



MLB Pitch Prediction



(Sprint 2)

Ty Painter, Connor Mignone, Dong Chen,
Daron Greenblatt

Objectives

- Create a model and eventually dashboard that will help the offensive team (the hitting team) in pitch anticipation during an at-bat.
- Help predict pitch type based on batter, pitcher, runners on base, count during the at-bat, pitch sequence, etc.
- This dashboard would help teams predict what will be thrown therefore hopefully increasing the chances the batter will be successful.

Data Summary

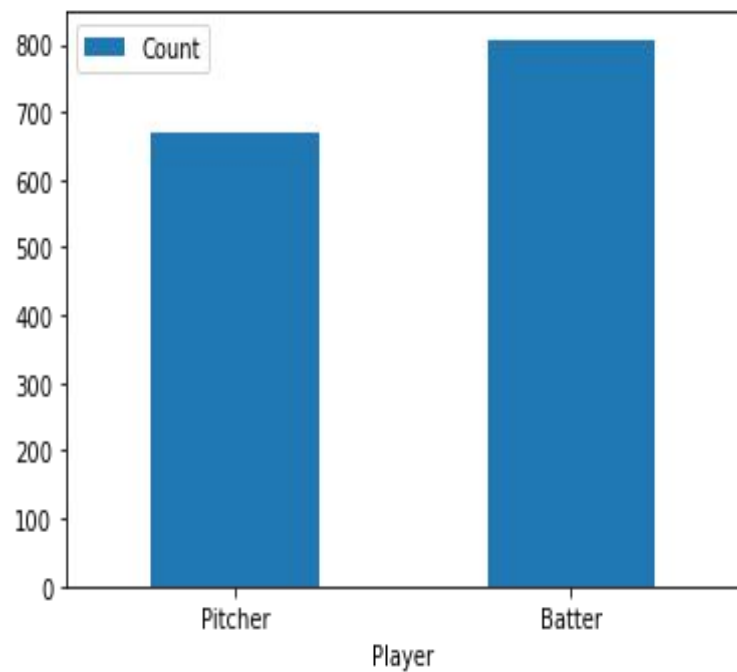
- Every pitch thrown during the Major League Baseball (MLB) 2019 season
- 8 total datasets, 1.06GB
- Merged 3 datasets together
 - Pitch data
 - At-Bat data
 - Player data
- Merged dataset
 - ~600K rows
 - 54 columns
 - Too big for GitHub - revised down to 25 columns

Data Cleaning

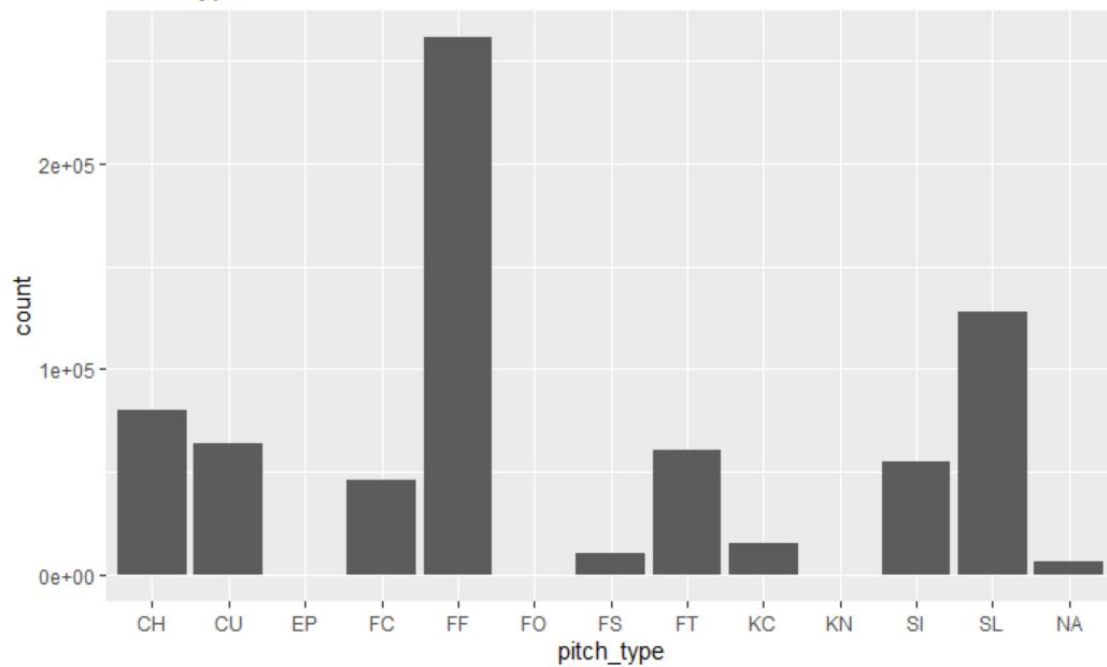
Seeing inaccurate or incomplete data:

