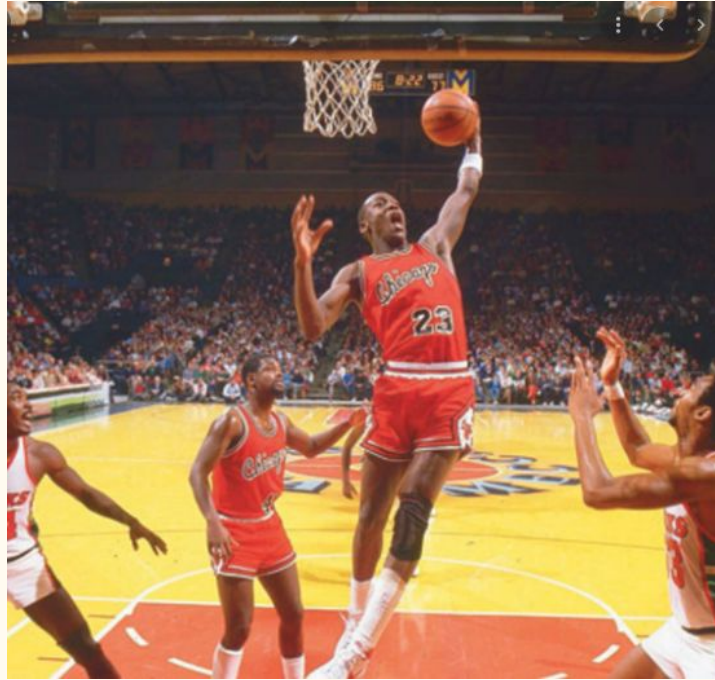


Predicting NBA Center Rebounds

Mario Hage



Introduction

This project is an attempt to answer the below questions, and to discover answers for currently unknown questions.

- Can compiling individual game data from all the NBA Centers in the league tell us something that can be used to our advantage whilst sports betting?
- What hints at players having big games?
- What hints at players falling short?
- Can I come up with data points that can be fed into a model that predicts accurate enough to profit whilst sports betting?
- What is the significance of a player or team “being on a roll?”
- What about the opposite of “being on a roll?”

Data

nba_api : https://github.com/swar/nba_api

Distance data: <https://www.rostrum.blog/2018/12/24/nba-travel/>

The nba_api contained player box score history, which was filtered for centers. It also contained the Match-Up column, which was split into Team & Opponent; later joined with the Distance data on Start/End Routes.

Data

(join process)

Nba_api original column

MATCHUP
LAC @ OKC



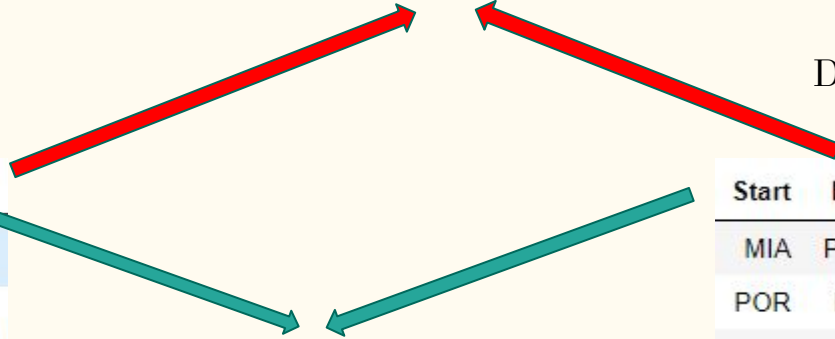
Nba_api (after data wrangling)

