

---

**CSE 519: Data Science**  
**Steven Skiena**  
**Stony Brook University**

---

Lecture 1: Introduction to Data Science

---

# What is Data Science?

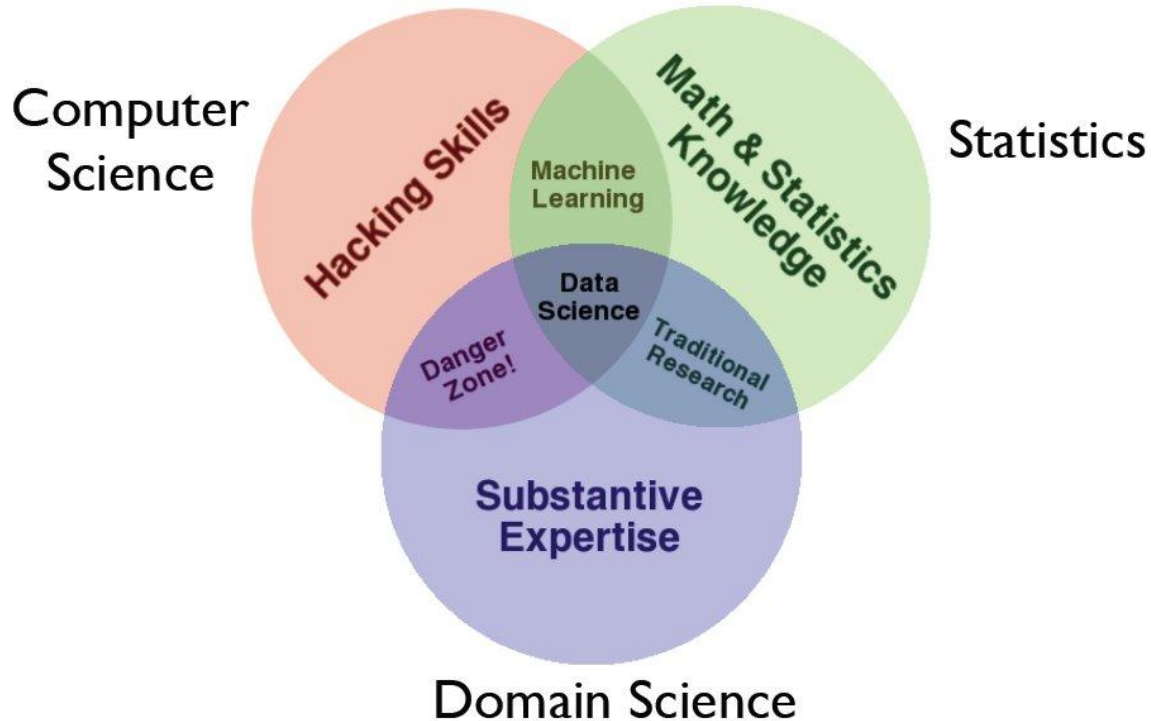
---

Like any emerging field, it isn't yet well defined, but incorporates elements of:

- Exploratory Data Analysis and Visualization
  - Machine Learning and Statistics
  - High-Performance Computing technologies for dealing with scale.
-

# Skill Sets for Data Science

---



# Appreciating Data

---

Computer Scientists do not naturally appreciate data: it's just stuff to run through a program.

The usual way to test algorithm performance is to run the implementation on “random data”.

But interesting data sets are a scarce resource, which requires hard work and imagination to obtain.

---

# Computer vs. Real Scientists (1)

---

- Scientists strive to understand the complicated and messy natural world, while computer scientists build their own clean and organized virtual worlds. Thus:
  - Nothing is ever completely true or false in science, while everything is either true or false in Computer Science / Mathematics.
-

# Computer vs. Real Scientists (2)

---

- Scientists are data-driven, while computer scientists are algorithm-driven.
  - Scientists obsess about discovering things, which computer scientists invent rather than discover.
  - Scientists are comfortable with the idea that data has errors; computer scientists are not.
-

# Genius vs. Wisdom

---

Software developers are hired to produce code.

Data Scientists are hired to produce insights.

Genius shows in finding the right answer!!!

Wisdom shows in avoiding the wrong answers.

