# Baseball's Offensive Evolution Throughout the Eras

By Jacob Adamsky

Link to website: https://jacobadamsky.github.io/final

## Overview and Motivation



Yes, that's me playing baseball. I may not make the best faces when I'm playing, but I've almost always enjoyed the game. It has been hard for the past couple years, missing out on both my senior year of high school baseball and freshman year of college baseball, so being able to do a project on baseball is kind of therapeutic in a way. As for why I chose to focus on batting statistics, I used to be a two-way player, meaning I pitched and played the field (when I wasn't pitching), so I was reminiscing on the good old days when I was still allowed to hit.

As for the three plots I made, I plotted batting average against year to show the progression and concentration of different batting averages as time passed, hits + walks to runs to see how a higher OBP (on base percentage) increased a player's potential to score, and strikeouts to walks to see if there is a correlation between walk rate and strikeout rate. The database used was Sean Lahman's Baseball Statistics downloaded as a collection of CSVs (there is a MySQL version, but I never had success with it).

# Related Work

One major related website that motivated me to do this was baseball-reference.com. I don't believe they have any visualizations for their data other than tables with player statistics, so I wanted to create some visualizations using the same data. Other than that, there wasn't anything that inspired me to do my final project on this subject.

# Questions

There were three questions that I based my plots on: "How much of a difference was there in batting average as baseball evolved?", "How much correlation is there between a high walk and hit rate with how often a player scores?", and "Is there a correlation between a higher strikeout rate and lower walk rate, and vice versa?". At first, I was just plotting the correlation between batting average and years passed, but this was far too simple and couldn't get all the project requirements in. I added the second question, and later the third as I looked at the statistics database more.

# Data

D3.js is surprisingly effective when it comes to filtering data, only taking a couple of seconds to filter out extraneous data (statistics from years not in the era being plotted), as well as plotting it. I did have some trouble searching for players at first as it was taking over a minute but realized that I had made a small mistake that was exponentially increasing the time taken to search for a player.

As for the data source, all of it comes from Sean Lahman's Baseball Statistics Database. It has several different access methods, but I wanted to have the ability to use GitHub pages, so I downloaded the CSV version. It contains tens of thousands of player entries, as well as hall of fame status, all-star statistics, and much more.

# Exploratory Data Analysis

I only looked at one visualization and didn't receive much inspiration from it. http://dgwartney.github.io/tableau/ plots a few statistics from three steroid-era players (Barry Bonds, Mark McGwire, and Sammy Sosa) who were all known steroid users. Aside from this, I came up with pretty much everything about my visualization from scratch.

# Design Evolution

At first, I was considering doing a line plot, but with the massive quantity of data, side-scrolling through all the results could take literal hours. I decided on a more condensed approach of using three scatter plots which, in my opinion, visualized the correlation between certain statistics much better than a line plot or another option would have. As for my proposal, I didn't have anything in it about what kind of plotting I was going to do, so I didn't deviate from it in any way.