Left Join On Start and End Columns

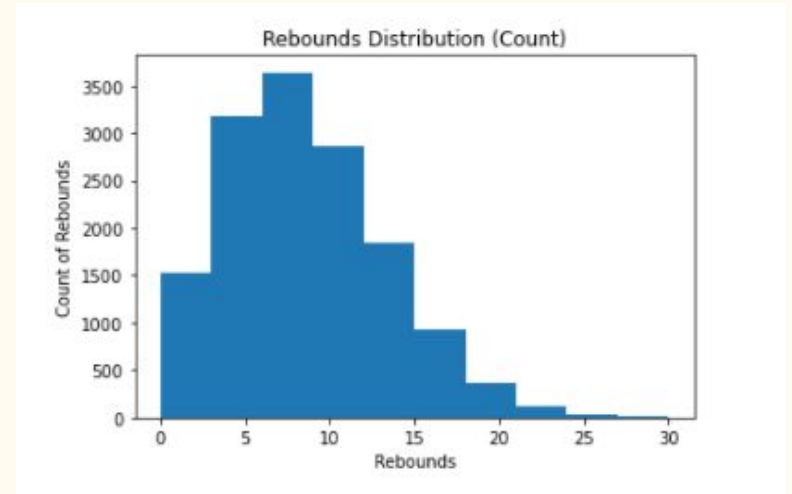
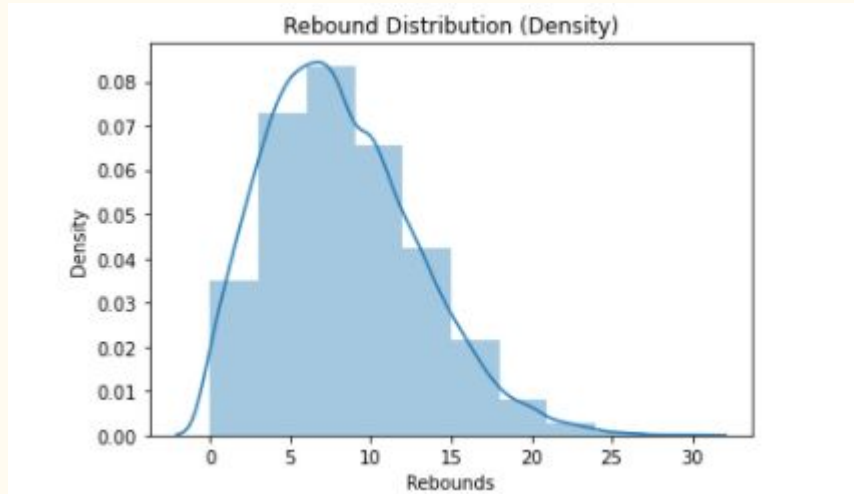
Distance Data

Home/Away	PreviousHome/Away	PreviousOpp	start	end
Away	Away	HOU	HOU	OKC
Away	Away	CHA	CHA	HOU
Away	Away	TOR	TOR	CHA
Away	Home	NYK	LAC	TOR

