# How External Factors Influence the Flow of Baseball Games
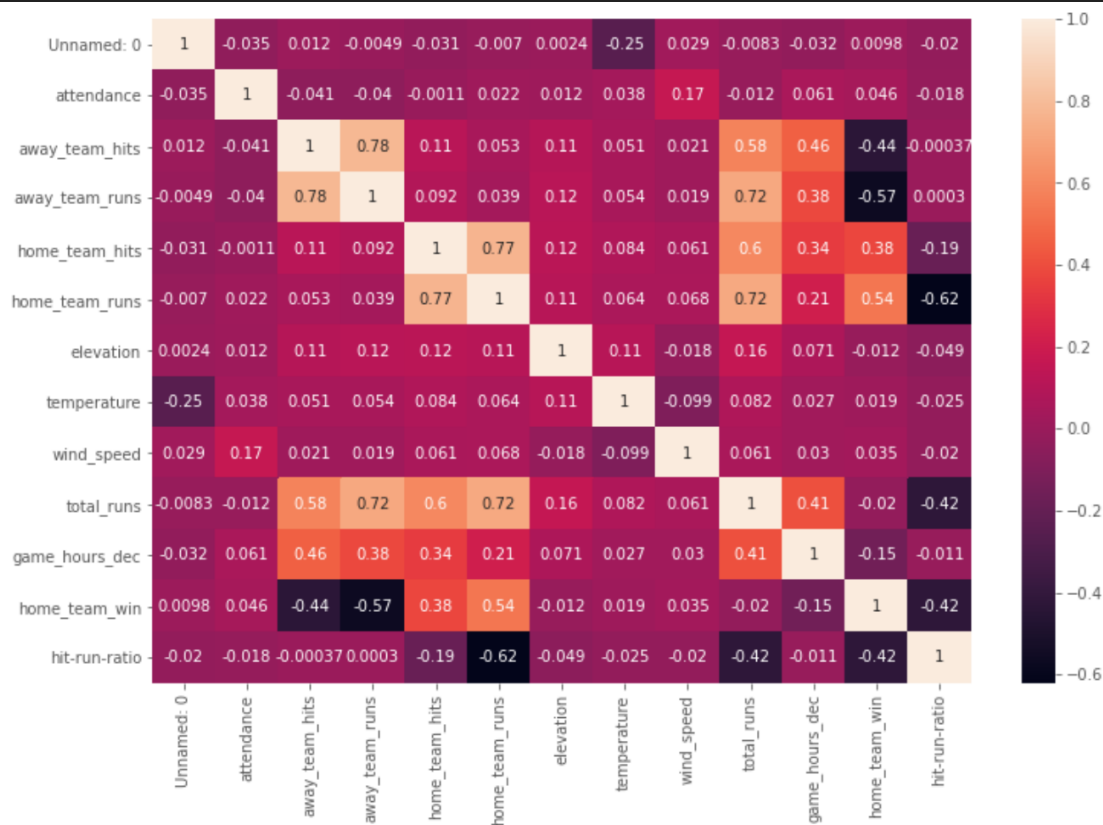
Josh De Leeuw, Michael Raybetz, and Dezmon Nash

# Tools and Technologies

- We received our data from Kaggle (and Baseball Reference), and used various Python libraries (Pandas, Numpy, Matplotlib) to produce the visualizations shown in the following slides

