

# Analyzing Baseball Fans and Predicting Game Attendance

---

Troy Hepper

DC-DSI 4

May 25, 2017

# AGENDA

- Introduction
- Exploring the Data
- Building Models
- Bayesian Inference
- Further Analysis

# PROBLEM STATEMENT

Predicting Major League  
Baseball game attendance  
based on historical data.

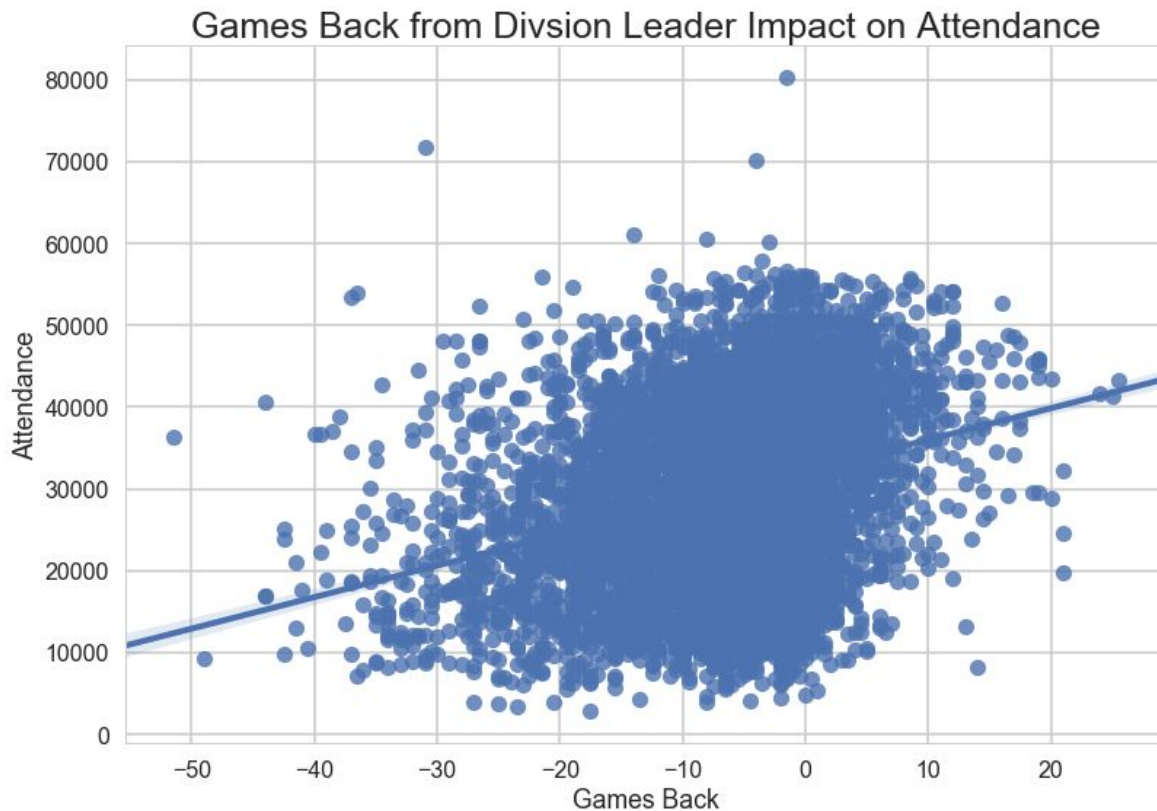
---

# COLLECTING DATA

- Collected game data from [baseball-reference.com](http://baseball-reference.com) for all MLB games from 1990 to 2016
- Engineered features that I thought might have an affect on attendance
- Also collected season summary information for each team

