

Above or Below a Certain Threshold?



Introduction

Being an avid NBA Fan, I wanted to see if I could build a model that could accurately predict whether or not a certain player in the NBA will achieve 8 or more rebounds in a game. The reason I did not want to try and predict the actual number of rebounds that a player might get is just because there are too many factors that could hinder the model; and the fact that I have friends who have a hobby of sports betting, and this would fit more in line with someone who wanted to use the model for those purposes.

Steps taken

I used the nba_api to pull data on any player of choice. (model accuracies vary by player,) as some things that cannot be tracked may affect the workings of the models I selected to use. Some examples of features that couldn't be implemented:

-
- Player Injury Tracking
 - Team Swaps and their its affects
 - Player rest
 - "Rhythm"
 - Mood

From there, I cleaned and transformed the data to separate Home vs Away, and get dummies for all teams. I also plotted several visualizations to help me better understand feature weights, and the data itself better.

Conclusion

Though accuracy scores did not exceed 0.635 (Random Forest Model) It is still enough to consistently beat odd spreads. Also, the models were better at predicting when a player *wouldn't* achieve 8 rebounds or more; as well as being more accurate with center players as they are more consistently going to be grabbing rebounds purely based on their respective position as players.