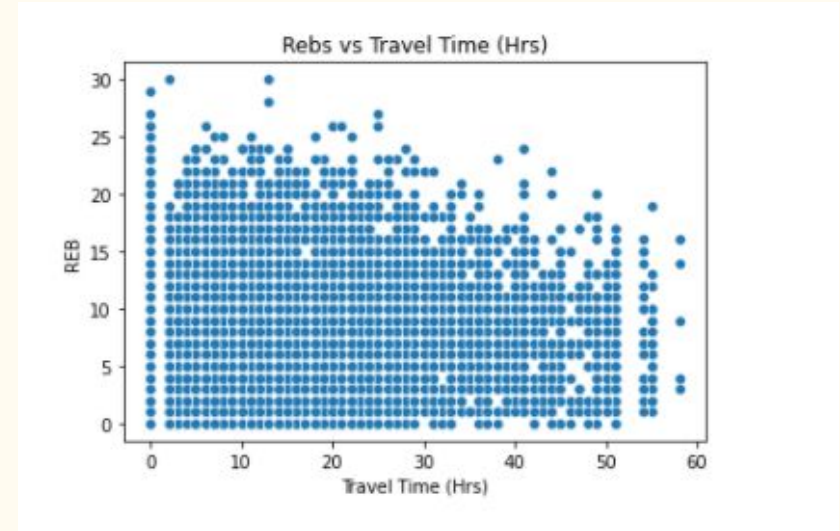
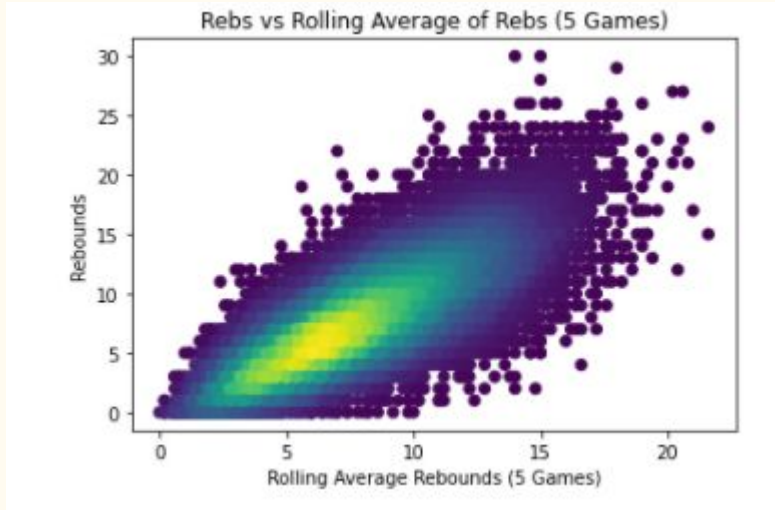
Start	End	Duration (mins)
MIA	POR	3459
POR	MIA	3457
GSW	BOS	3311
BOS	GSW	3307
POR	BOS	3302



