

Stat 340 Group Progress Report

11/11/2021

Matthew Chiang: 907 723 8120 (Mchiang7)

Ishaan Backliwal: 908 134 7719 (backliwal)

Jordan Livingston: 908 132 1151 (jlivingston4)

Eric Dietze: 907 935 8843 (edietze)

Vu Pham: 907 808 5595 (vmpham2)

Description of Data

Baseball is an American sport played between two teams. One team plays defense and has a player (pitcher) try to throw the ball past the other team's batter who tries to hit the ball. If the batting team successfully hits the ball enough times they can score runs (points). We will use data sets from baseball-reference.com which contain team and individual statistics. This database also offers data 'splits', which show comparisons between home and away games. For our primary analysis, we will be focusing on the Houston Astros 2017 Team Batting Splits and the 2017 Individual Player Batting Splits. These data sets outline various different statistics of the Houston Astros 2017 season in which it was confirmed that they had cheated in home games. We will use the 'splits' data sets to determine a difference of statistics between home and away games. This database also contains data on the Astros for different seasons and different teams. Depending on the results of our initial analysis, we will compare the players on the roster of the 2017 Astros to previous seasons to determine a difference in performance before and after cheating. Our statistical question is whether or not the 2017 Houston Astros cheated. Two years after the alleged cheating occurred, reports from a player who was on the 2017 Astros alleged that there were cameras positioned in the Astros home field stadium, Minute Maid Park in Houston, that could see the relayed signs between the pitcher and catcher. Note: the catcher communicates to the pitcher what pitch the pitcher should throw based on a series of hand movements in between pitches. The Houston Astros camera could see these hand signals and people were in charge of watching the film and decoding the symbols to determine what pitch was coming. Once a pitch was shown, members of the Astros would relay what pitch was coming to the batter through their own series of symbols (most infamously a banging of a garbage can).

The Astros only had this technology at their Home Stadium, so they should have higher performance at home compared to on the road. Note: this sign decoding should reason that only their batting improved at home and not their pitching. We used several stats for measuring their performance.

Statistical Questions

- Is there a noticeable difference between the batting splits for home and away games of the Houston Astros team during the 2017 season?
- How does this difference, if any, compare to other teams' performance in the MLB during the 2017 season?

- Is there sufficient statistical evidence to suggest the Houston Astros benefited from cheating in the 2017 season?
 - If there is, how much may the cheating have affected their season? And what would their season have looked like if they didn't cheat?

Why We Chose This Dataset

- This has been a controversy within the baseball community regarded to be one of the biggest cheating scandals in the sport's history.
- In 2017 the Houston Astros won the world series. However, it came out in November of 2019 that they were using technology to steal signs during home games.
- To us, this data set is interesting because it contains (in detail) records of each game, each player, and almost all variables that could take place in this event. So, this dataset is very in depth and useful for testing hypotheses thoroughly regarding the incident. Like they said "Numbers don't lie"!

Variables

Below is a list of some important variables in our dataset:

Name	Abbr.	Description
On base and slugging percentage	OPS	Measures a players On base Percentage (percentage of At bats a player has gotten on base) and a players slugging percentage (a weighted batting average)
On base and slugging percentage (Player)	tOPS	This is adjusted so that 100 is the team average, so if tOPS is less than 100, the batters did worse, and if it is higher than 100, the batters did better.
On base and slugging percentage (Team and Player)	sOPS	This is adjusted so that 100 is the league average, so if sOPS is less than 100, the batters did worse than the league average and if it is higher than 100, the batters did better.
Runs Batted In	RBI	The number of runs a player has created by hitting a player home. Weights a players higher
Hits	H	When a batter reaches base without doing so via error or fielder's choice.
Homeruns	HR	Scored when the ball is hit in such a way that the batter is able to circle the bases and reach home safely.
Walks	W	Occurs when a pitcher throws four pitches outside of the strike zone, none of which are swung at by the batter.
Batting Average	BA	Percentage of At Bats a player gets a hit.

