

Phuc Ton Nguyen

Prediction of salary for MLB player

Guided Capstone Project Report

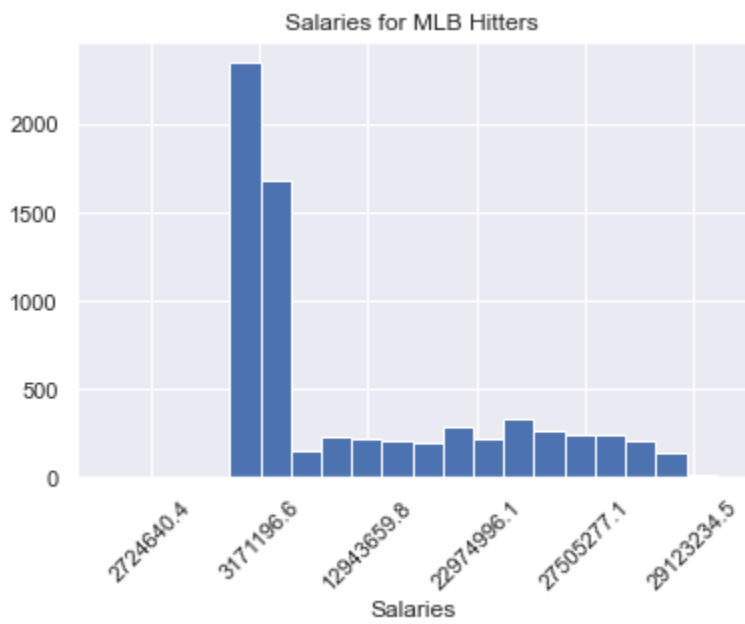
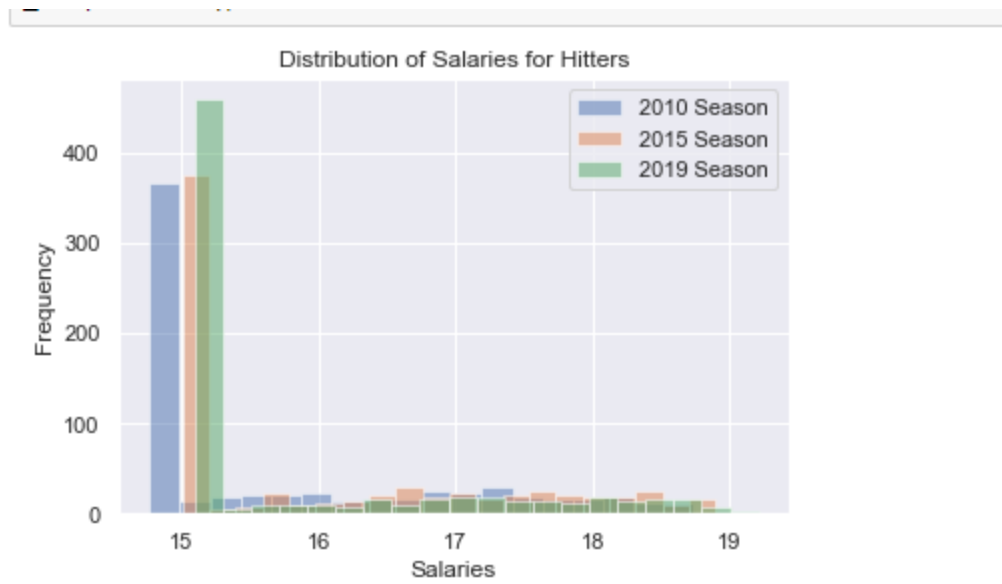
Introduction

