



Predicting Board Game Success

Dan Feldman - 9/4/2018

Goals

- ◆ *Identify trends to gain insight into the current board game market.*
- ◆ *Build a predictor for board game success.*



Potential Clients



- ◆ *Board game creators: how can they create better games with more effective market strategy?*
- ◆ *Board game sellers/stores: which games should they sell and emphasize to their buyers?*

Table Toys shop in Indonesia (<https://www.tokopedia.com/tabletoys>)

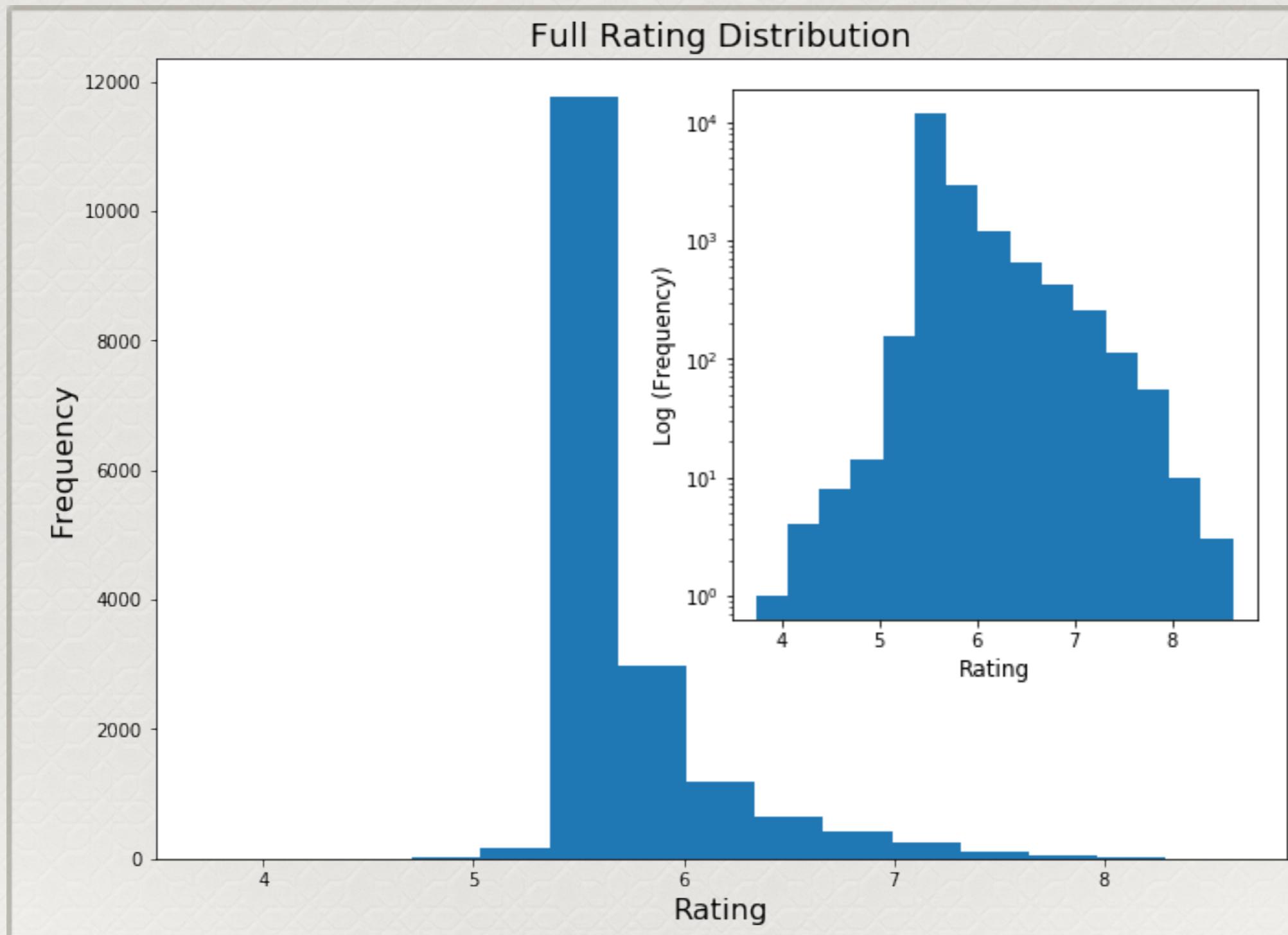
The Dataset: Board Game Geek

- ◆ *boardgamegeek.com (BGG)*
- ◆ *17,625 total board games*
- ◆ *Likely selection bias against casual gamers*
- ◆ *Data wrangled from XML API*