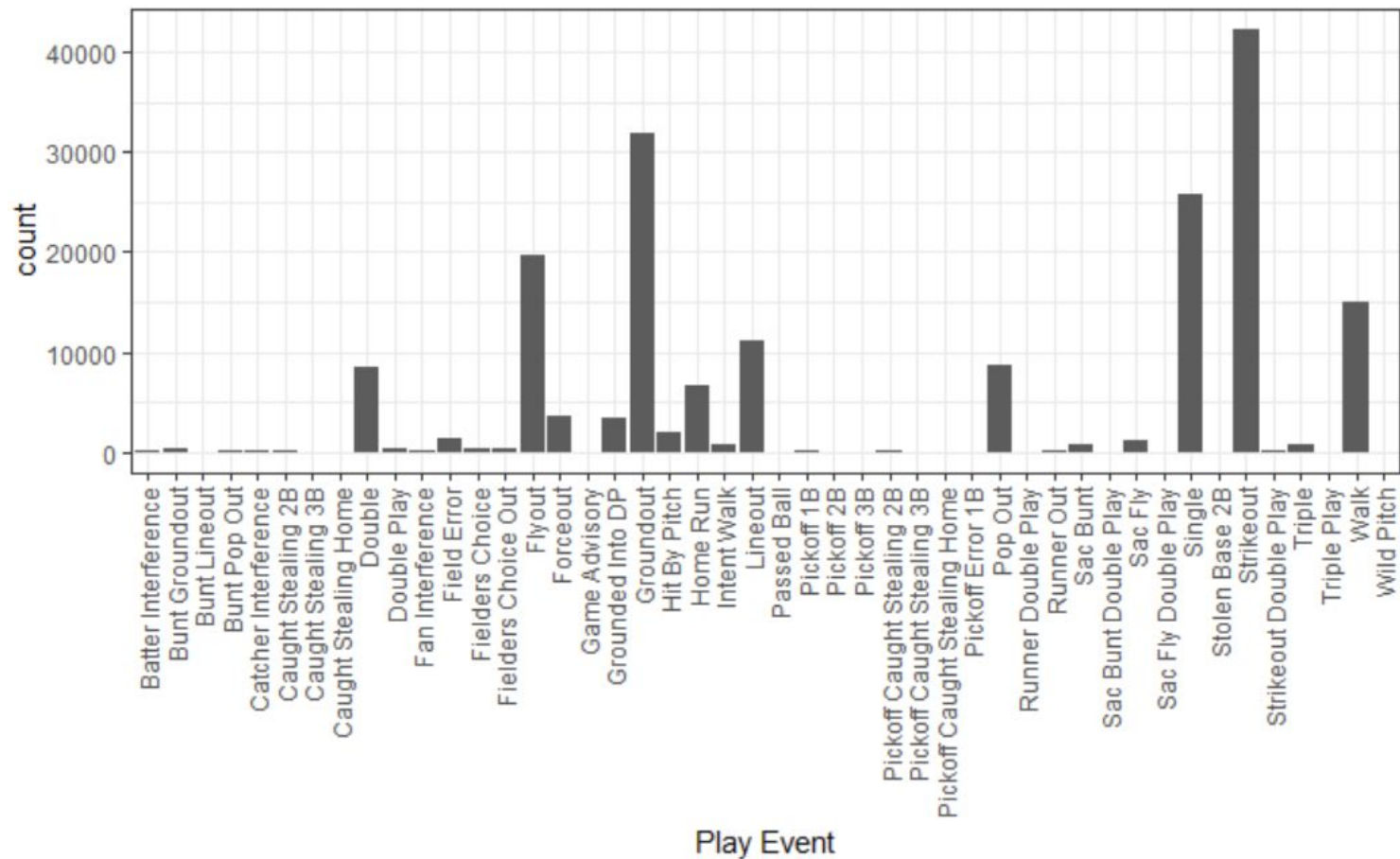
- Pitch count often not changing after a pitch is thrown - will be difficult to include this important feature in our model if it is unreliable
- Many other columns have questionable data but we will likely use filler values such as mean or mode

Batter and Pitcher Counts



Pitch Type





Next Steps

- Feature engineering
 - Batter
 - Pitcher
 - Count
 - Base runners
 - Score
- More EDA
- Split train/test data
 - Eliminate pitchers/hitter with below a threshold of pitches thrown