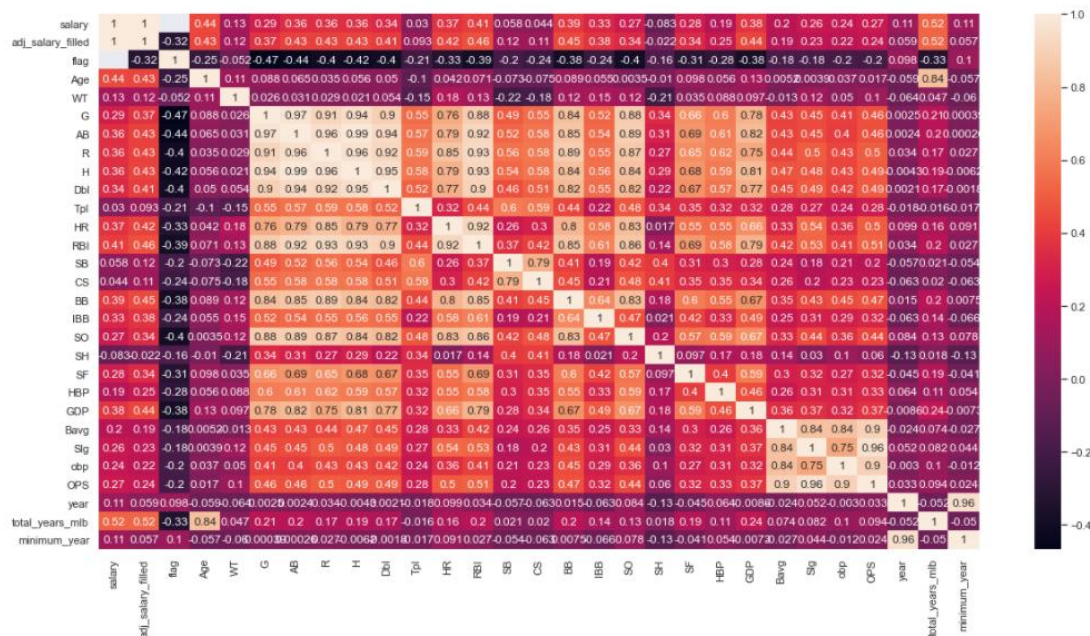
San Francisco Giants wants to buy player and they are asking the evaluation of the players on the market to buy.

Exploratory Data Analysis

The cleaned and merged dataset contained 6928 observations (players) each with 38 features. I did univariate, bivariate and multivariate analyses along with visualizations. I also performed some hypothesis testing for more statistically rigorous statements. The followings are the summary of what I found through exploratory data analysis (EDA) and inferential statistics.

Hitter



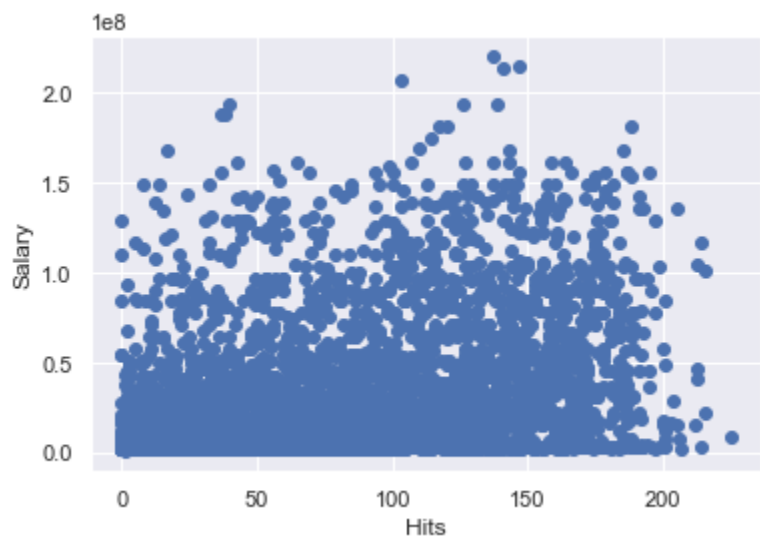
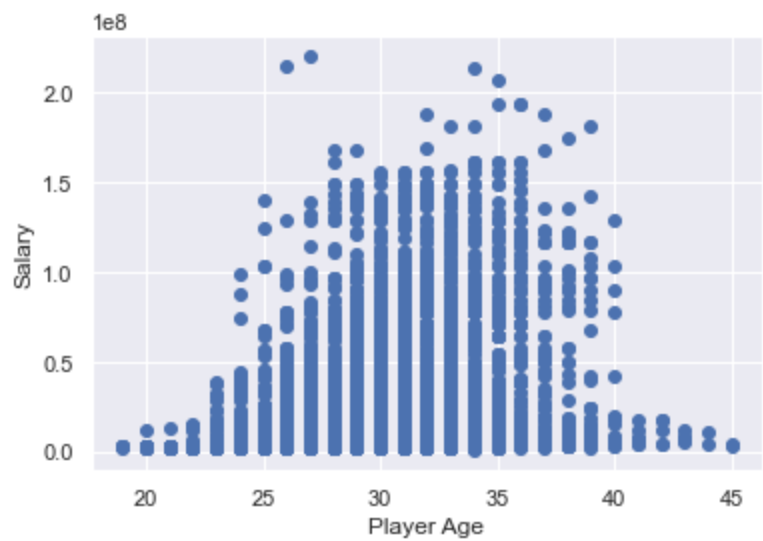