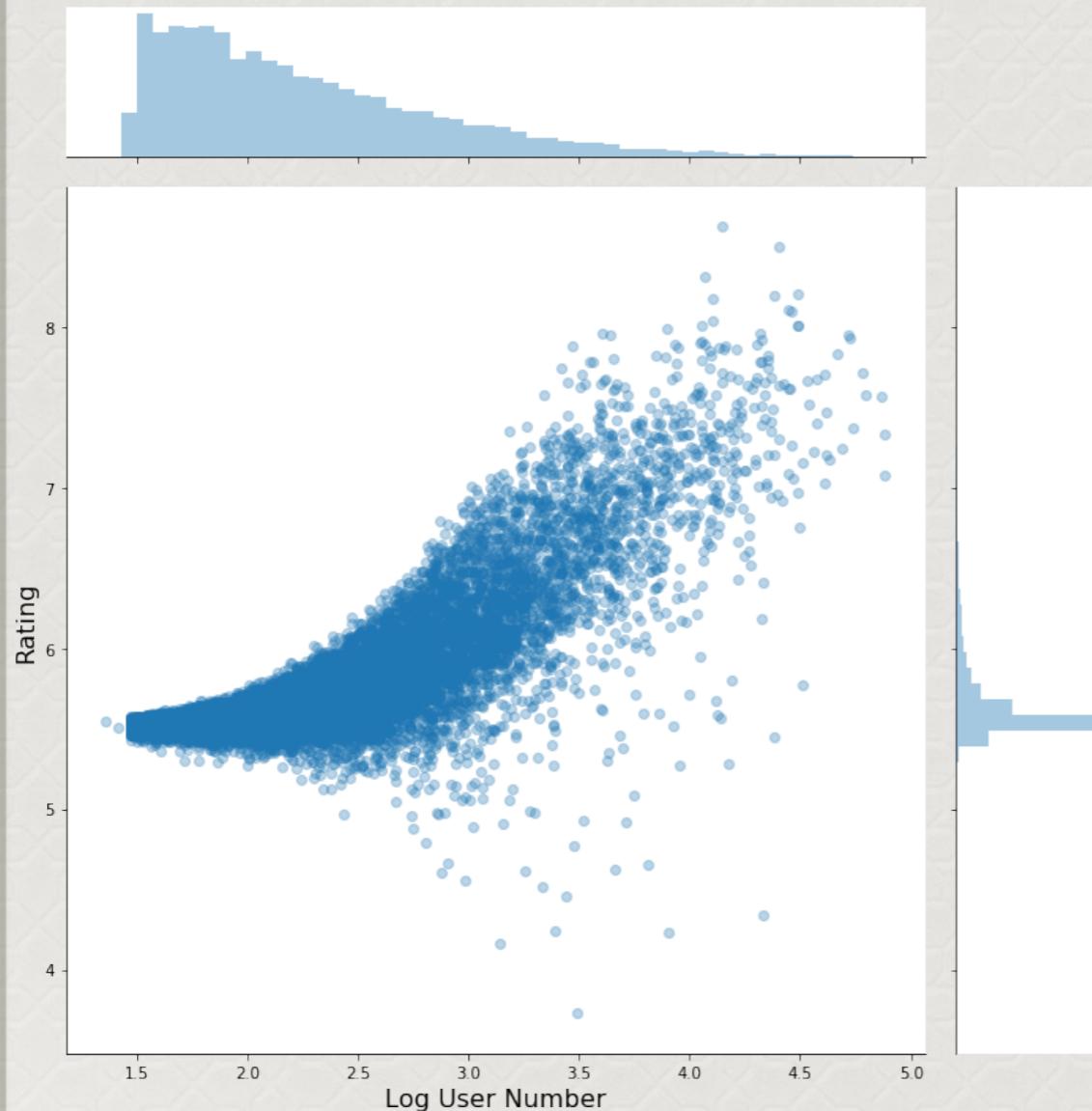


BGG Logo