Loading Data

```
# will store the data for all teams in the 2017 season
season_2017 <- NA
data_empty <- TRUE

# Getting initial path data
data_path <- "../data"
data_dirs <- list.files(data_path)
data_dirs <- data_dirs[!str_detect(data_dirs, ".csv")]

# iteratively get all csv file paths and store them
for(path in data_dirs){
```

```

# gets the directory for a specific teams csv files
team_dir <- paste(data_path, path, sep="/")
csv_list <- list.files(team_dir)

# steps through csv files, loads and edits them as dataframes
for(csv in csv_list){
  # get full path, and other information from csv name
  full_path <- paste(team_dir, csv, sep="/")
  split <- ifelse(str_detect(csv, "home"), "HOME", "AWAY")
  team_name <- str_split_fixed(team_dir, pattern="data/", 2)[2]

  # Read in data
  df <- read.csv(full_path)%>%
    mutate(Name = str_split_fixed(Name, "\\\\", 2)[,1]) %>%
    mutate(split = split, team = team_name)
  df <- df[ -length(df[,1]), ]

  # binds all teams data together
  if(data_empty){
    season_2017 <- df
    data_empty <- FALSE
  }else {
    season_2017 <- rbind(season_2017, df, make.row.names=TRUE)
  }
}
}
astros_2017 <- season_2017 %>%
  filter(team=="astros", AB != 0)

# Loading Data pt 2
# btw this sucked
wins_and_losses_2017 <- read.csv("./data/win_loss_2017.csv") %>%
  mutate(split = gsub(" ", "", split, fixed = TRUE), team = gsub(" ", "", team, fixed = TRUE)) %>%
  mutate(split = toupper(split))

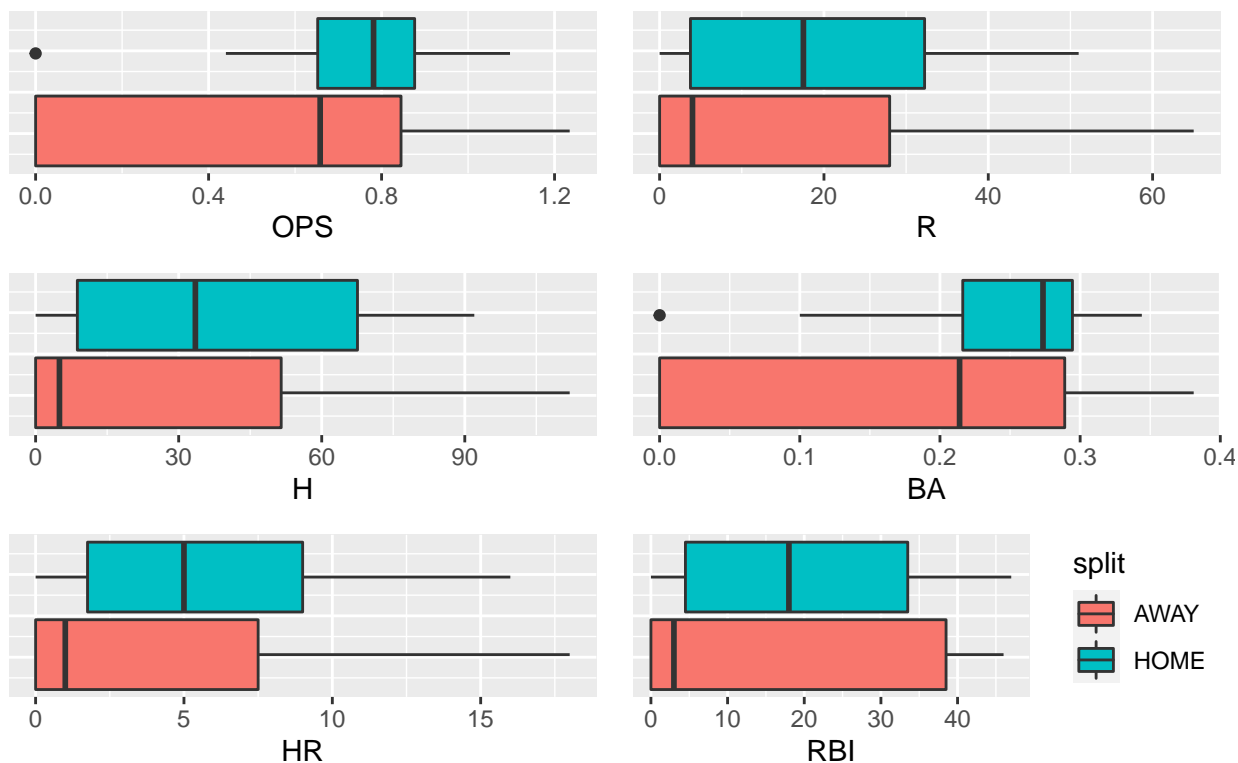
```

Preliminary Plots

Comparing The Astros 2017 Home and Away Data

to start we will look at some of the variables we take a look some comparisons of specific variables between home and away games.

