

Baseball Sabermetrics

A Closer Look at Advanced Baseball Statistics and Their Change Over Time

Nathan Shaver
109410328
University of Colorado Boulder
CSCI 4502
nash4637@colorado.edu

Thomas Hawley
109706364
University of Colorado Boulder
CSCI 4502
thomas.hawley@colorado.edu

ABSTRACT

Sabermetrics is the empirical analysis of baseball, especially baseball statistics that measure in-game activity. These activities measure pitching, batting, offensive and defensive metrics. Analysis of these statistics has resulted in more advanced statistical measurements that coaches and team front offices use to analyze players that simpler statistics wouldn't offer. At its foundation, sabermetrics raises questions about how baseball is played and the most effective strategies to achieve, then attempts to answer them empirically. We used pybaseball, a Python library for storing baseball data, to analyze various relationships between different baseball statistics and extrapolate that relationship by stressing the data with multiple variables. As such, our general null hypothesis is as follows: various events and/or rule changes have no impact on the overall statistical trends throughout Major League Baseball.

INTRODUCTION

To start, let's take a look at several classic pitching statistics. These metrics have been used in the sport for over a century.

- ERA (Earned Run Avg)
- IP (Innings Pitched)
- BB (Base on Balls)
- K (Strikeouts)
- H (Hits Allowed)

These statistics follow a fairly consistent pattern. How many walks did the pitcher give up in a season? How many batters did the pitcher strike out in a season?

They are all cumulative metrics. Now let's take a look at several advanced sabermetric statistics that have appeared in the last 20 years.

- WHIP (Walks plus hits per Innings Pitched)
- BABIP (Batting average on balls in play)
- FIP (Fielding Independent Pitching)
- SIERA (Skill-interactive Earned Run Avg)
- BQR (Bequeathed Runners Scored)

These metrics take the classic statistics and convert them into measurable data. For example, WHIP takes into account a pitcher's ability to keep batters off the basepaths. It combines walks and hits allowed per inning pitched.

Through the use of these metrics and tools described in the background, we examined various relationships between different baseball statistics and various rule changes. We graphed and analyzed the relationships on various plots. From there we drew conclusions on the results we found.

RELATED WORK

Sports teams would originally evaluate athletes based on gut instincts and arbitrary attributes such as fitness level, motivation, appearance, and body type. However, data scientists and statisticians have entered the sport with advanced measures over the previous twenty years. The 2002 Oakland Athletics were the first and most successful team to employ this strategy, taking a statistical approach to their season by using data analytics to draft, trade, and manage its players. Steinberg from [1] stated that the 2002 Oakland Athletics "sabermetrics system did not

Mid 1970s: Various MLB teams moved to larger stadiums with longer distances from home plate to the outfield fence.

In 2003: MLB began testing players for steroids.

After graphing the data between 1900-2022, it was necessary to begin data cleaning. We didn't want extraneous factors impacting our results. As such, we decided to clean the data to only include the year of the rule change plus or minus 5-7 years, depending on the experiment. We then analyzed the null hypothesis through a difference of means t-test. Our decision on whether to reject or to fail to reject the null hypothesis comes down to if our t-statistic is inside or outside the t-critical value bounds. From there, we looked at the key findings and implications of our results.

PRELIMINARY WORK

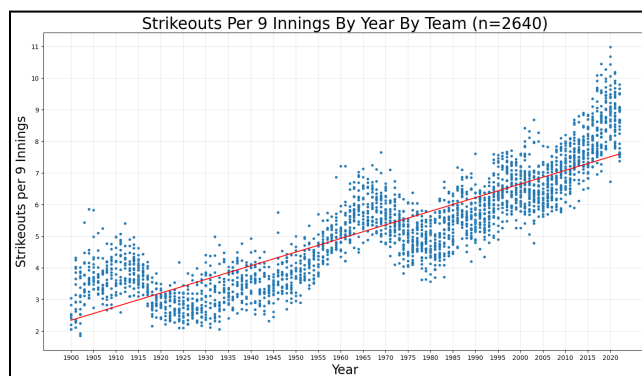
For our initial setup, we drew up code in Python to analyze statistics across a broad range (1900-2022). We wanted to analyze specific statistics that we thought could be impacted by the rule changes outlined in our methodology. As such, we plotted and analyzed K/9 (strikeouts per nine innings), SB (stolen bases), and HR (home runs). Each data point was averaged by team per year. For example, the 2022 Texas Rangers stole 128 bases as a team. That data is represented as a single point. The summary data collected for each team and each year are as follows.

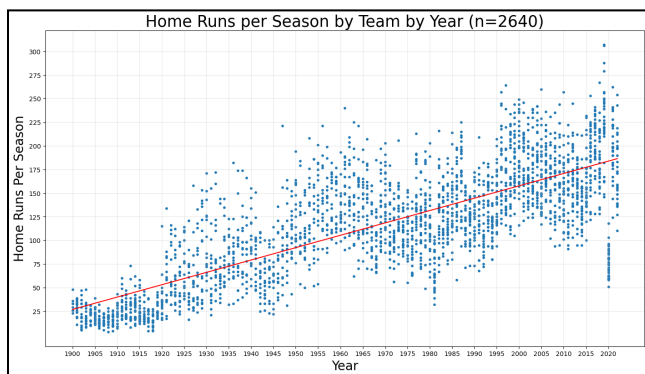
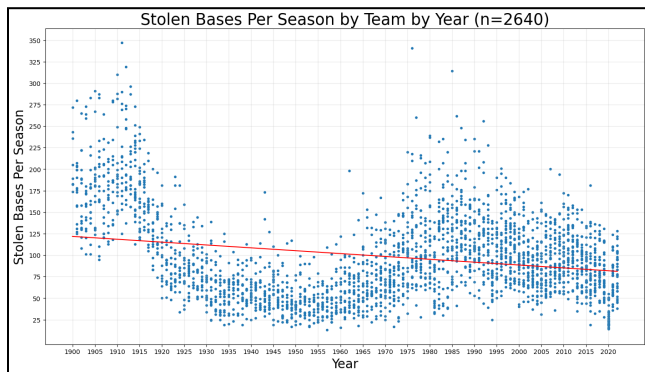
K/9 (1900-2022)	
n	2640
mean	5.3500
standard deviation	1.7571
min	1.8500
max	10.9800
variance	3.0873

SB (1900-2022)	
n	2640
mean	98.5341
standard deviation	51.3285
min	13
max	347
variance	2634.6241