- The data was organized and the graphs were made in a Jupyter Notebook
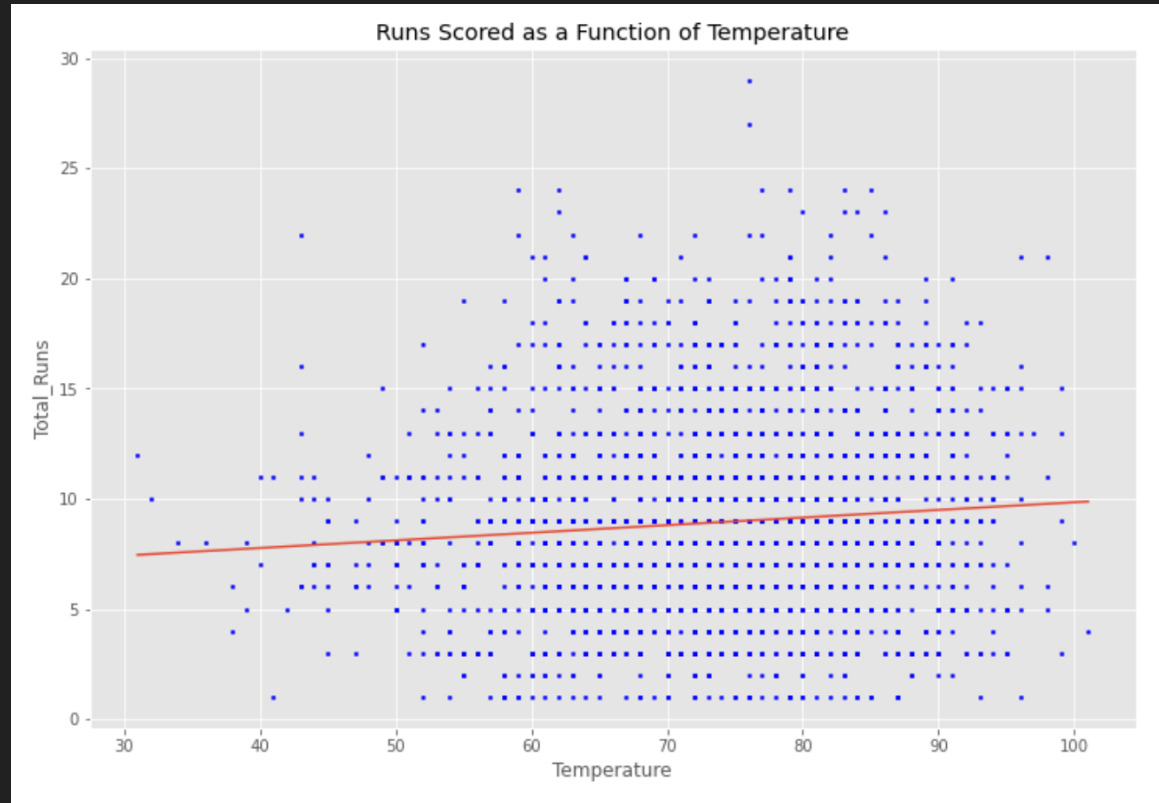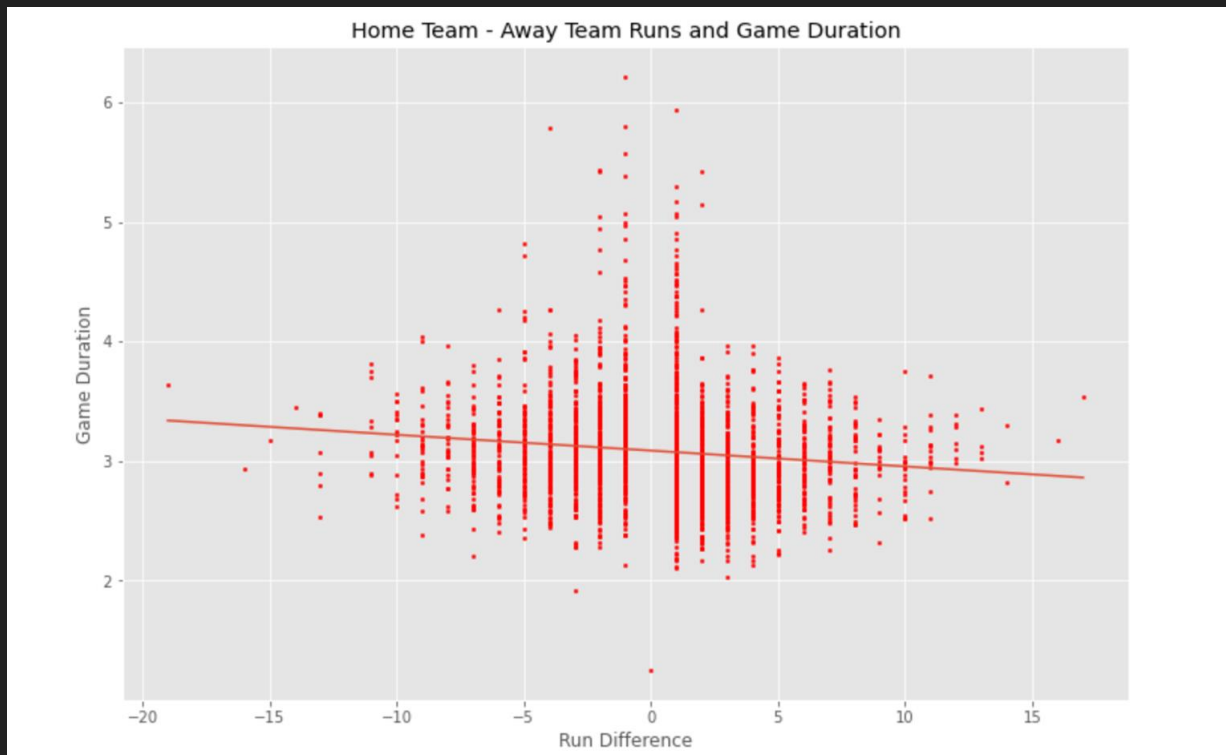
# Correlation Matrix

