

Predicting Baseball Hall Of Fame Inductions

Michael Hirsch

ILLC, University of Amsterdam

michaelahirsch@gmail.com

Abstract

Every year, the Baseball Writers Association of America votes on a new Hall of Fame class. Each ballot consists of 10 votes, and players need to appear on 75% of ballots in order to be inducted into the hall. Players have 15 years to get inducted, and are no longer eligible if that time period has passed. Here we attempt to classify hall of fame batters based on their career statistics using decision trees and an ensemble method, random forests.

1. Introduction

Major League Baseball (MLB) has been keeping thorough records of batting, pitching, and fielding statistics since its inaugural season in 1869. Recently, with the advent of sabermetrics by the Society for American Baseball Research (SABR), many new metrics building on traditional statistics were created. Such metrics caught the eye of statisticians, as these new metrics allowed for the creation of even more powerful predictive models.

In this paper, we are investigating what makes a Hall Of Fame (HOF) batter. There have been over 20,000 Major League baseball players in the history of the organization, and only 211 of them have been inducted into the HOF. There are several ways in which a player can be inducted, and we only concern ourselves with players inducted from Baseball Writers Association of America (BBWAA) ballots, as this process is regulated. Often, fans and players feel that a worthy candidate is unfairly denied entry into the HOF because of voter bias, stacked ballots, or negative associations with performance enhancing drugs. The model proposed in this paper should be able account for these cases.

2. Data

Data was collected from the Lahman Baseball Database, a freely available database that has been continuously updated since 1994 with the help of SABR and many individual researchers. There are many data tables available for download, but the ones that we are focusing on are *Batting.dat* and *HallOfFame.dat*. Batting statistics were supplied on a by-year basis, so a player's career statistics were computed by aggregating his yearly results. The batters' data was then merged with hall of fame voting results using the player's ID string. Attention is paid to 18 different batting statistics over the course of the player's career:

Below is a list of features used

1. G: games played
2. AB: at-bats
3. R: runs
4. H: hits
5. X2B: doubles
6. X3B: triples
7. HR: homeruns
8. RBI: runs batted in
9. SB: stolen bases
10. CS: caught stealing
11. BB: walks
12. SO: strike outs
13. IBB: intentional walks
14. HBP: hit by pitches
15. SH: sacrifice hits
16. SF: sacrifice flies
17. GIDP: grounded into double plays
18. Inducted: classifier for HOF induction

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Pellentesque quis interdum velit. Nulla tincidunt sem quis nisi molestie nec hendrerit nulla interdum. Nunc at lectus at neque dapibus dapibus sit amet in massa. Nam ut nisl in diam consectetur dignissim. Sed lacinia diam id nunc suscipit vitae semper lorem semper. In vehicula velit at tortor fringilla elementum aliquam erat blandit. Donec

pretium libero et neque vehicula blandit. Curabitur consequat interdum sem at ultrices. Sed at tincidunt metus. Etiam vulputate, lacus eget fermentum posuere, ante mi dignissim augue, et ultrices felis tortor sed nisl.

Maecenas [1] fermentum [2] urna ac sapien tincidunt

[1] J. M. Smith, A. B. Jones, Book Title, Publisher, 7th edition, 2012.

[2] A. B. Jones, J. M. Smith, Article Title, Journal Title 13 (2013) 123–456.