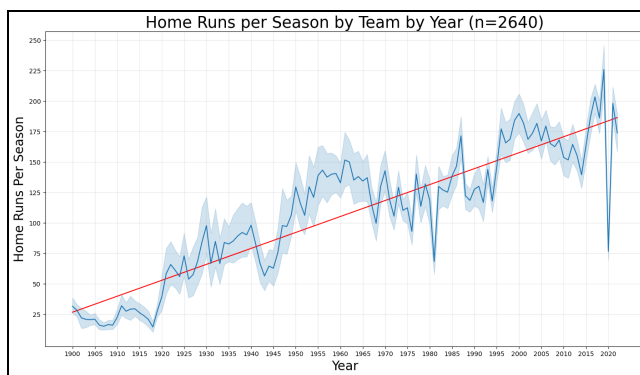
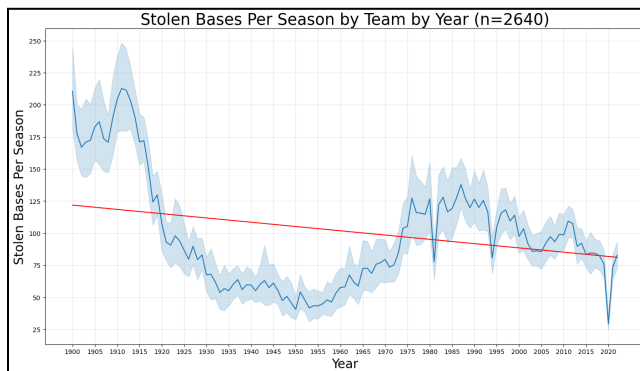
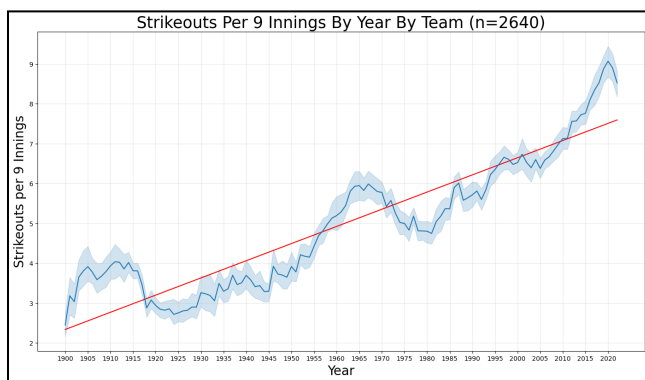
HR (1900-2022)	
n	2640
mean	118.3057
standard deviation	59.4055
min	3
max	307
variance	3529.0081

The scatterplots of each statistic vs year are as follows. Blue dots represent a single data point, the red line represents the best fit line of the entire dataset.





For visual purposes, we wanted to clean the data and contain the line plot to a 95% confidence level. Below are the plots for each statistic. The dark blue line represents the yearly mean, the shaded blue area represents the 95% confidence interval of all data points for each year.



Now that our preliminary data was collected and displayed, it was time to move on to the evaluation part of our three experiments.

EVALUATION

With the established framework for the project, we move on to the major topic. We formed the general null hypothesis: various events and/or rule changes have no impact on the overall statistical trends throughout Major League Baseball. The experiments tested are the following rule changes: MLB lowered the height of the pitching mound from 15 in to 10 in in 1968, various MLB teams moved to larger stadiums with longer distances from home plate to the outfield fence in the 1970s and MLB began testing players for steroids in 2003. Diving into the rule changes, we start with MLB lowering the height of the pitching mound.

Experiment 1

Null hypothesis: MLB lowering the pitching mound in 1968 from 15 inches to 10 inches had no impact on team strikeouts per 9 innings ratio.

Null hypothesis: $\mu_1 - \mu_2 = 0$

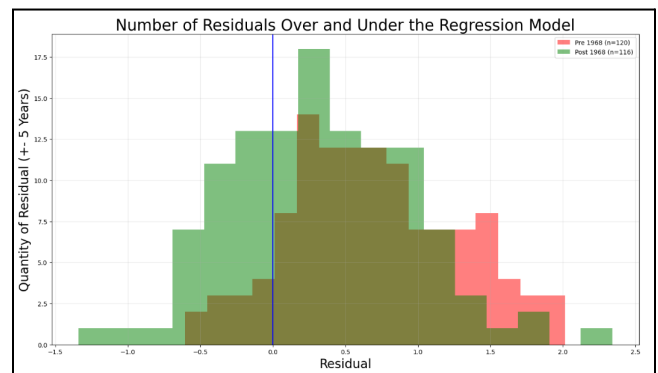
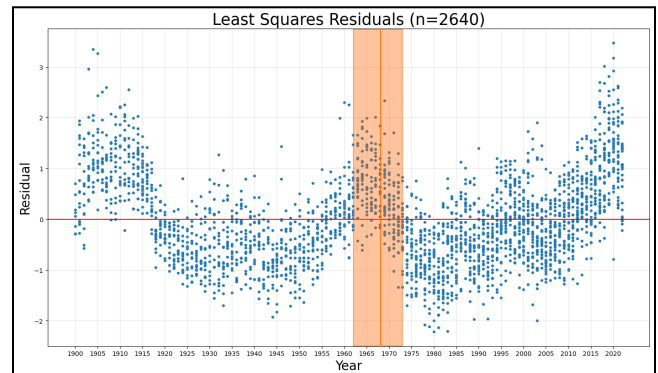
Alternative hypothesis: $\mu_1 - \mu_2 \neq 0$

As written above, we used pyBaseball to extract team pitching data from 1900 to 2022. We then mined the K/9 column, averaged each team by year, and plotted them by year. However, we wanted to clean the data further by only analyzing the years surrounding 1968, when the rule was put into effect. This way, no extraneous variables in the game would impact our results. The summary data collected for each team and each year are as follows.

K/9 (1962-1967)	
n	120
mean	5.8236
standard deviation	0.5961
min	4.4900
max	7.2400
variance	0.3553

K/9 (1968-1973)	
n	116
mean	5.6853
standard deviation	0.6026
min	4.100
max	7.6500
variance	0.3631

We then graphed the least squares residuals for this rule change. The rule went into effect in 1968, so we wanted to look at what percentage of residuals were over/under the lifelong trendline, which is highlighted in orange on the histogram. Each histogram was separated into years before and after the rule went into effect.



Of the 120 residual points between 1962-1967 (red), 90.00% of them were at or above the lifelong trendline. Of the 116 residual points between 1968-1973 (green), 68.96% of them were at or above the lifelong trendline.

To test our null hypothesis, we used a difference of means hypothesis test (two tailed, $\alpha < .05$). This decision comes down to the t-statistic value. If it is outside the critical values from the t-table below, we reject the null. If it is inside, we fail to reject the null hypothesis. Using the equation: $t = [(\bar{x}_1 - \bar{x}_2) / SE]$, we found the test statistic to be $t = 1.7723$. Using a T-table, $P(t \leq 1.7723) = 0.96118$. Therefore, $P(t > 1.7723) = 0.0388$. The two tailed P value equals 0.776. Since $\alpha = .05$, this P value is not considered to

be significantly significant, and we fail to reject the null hypothesis.

Experiment 2

Null hypothesis: Several baseball teams moving into stadiums with larger fields in the 70s had no impact on team successful stolen bases.