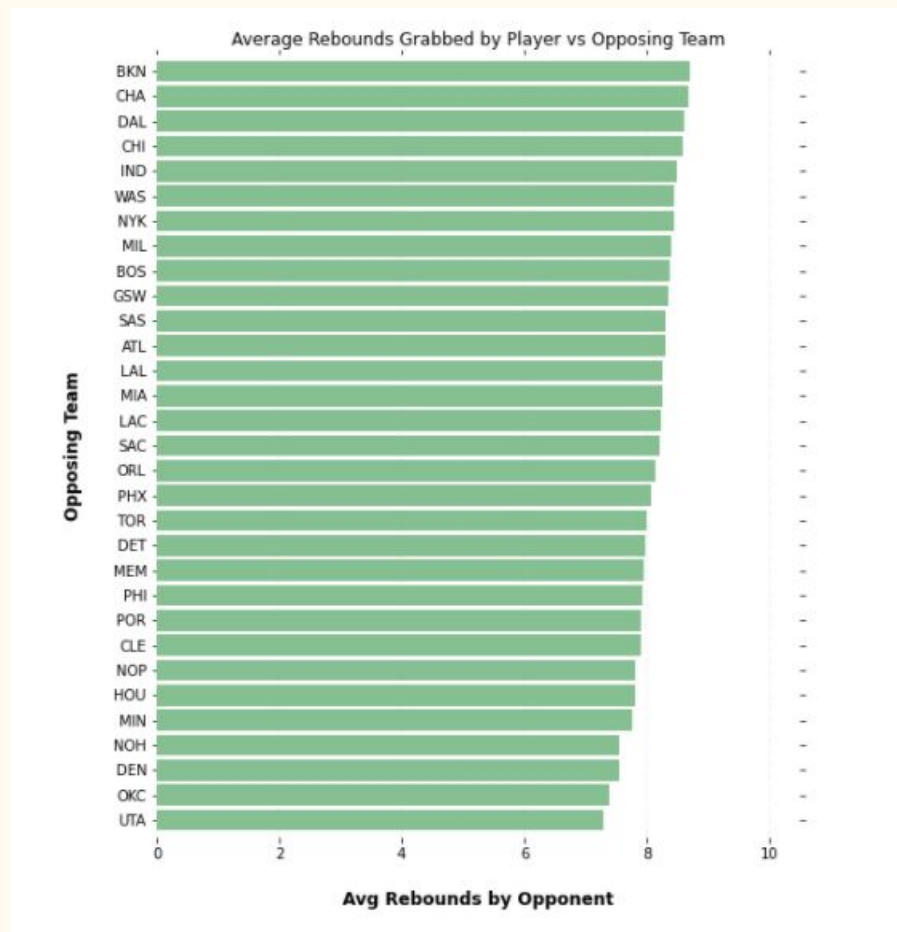
Rebounds Distribution



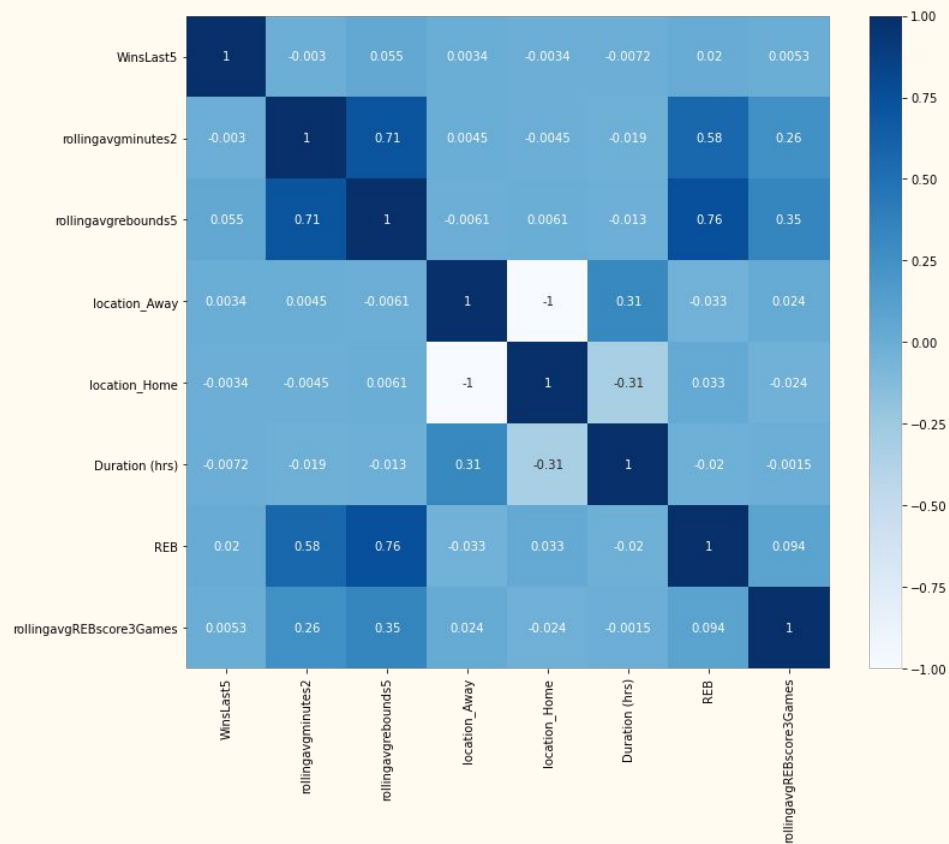
Rebounds vs Rolling Average Rebounds & vs Travel Time