# Implementation



*Figure 1: Home page containing database of all player information (not statistics)*
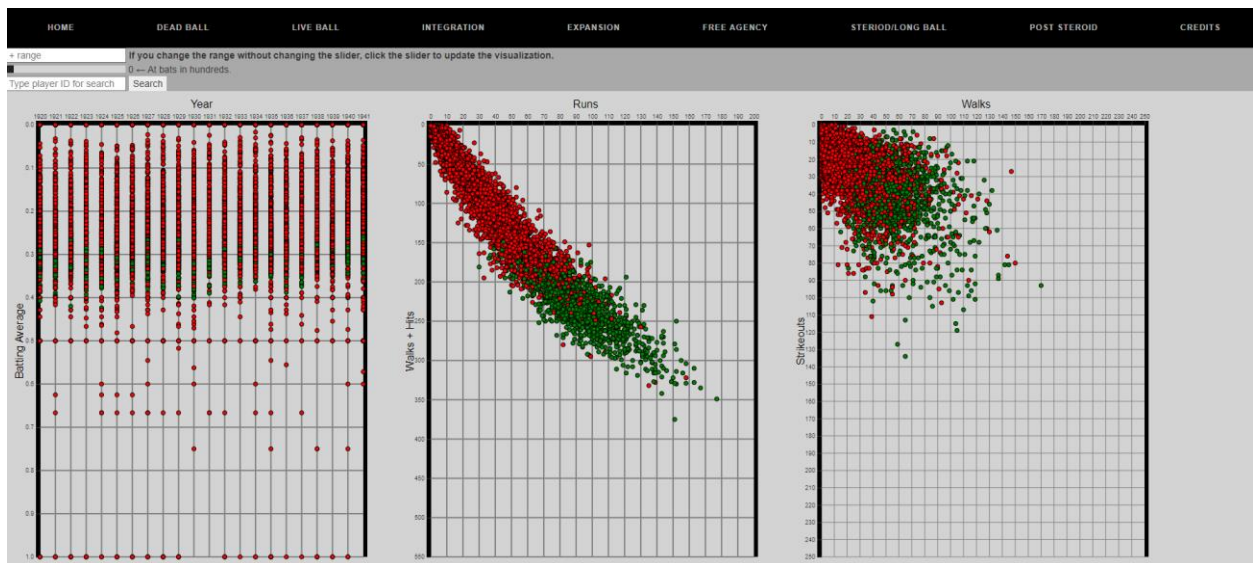
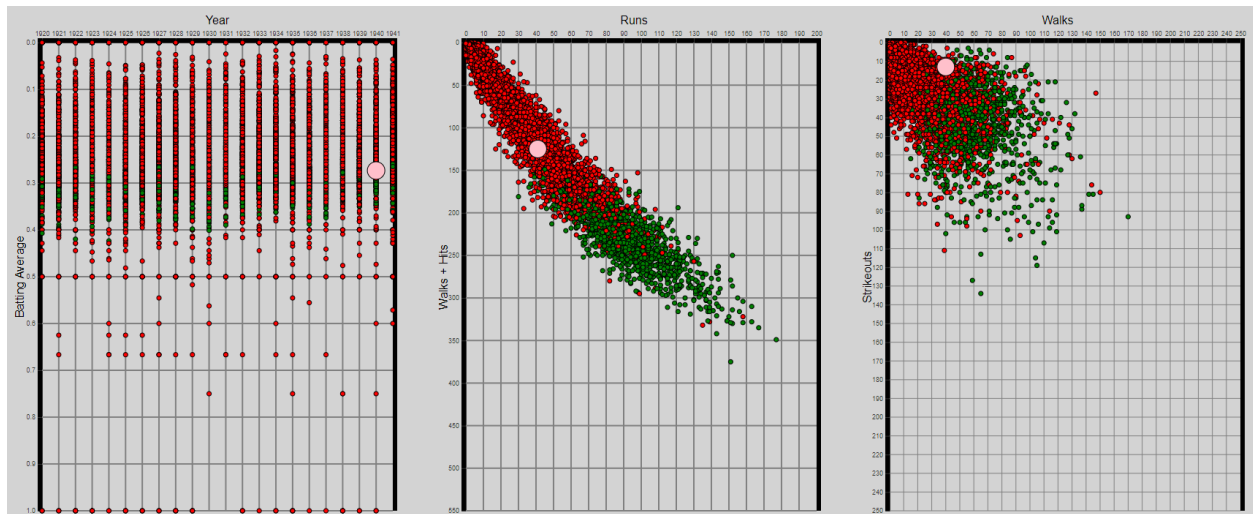*Figure 2: Initial view of visualization*
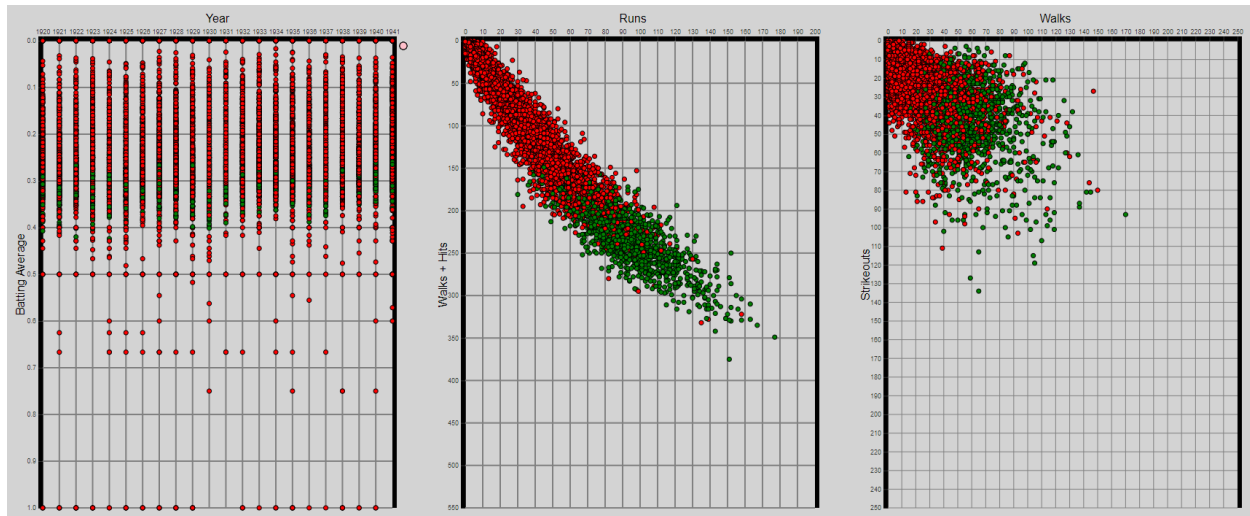