Null hypothesis: $\mu_1 - \mu_2 = 0$

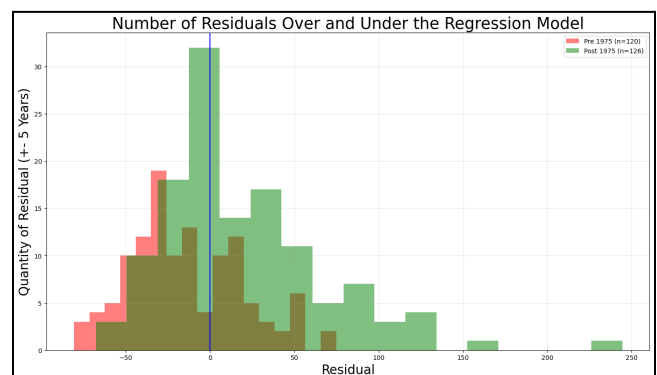
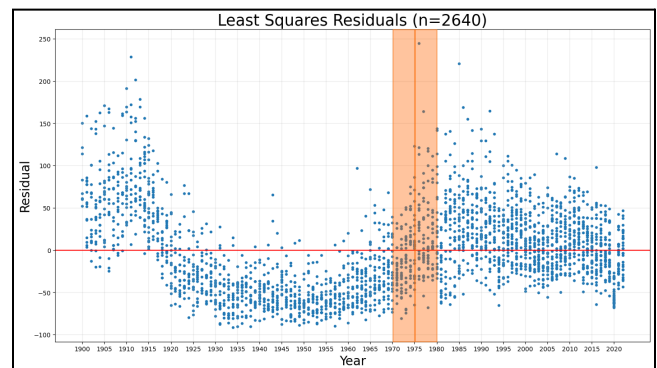
Alternative hypothesis: $\mu_1 - \mu_2 \neq 0$

As written above, we used pyBaseball to extract team batting data from 1900 to 2022. We then mined the SB column, averaged each team by year, and plotted by year. However, we wanted to clean the data further by only analyzing the years surrounding 1975. Unlike before, when there was a rule put into effect in a single year, teams moving into larger stadiums happened over the course of several years. We decided 1975 would be a good year to separate and clean our data. This was after the Phillies, Reds and Pirates switched stadiums. This way, no extraneous variables in the game would impact our results. The summary data collected for each team and each year are as follows.

SB (1970-1974)	
n	120
mean	83.3583
standard deviation	33.0716
min	17.0000
max	172.0000
variance	1093.7277

SB (1975-1979)	
n	126
mean	115.7540
standard deviation	48.6203
min	28.0000
max	341.000
variance	2363.9310

We then graphed the least squares residuals. The stadium change happened roughly in 1975, so we wanted to look at what percentage of residuals were over/under the lifelong trendline, which is highlighted in orange on the histogram. Each histogram was separated into years before and after 1975.



Of the 120 residual points between 1970-1974 (red), 33.33% of them were at or above the lifelong trendline. Of the 126 residual points between

1975-1979 (green), 61.91% of them were at or above the lifelong trendline.

To test our null hypothesis, we used a difference of means hypothesis test (two tailed, $\alpha < .05$). This decision comes down to the t-statistic value. If it is outside the critical values from the t-table below, we reject the null. If it is inside, we fail to reject the null hypothesis. Using the equation: $t = (\bar{x}_1 - \bar{x}_2) / SE$, we found the test statistic to be $t = -6.1358$. Using a T-table, $P(t \leq -6.1358) \approx 0$. Therefore, $P(t > -6.1358) = 99.999$. The two tailed P value is < 0.0001 . Since $\alpha = .05$, this P value is considered to be significantly significant, and we reject the null hypothesis.

Experiment 3

Null hypothesis: MLB testing players for steroids in 2003 onward had no impact on team home runs throughout a season

Null hypothesis: $\mu_1 - \mu_2 = 0$

Alternative hypothesis: $\mu_1 - \mu_2 \neq 0$

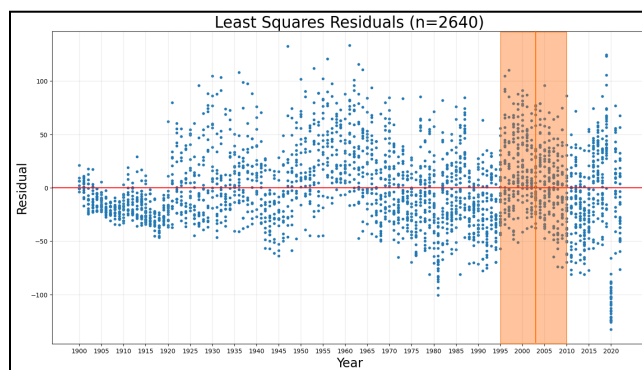
As written above, we used pyBaseball to extract team batting data from 1900 to 2022. We then mined the HR column, averaged each team by year, and plotted them by year. However, we wanted to clean the data further by only analyzing the years surrounding 2003, when the rule was put into effect. This way, no extraneous variables in the game would impact our results. The summary data collected for each team and each year are as follows.

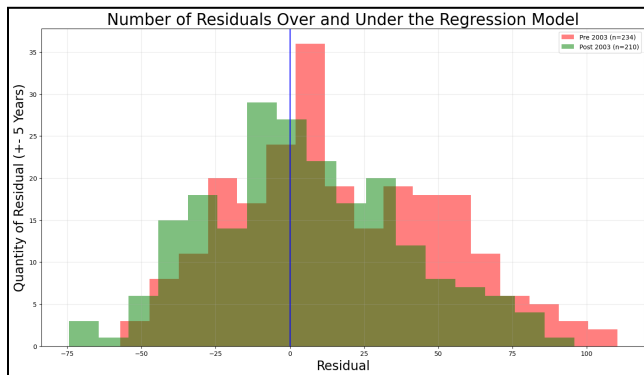
HR (1995-2002)	
n	234
mean	173.0128
standard deviation	35.3272

min	94.0000
max	264.0000
variance	1248.0127

HR (2003-2009)	
n	210
mean	171.1333
standard deviation	32.7115
min	94.0000
max	260.0000
variance	1070.0396

We then graphed the least squares residuals for this rule change. The rule went into effect in 2003, so we wanted to look at what percentage of residuals were over/under the lifelong trendline, which is highlighted in orange on the histogram. Each histogram was separated into years before and after the rule went into effect.





Of the 234 residual points between 1995-2002 (red), 66.24% of them were at or above the lifelong trendline. Of the 210 residual points between 2003-2009 (green), 54.76% of them were at or above the lifelong trendline.

To test our null hypothesis, we used a difference of means hypothesis test (two tailed, $\alpha < .05$). This decision comes down to the t-statistic value. If it is outside the critical values from the t-table below, we reject the null. If it is inside, we fail to reject the null hypothesis. Using the equation: $t = [(\bar{x}_1 - \bar{x}_2) / SE]$, we found the test statistic to be $t = 0.5820$. Using a T-table, $P(t \leq 0.5820) = 0.71957$. The two tailed P value equals 0.5625. Since $\alpha = .05$, this P value is not considered to be significantly significant, and we fail to reject the null hypothesis.

DISCUSSION

One critical challenge we encountered when analyzing sabermetrics is that many have not been recorded for more than seven years. The 'Statcast-Era' began in 2015 when Major League Baseball began tracking several advanced analytics. These include pitching velocities, baserunning speeds, fielding percentages and much more. This made it difficult to conduct a proper analysis of the data when there is little to begin with.

Additionally, pybaseball limits the amount of pull requests by source for a certain time period. This posed a minor problem for making rapid adjustments or quick changes. Often, we had to wait for our time period to reset and we could pull more data. It may be

helpful to store the data locally and then run tests. This can only work if the dataset is small enough to be reasonably stored locally.

CONCLUSION

The three experiments had the following results:

Experiment 1:

Fail to reject the null hypothesis that MLB lowering the pitching mound in 1968 from 15 inches to 10 inches had no impact on team strikeouts per 9 innings ratio.

Experiment 2:

Reject the null hypothesis that several baseball teams moving into stadiums with larger fields in the 70s had no impact on team successful stolen bases.