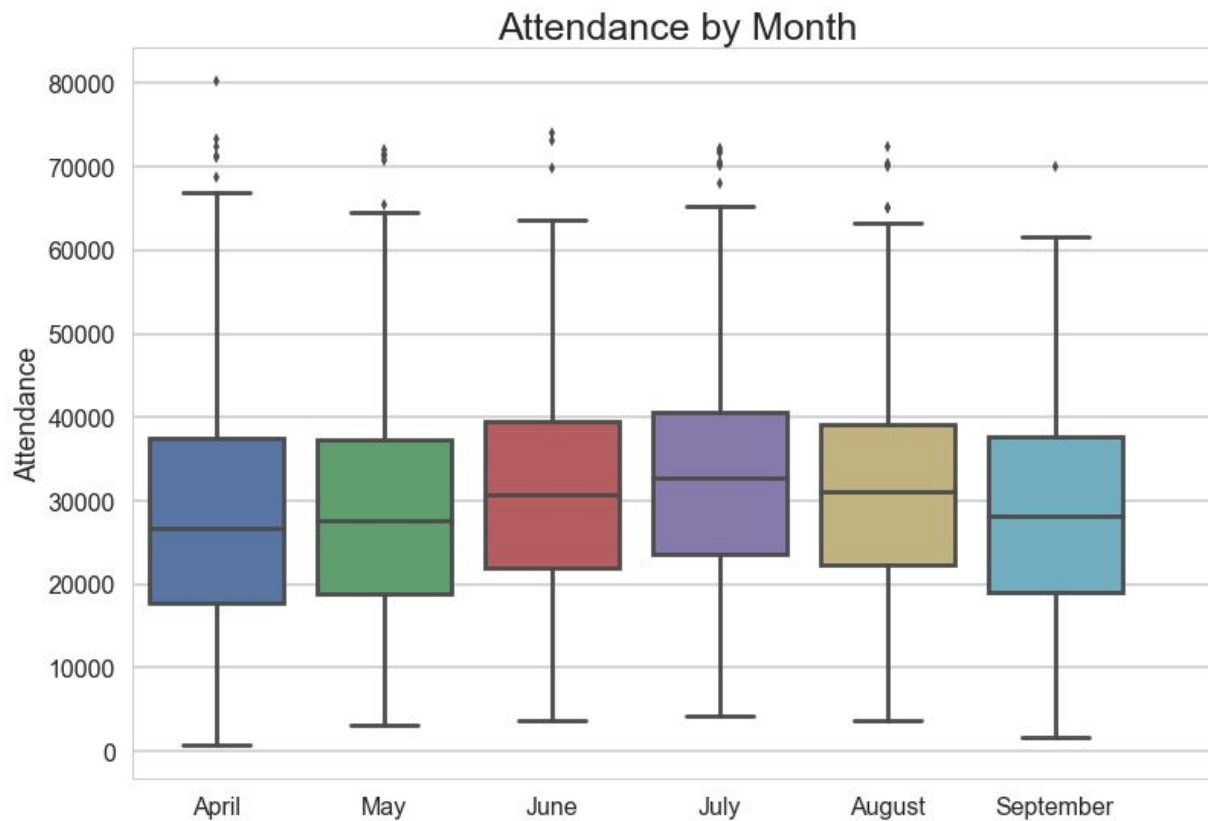
# EXPLORING THE DATA

Team performance  
is directly correlated  
with fan attendance

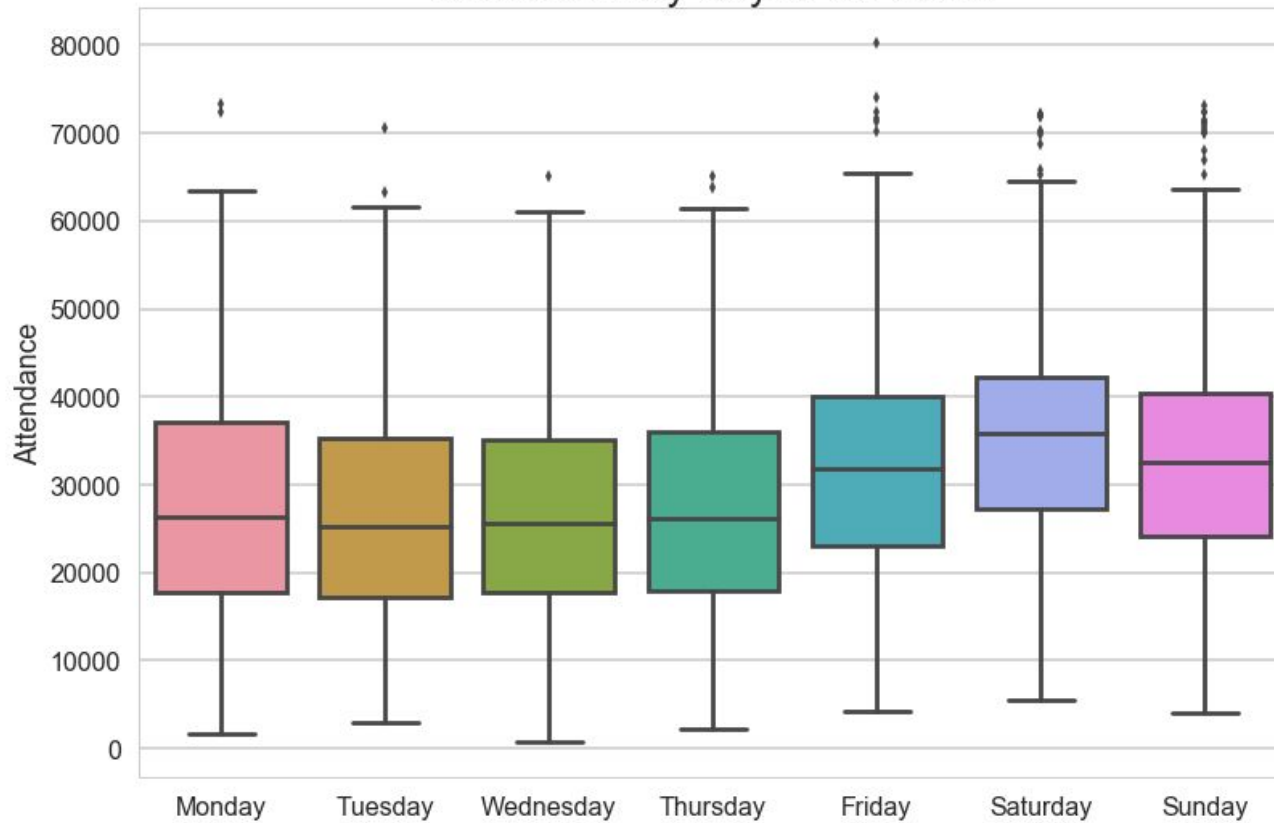


# EXPLORING THE DATA

Time of year and  
day of the week  
certainly affects  
fan turnout

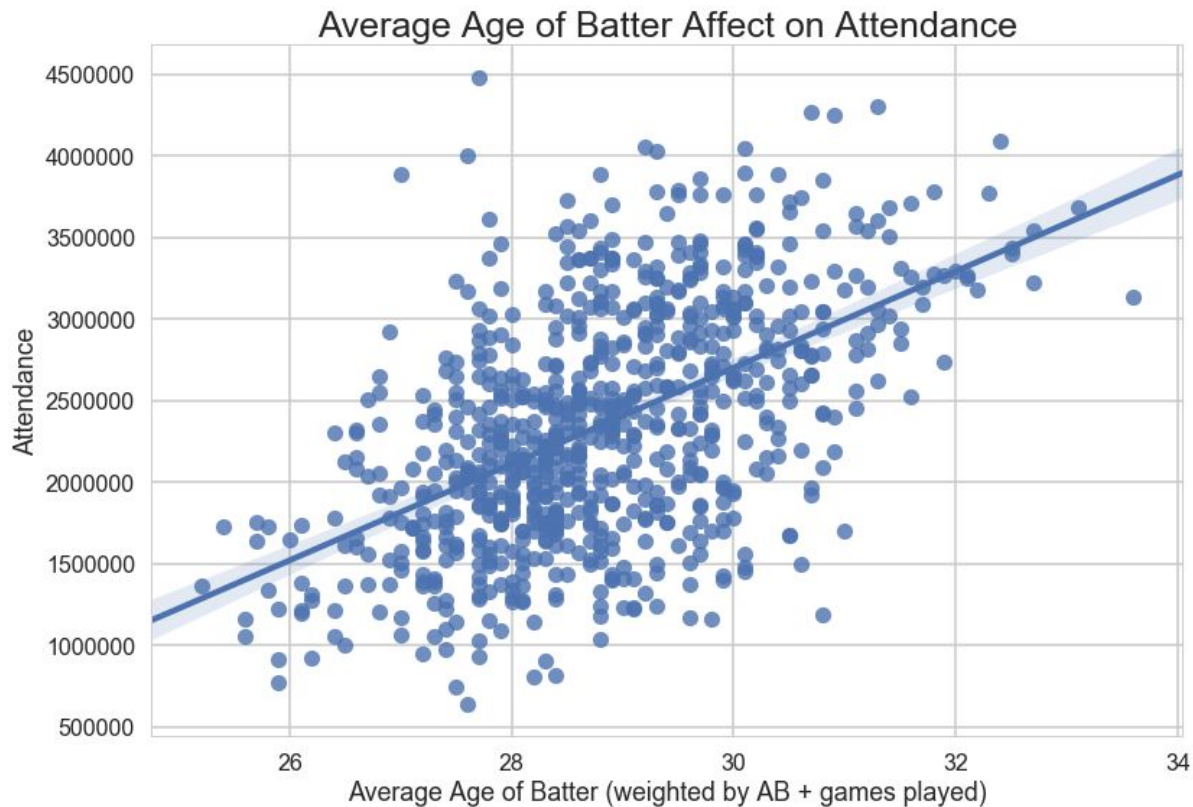


Attendance by Day of the Week



# LESS OBVIOUS FACTORS

The yearly summary data showed that the average age of the players on the team were highly correlated with attendance