Astros 2017 data – Home vs Away



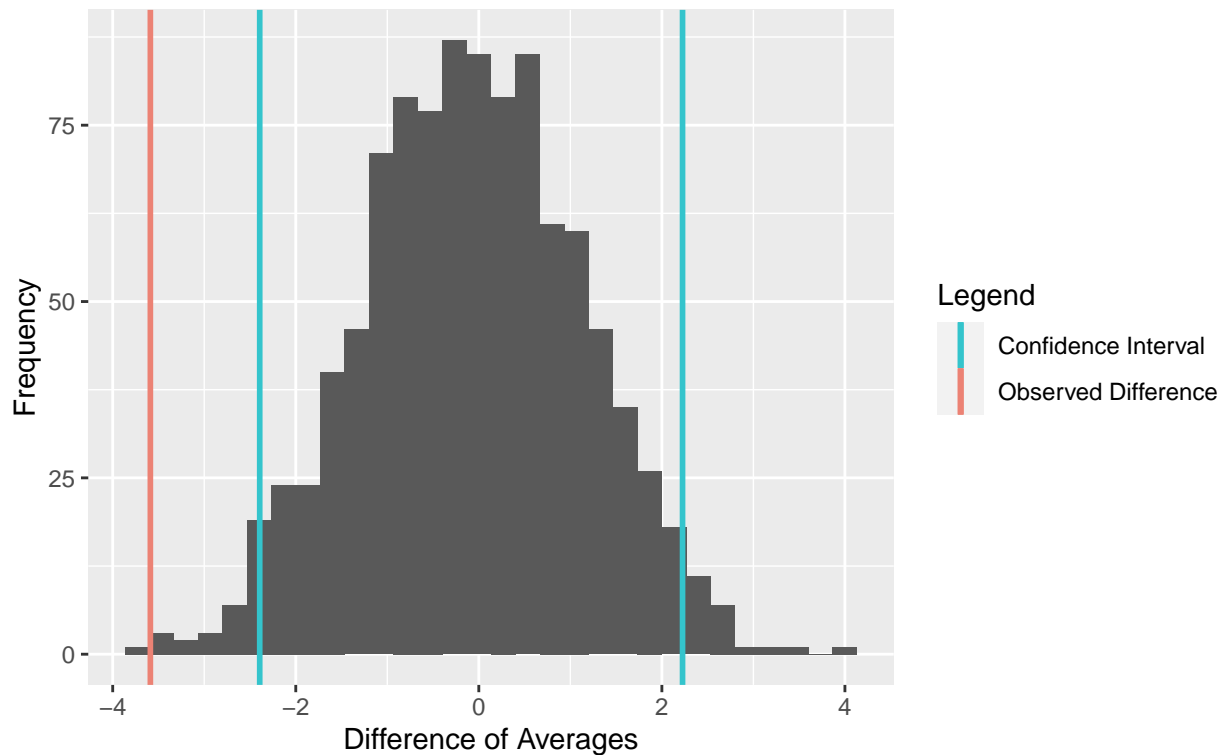
These box plots show how there is a clear difference between home and away games. The largest difference is in average Hits (H) per season that a player gets between home and away games.

Runs analysis - Home vs. Away

To do this analysis, we will be doing a monte carlo simulation on the difference of average runs a player gets per season between home and away games. Runs will be modeled by a poisson random variable. Under the assumption that Runs from home and away games comes from the same distribution, we will set lambda as the mean of runs for home and away games combined.

Difference of Average Player Runs

Home vs. Away – Astros 2017



After performing a 95% confidence interval on our data, we see that our observed difference falls outside of the confidence interval, leading us to reject our initial assumption. Thus, we can say with 95% confidence that there is statistical significance in the difference of Run means between home and away games.

However, because Away game Runs are lower than Home game Runs, this implies that the Astros cheating had a positive effect on their performance.

Is It Just a Home-Field Advantage

Since our original statistical question is basically “can we prove that the Astros cheated in the 2017 season?” we want to determine if the results we found in the confidence interval above just prove that a home field advantage exists or if the Astros have a better home field advantage due to the existence of cheating.