*Figure 3: Clicking on a specific point*

*Figure 4: Clicking on a point you already clicked on moves it to the right side of the plot (limit of 250 per, can only move to 3 plots before space runs out)*
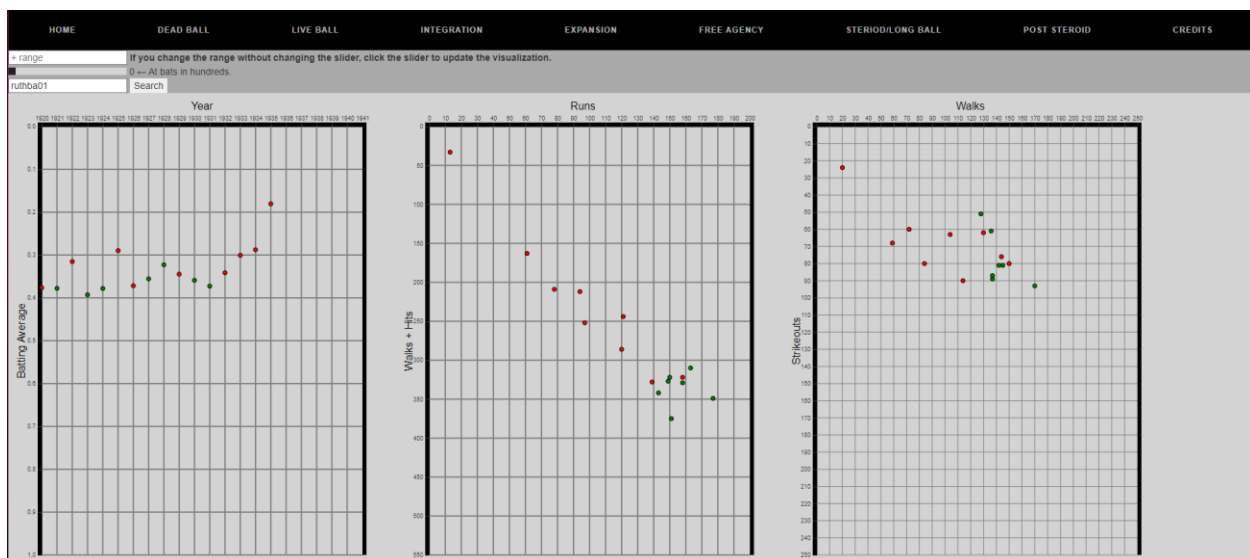