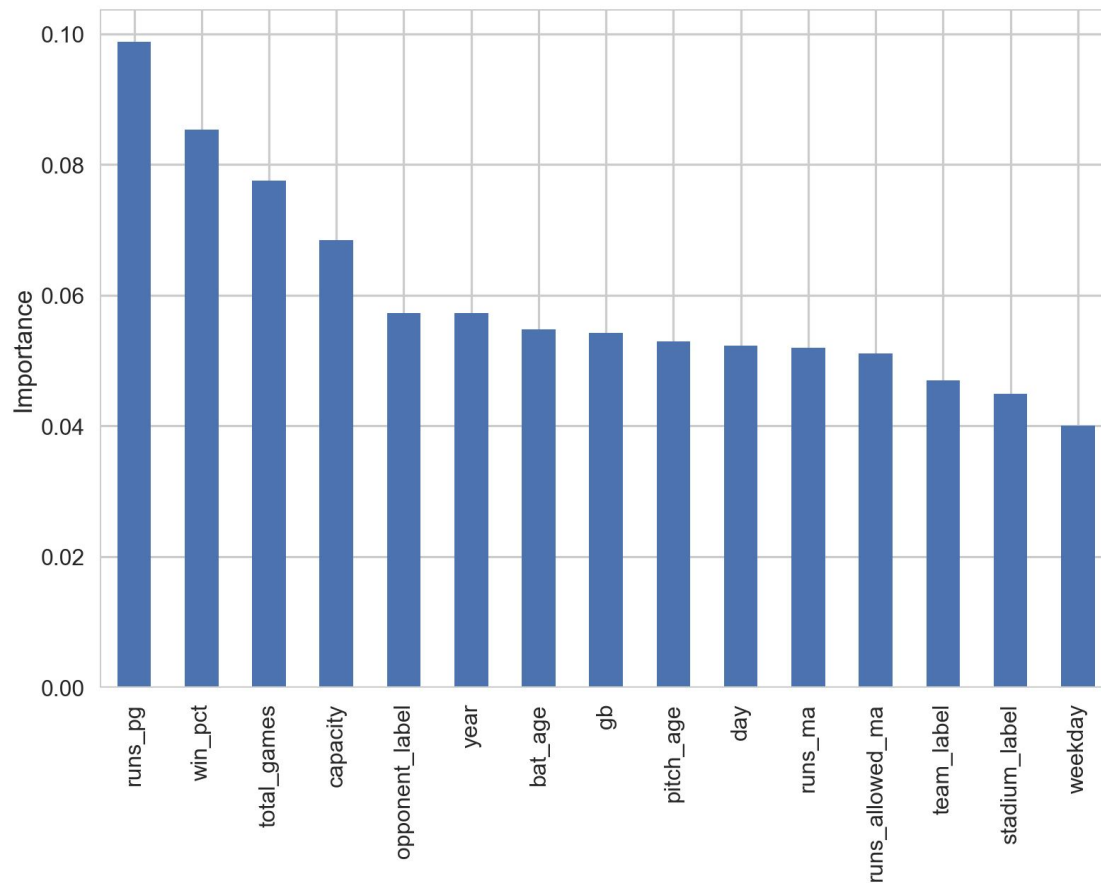


# BUILDING MODELS

# 0.833

The  $R^2$  value of my Gradient Boosting Regression model, which represents how close the actual attendance values are to the fitted regression line.

# FEATURE IMPORTANCES



# APPLYING THE RESULTS

A team may find these results useful  
in order to develop effective  
marketing/promotional strategies

