started	Yes	Yes	Yes	No	Yes	2
Games_Played	Yes	Yes	Yes	No	Yes	2
Total_Rebounds	Yes	Yes	Yes	No	Yes	2
Field_Goals_Attempted	Yes	Yes	Yes	No	No	2
Field_Goals	Yes	Yes	Yes	No	Yes	2
Defensive_Box_Plus_Minus	Yes	Yes	Yes	Yes	Yes	2
Defensive_Rebound_Pct	Yes	No	Yes	Yes	Yes	2

Usage_Pct	Yes	Yes	Yes	No	Yes	2
Free-Throws	Yes	Yes	Yes	No	Yes	2

Table 1. Feature selection summary.

Based on the table summaries, and on previous discussions in this paper, we have selected 14 variables to be used in our initial model. They are: Age, Player Efficiency Rating, Blocks, Turnover Pct, Steals, Assists Pct, Games Started, Games Played, Total Rebounds, Field Goals, Defensive Box Plus Minus, Defensive Rebound Pct, Usage Pct, and Free Throws.

- **Model Build**—we used an iterative modeling approach to determine the best model for our predictions. We ended up using five different models. We divided the data into a test and training subset before each model fit. We used a 70/30 test to train ratio for each model run. The first four models were run using the sklearn Python package. The fifth model was run using the xgboost Python package. A description of each model considered is discussed below:
  - Model 1: Ordinary Least Squares (OLS)—our reason for selecting this model is that it is one of the most commonly used models and would serve as a good baseline with which we could compare the subsequent models. Table 2 shows the results from this model.



Model	Train RMSE	Train R2	Test RMSE	Test R2
OLS	4398500	0.63	5024129	0.57

Table 1. OLS Model Fit Results.

- The R2 values for the models indicate that we are able to account for about 60% of the salary prediction variability with the features we selected.
  - The root mean square error (RMSE) values will serve as a reference point for the other linear models we are going to use.
- Model 2: Ridge Regression (RR)—the ridge regression model is an extension of linear regression where the loss function is modified to minimize the complexity of the model. Ridge regression includes an  $\alpha$  variable that we can adjust to help improve the fit. The higher the  $\alpha$  value, the more likely our model is going to suffer from under-fitting. A lower  $\alpha$  value can lead to under-fitting. Table 3 shows the results from the RR model for two different  $\alpha$  values.

Model	Train RMSE	Train R2	Test RMSE	Test R2
RR, $\alpha=0.01$	4398504	0.63	5024400	0.57
RR, $\alpha=10$	4608385	0.59	5203120	0.54

Table 2. RR Model Fit Results.

- The R2 values for the models indicate that we are able to account for about 60% of the salary prediction variability with the features we selected.
  - The RMSE values for the RR method are slightly higher than the OLS values. This indicates that the OLS predictions are more accurate.
  - Increasing the  $\alpha$  value decreased the model accuracy.
- Model 3: Lasso Regression (LR)—Lasso regression is another modification of linear regression modeling. In Lasso, the loss function is modified to minimize the complexity of the model by limiting the sum of the absolute values of the model coefficients. We will select different  $\alpha$  values to tune the model for the best fit. Table 4 shows the results from the LR model for two different  $\alpha$  values.

Model	Train RMSE	Train R2	Test RMSE	Test R2
LR, $\alpha=0.01$	4398500	0.63	5024129	0.57
LR, $\alpha=100$	4398502	0.63	5023978	0.57

Table 3. LR Model Fit Results.

- The R2 values for the models indicate that we are able to account for about 60% of the salary prediction variability with the features selected.
  - The RMSE values for the LR method are essentially identical to the OLS values. This indicates that the OLS and Lasso predictions are basically the same.
  - Increasing the  $\alpha$  value had little impact on the model accuracy.
- Model 4: ElasticNet Regression (ENR)—ElasticNet regression combines the properties of both the Ridge and Lasso regression methods. ENR also utilized an  $\alpha$  term to help improve fit. Table 5 shows the results from the ENR model for two different  $\alpha$  values.

