

# Predicting MLB Games - Midterm Report

Brian Won (bw436) and Jiacheng Wu (jw2545)

November 7, 2019

## **Our Data**

Our current dataset only consists of 1615 games from the 2018 season, with each row representing a game. For each game, we have hitting and pitching stats for the two teams. We decided to use stats on a rolling window basis; for example, for hitting, we only look at how teams have been hitting over their previous 10 games. Our reasoning is that recent performance and trends are more predictive than past season stats for predicting the outcome of a single game. As of now, the only pitching info we have are the rolling window stats of the two starting pitchers of a game.

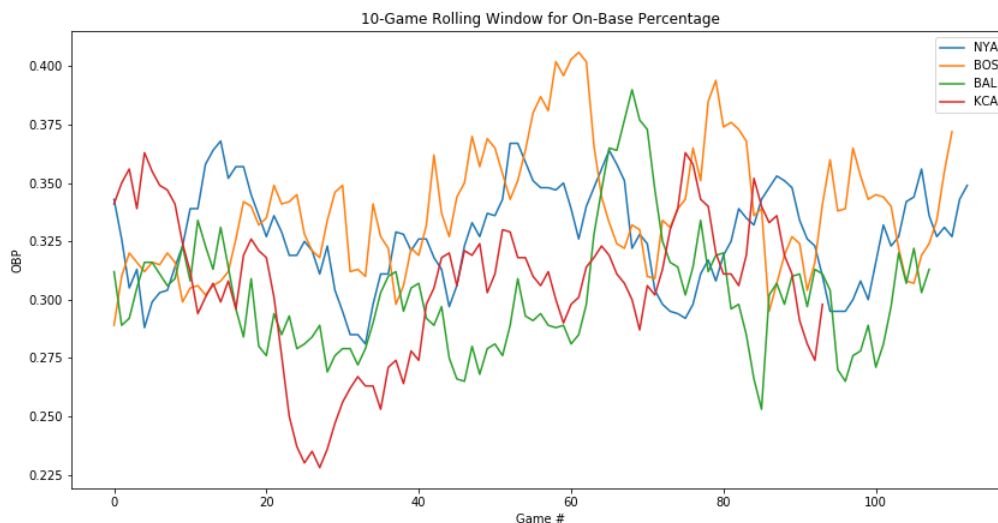
## **Initial Analysis**

### **Class Balance**

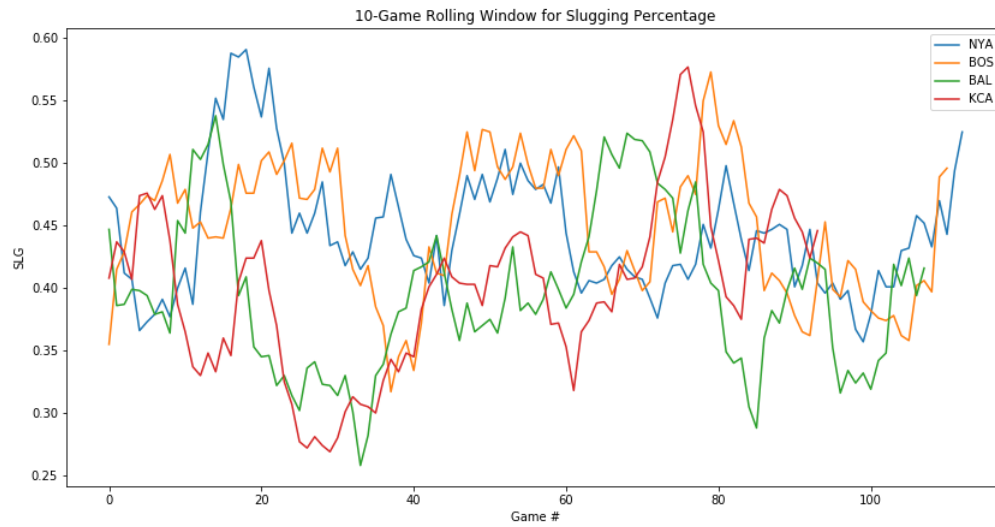
Of the 1615 games, the home team won about 53% of the time. It makes sense that home teams tend to win slightly more and it's great that we don't have to deal with much class imbalance for this classification problem.

## Plots

Below is a plot for rolling on-base percentages (OBP) for a few teams:



Here we've plotted rolling OBP for the Yankees, Red Sox, Orioles, and Royals. OBP is simple: it's the number of times a hitter reaches base safely divided by total number of plate appearances. OBP is a "good" stat because it represents efforts towards a goal in baseball and only accounts for the skill of the featured player or team. The goal in baseball is to score runs. Teams score runs by getting people on base. The plot above shows that OBP can be used to differentiate good and bad teams. We can see that for most of the 2018 season, the Yankees and Red Sox (really good teams) had better rolling OBP than those of the Orioles and Royals (really bad teams). The "good" stat criteria is the reason that we didn't include traditional stats such as runs batted in or RBI. Granted, RBI directly represents scoring runs, but it doesn't only account for the skill of the featured hitter. For example, a hitter can have inflated RBI numbers if his teammates who hit ahead of him in the lineup get on base often, or if they're fast runners so they can score from second on a mere single.



The plot above shows that slugging percentage (SLG) can also be a good stat to differentiate teams.

## Moving Forward