The Target Variable: Ratings

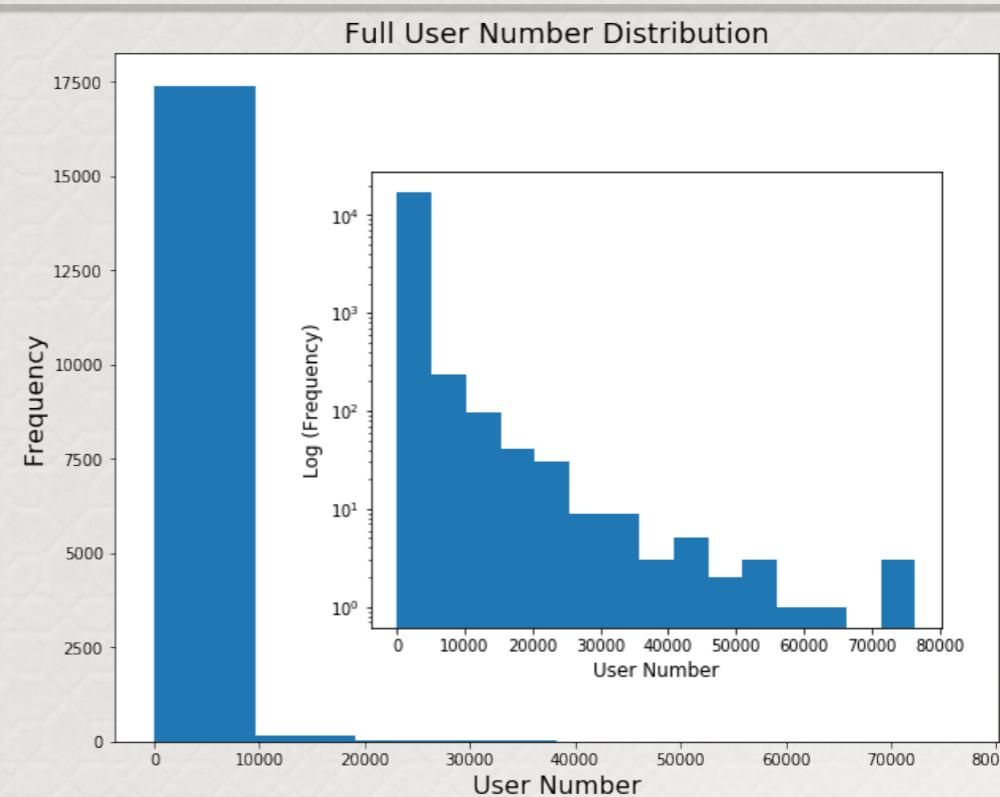