Model	Train RMSE	Train R2	Test RMSE	Test R2
ENR, $\alpha=0.01$	4432638	0.62	5070003	0.56
ENR, $\alpha=10$	7010747	0.06	7488845	0.05

Table 4. ENR Model Fit Results.

- The R2 values for the models were still around 60%, but were slightly lower than what we saw with the other models.
  - The RMSE values for the ENR method were the highest of all the models. This indicates the least accurate prediction model. This is likely due to the fact that we removed features that were highly correlated with one another.
  - Increasing the  $\alpha$  value had a negative impact on our prediction accuracy.
- Model 5: Extreme Gradient Boosting Regression (XGBR)—XGBR is an ensemble modeling technique where new models are added to correct the errors made by existing models. This method uses the gradient boosting decision tree algorithm. Table 6 shows the results from the XGBR model.

Model	Train RMSE	Train R2	Test RMSE	Test R2
XGBR	1818766	0.92	4531190	0.65

Table 5. XGBR Model Fit Results.

- The R2 values for the XGBR method were the best of all models evaluated. Our training set R2 had a value of 92%, which is excellent. The results diminished to only 65% when applied to the test set, but this about 10% better than our other models. The reason for the lower R2 value for the test set is likely due to the lower number of data points. Remember that the test set contained only one-third of the original data.

- The RMSE values for the XGBR method were also the lowest of all models evaluated.

This indicates the most accurate prediction model.

- Because the XGBR model provided the best overall fit, this is the model we selected as our final model.

## Predictive Models

- The prediction results for 25 players are shown in Figure 11. The predictive model is random forest classifiers.