Experiment 3:

Fail to reject null hypothesis that MLB testing players for steroids in 2003 onward had no impact on team home runs throughout a season

In Experiment 1's case, there are several reasons why we failed to reject the null hypothesis. Our leading theory is that strikeouts are not an absolute performance measure when it comes to pitching. Pitchers can have very good strikeout rates and still perform poorly by giving up home runs and hits. There are a variety of other possible statistics that could have resulted in a rejection of the null hypothesis. In a future experiment, we would like to look at other sabermetrics, such as BABIP (batted balls in play) or OBA (opponent batting average). These metrics could offer a more absolute indication of pitching performance and result in a rejection of the null hypothesis. However, experimenting on these metrics would require an alteration in the experiment itself, as these metrics haven't been recorded as far back as 1900.

In Experiment 2's case, it remains fairly clear through both the graphs and analysis that there was a significant difference in the means. However, we remained curious about these results. It's possible that teams moving into larger stadiums with bigger fields were not the direct cause of the increase in stolen bases. We decided to do some follow up research on if this was partly the case. According to John McMurray at SABR, it appears that it was. "With

the opening of several large stadiums in the early 1970s...it is easy to see why teams emphasized speed over home run hitting power. Further, the attraction of scoring runners from second base on singles hit on Astroturf surely influenced teams to place an additional premium on fast players.” [3]. The switch from grass to astroturf was not something we had considered in our research, and we feel confident saying that it could have been another significant factor.

For Experiment 3, we failed to reject the null hypothesis and have a leading theory as to why. The steroid era was a time filled with infamous players playing with performance enhancing drugs. However, the number of people taking these PEDs remained relatively low. We believe that when steroid testing began in 2003, only a handful of players were impacted. On a team based scale, like the one we measured, there wasn’t a significant decrease in home runs because most people weren’t on PEDs. However, it will be impossible to prove this as the number of players who were on PEDs has yet to be proven. All we can rely on is a speculating theory. If a statistic like home runs per player existed in the pyBaseball library, it would certainly be a focus for a future experiment.

Baseball remains unique in the fact that it has one of the most documented histories in sports. Through the past 122 years, almost every player and team has recorded statistics. Trends in these statistics have fluctuated and been influenced by many events and rule changes as the sport evolves. Often, these rule changes are applied to have a direct impact on trends of certain statistics. However, as we have proven in our experiments, it’s not always the case that a statistic is impacted. This will have to be something MLB takes into consideration as they continue to modify and change the sport.

APPENDIX

Division of Work:

Python Setup - Nathan

Project Proposal Report - Nathan

Project Proposal Slides - Nathan

Coding Experiments - Nathan

Project Checkpoint Report - Nathan/Thomas

Project Checkpoint Slides - Nathan/Thomas

Coding Evaluation and Plots - Nathan
Project Final Slides - Nathan/Thomas
Project Final Report - Nathan/Thomas
Code Architecture Submission - Nathan

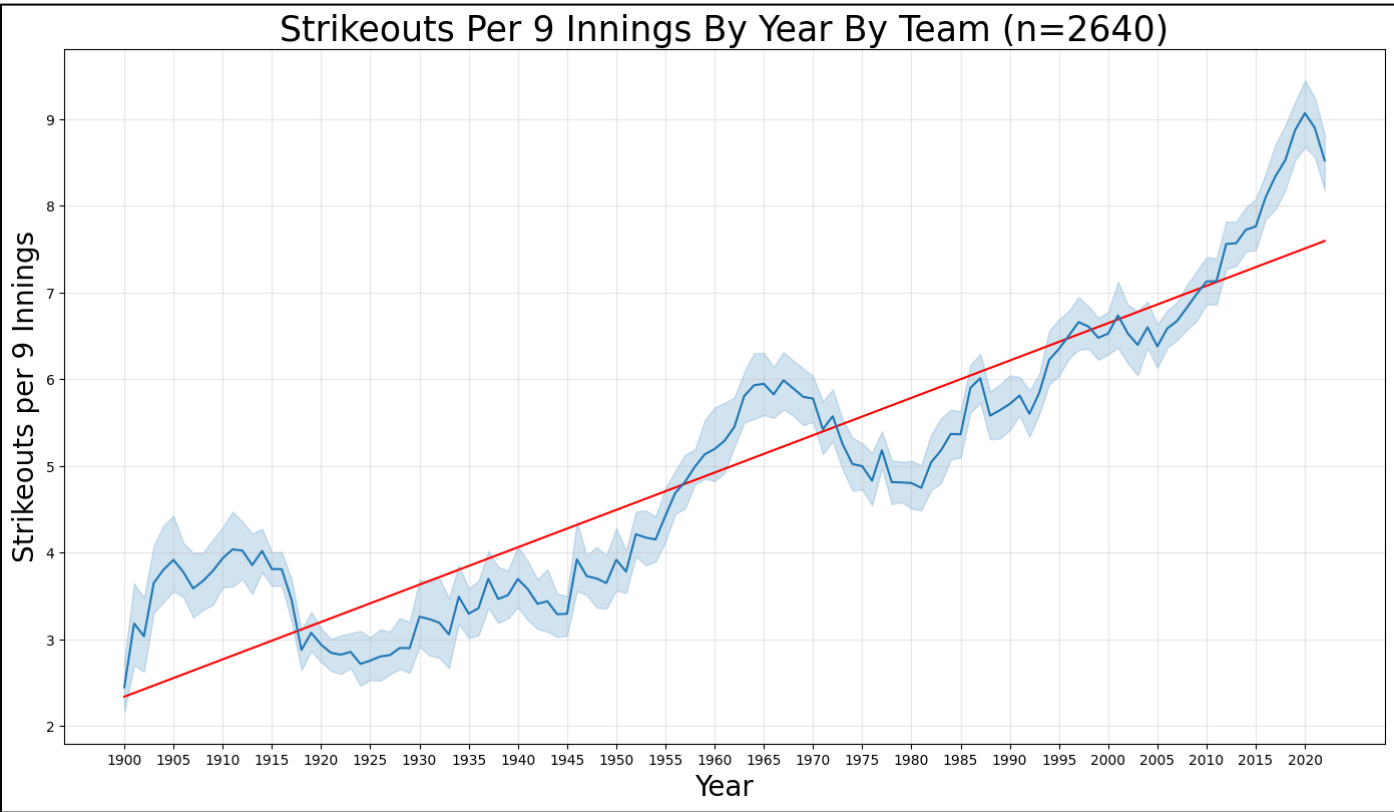
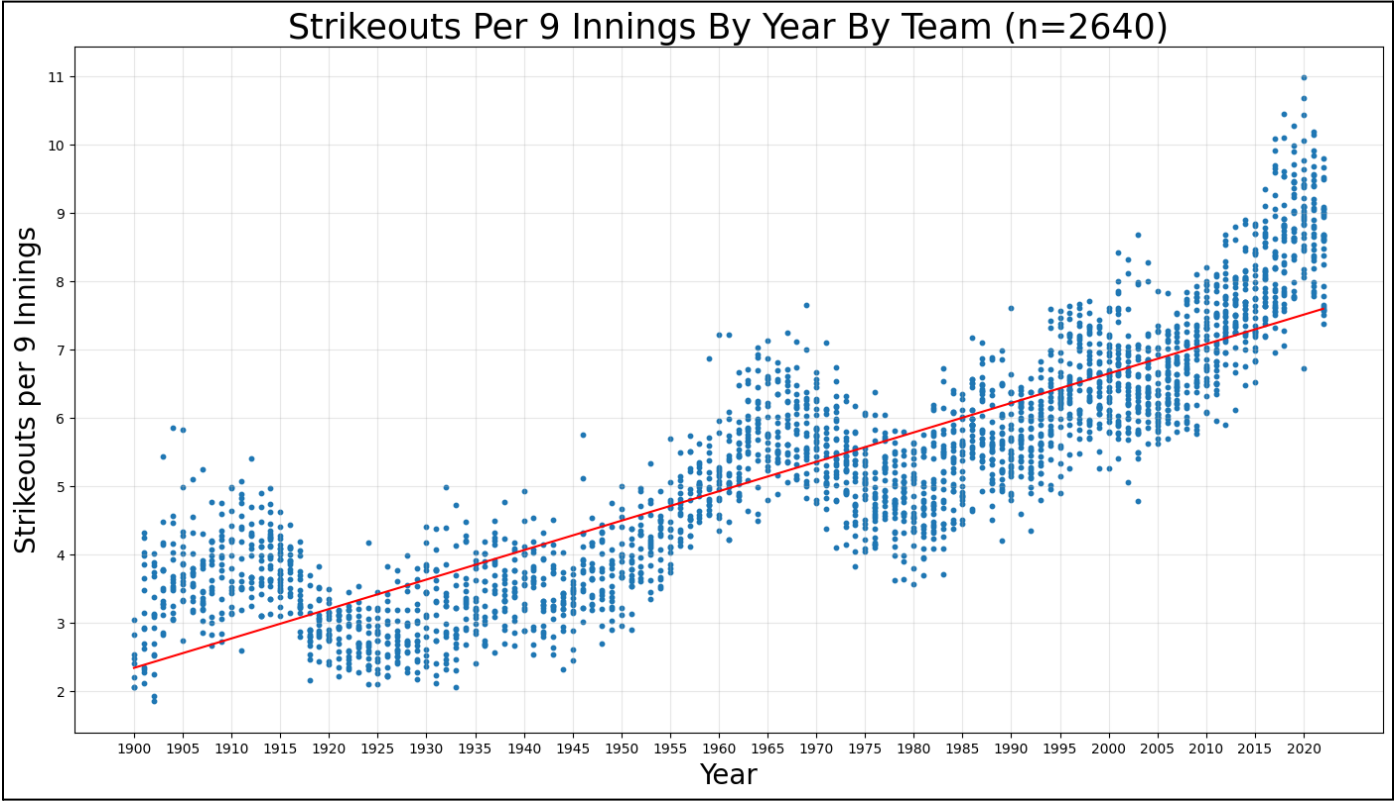
On my honor, as a University of Colorado Boulder student, I have neither given nor received unauthorized assistance.