# Interesting Takes from the Matrix

- Most factors regarding external factors (e.g. temperature, wind speed, crowd size, etc…) were not particularly correlated
  - This does not include things like "home team runs" vs "home team hits", which of course have strong dependencies
- Nevertheless, there were some interesting trends (and non-trends), that we will take you through in the following slides
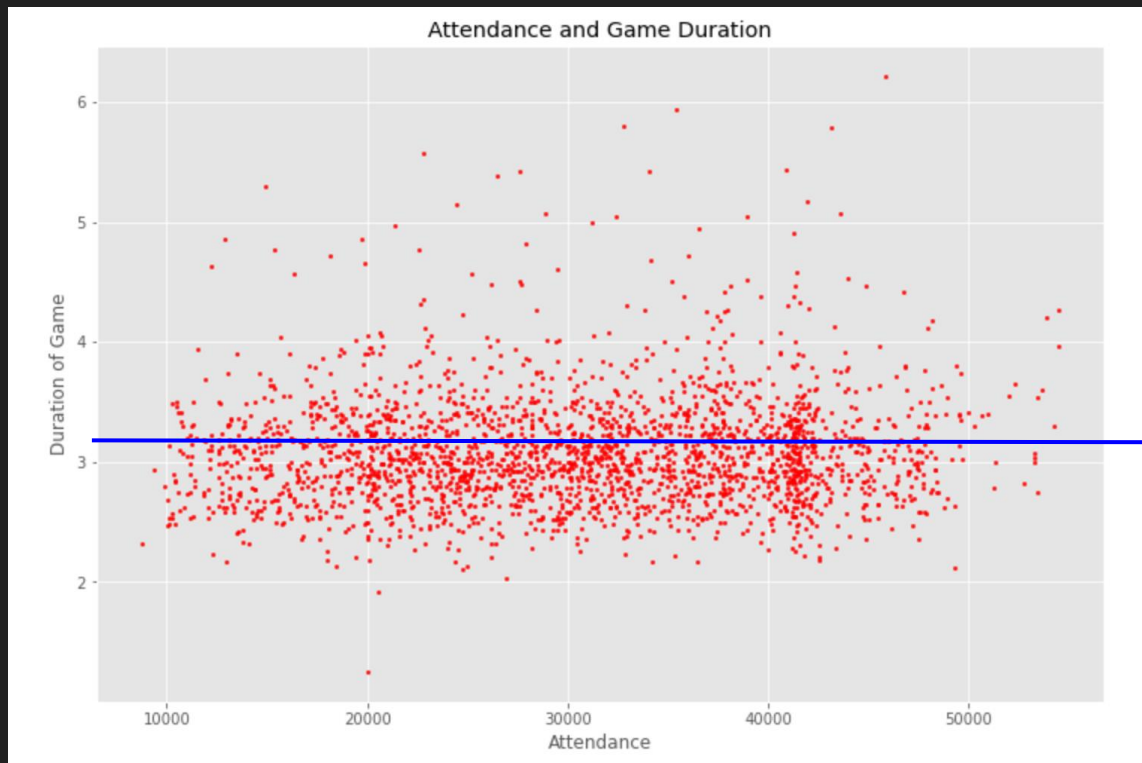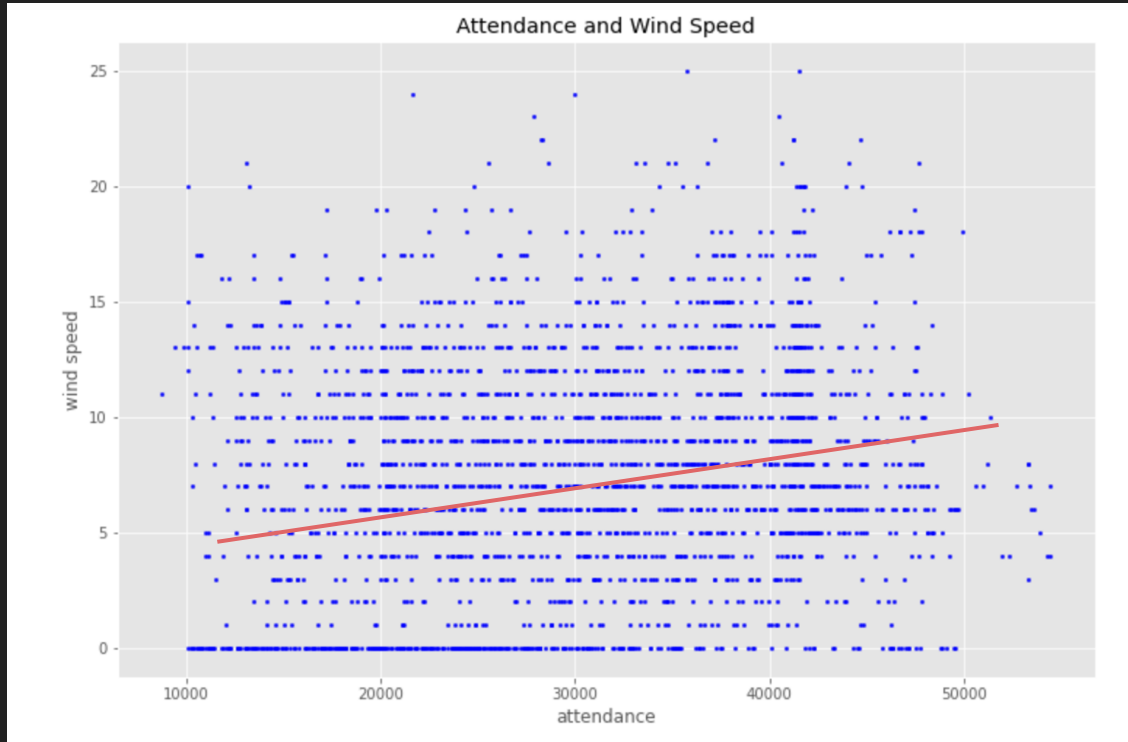
# Runs Scored vs. Temperature

# How Run Difference Affects the Length of Game

# Do more attended games last longer?

# Surprising Correlation: Attendance vs Wind Speed

# Conclusions From Graphs

- Even when correlations were present, it can be hard to prove causation
  - Either factor could be causing the other, or a separate factor may be causing both of them
  - My favorite statistical example: ice cream consumed per month and drownings per month are correlated, even though eating ice cream clearly does not cause one to drown
    - Hot weather causes both of them!
- However, some correlations we found were runs scored vs. temperature, attendance vs. wind speed, and run difference vs. the length of game
- Further analysis could involve working with different Kaggle datasets, or trying to dig further into this one