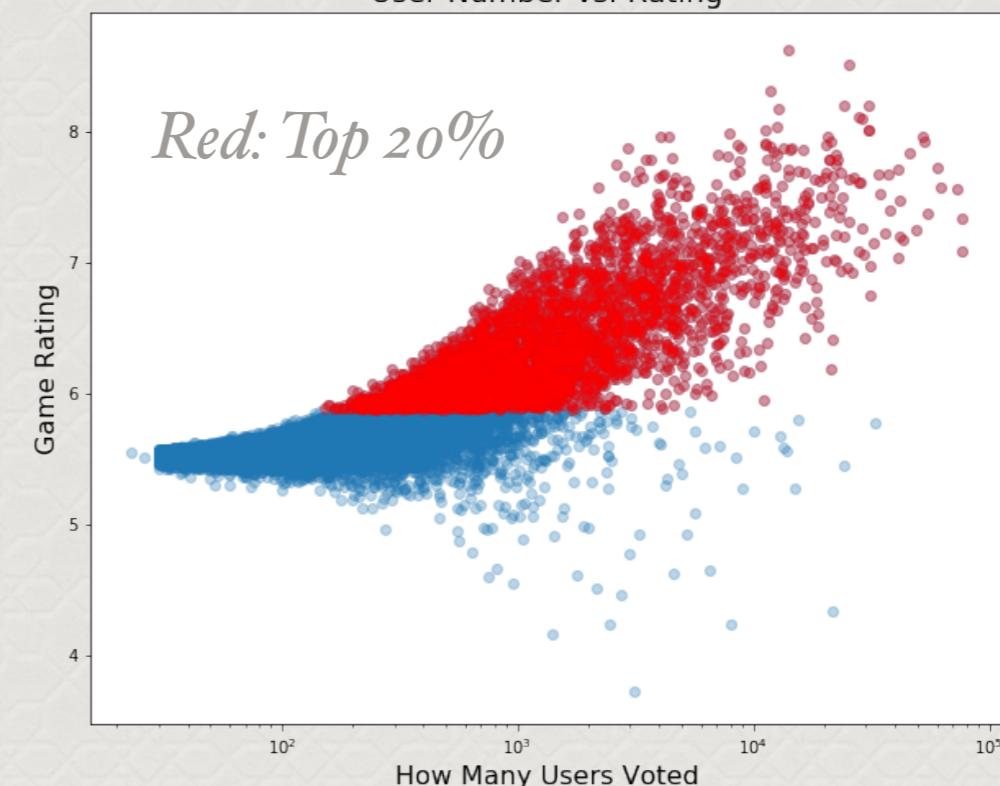