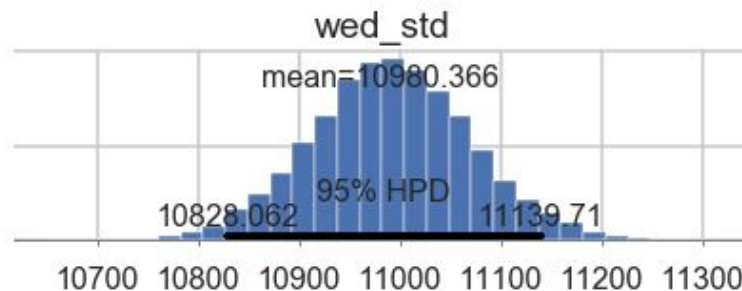
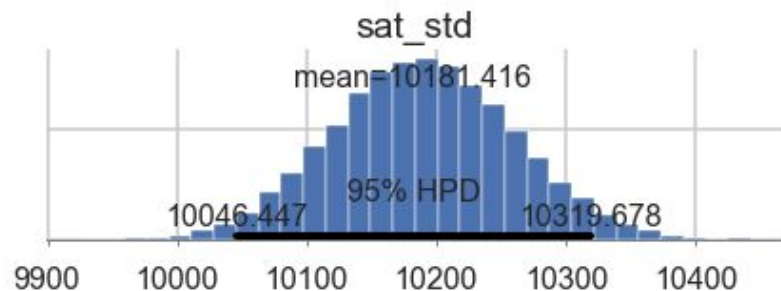
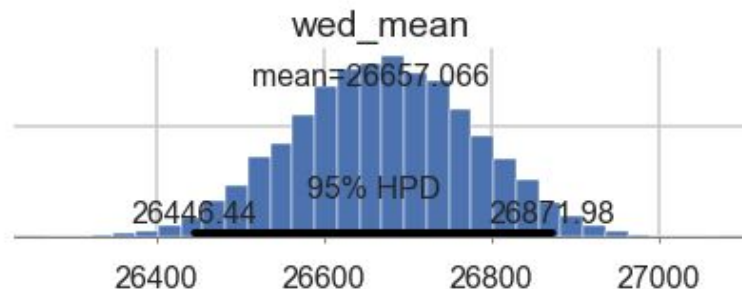
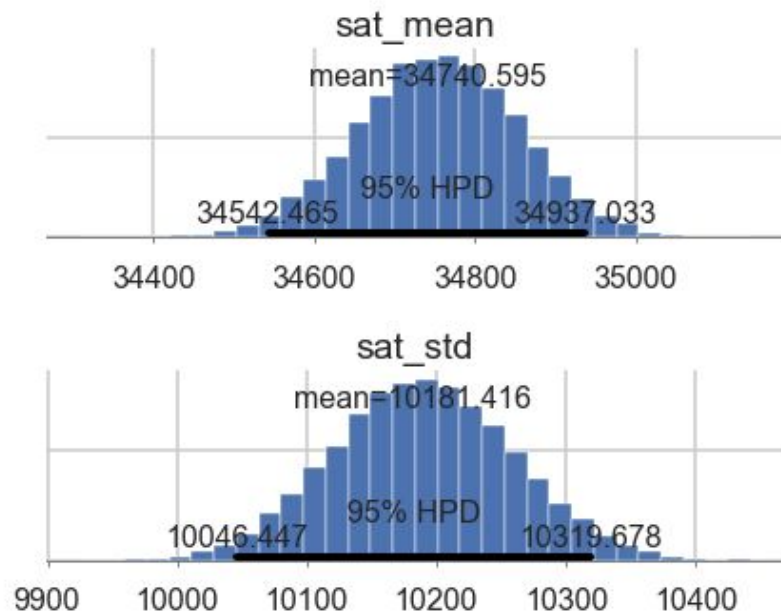
---

# BAYESIAN INFERENCE

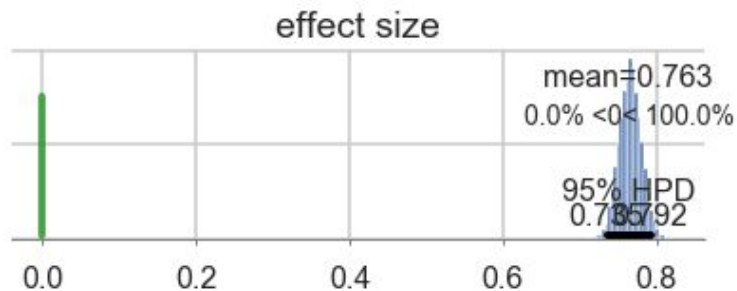
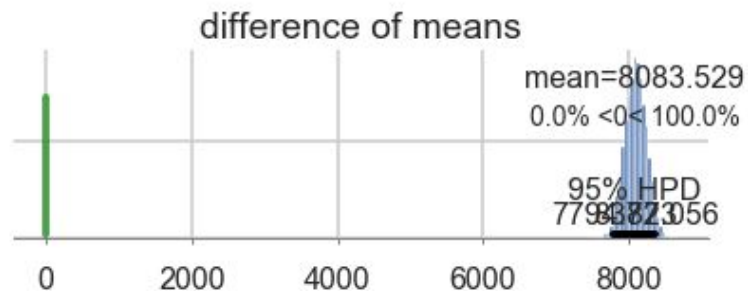
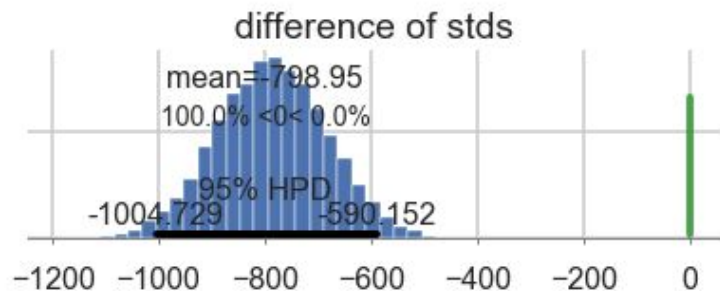
# COMPARING TWO DIFFERENT DAYS OF THE WEEK

- I decided to compare the average attendance between Wednesday and Saturday games
- As a prior, I used all other days of the week

# PLOTTING POSTERIOR DISTRIBUTIONS



# PLOTTING POSTERIOR DISTRIBUTIONS





# COMPARING TWO FAN GROUPS

## Washington Nationals

Stadium: Nationals Park

Opened: 2008

Capacity: 41,500

Avg Win Pct: 0.498

## New York Mets

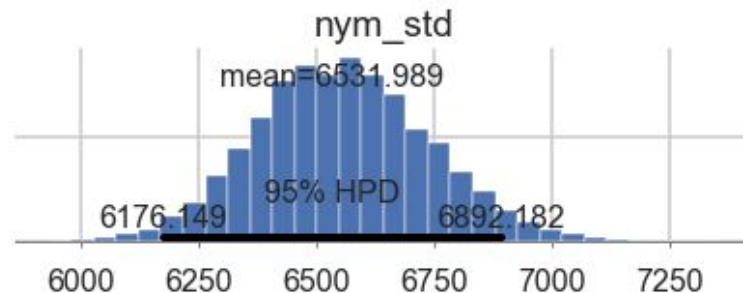
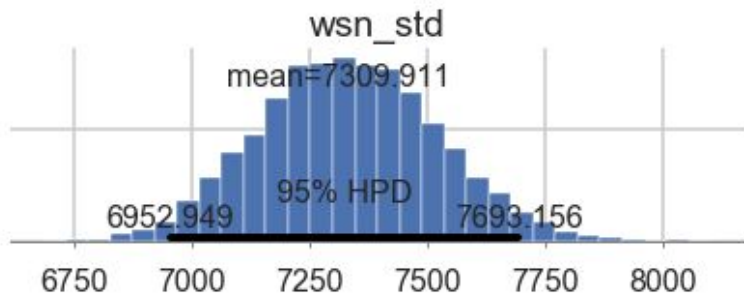
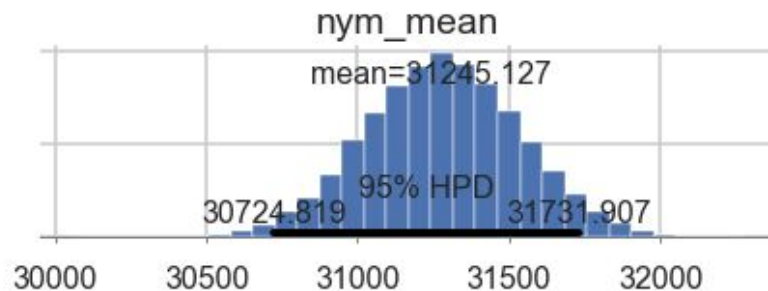
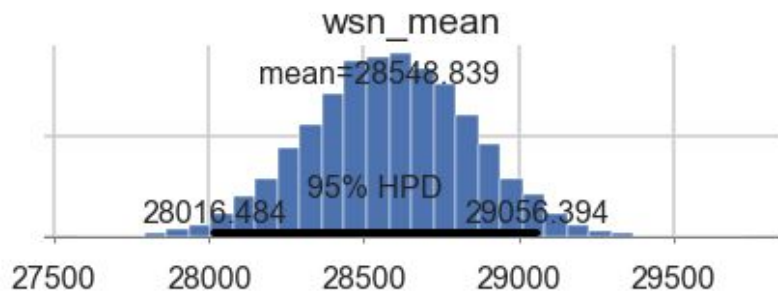
Stadium: Citi Field