To do this we will perform some two sample t tests to determine a difference between the average statistics between Astros at home and other teams at their home stadium. We have realized that sign stealing should have a higher effect on players batting averages because it allowed the astros players to know which kind of pitch was coming in and thus have a better chance at hitting the ball. Batting Average is calculated as $BA = Hits/AtBats$.

```
##
## Welch Two Sample t-test
```

```
##
## data:  astros_home$BA and astros_away$BA
## t = 2.2676, df = 48.977, p-value = 0.0278
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  0.007848595 0.130125598
## sample estimates:
## mean of x mean of y
## 0.2466000 0.1776129
```

Above we see that the batting average of the Astros is significantly better at home than away. This again does not show conclusive evidence that they cheated, but instead adds strength to our earlier analysis that the Astros definitely have some sort of a home field advantage. Next we will compare the Astros with the rest of the league, by performing another two sample t test between the astros home games and the home games of all other teams in the league.

```
##
## Welch Two Sample t-test
##
## data:  astros_home$BA and c$BA
## t = 1.6987, df = 21.411, p-value = 0.1039
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.007287044 0.072700567
## sample estimates:
## mean of x mean of y
## 0.2466000 0.2138932
```

Here the results do not lead us anywhere conclusive as our p-value is too high to reject the hypothesis that the true mean of the distribution of the Astros batting average is the same as the true mean of the batting average of all other teams. We have also tried this in other variables in our data set as well, however the results have all been similar. There is no distinction between the astros home statistics versus other teams home statistics.

A New Approach

So far it seems that our data is too general for our statistical question. There are a myriad of different factors that play into how a team does in a season, differing game to game. We believe that to tackle this question we need to take a look at more in-depth statistics than just player hits, batting average, etc. We have found a new dataset from “Baseball Savant” that offer’s player data for the percent of pitches that are out of the strike zone that a player swings at and the percent of pitches that are in the strike zone that a player swings at.