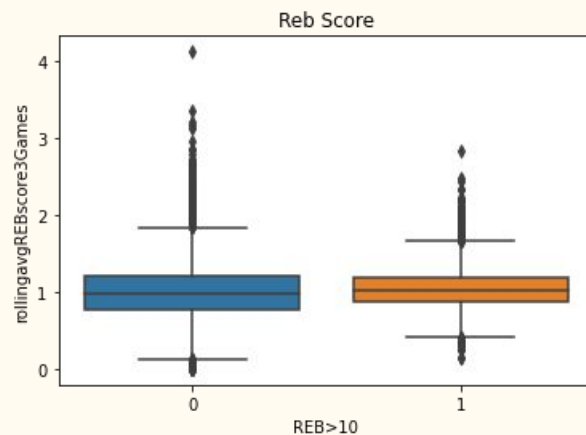
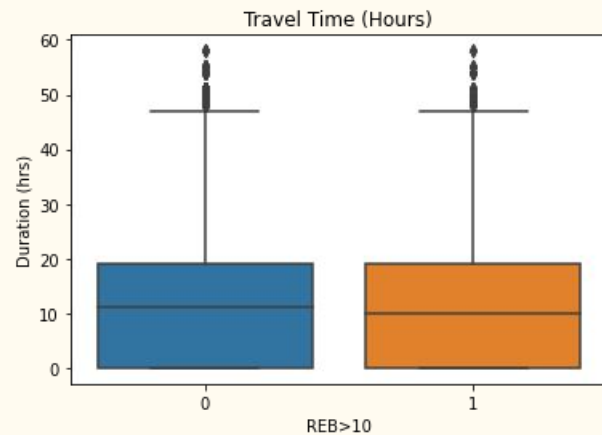
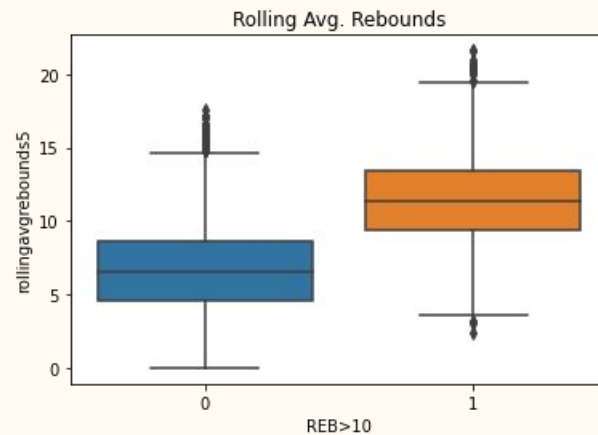
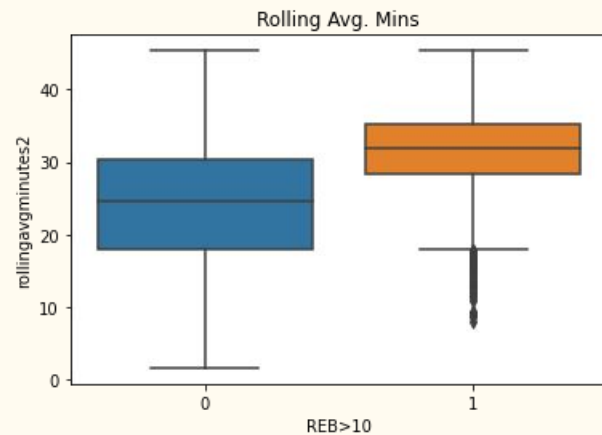
Average Player Rebounds by Opponent