Opened: 2009

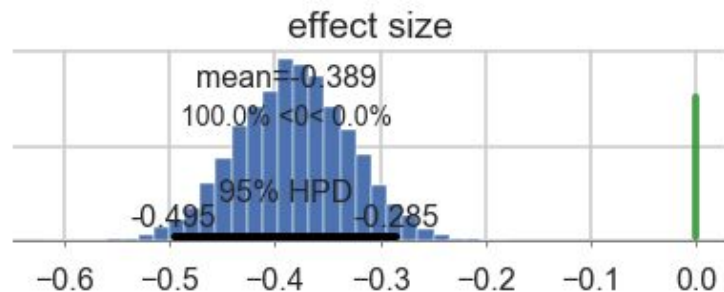
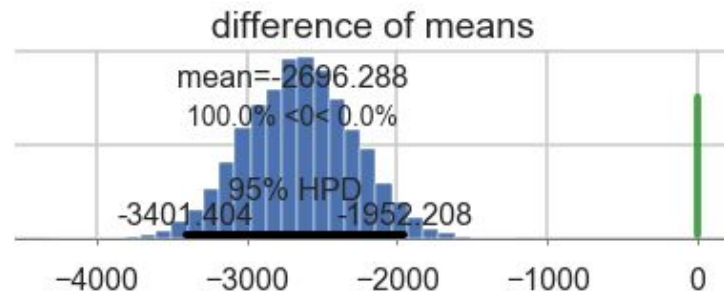
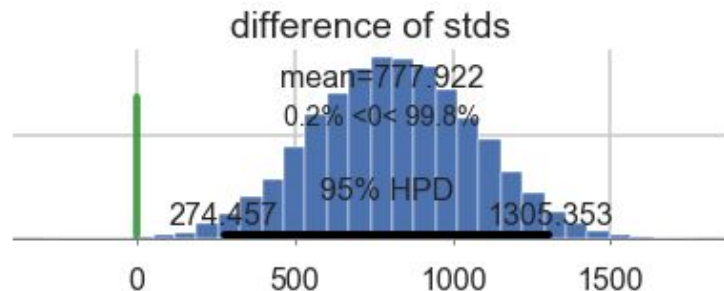
Capacity: 41,800