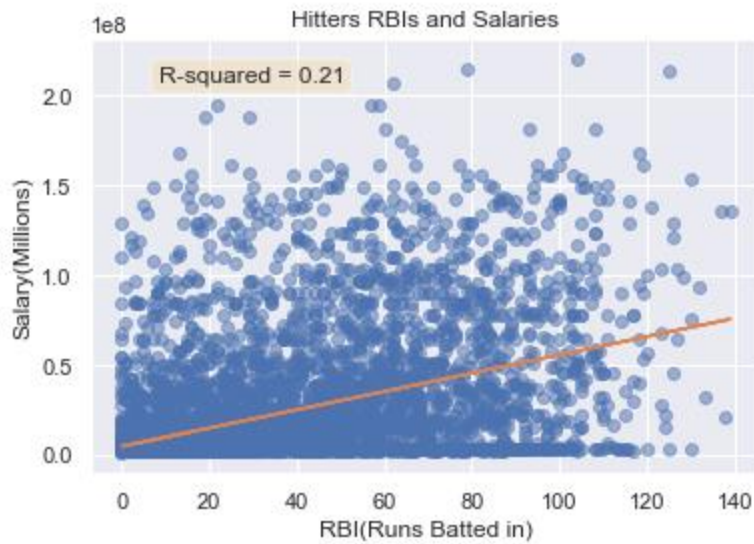
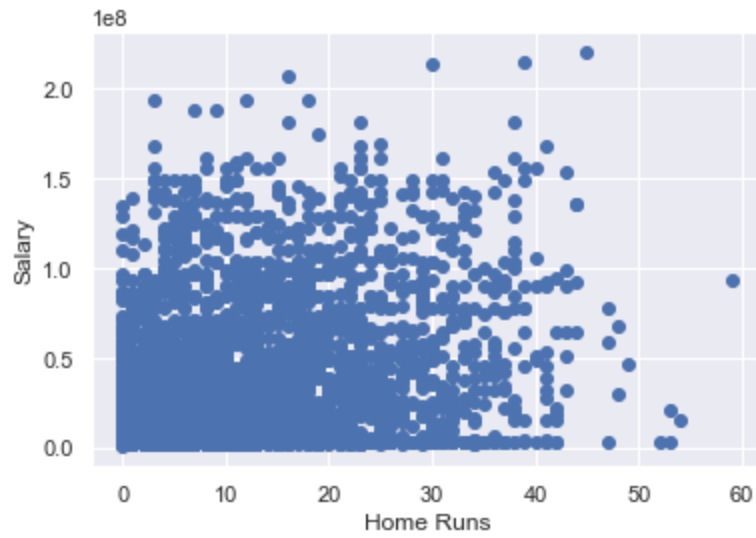
```

salary          1.000000
adj_salary_filled 1.000000
flag            0.315781
Age             0.433457
G               0.374448
AB              0.426127
R               0.426275
H               0.429633
DBI             0.407866
HR              0.423042
RBI             0.463561
BB              0.447756
IBB             0.377933
SO              0.343339
SF              0.344241
GDP             0.438088
total_years_mlb 0.524767
Name: adj_salary_filled, dtype: float64

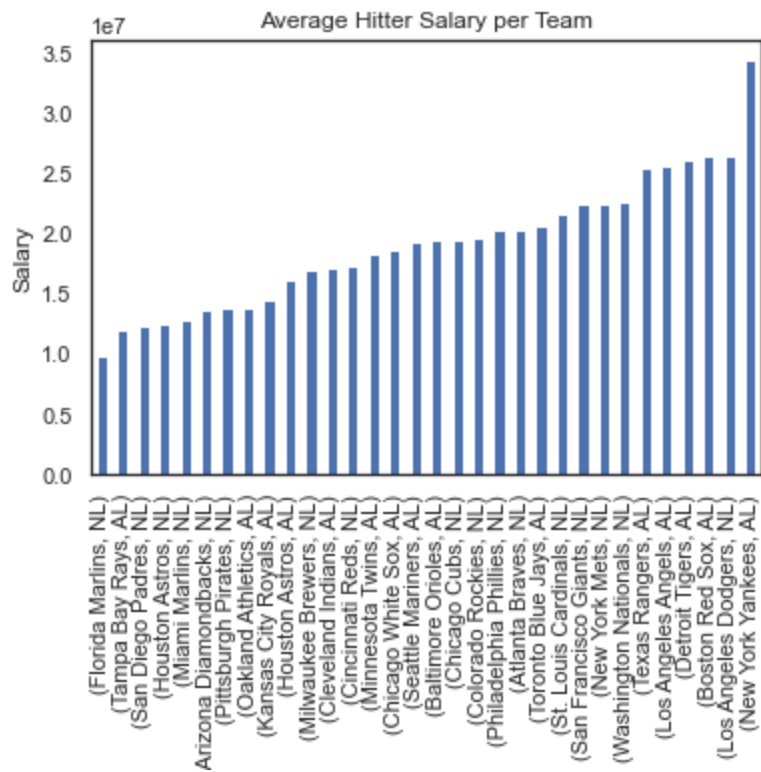
```