Data science (like most things) benefits more from wisdom than from genius.

---

# Developing Wisdom

---

- Wisdom comes from experience.
- Wisdom comes from general knowledge.
- Wisdom comes from listening to others.
- Wisdom comes from humility, observing how often you have been wrong and why/how.

I seek pass on wisdom, by providing experience on the difficulty of making good predictions.

---



# Developing Curiosity

---

- The good data scientist develops a curiosity about the domain/application they are working in.
  - They talk shop with the people whose data they are working on.
  - They read the newspaper every day, to get a broader perspective on the world.
-

# Asking Good Questions

---

Software developers are not encouraged to ask questions, but data scientists are:

- What exciting things might you be able to learn from a given data set?
  - What things do you/your people really want to know?
  - What data sets might get you there?
-

# Let's Practice Asking Questions!

---

Who, What, Where, When, and Why on the following datasets:

- [Baseball-reference.com](http://baseball-reference.com)
  - International Movie Database (IMDb)
  - Google ngrams
  - NYC taxi cab records
-

# Baseball-Reference.com: biosketch



play index **players** teams seasons managers leaders awards postseason boxes japan nlb minors draft

Mobile Site You Are Here > Home > Encyclopedia of Players > R Listing > Babe Ruth Statistics and

News: s-r blog:KBO Stats back to 1999 - Baseball-Reference.com

Babe Ruth Player Page > Batting Pitching Fielding Minors News Archive (1456) Bullpen Oracle



## Babe Ruth

Like 1,213 people like this.

+25 Recommend this

George Herman Ruth (Babe, The Bambino or The Sultan Of Swat)

Positions: Outfielder and Pitcher

Bats: Left, Throws: Left

Height: 6' 2", Weight: 215 lb.

**Born:** February 6, 1895 in Baltimore, MD

**High School:** St. Mary's HS (Baltimore, MD) (All Transactions)

**Debut:** July 11, 1914 (Age 19.155)

**Rookie Status:** Exceeded rookie limits during 1915 season [\*]

**Teams** (by GP): Yankees/RedSox/Braves 1914-1935

**Final Game:** May 30, 1935 (Age 40.113)

**Inducted** into the Hall of Fame by BBWAA as Player in 1936 (215/226 ballots). Induction ceremony in 1

View [Babe Ruth Page](#) at the Baseball Hall of Fame (plaque, photos, videos).

**Died:** August 16, 1948 in New York, NY (Aged 53.192)

**Buried:** Gate of Heaven Cemetery, Hawthorne, NY

**View Player Bio** from the [SABR BioProject](#)

[About biographical information](#)



S-R: M

## Transactions

**July 9, 1914:** Purchased with [Ernie Shore](#) and [Ben Egan](#) by the [Boston Red Sox](#) from Baltimore (International) for more than \$25000. more than \$25000

**December 26, 1919:** Purchased by the [New York Yankees](#) from the [Boston Red Sox](#) for \$100,000.

**February 26, 1935:** Released by the [New York Yankees](#).

**February 26, 1935:** Signed as a Free Agent with the [Boston Braves](#).

The transaction information used here was obtained free of charge from and is copyrighted by [RetroSheet](#). We attempt to update transactions throughout the season.

## Salaries

Convert to YYYY \$'s

Salaries may not be complete (especially pre-1985) and may not include some earned bonuses