Feature Correlation



Feature Target Analysis



Feature Importance (RF Classifier)

```
'WinsLast5': 0.050758547749473554,  
'rollingavgminutes2': 0.17884414268054308,  
'rollingavgrebounds5': 0.3788050152290544,  
'OPP_ATL': 0.004329510366784768,  
'OPP_BKN': 0.004164819770918081,  
'OPP_BOS': 0.004982057323514048,  
'OPP_CHA': 0.005097794617651523,  
'OPP_CHI': 0.0049598639777465765,  
'OPP_CLE': 0.004932042569220045,  
'OPP_DAL': 0.005010166652688406,  
'OPP_DEN': 0.004459970051046029,  
'OPP_DET': 0.005621476298267344,  
'OPP_GSW': 0.004813936077680448,  
'OPP_HOU': 0.004991574361358178,  
'OPP_IND': 0.004459989762137683,  
'OPP_LAC': 0.005445253174695959,  
'OPP_LAL': 0.005305021186284147,  
'OPP_MEM': 0.004614986137080766,  
'OPP_MIA': 0.004576440615629907,  
'OPP_MIL': 0.0044083132059503614,  
'OPP_MIN': 0.005070994549088146,  
'OPP_NOH': 0.00036011158931877735,  
'OPP_NOP': 0.004904014673263913,  
'OPP_NYK': 0.004649399751365807,  
'OPP_OKC': 0.004486830406000347,  
'OPP_ORL': 0.005110949092651474,  
'OPP_PHI': 0.004820935894873621,  
'OPP_PHX': 0.005306320080746615,  
'OPP_POR': 0.004710643984307257,  
'OPP_SAC': 0.005127055919624439,  
'OPP_SAS': 0.005405036328956066,  
'OPP_TOR': 0.00467677519648686,  
'OPP_UTA': 0.004348827649502595,  
'OPP_WAS': 0.004698840051536204,  
'location_Away': 0.00649726837373837,  
'location_Home': 0.006214753874765955,  
'Duration (hrs)': 0.06870002340418552,  
'PreviousHome/Away_Home': 0.009368823598495439,  
'PreviousHome/Away_Away': 0.009395654109252465,  
'rollingavgREBscore3Games': 0.14556581966411472}
```

Note: Duration (hrs) *not* statistically significant (p value = 0.7)

Modeling

Gridsearch CV was used on all models to find optimal parameters

Models used:

- K-Nearest Neighbors

```
n_neighbors: 24  
0.8700222328206803
```

- Logistic Regression

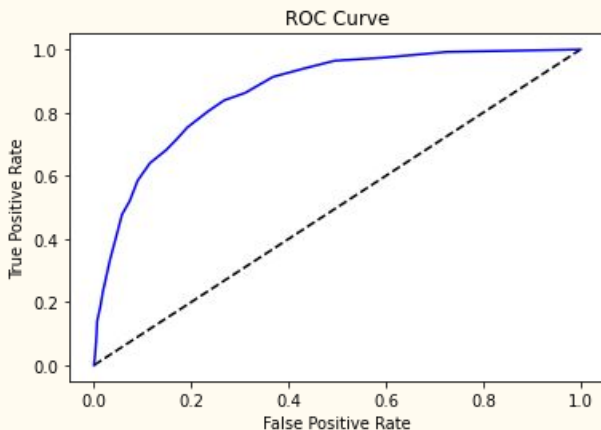
```
Tuned Logistic Regression Parameters: {'C': 0.05179474679231213}  
Best Score: 0.8917851075458696
```

- Random Forest

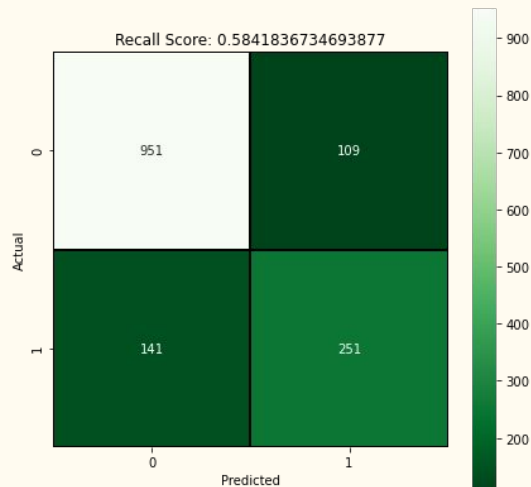
```
{'criterion': 'entropy',  
 'max_depth': 8,  
 'max_features': 'auto',  
 'n_estimators': 500}
```