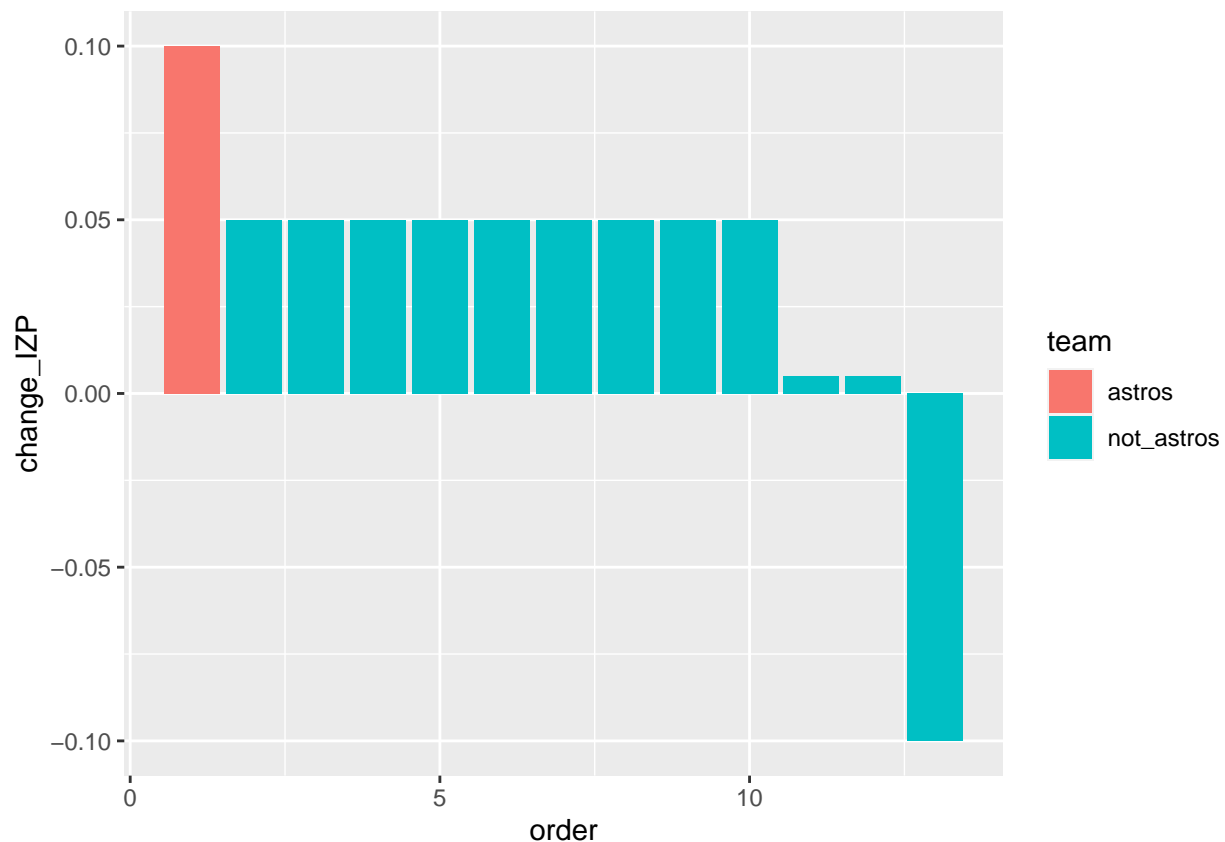
```
test = data.frame(
  player= c("a", "b", "c", "d", "e", "f", "g", "h", "i", "j", "k", "l", "m",
            "a", "b", "c", "d", "e", "f", "g", "h", "i", "j", "k", "l", "m"),
  IZP = c(.1, .2, .3, .4, .5, .6, .7, .8, .9, .8, .14, .35, .66,
          .2, .25, .35, .45, .55, .65, .75, .85, .95, .85, .145, .355, .56),
  OZP = c(.15, .25, .35, .45, .55, .65, .76, .86, .96, .86, .146, .356, .65,
          .2, .28, .38, .48, .58, .68, .75, .85, .95, .85, .145, .355, .56),
  year = c("2017", "2017", "2017", "2017", "2017", "2017", "2017", "2017", "2017", "2017", "2017", "2017", "2017",
            "2016", "2016", "2016", "2016", "2016", "2016", "2016", "2016", "2016", "2016", "2016", "2016"),
```

```

team = c("A", "B", "C", "D", "E", "F", "G", "H", "I", "J", "K", "L", "M",
        "A", "B", "C", "D", "E", "F", "G", "H", "I", "J", "K", "L", "M")
)

IZP_diffs <- test %>%
  drop_na() %>%
  group_by(player) %>%
  mutate(change_IZP = (IZP[2] - IZP[1]), change_OZP = (OZP[2] - OZP[1])) %>%
  select(player, team, change_IZP, change_OZP) %>%
  unique() %>%
  arrange(desc(change_IZP)) %>%
  mutate(team = ifelse(team == "A", "astros", "not_astros"))
IZP_diffs$order <- 1:length(IZP_diffs$change_IZP)
IZP_diffs %>%
  ggplot() +
  geom_col(aes(x=order, y = change_IZP, fill = team))

```



```

zones <- read.csv('./data/zone_swings.csv')
zones_clean <- zones %>%
  unite('Name', first_name:last_name, remove = TRUE) %>%
  mutate(Name = str_replace(Name, " ", "")) %>%
  mutate(Name = str_replace(Name, "_", " ")) %>%
  select(Name, out_zone_swing, in_zone_swing, year) %>%
  arrange(Name)

```

```

astros_zones_2017 <- merge(astros_2017, zones_clean, by = "Name", all = TRUE)
# astros_zones_2017
# zones_clean
z_2016 <- zones_clean %>%
  filter(year == "2016")

z_2017 <- zones_clean %>%
  filter(year == "2017")

# Zone data for Astros 2017
astros_z_2017 <- left_join(z_2017, astros_2017, by = "Name") %>%
  select(Name, out_zone_swing, in_zone_swing, year, team) %>%
  drop_na() %>%
  unique()

# Zone data for Astros 2016
astros_z_2016 <- left_join(z_2016, astros_2017, by = "Name") %>%
  select(Name, out_zone_swing, in_zone_swing, year, team) %>%
  drop_na() %>%
  unique()

# Zone data for Non-Astros teams 2017
not_astros_z_2017 <- left_join(z_2017, not_astros, by = "Name") %>%
  select(Name, out_zone_swing, in_zone_swing, year, team) %>%
  drop_na() %>%
  unique()

# Zone data for Non-Astros teams 2016
not_astros_z_2016 <- left_join(z_2016, not_astros, by = "Name") %>%
  select(Name, out_zone_swing, in_zone_swing, year, team) %>%
  drop_na() %>%
  unique()