Independent Variables Vs. Dependent Variable Scatter Plots

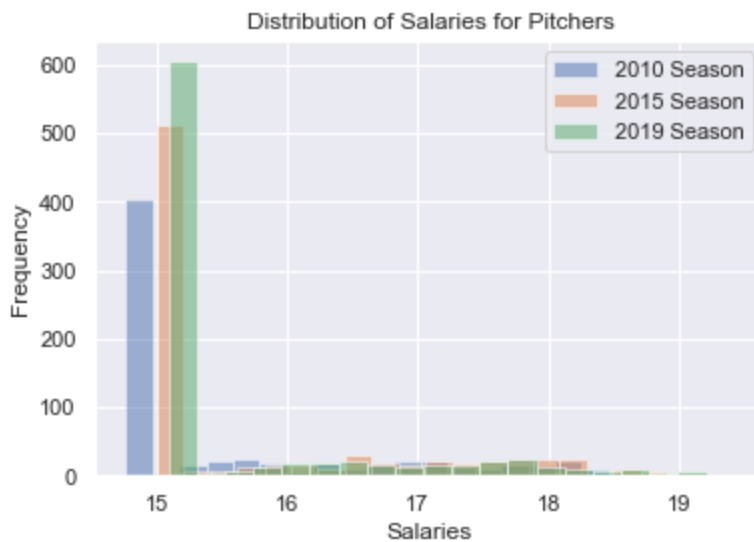


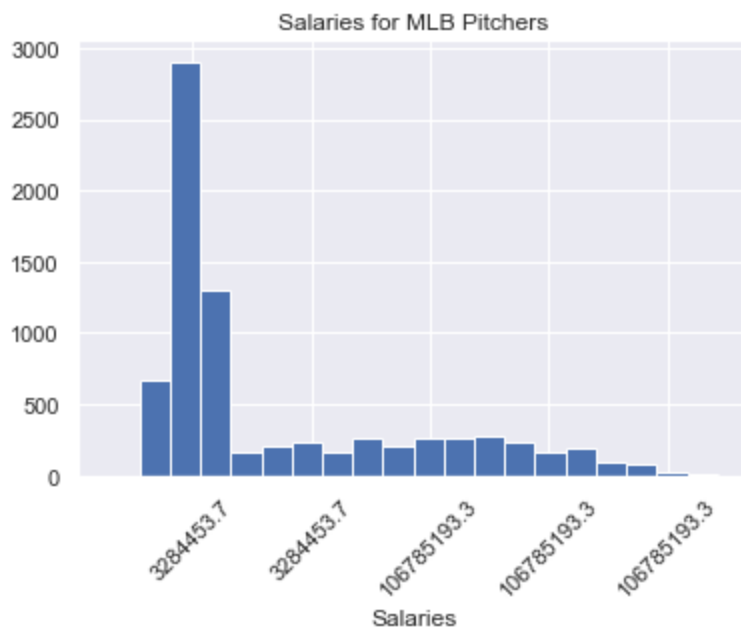


From a hitters perspective what team pays the most?