Avg Win Pct: 0.486

# PLOTTING POSTERIOR DISTRIBUTIONS

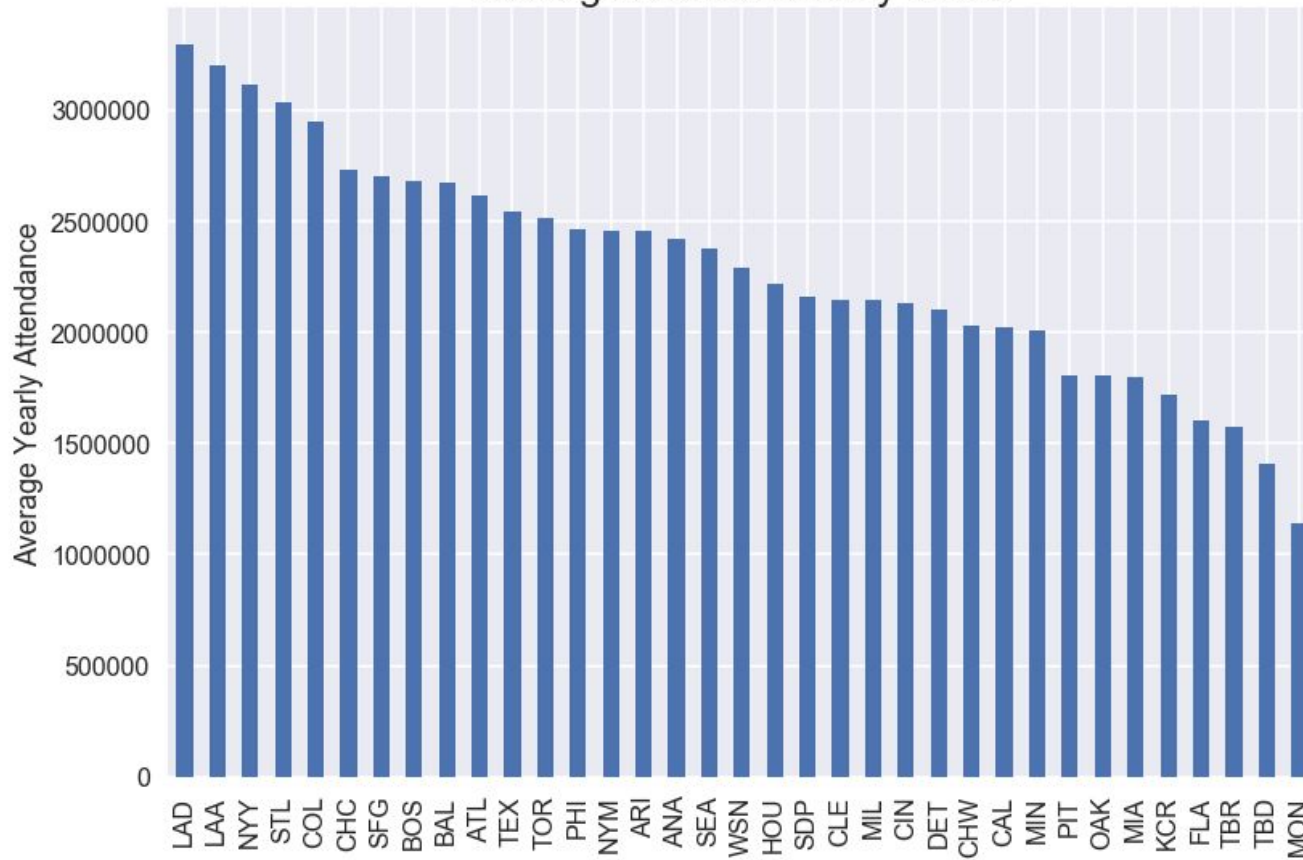


# PLOTTING POSTERIOR DISTRIBUTIONS

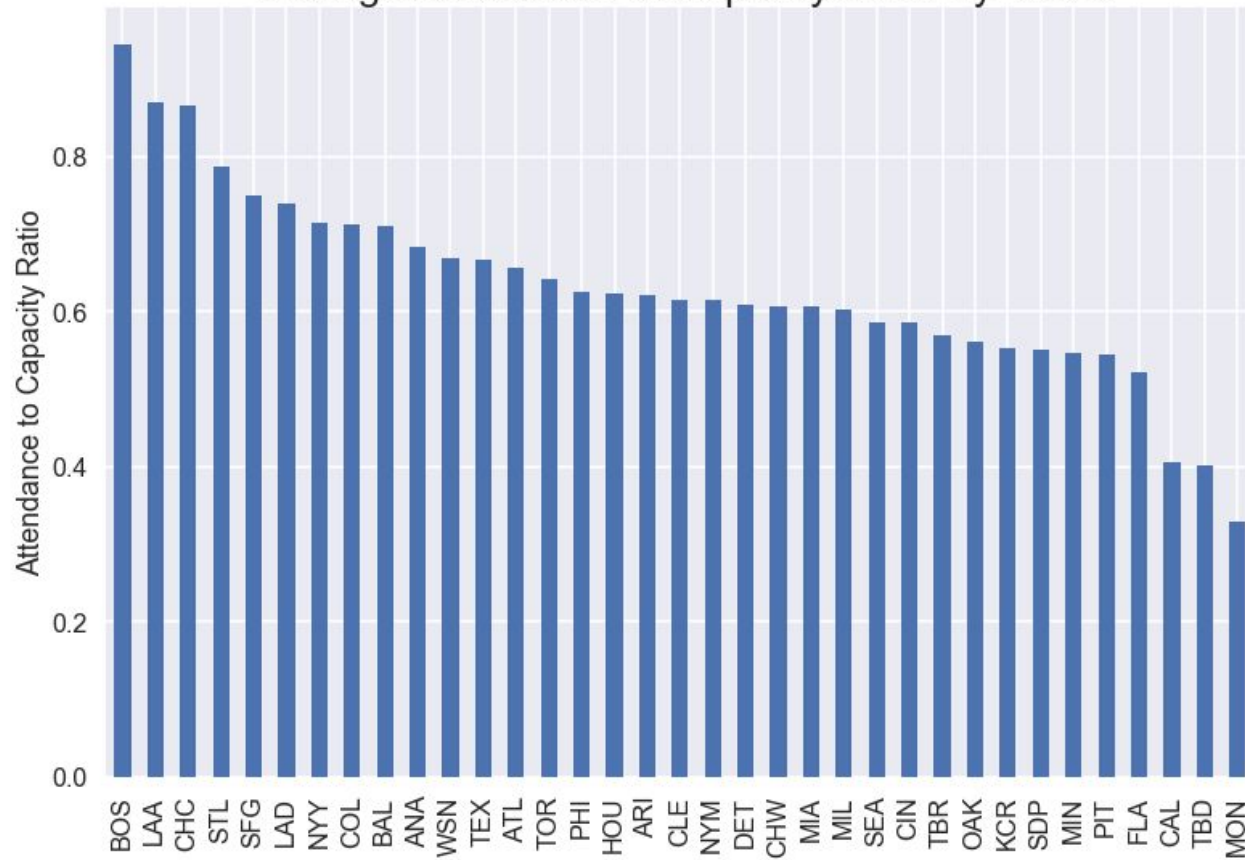