Modeling - K-Nearest Neighbors

Gridsearch CV was used on all models to find optimal parameters



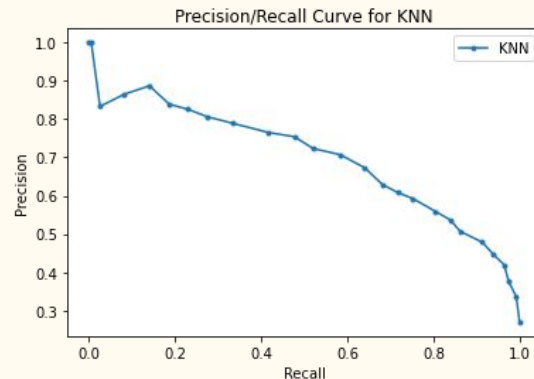
$$ROC_AUC = 0.86$$



$$fbeta = 0.82$$

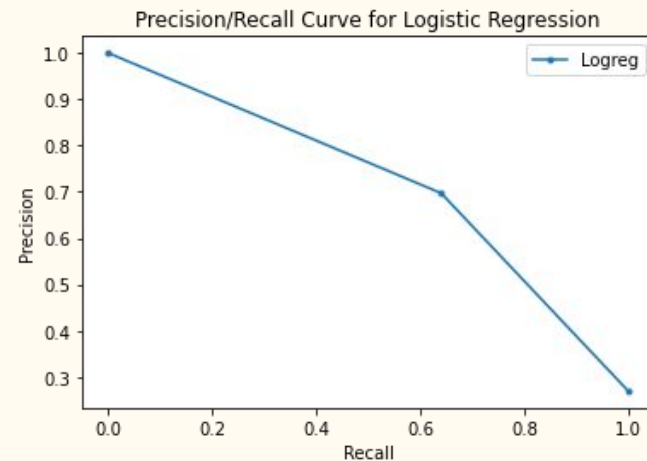
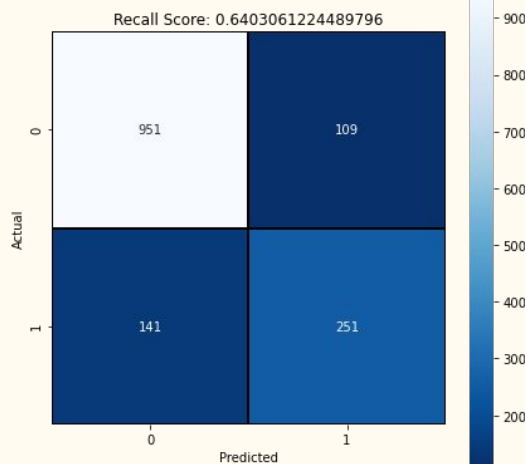
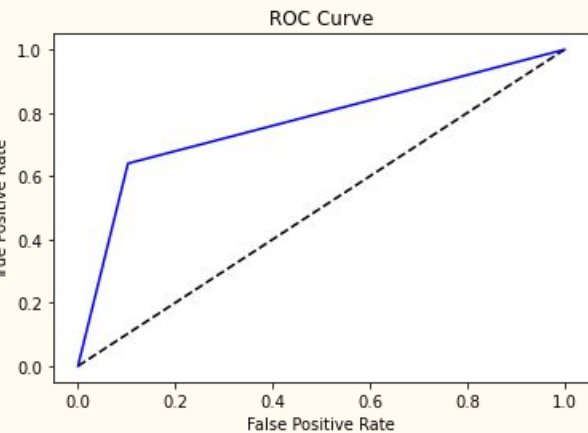
Classification Report

	precision	recall	f1-score	support
0	0.86	0.91	0.88	1060
1	0.71	0.58	0.64	392
accuracy			0.82	1452
macro avg	0.78	0.75	0.76	1452
weighted avg	0.82	0.82	0.82	1452