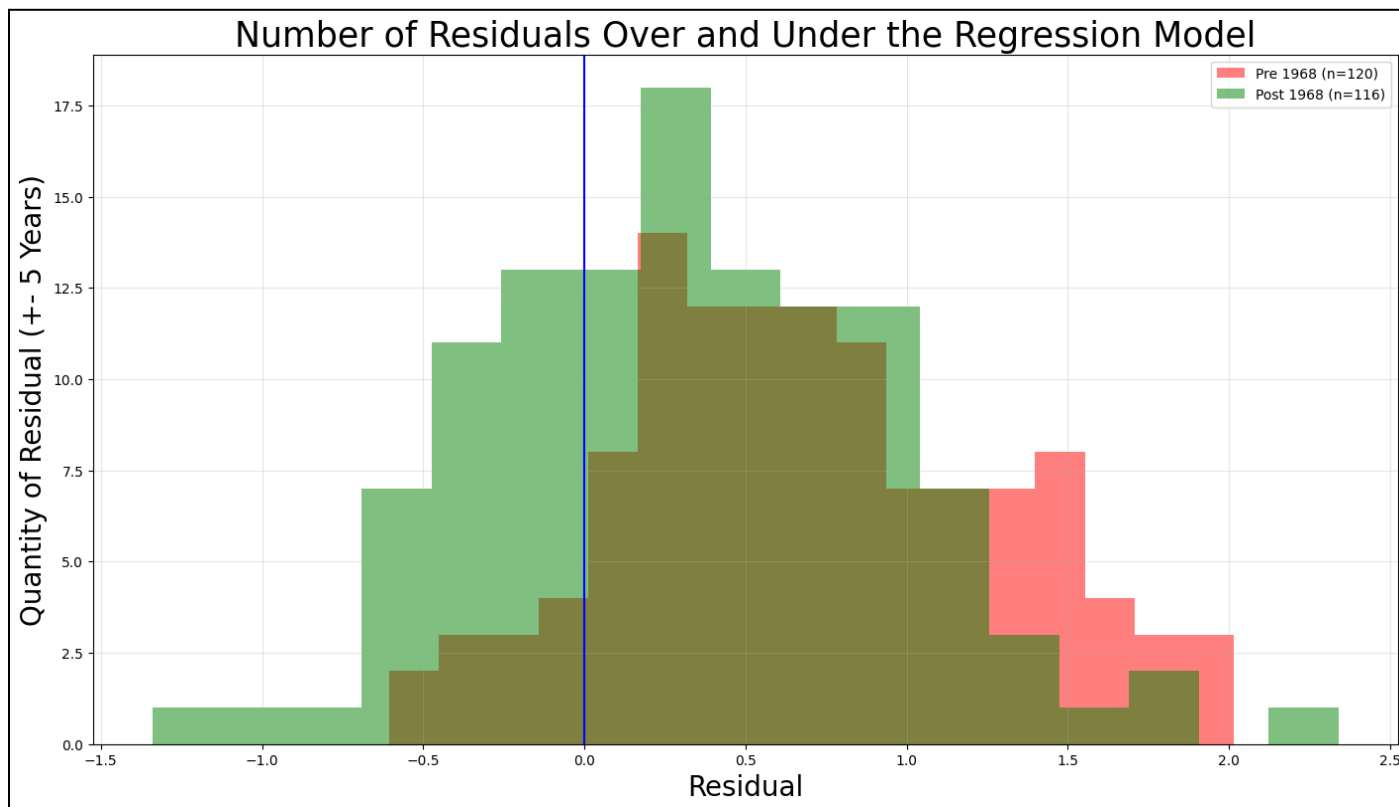
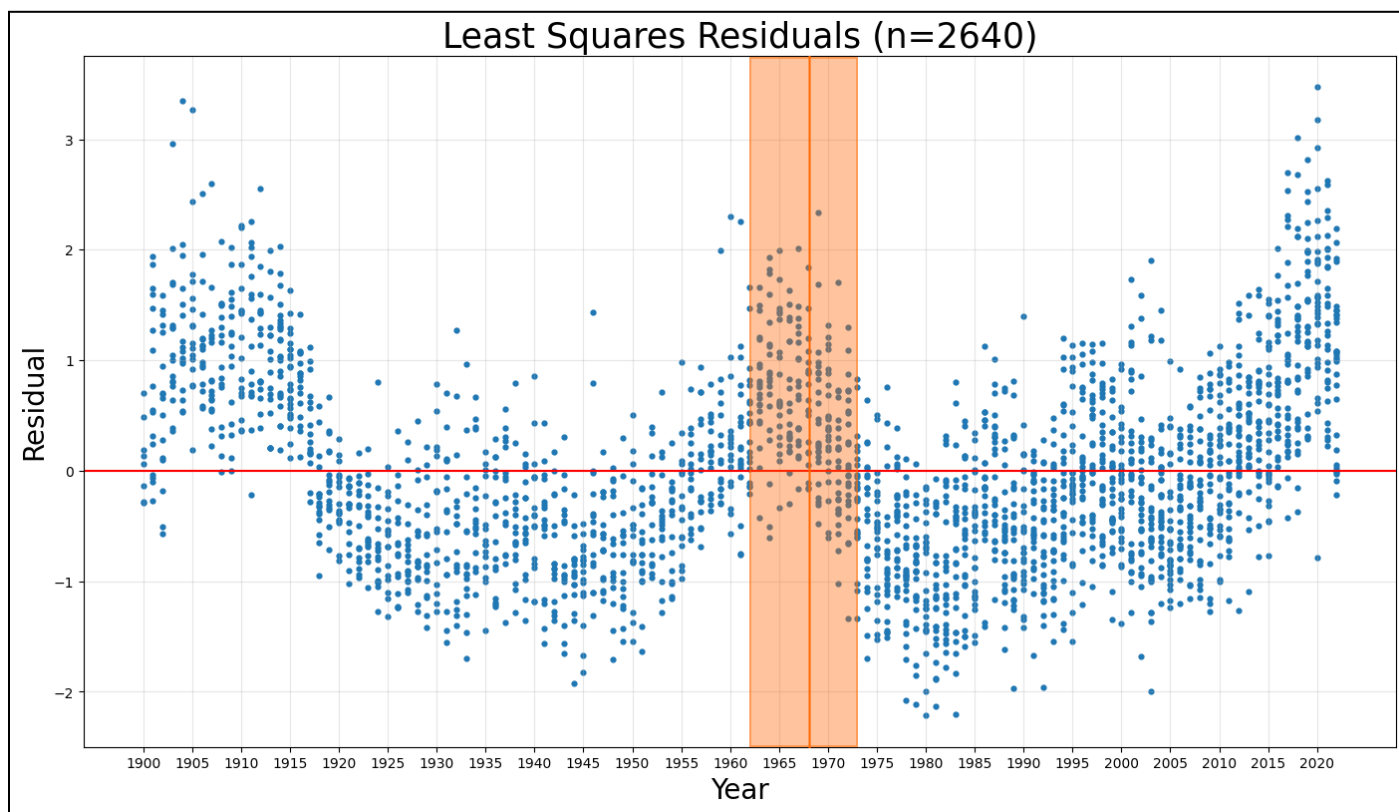
See images below for full sized renders.