User Number vs. Rating



*Highly peaked but
strongly correlated*

User Number vs. Rating



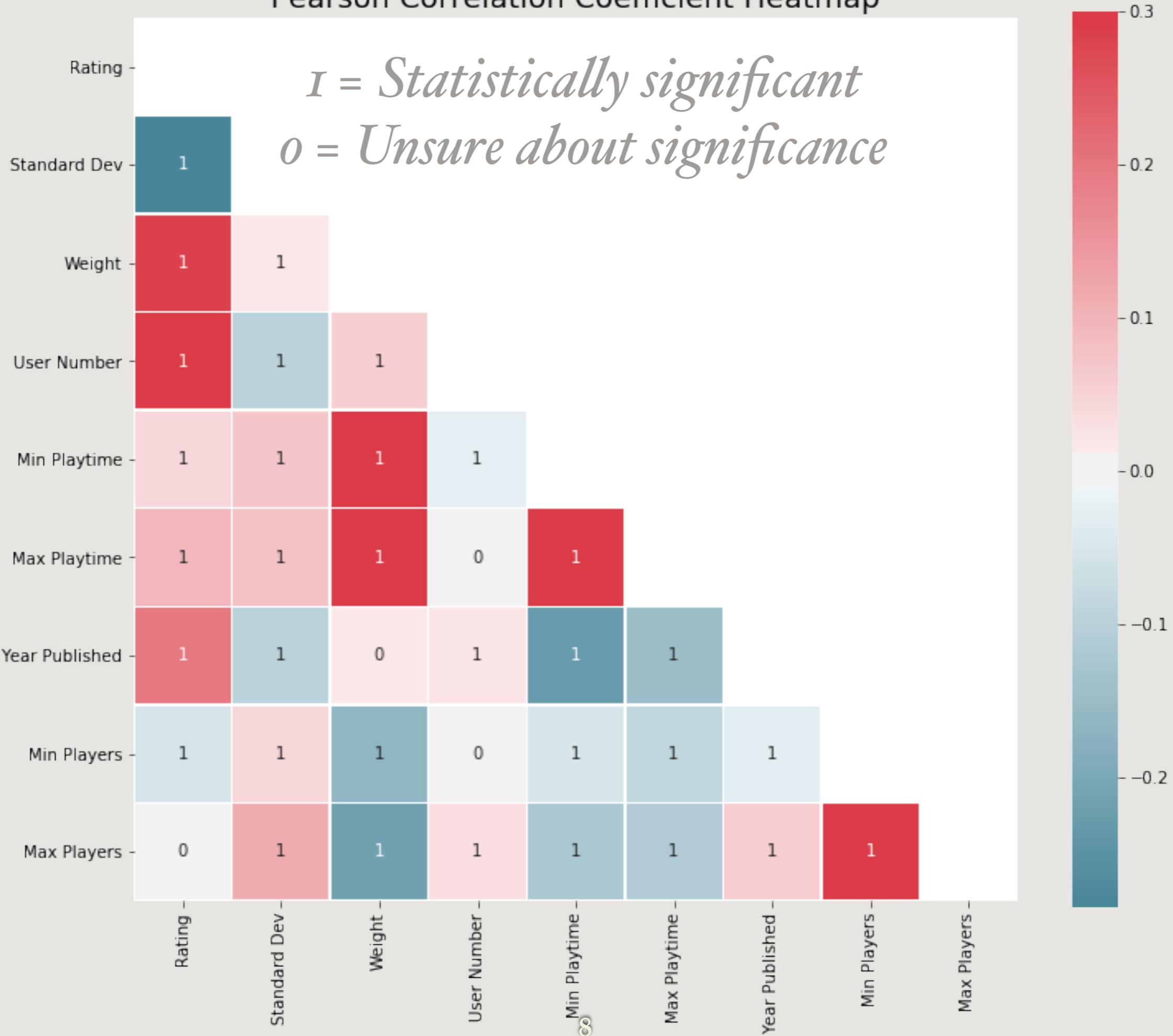
*Scatter
plots
at a
glance*

Rating



Pearson Correlation Coefficient Heatmap

I = Statistically significant
o = Unsure about significance



Modeling: Binary Classification