Pitcher



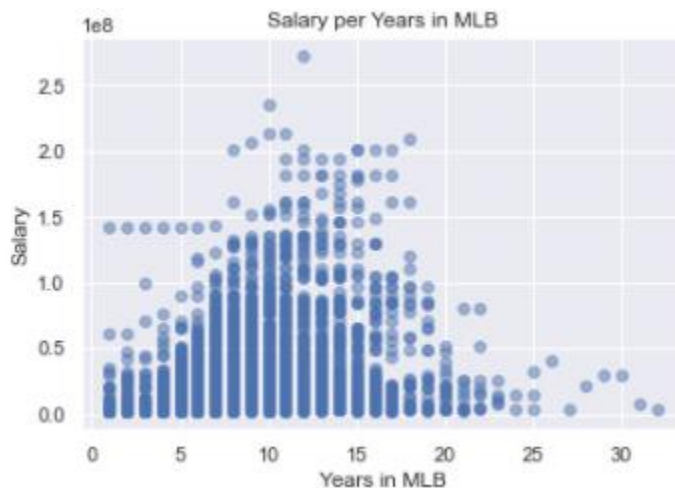


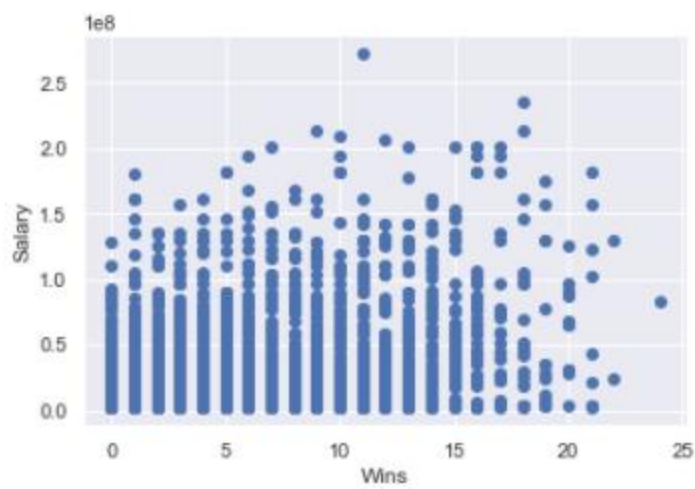
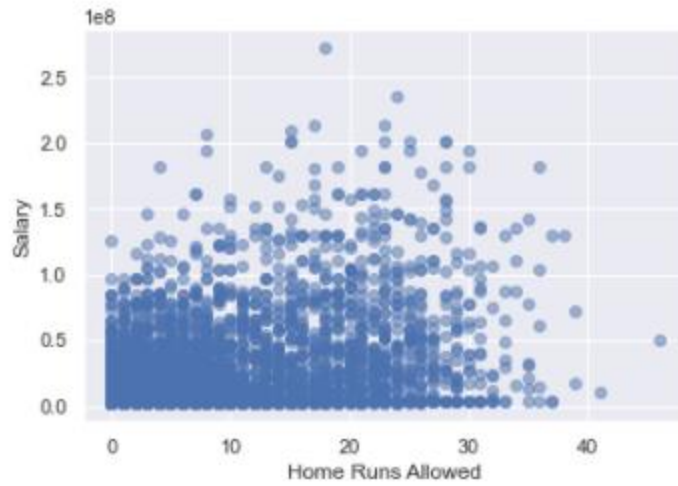
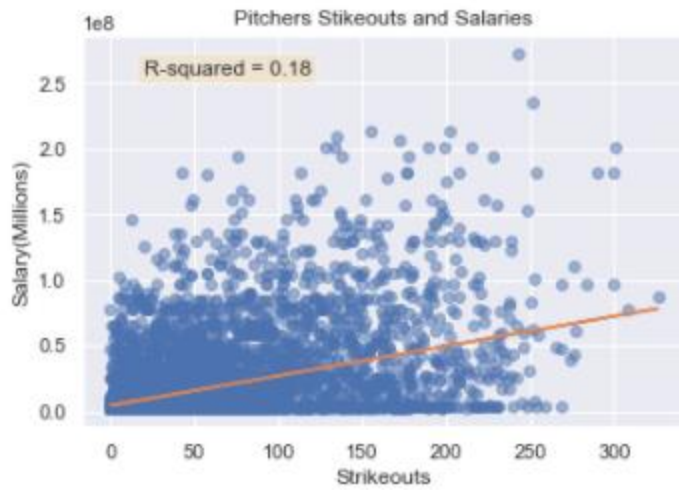
```

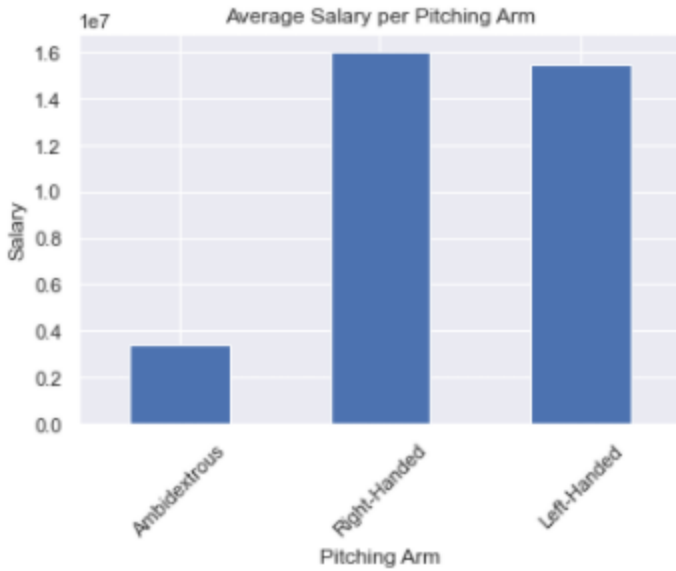
salary          1.000000
adj_salary_filled 1.000000
flag            0.310272
Age             0.391690
W               0.406481
L               0.355426
GS              0.407395
IP              0.420481
H               0.404449
HR              0.392769
R               0.385521
ER              0.385179
BB              0.333458
SO              0.426906
total_years_mlb 0.515814
Name: adj_salary_filled, dtype: float64