# FURTHER ANALYSIS

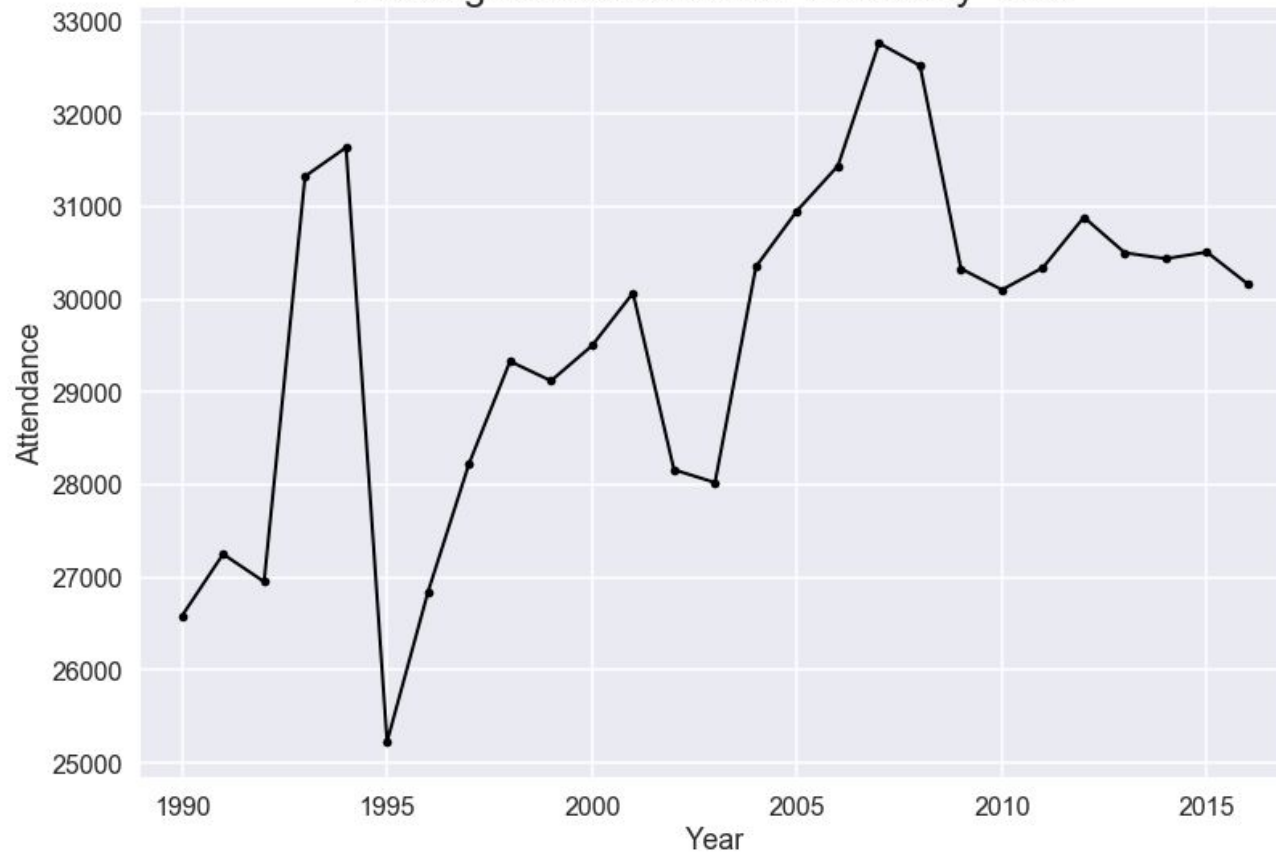
Average Attendance By Team



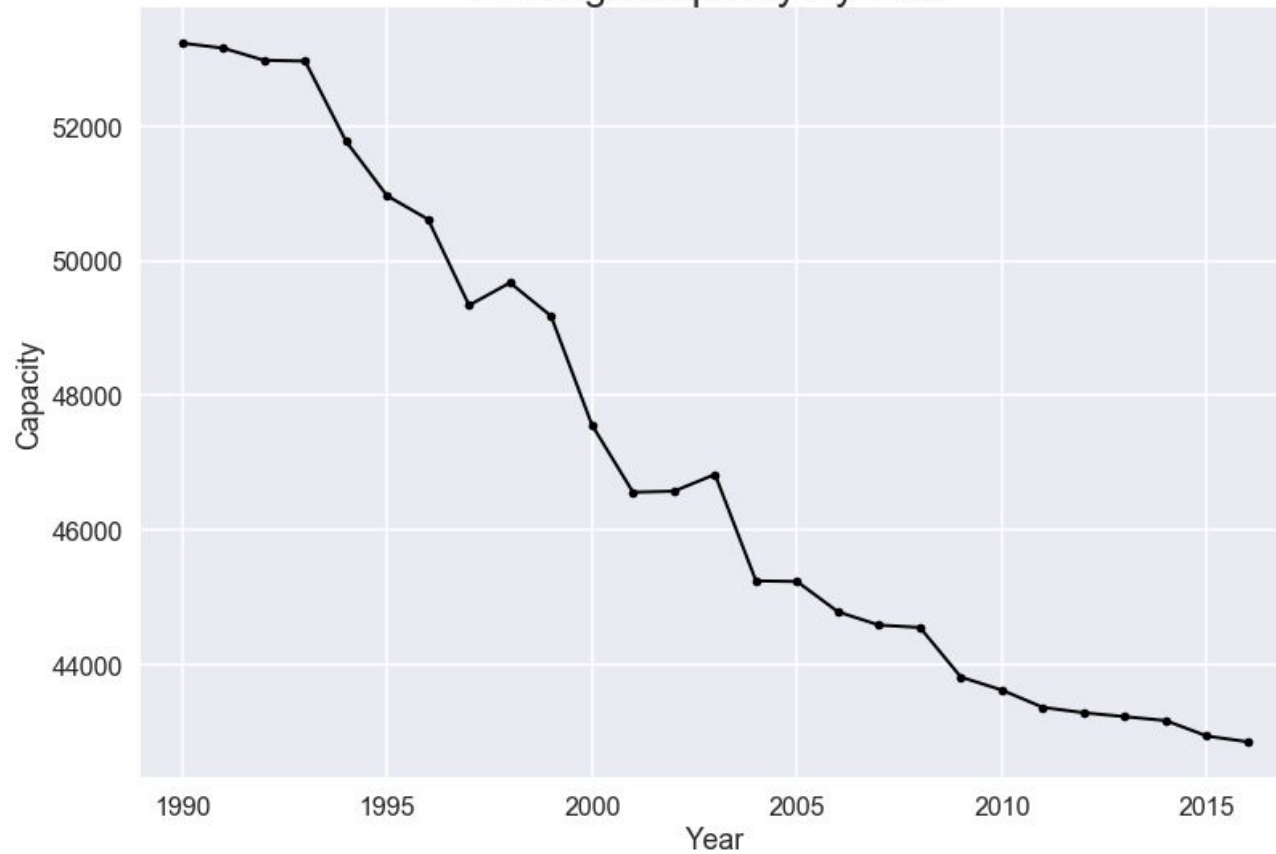
Average Attendance to Capacity Ratio By Team