REFERENCES

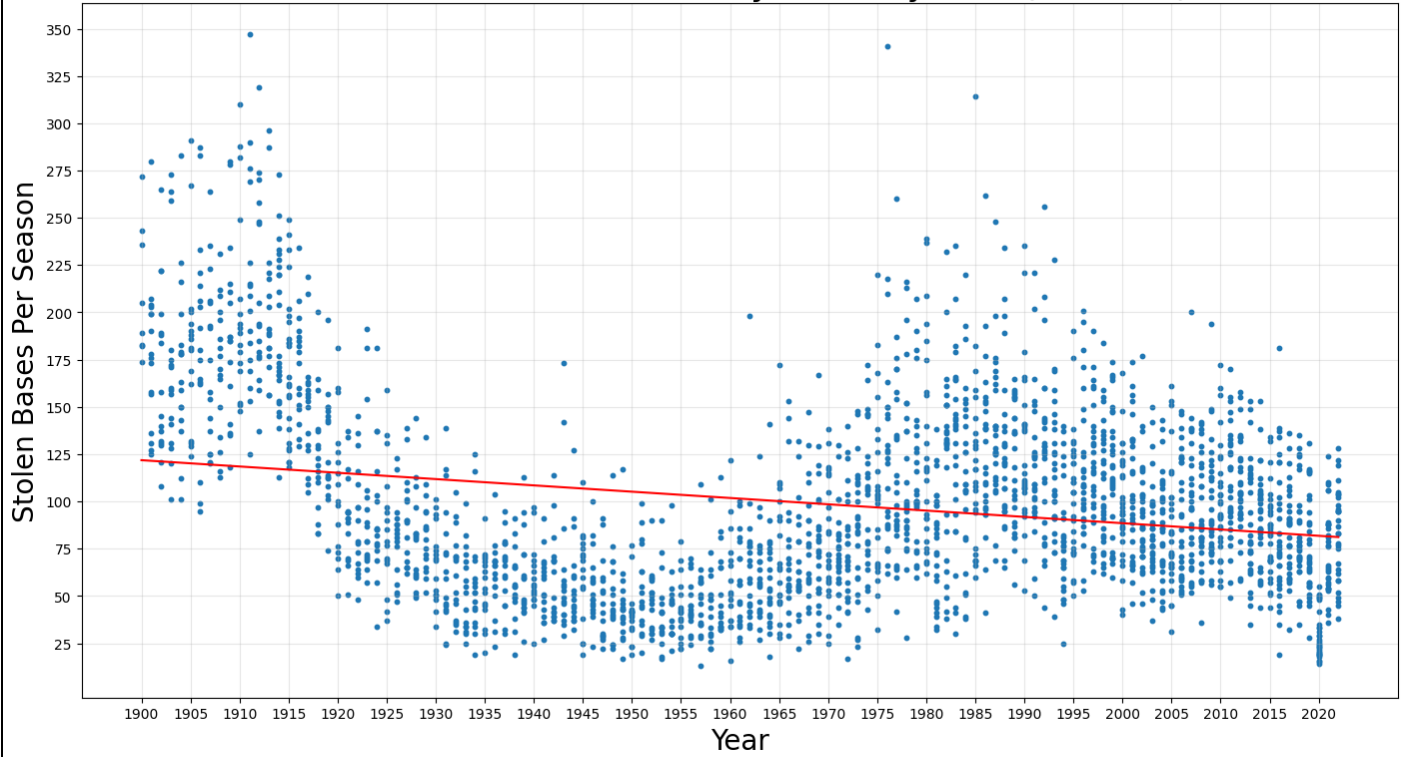
- [1] Steinberg, Leigh. “Changing the Game: The Rise of Sports Analytics.” *Forbes*, Forbes Magazine, 18 Aug. 2015,
- [2] Beneventano, Philip, Paul D. Berger, and Bruce D. Weinberg. "Predicting run production and run prevention in baseball: the impact of Sabermetrics." and *Int J Bus Humanit Technol* 2.4 (2012): 67-75.
- [3] McMurray, John. “Examining Stolen Base Trends by Decade from the Deadball Era through the 1970s.” Society for American Baseball Research, Society for American Baseball Research, Oct. 2015, <https://sabr.org/journal/article/examining-stolen-base-trends-by-decade-from-the-deadball-era-through-the-1970s/>.