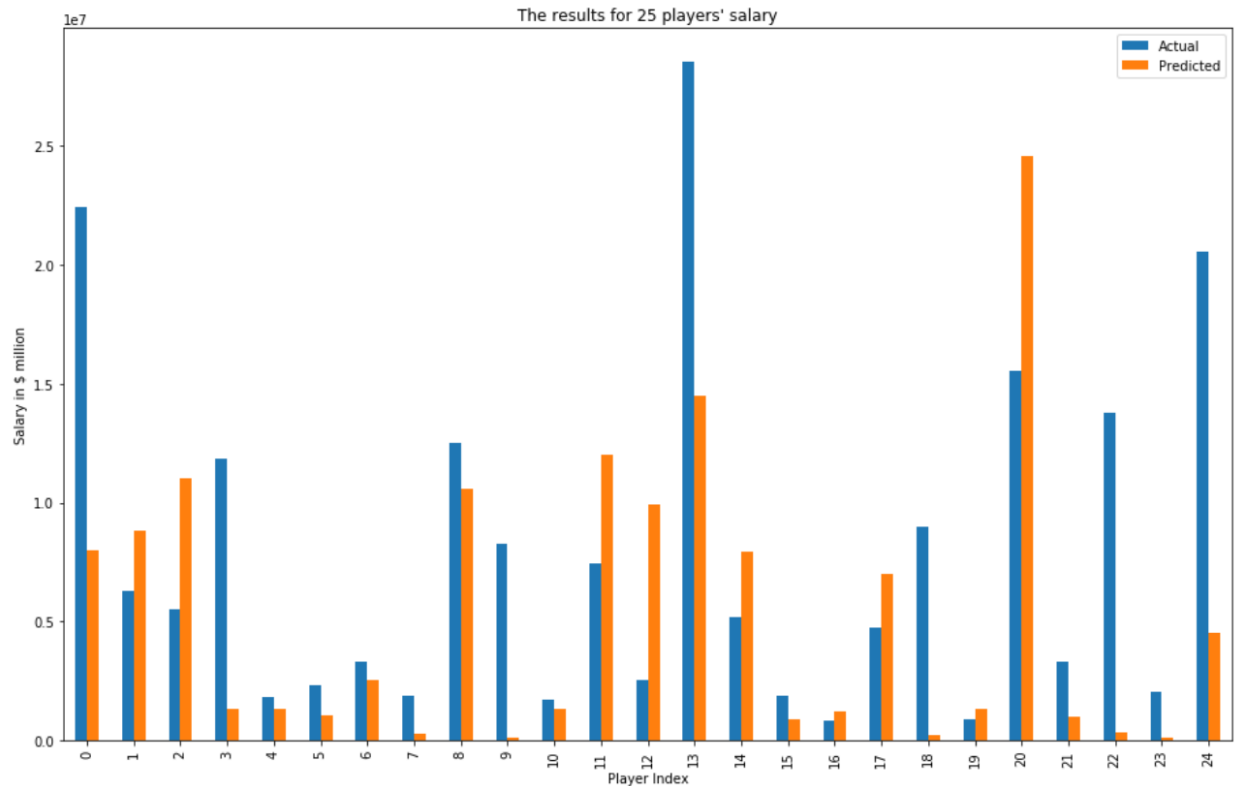


Figure 11. Predictions by Random Forest Classifier.

- The prediction results for 25 players are shown in Figure 12. The predictive model is linear regression.

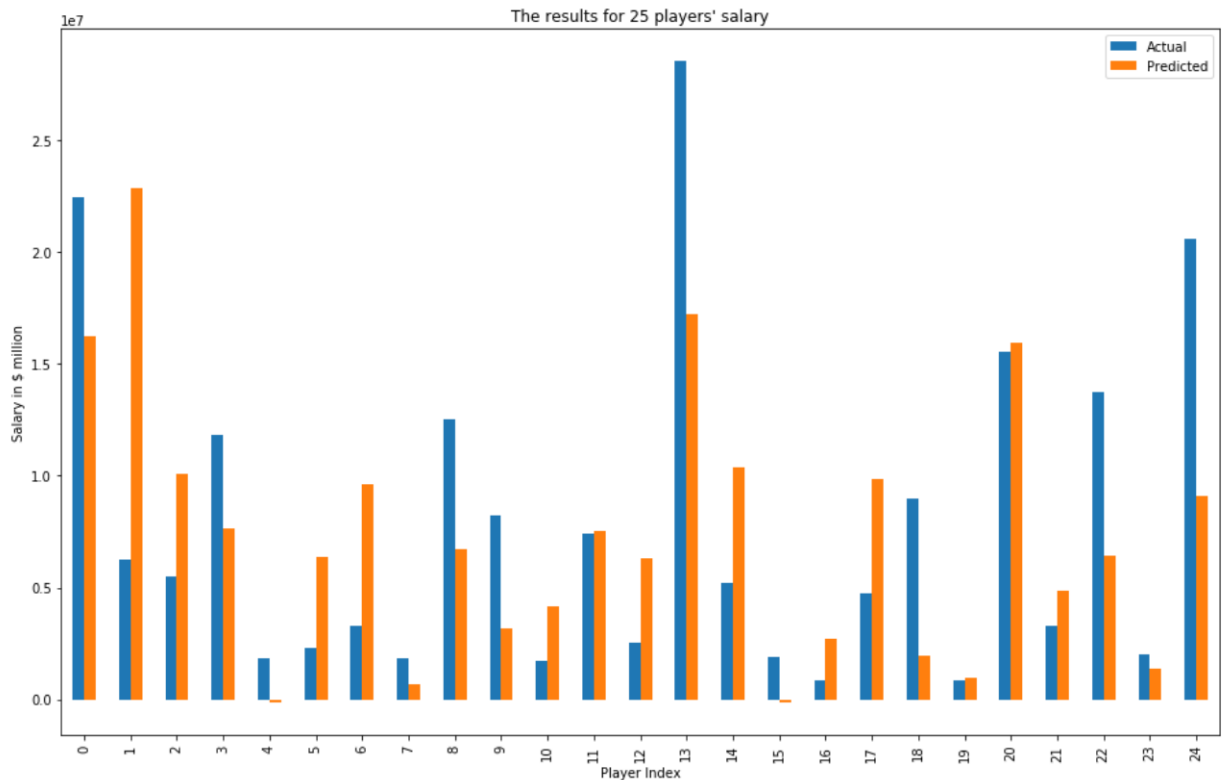


Figure 12. Predictions by Linear Regression.