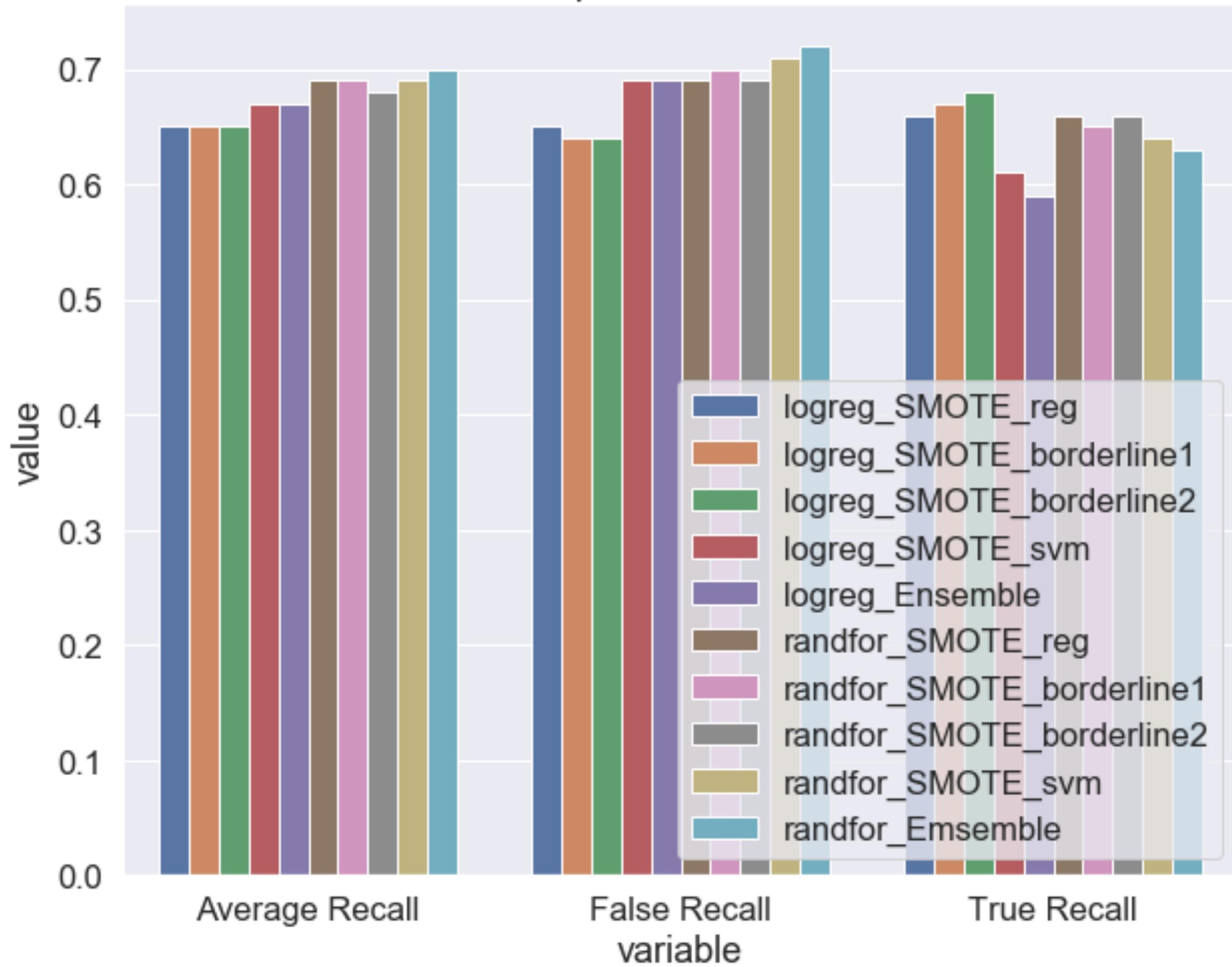
- ◆ *To make prediction easier, changed from continuous target variable (ratings) to binary classification*
- ◆ *If true: successful game, if false: unsuccessful game*
- ◆ *Chose a cutoff of the 80th percentile (rating ~ 5.88 or above considered successful)*
- ◆ *Began with Logistic Regression classifier*