FULL SIZE PLOT IMAGES

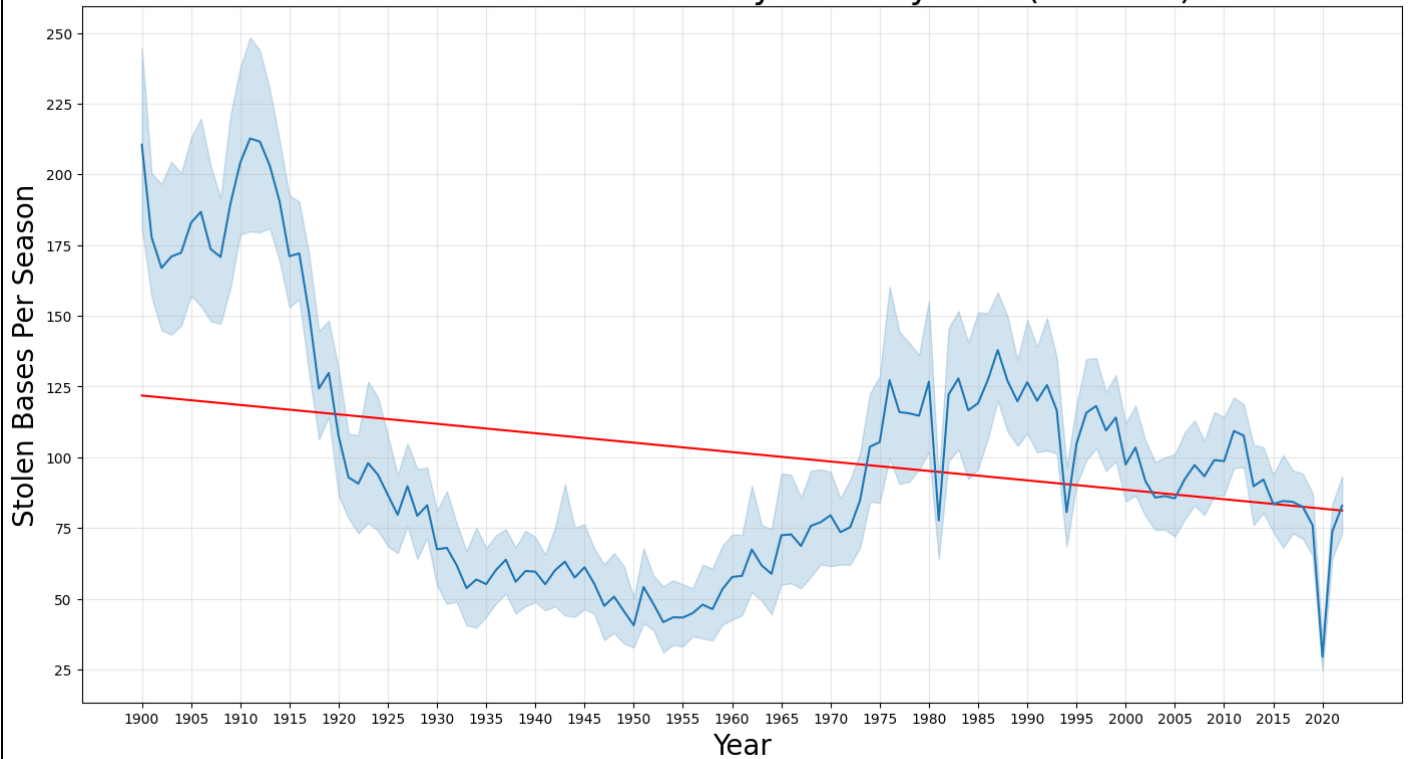


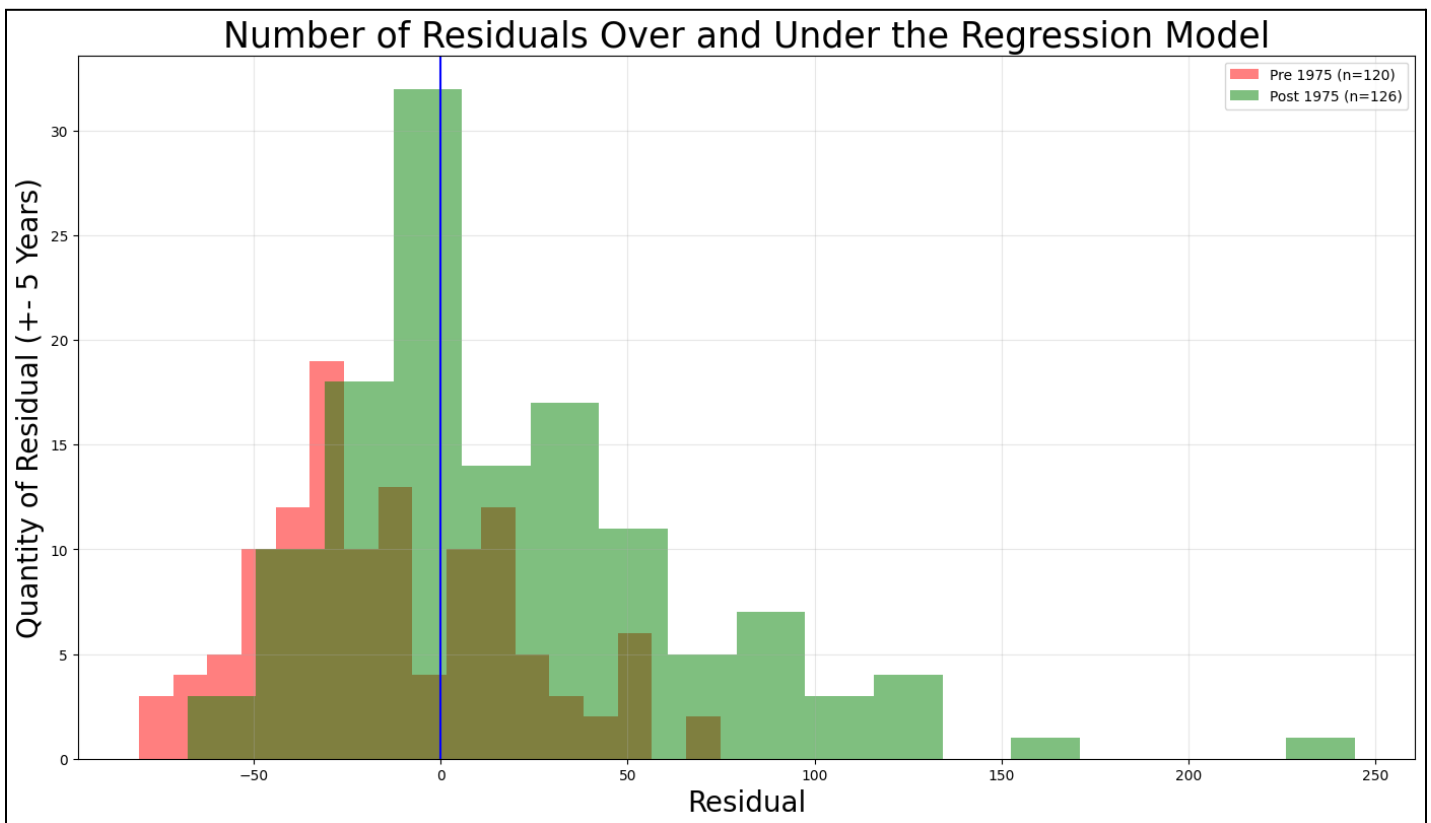
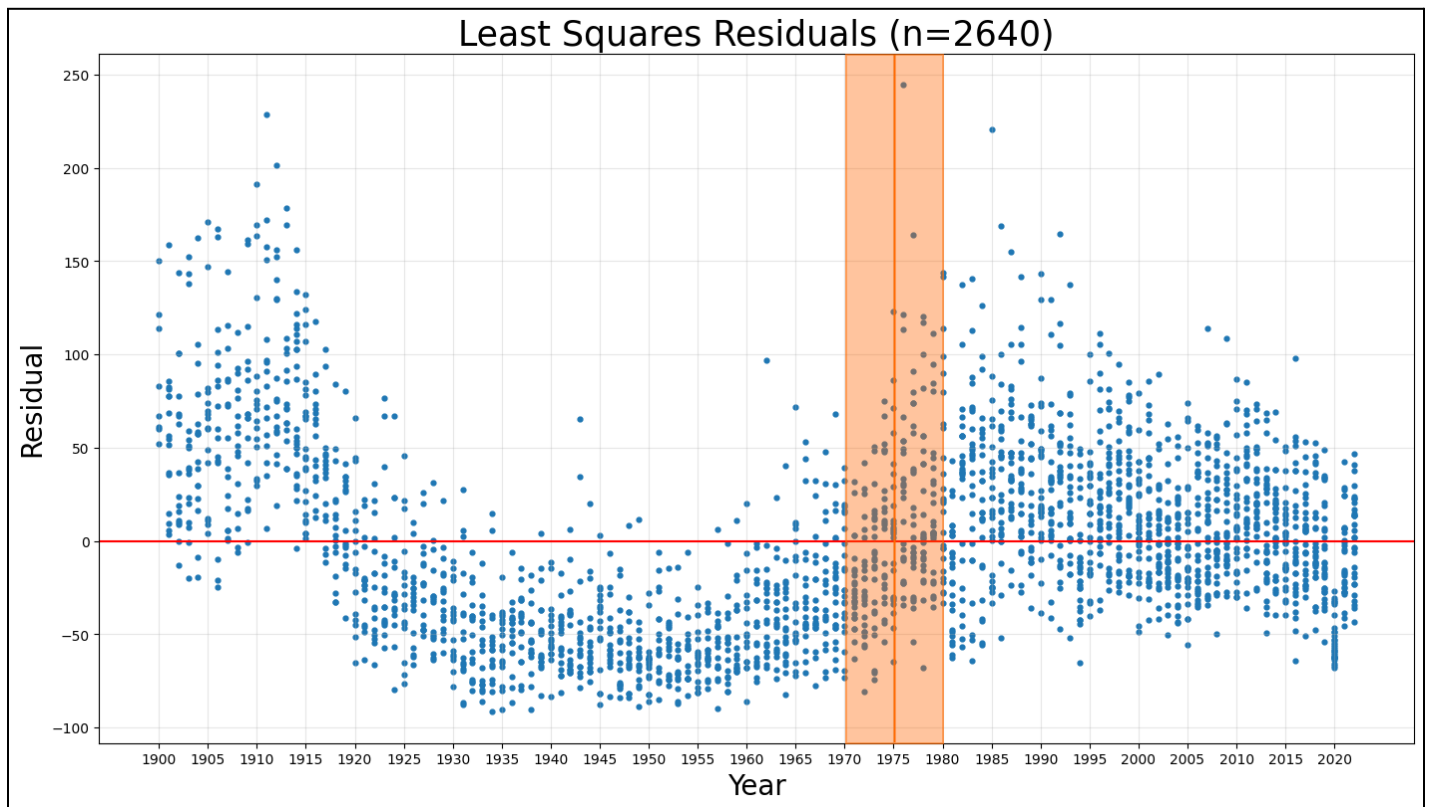