*Figure 5: Visualization after searching for a player. In this case, ruthba01 (Babe Ruth)*
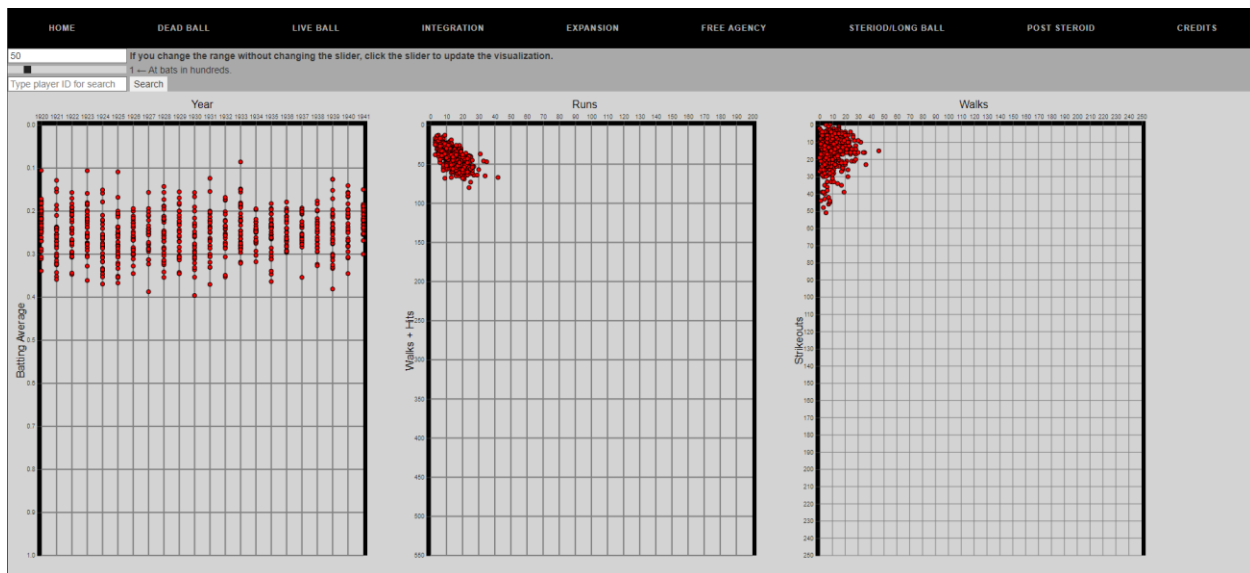
*Figure 5: Only showing players who have 100-150 at balls (inclusive).*

Each picture here shows a different step in the interaction process. At first, the visualization is displaying all the data. You can then search for a specific player, click on random players in the visualization to highlight them, or filter players based on number of at bats. All three plots are linked, so clicking on a point from plot 1 will highlight that same point on the other plot, and the same goes for plots 2 and 3.

# Evaluation

As I was already familiar with the correlation between different statistics and the difference in statistics based on era, I didn't necessarily learn much about the data. Though that may be the case, plotting the data gave me a better appreciation for what it represented and how different each era of baseball was offensively. There are dozens of visualizations that could be created with the database, so I'm sure there is so much more that I could learn and understand from visualizing it.

As for functionality, I feel that my visualization works well. Considering that each plot has 20,000+ points in it, it only takes about a second for the data to load initially. Searching for players is fast as well, only taking about a couple of seconds to search through all three plots and hide all extraneous points. As for the transitions, they can be a little bit laggy, but I'm not sure if that is a browser issue or a problem with my code as I wasn't able to do anything about it.