```

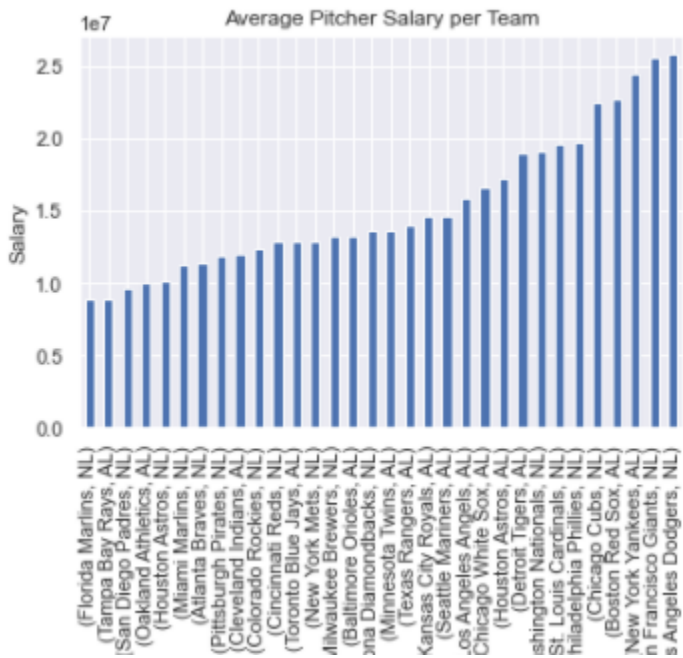
Independent Variables Vs. Dependant Variable Scatter Plots







From a pitchers perspective what team pays the most?

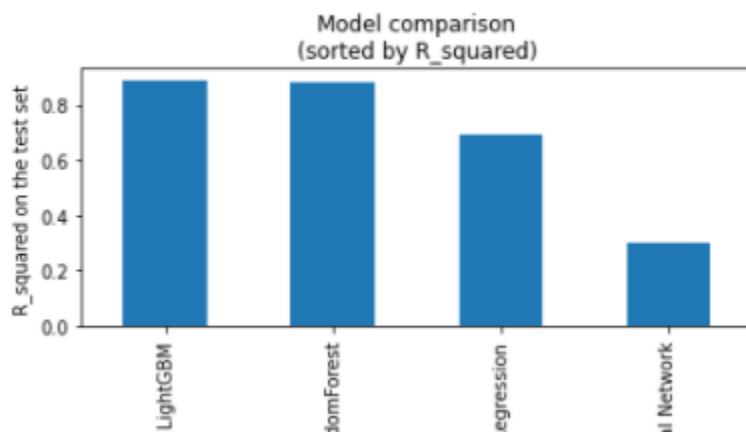


Model

I used 4 models Linear Regression, Random Forest and Light Gradient Boosting, and Neural Network to evaluate the models.