Average Attendance Per Game By Year

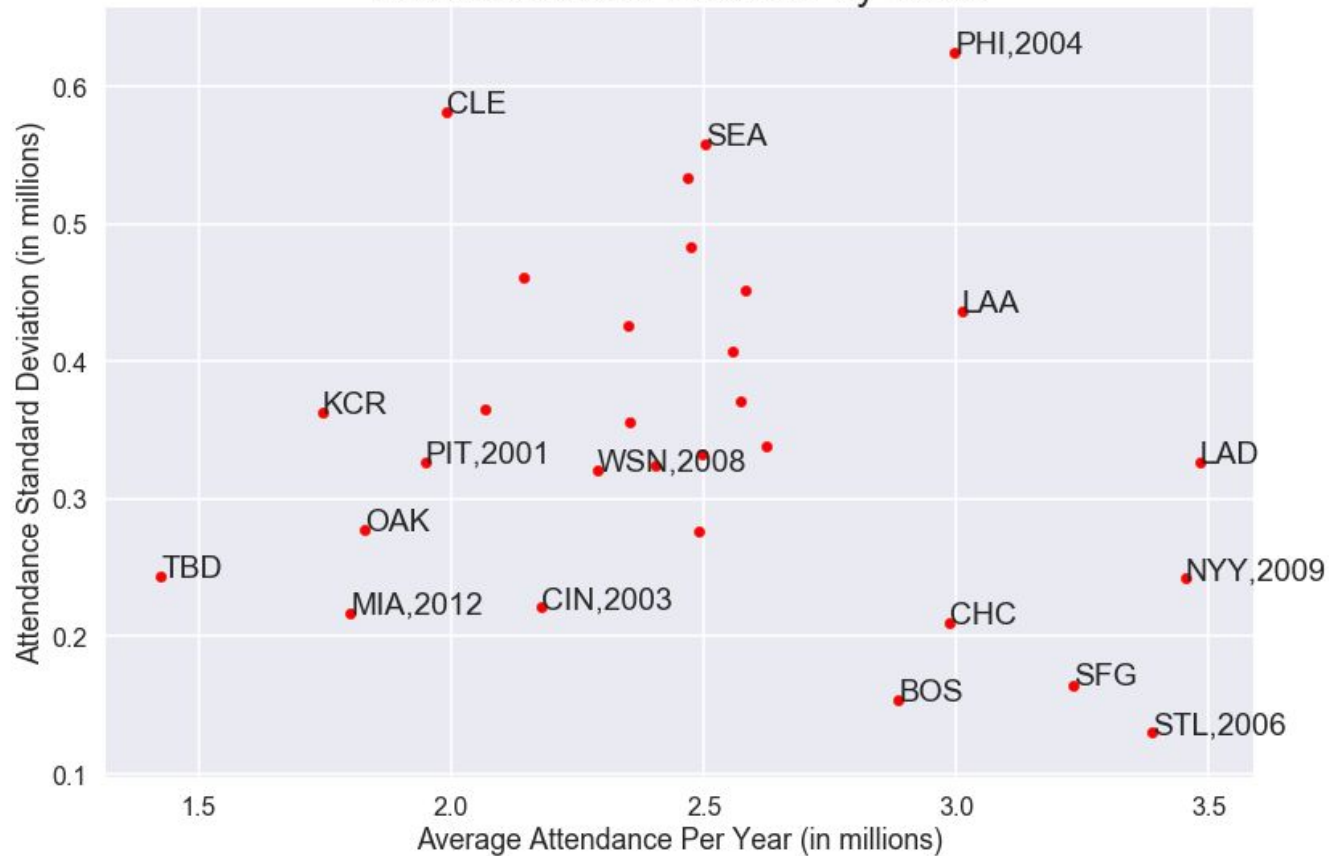


Average Capacity By Year

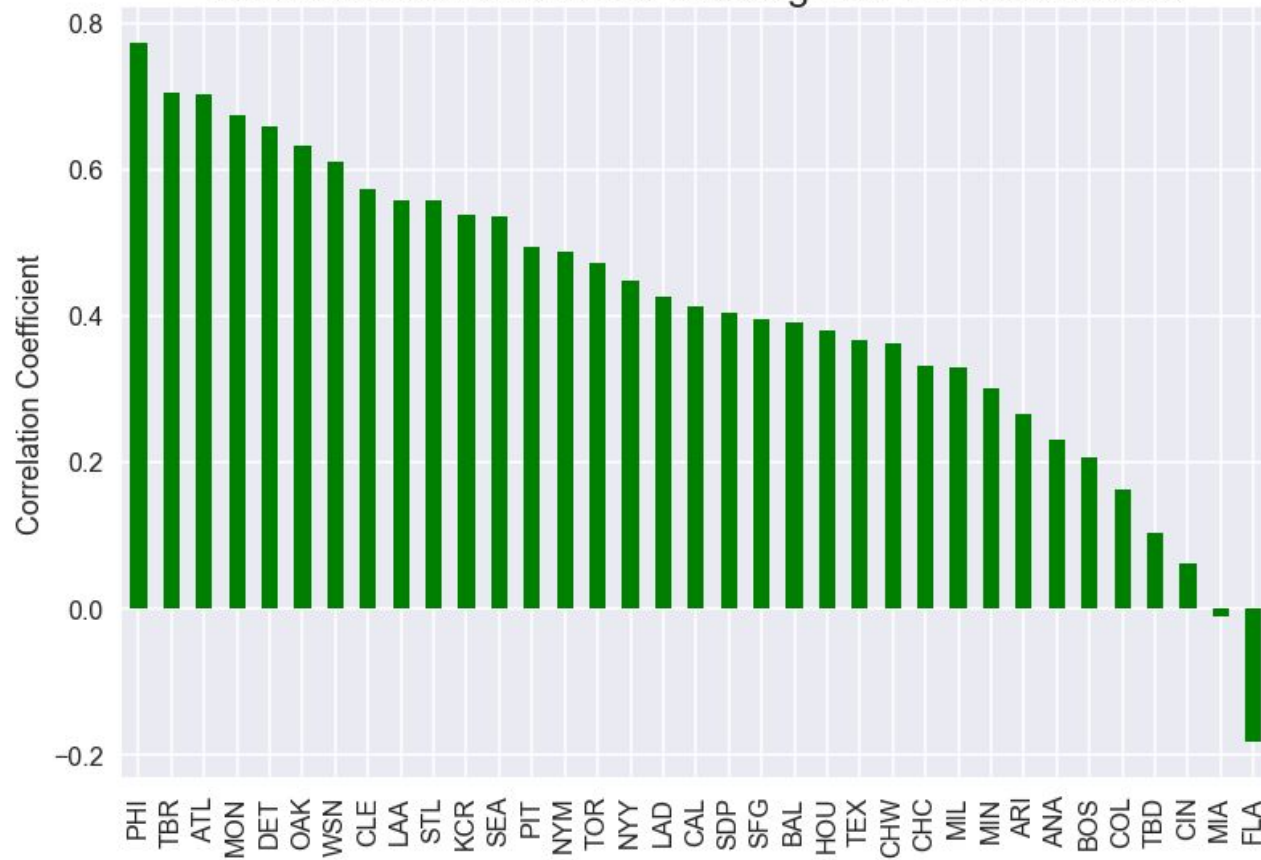




Fan Attendance Variance by Team



The Correlation Between Winning and Fan Attendance



RECAP

QUESTIONS?

---