Initial Modeling: Imbalanced Dataset

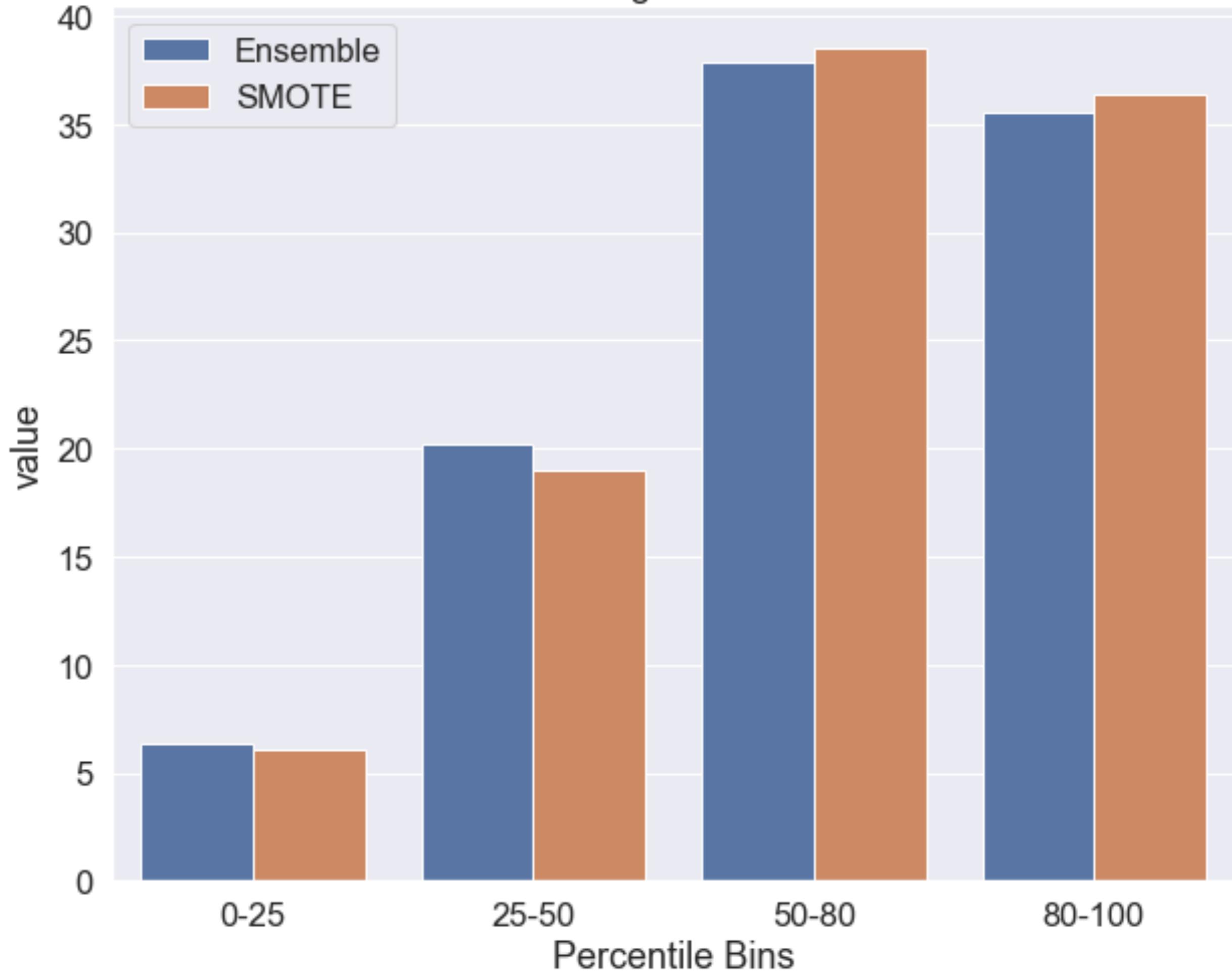
[Test Classification Report:]				
	precision	recall	f1-score	support
False	0.80	1.00	0.89	4238
True	0.35	0.01	0.02	1050
avg / total	0.71	0.80	0.72	5288

- ◆ *Because the top 20% of games are successful by definition, the ratio of false to true is 4:1*
- ◆ *Classifier maximizes accuracy, classifying all games as unsuccessful (false)*
- ◆ *Fix Imbalance using over-sampling with SMOTE, and under-sampling with random ensemble*