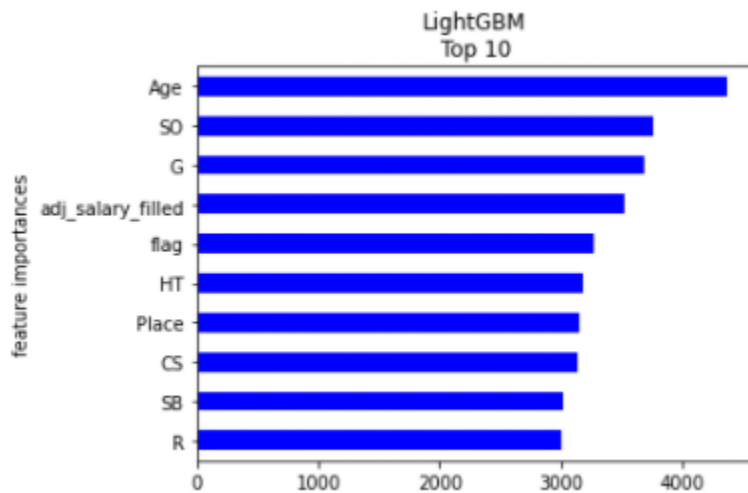
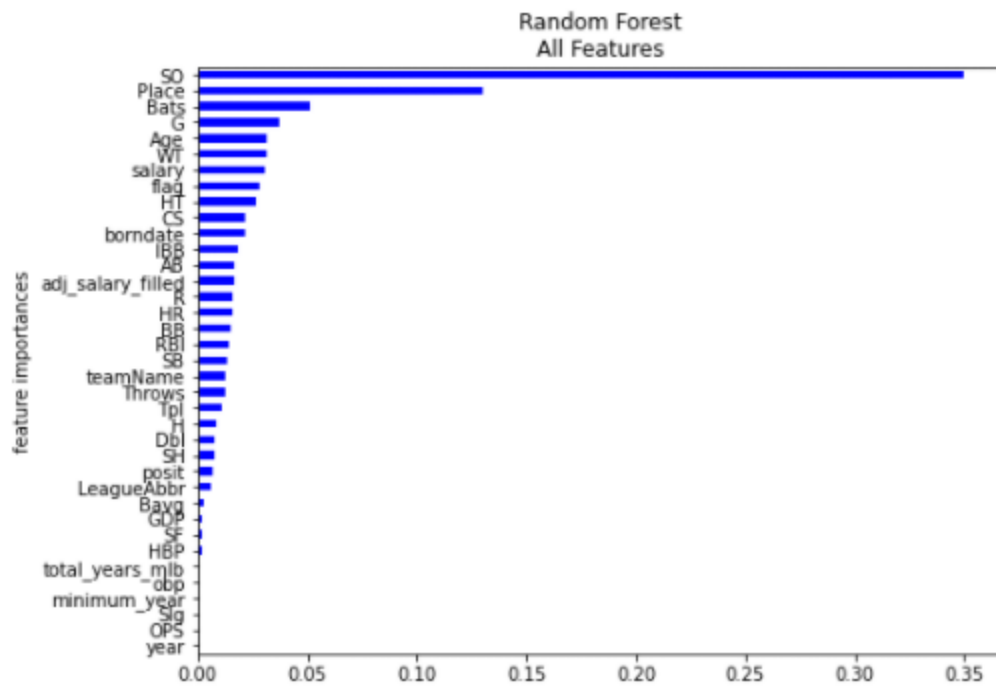
Result

	RMSE_val	RMSE_test	R_squared	Time	Scaling
LightGBM	120472	120537	0.8883	10min 31s	True
RandomForest	131979	124013	0.8818	3min 6s	False
LinearRegression	202941	200704	0.6904	5.98 s	True
Neural Network	26461500	27301891	0.2996	36.1 s	False



LightGBM is best followed by RandomForest followed by LinearRegression

Feature Importances



Conclusion

I have found the best MLB salary prediction models is the LightGBM models which showed high speed and best performance in RMSE. If one model should be selected I would recommend to use the LightGBM model since it is faster and it makes fewer outliers.