Year	Age	Team	Salary	ServTm (OpnDay)	Sources	Notes/Other Sources
1914	19	Boston Red Sox	\$2,500		? Bill James Historical Abstract	Annualized rate; came up late in season
1915	20	Boston Red Sox	\$3,500		? Bill James Historical Abstract	
1916	21	Boston Red Sox	\$3,500		? Contract at HOF	
1917	22	Boston Red Sox	\$3,500		? Contract at HOF	BJHA: \$5,000; Baseball Timeline \$7,000
1918	23	Boston Red Sox	\$9,000		? Allan Wood, 1918, at 183	Includes \$1,000 midseason raise, \$1,000 WS bonus
1919	24	New York Yankees	\$10,000*		? Michael Haupert research of HOF contracts	Contract at HOF:10000.00.
1920	25	New York Yankees	\$20,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:20000.00.
1921	26	New York Yankees	\$20,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:30000.00,Plus \$5K for '20 and '21 exhibitions, \$50/HR (\$9)m
1922	27	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00.
1923	28	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00.
1924	29	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00.
1925	30	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00.
1926	31	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00.
1927	32	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	S/23/27 AL letter:70000.00.
1928	33	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	S/23/27 AL letter:70000.00.
1929	34	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	S/23/27 AL letter:70000.00.
1930	35	New York Yankees	\$70,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:80000.00.
1931	36	New York Yankees	\$70,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:80000.00.
1932	37	New York Yankees	\$70,000*		? Michael Haupert research of HOF contracts	M. Smelser, Life That Ruth Built, p. 441:75000.00,Plus 25% of all exhibition-game profits
1933	38	New York Yankees	\$80,000*		? Michael Haupert research of HOF contracts	M. Smelser, Life That Ruth Built, p. 456:52000.00,Plus 25% of revenue from in-season exhibitions
1934	39	New York Yankees	\$80,000*		? Michael Haupert research of HOF contracts	1/16/36 TSN, per government report:36696.00,\$35,000 salary plus 25% of exhibition profits
1935	40	New York Yankees	\$75,000*		? Michael Haupert research of HOF contracts	Bill James Historical Abstract:35000.00,Annualized rate; retired early in season
1936	41	New York Yankees	\$52,000*		? Michael Haupert research of HOF contracts	
1937	42	New York Yankees	\$35,000		? Michael Haupert research of HOF contracts	

Career to date (may be incomplete) **\$1,020,000**



# Baseball Questions

---

- How to best measure individual player's skill, value or performance?
  - How fair do trades between teams work out?
  - What is the trajectory of player's performances as they mature and age?
  - To what extent does batting performance correlate with the position played?
-

# Demographic Questions

---

- Do left-handed people have shorter lifespans than right-handers?
  - How often do people return to where they were born?
  - Do player salaries reflect past, present, or future performance?
  - Are heights and weights increasing in the population?
-



# IMDb: Movie Data

All

[Movies, TV & Showtimes](#) [Celebs, Events & Photos](#) [News & Community](#) [Watchlist](#)



**It's a Wonderful Life** (1946) Top 5000

Approved 130 min - Drama | Family | Fantasy -  
7 January 1947 (USA)

**8.7** Your rating: ★★★★★★☆☆ -/10  
Ratings: **8.7/10** from 202,743 users  
Reviews: 632 user | 187 critic

An angel helps a compassionate but despairingly frustrated businessman by showing what life would have been like if he never existed.

**Director:** Frank Capra  
**Writers:** Frances Goodrich (screenplay), Albert Hackett (screenplay), 4 more credits »  
**Stars:** James Stewart, Donna Reed, Lionel Barrymore | See full cast and crew »

[+ Watchlist](#) [Watch Trailer](#) [Share...](#)

[More at IMDbPro »](#)

## Details

[Edit](#)

**Country:** USA

**Language:** English

**Release Date:** 7 January 1947 (USA) [See more »](#)

**Also Known As:** The Greatest Gift [See more »](#)

**Filming Locations:** California, USA [See more »](#)

## Box Office

**Budget:** \$3,180,000 (estimated)

**Opening Weekend:** £49,845 (UK) (19 December 2008)

**Gross:** £682,222 (UK) (24 December 2010)

[See more »](#)

## Company Credits

**Production Co:** Liberty Films (II) [See more »](#)

[Show detailed company contact information on IMDbPro »](#)

## Technical Specs

**Runtime:** 130 min | 118 min (DVD edition)

**Sound Mix:** Mono (RCA Sound System)

**Color:** Color (colorized) | Black and White

**Aspect Ratio:** 1.37 : 1

[See full technical specs »](#)



# IMDb: Actor Data



**James Stewart** (I) (1908–1997)

Actor | Soundtrack | Director



James Maitland Stewart was born on 20 May 1908 in Indiana, Pennsylvania, where his father owned a hardware store. He was educated at a local prep school, Mercersburg Academy, where he was a keen athlete (football and track), musician (singing and accordion playing), and sometime actor. In 1929 he won a place at Princeton, where he studied ... [See full bio](#) »