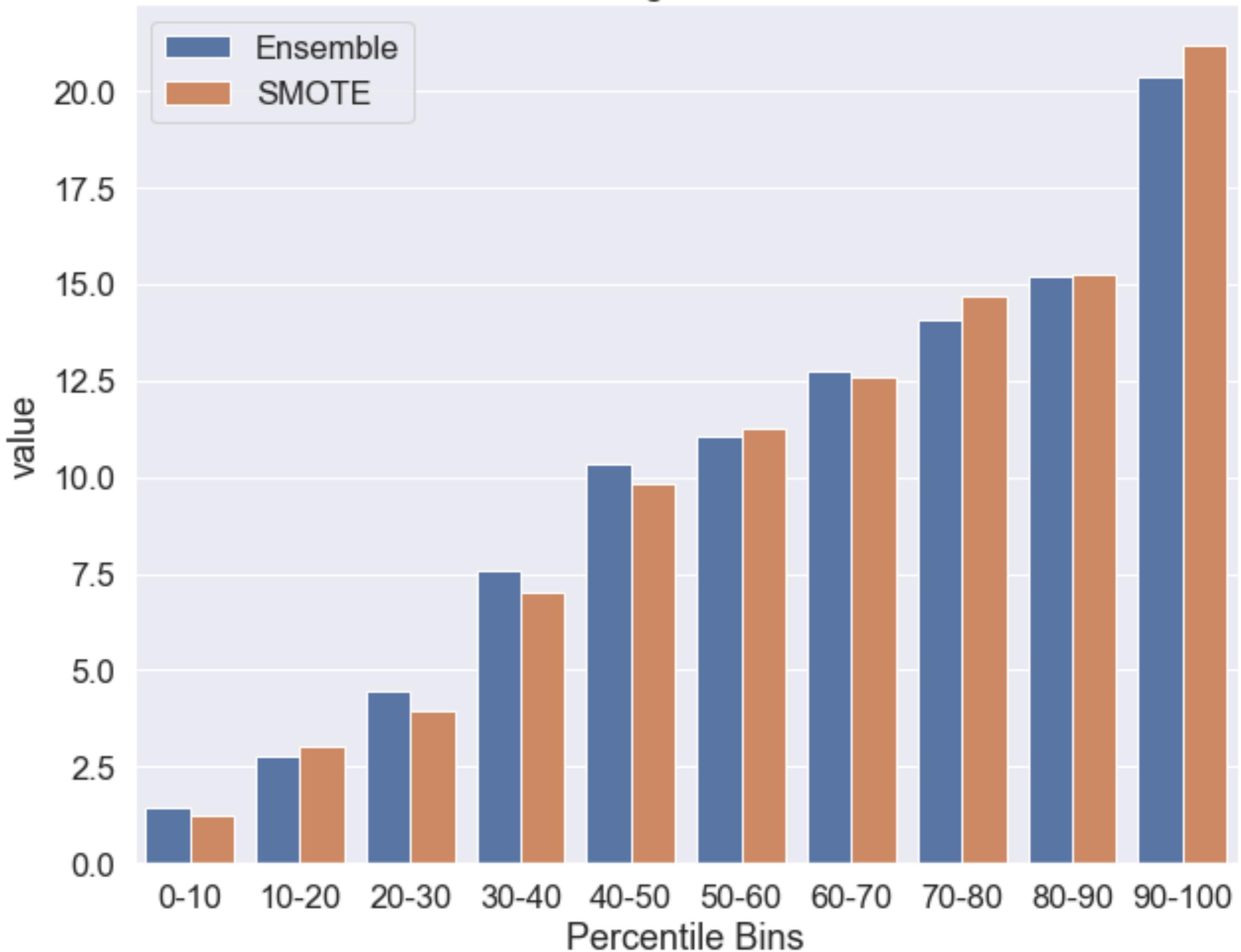
Comparison of Metrics



Distribution of Ratings for Predicted Success



Distribution of Ratings for Predicted Success



Findings/Results Summary

- ◆ *Ratings are correlated with board game complexity, number of users voting (popularity proxy?), year published; anti-correlated with standard deviation (controversial proxy?)*
- ◆ *Random forest performed better than logistic regression, logistic regression coefficients reinforce EDA correlations*
- ◆ *Over-sampling with SMOTE and under-sampling ensemble techniques roughly same effectiveness, SMOTE slightly better for both computation time and profitability/risk aversion*
- ◆ *Modeling can predict strong board game success approximately 69% of the time*