Stolen Bases Per Season by Team by Year (n=2640)

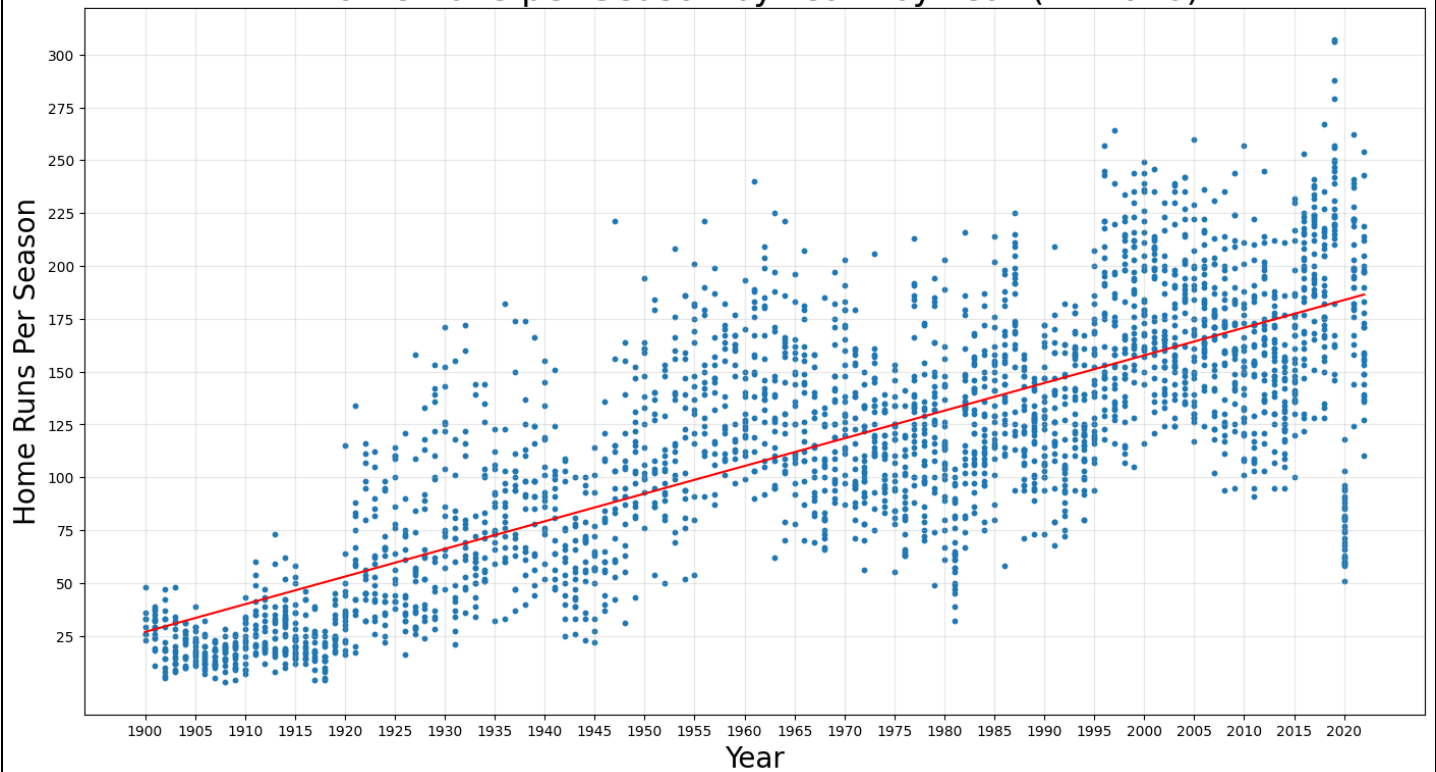


Stolen Bases Per Season by Team by Year (n=2640)





Home Runs per Season by Team by Year (n=2640)



Home Runs per Season by Team by Year (n=2640)