## DISCUSSION & CONCLUSION

As the exploratory data analysis was performed, the dataset and its variables were better understood. Using various data analysis techniques such as plotting histograms, generating a confusion matrix, irrelevant features were identified and removed. Predictive models were ready to be built by meaningful features. Furthermore, more feature selection techniques such as random forest classifier, RFE, extra trees classifier were used to ensure the

features were truly important. From 52 variables, only 14 variables were finalized and selected to analyze player salaries.

The results of two of the most accurate models were displayed by random forest classifier and linear regression. In addition, the results of these two models were very close. Minor differences could be used for further examination on whether a player was actually overpaid or underpaid according to their performance. Due to the limited space of the Jupyter notebook, only 25 players are shown, the result for the rest of the player salaries could be accessed from the source codes.

## ACKNOWLEDGEMENTS

We want to thank the data science community from Kaggle for the dataset, instructions and help on the project, and also thank to many detailed answers from StackOverFlow.

## REFERENCES

1. Wikipedia. *National Basketball Association*. Retrieved from [https://en.wikipedia.org/wiki/National\\_Basketball\\_Association](https://en.wikipedia.org/wiki/National_Basketball_Association)
2. Casalan, A. (2018, October 20). *NBA player stats 2017-18*. Retrieved from <https://www.kaggle.com/acasalan/nba-player-stats-201718>
3. Abbott, D. (2014). *Applied predictive analytics: principles and techniques for the professional data analyst*. Indianapolis, IN: Wiley.
4. Davis, S. (2019, June 19). *WHERE ARE THEY NOW? The biggest NBA Draft busts of all time*. Retrieved from <https://www.businessinsider.com/nba-draft-busts-2015-6>.
5. Ye, D. (2019, July 4). *Top 10 Most Underrated NBA Players of the 2018-19 Season*. Retrieved from <https://howtheyplay.com/team-sports/Top-10-Most-Underrated-NBA-Players-of-the-2018-19-Season>

6. Gleeson, S. (2019, June 19). *NBA draft: Ranking the five biggest busts as the No. 1 pick*. Retrieved from <https://www.usatoday.com/story/sports/nba/draft/2019/06/19/nba-draft-ranking-five-biggest-busts-no-1-pick/1488407001>
7. These are the 2016/17 salaries of all NBA teams. Retrieved from <https://hoopshype.com/salaries/2016-2017>
8. Krishnan, S. (2019, December 20). *Variable Selection using Python - Vote based approach*. Retrieved from <https://medium.com/@sundarstyles89/variable-selection-using-python-vote-based-approach-faa42da960f0>
9. Dubey, A. (2018, December 15). *Feature Selection Using Random forest*. Retrieved from <https://towardsdatascience.com/feature-selection-using-random-forest-26d7b747597f>
10. Rade, D. (2019, September 2). *Feature Selection in Python - Recursive Feature Elimination*. Retrieved from <https://towardsdatascience.com/feature-selection-in-python-recursive-feature-elimination-19f1c39b8d15>
11. Chi Square test for feature selection. (2018, July 5). Retrieved from <http://www.learn4master.com/machine-learning/chi-square-test-for-feature-selection>
12. Dubey, A. (2019, February 4). *Feature Selection Using Regularization*. Retrieved from <https://towardsdatascience.com/feature-selection-using-regularisation-a3678b71e499>
13. Singh, D. (2019, May 17). *Linear Lasso Ridge Regression in Scikit-learn*. Retrieved from <https://www.pluralsight.com/guides/linear-lasso-ridge-regression-scikit-learn>
14. Brownlee, J. (2019, August 21). *A Gentle Introduction to XGBoost for Applied Machine Learning*. Retrieved from <https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning>
15. Gupta, A. (2019, July 26). *ML: Extra Tree Classifier for Feature Selection*. Retrieved from <https://www.geeksforgeeks.org/ml-extra-tree-classifier-for-feature-selection>