# Baseball Salary Analysis

Edgar Hernandez, Ryan Sims, Jessica Tomas

2022-08-03

## Contents

## Introduction

We will be examining what statistics related to a baseball player's performance have a significant effect on the salary of a player. This is particularly interesting as it is sort of the flip-side of the coin to the traditional Sabermetrics employed by the Oakland Athletics in the 1990's to analyze player statistics to try to assemble a winning team. We will be instead be leveraging the dataset to instead work for the players, hopefully determining which statistic(s) a player should focus on improving in order to increase their compensation.

Before we begin our analysis, we must first load the data using the `Lahman` library:

```
library(Lahman)
```

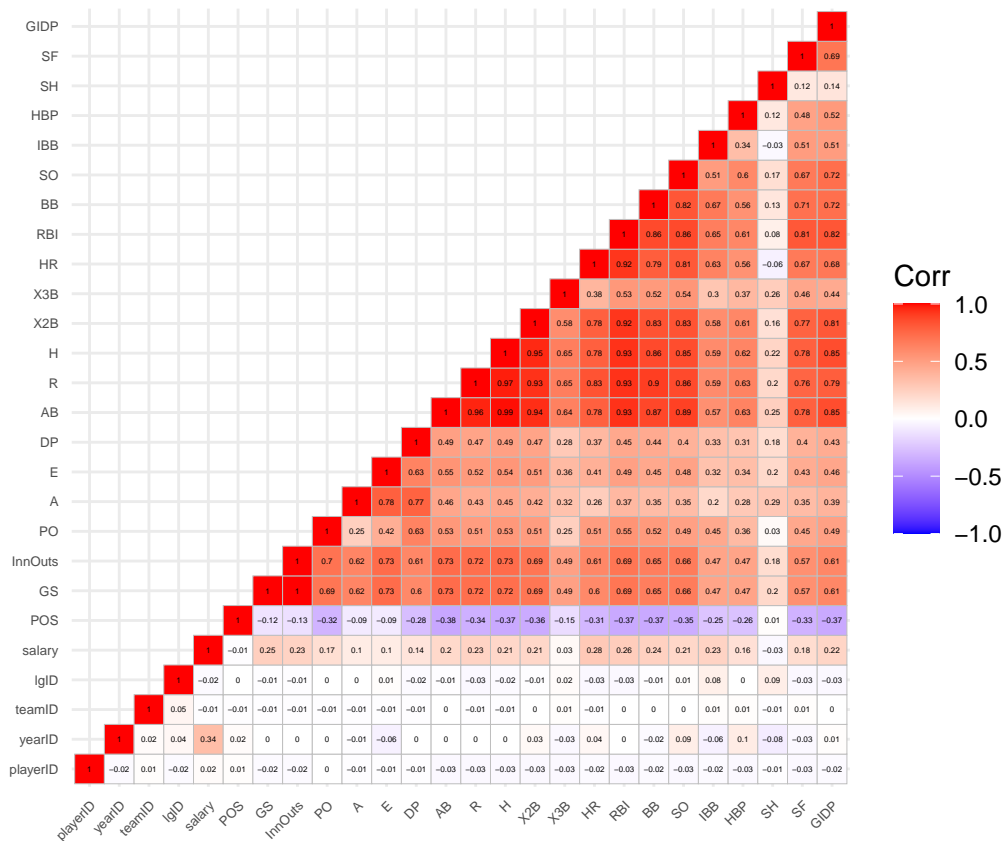We then join the Salaries, Fielding, and Batting tables:

### Methods

First, we visualize pairwise correlations between `Salary` and other variables:

- Calculating the correlations:

```
salary_data = as.data.frame(lapply(salary_data, as.integer))
salary_data_cor = cor(as.data.frame(salary_data))
```

- Plotting the relationships:

```
library(ggcorrplot)
my_plt = ggcorrplot(salary_data_cor, lab = TRUE, lab_size = 1, show.diag = TRUE, type = "lower")
my_plt + theme(axis.text.x = element_text(size = 5), axis.text.y = element_text(size = 5))
```

```r
salary_data = as.data.frame(salary_data)
baseball_fit = lm(salary ~ ., data = salary_data)
baseball_aic = step(baseball_fit, direction = "backward", trace = 0)

range(salary_data$yearID)

## [1] 1985 2016
```