95.0 confidence interval Recall: 54.1% and 58.4%
95.0 confidence interval Precision: 70.4% and 74.6%

Modeling - Logistic Regression

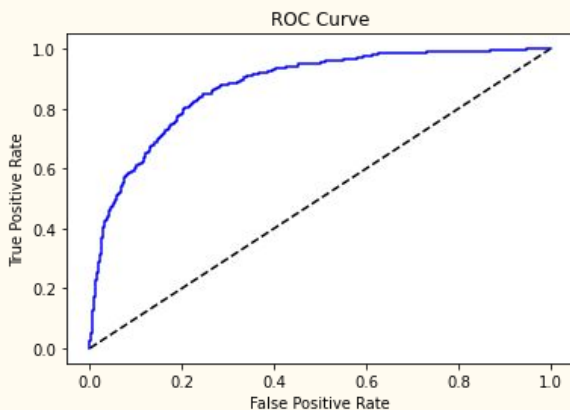


$$ROC_AUC = 0.77$$

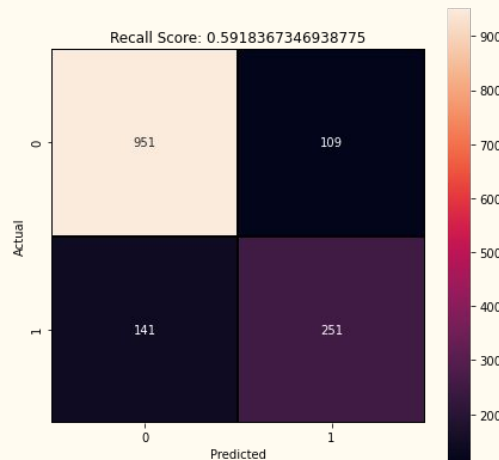
$$fbeta = 0.83$$

95.0 confidence interval Recall: 59.3% and 63.7%
95.0 confidence interval Precision: 71.7% and 75.6%

Modeling - Random Forest

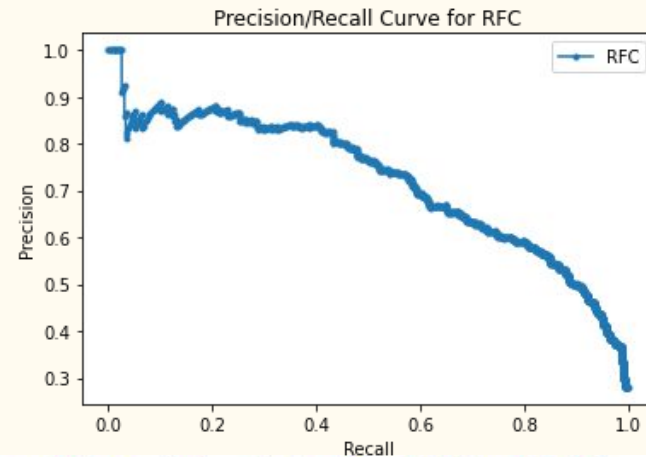


$$ROC_AUC = 0.87$$



$$fbeta = 0.82$$

Classification Report				
	precision	recall	f1-score	support
0	0.86	0.91	0.88	1060
1	0.70	0.59	0.64	392
accuracy			0.82	1452
macro avg	0.78	0.75	0.76	1452
weighted avg	0.82	0.82	0.82	1452



95.0 confidence interval Recall: 54.1% and 58.4%
95.0 confidence interval Precision: 70.4% and 74.6%

Conclusion

Mario Hage

Takeaways

- Random Forest Model is the optimal choice between the 3 models. All 3 struggle at predicting positive classifications, however the RF Model can accurately predict Negative classifications (91%)
- Model can be utilized to bet under scenarios
- EDA portion provided useful notifications/alerts

Looking Forward

- Spend time collecting data/quantifying more features
- Integrate betting odds
- Test and optimize