```

Discussion on Progress/Challenges

Challenges:

- On our last project report, the results that we gathered from the data were not what we expected. We originally cleaned the dataset wrong and it affected the outcome, giving us reason to believe the Astros didn't cheat. After going back and changing how we read in the data, our graphs prove what we were originally expecting, and show initial evidence that they were cheating. Our original findings were very surprising because this cheating scandal wouldn't have been a scandal if it weren't for some correlation and evidence that suggest they did. This was a big challenge because after our last project report, we thought we were gonna have to go back and change our whole project. But now we are back on track and can continue the analysis of our questions.
- Another challenge we faced during this project report was trying to figure out how to use a linear model in our analysis. We originally wanted to use a regression model on all of the variables in our dataset to predict wins and then we could see if home games (the supposed cheating) was a statistically significant variable. After initially trying this, we realized this approach was not a very good/strong indication of how home games helped them cheat. Then, we thought it would be beneficial to compare individual Astro players pre-season stats, to that of the regular season, and compare home and away.

This also came with its issues though because there wasn't individual player data for preseason. So we moved onto our next thought process

- From the data site Kieth gave us from the last project report, we found a perfect data set split up into the home and away games, but there were no column names for 161 variables. So we used python to assign column names from the reference sheet online.
- In doing this, we found out how hard it is to try and piece together all the datasets. This method also increases the margin of error, which is something we can compute and analyze for the final report.
- After discussing our project in office hours and realizing regression would not really prove anything, we decided to just keep going down the same path of comparing home and away splits. This led us to multiple t-tests (**fill in what VU did**) and continued along with our analysis. The challenge with this was that no matter how we computed home and away splits (we tried many different ways), there was nothing really jumping out at us. We attempted bootstrap but the results again were no good.
- This led to see that our data just wasn't specific enough. We couldn't find any real evidence that the Astros had a advantage at home do to cheating. If we want to see if the actual act of sign stealing had an impact on the success of the Astros at home, and not just do to home field advantage, we had to find more specific data. Although we were told we shouldn't be gathering any more data, we thought if we wanted be able to answer our statistical questions in any way, we would need to gather more data on in (strike) swing percentage and out of (strike) zone percentage. This way, we can actually see if the batter having the knowledge to what pitch he was gonna be thrown actually impacted his hits at home in comparison to away.
- If we had more time, we would spend much more time searching for a data set with in zone and out zone swing percentage for home and away. The only easily available data set to grab without too many cleaning issues combined home and away zone swing percentages, making it impossible for us to fully answer our question. We were able, however, to see a significant difference in the zones swing percentage from the previous year, (2016) to the cheating year, showing us that there is definitely more to be explored on this path.

Progress:

- We came to the realization that some of our data cleaning from our last project report was skewing the results. We fixed that problem, re-fit the previous graphs, and re-analyzed them
- We spend a tremendous amount of time reading in new data and cleaning it, in order to form a regression model. After forming a linear model to predict win_percentage based on Runs, Hits, Runs Batted In, Home Runs, and Split (Home vs, Away), we found that the model sucks, like a lot. So we will need to adapt our model to work better in the future, but we did not have enough time to finish it now.
- In experimenting with different ways to analyze home and away splits, we found that creating our own batting average would be best because we could create one for each home and away game according to the player. In doing this we also realized that teams rosters are different for home and away and this, along with players who never got up the the plate once in a whole season were skewing our results. In an attempt to resolve this, we only analyzed players who got up to bat once per game, or 162 times for all the regular season games in the pre-covid 2017 season.

Next Steps

- We want to figure out why the new regression models we created were so bad. Before the final project is due, we will work on how to improve our model, and attend office hours to gain some insight as why they currently are the way they are. We can then plot the models to visually see how they are doing

- We also might want to look into the performance of teams throughout the league at the Astros' home stadium to see if the conditions under which teams play at that location (called park factors) have any affect on performance, indicating a possible reason for the Astros lower performance at home.
- We are also considering looking into things such as individual Astros players performance and test to see if there were any noticeable, improbable differences between years, comparing it to the change in performance for other teams' top performers.