**Born:** James Maitland Stewart  
May 20, 1908 in Indiana, Pennsylvania, USA

**Died:** July 2, 1997 (age 89) in Los Angeles, California, USA



230 photos | 42 videos | 1180 news articles »

**Won 1 Oscar.** Another 25 wins & 19 nominations. [See more awards](#) »

## Cast

Edit

Cast overview, first billed only:

	James Stewart	...	George Bailey
	Donna Reed	...	Mary Hatch
	Lionel Barrymore	...	Mr. Potter
	Thomas Mitchell	...	Uncle Billy
	Henry Travers	...	Clarence
	Beulah Bondi	...	Mrs. Bailey
	Frank Faylen	...	Ernie
	Ward Bond	...	Bert
	Gloria Grahame	...	Violet
	H.B. Warner	...	Mr. Gower

# Movie Questions

---

- Can we predict how well people will like a movie? What about its gross?
  - What does the social network of actors look like? (Six degrees of Kevin Bacon)
  - What is the age distribution of actors and actresses in film?
  - Do stars live longer or shorter lives than the bit players or public?
-

# Google Ngrams

---

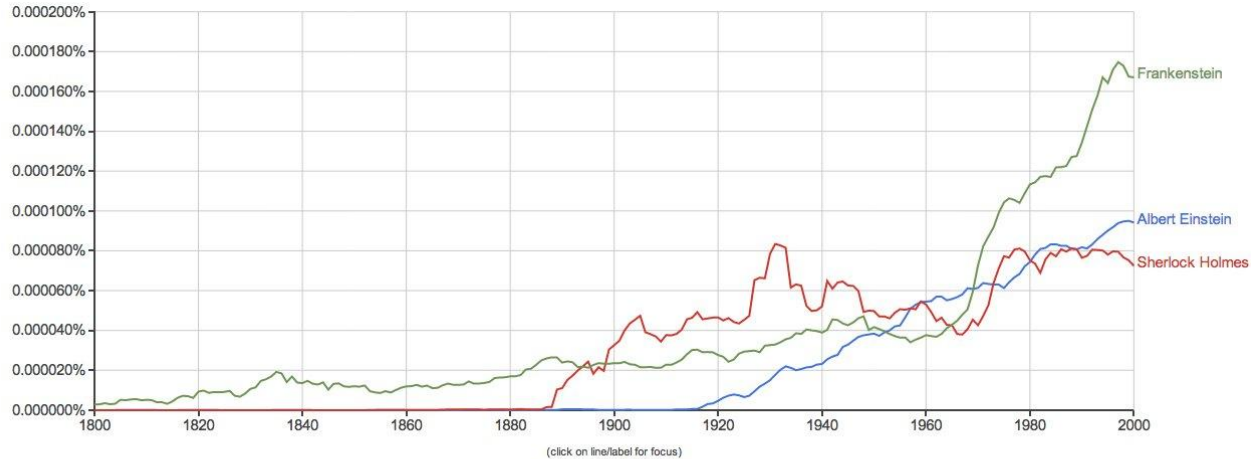
- Presents an annual time series of the frequency of every “popular” word/phrase with 1 to 5 words occurs in scanned books.
  - ‘Popular’ means appears >40 times in total.
  - Google has scanned about 15% of all books ever published, making this resource quite comprehensive.
-

# Google Ngram Viewer

Google books Ngram Viewer

Graph these comma-separated phrases:   case-insensitive

between  and  from the corpus  with smoothing of



Run your own experiment! Raw data is available for download [here](#).

# Ngram Questions

---

- How has the amount of cursing changed over time?
  - What is the lifespan of fame and technologies? Is it increasing/decreasing?
  - How often do new words emerge? Do they stay in common usage?
  - What words are associated with other words, i.e. can you build a language model?
-



# Taxicab Questions

---

- How much do drivers make each night?
  - How far do they travel?
  - How much slower is traffic during rush hour?
  - Where are people traveling to/from at different times of the day?
  - Do faster drivers get tipped better?
  - Where should drivers go to pick up their next fare?
-