Take me out to Big Data:

analyzing professional

baseball data online

Jordan Raddick IDIES, SciServer outreach team

with SciServer



I use SciServer to analyze and share game data from more than 100 years of Major League Baseball. This dataset provides an amazing opportunity to teach data science and statistics to an engaged global audience.

Box scores from game data Average MLB box score

All-ti	ime M	LB box	score	Avera	age M.	LB box	score
Team	Runs	Hits	Errors	Team	Runs	Hits	Errors
Visitor	794,733	1,637,837	167,352	Visitor	4.36	8.99	0.92

	Runs	Hits	Errors
r	794,733	1,637,837	167,352
,	822,947	1,621,067	172,440

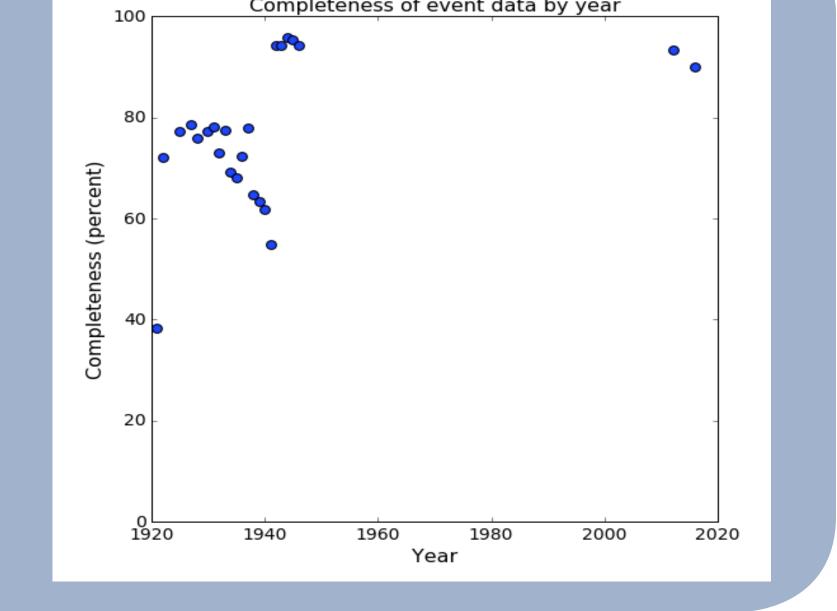
Iome 822,947 1,621,067 172,440 Hom	ne 4.52	8.90	0.95
Team		Hits	Errors
Median MLB box score Visito	or 4	9	1
Hom	e 4	9	1

Event data

Slowly processing event data in Compute

Loading progress

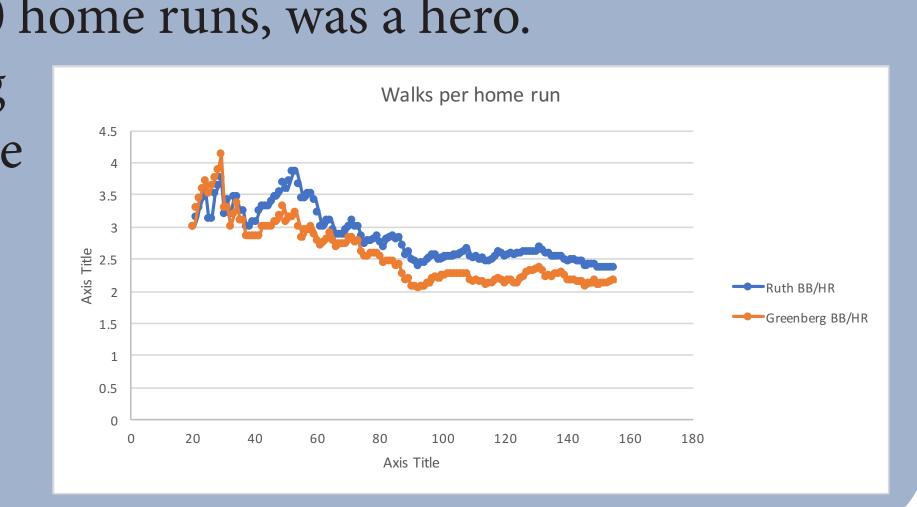
Zears .		Events							
921-1946		1,627,509							
012		177,720							
016		171,522							
		171,32							



Event data: the home run chase

1927: Babe Ruth hit 60 home runs, was a hero.

1938: Hank Greenberg hit 58 home runs. Were anti-Semitic pitchers avoiding him?



Ditab data

event	game	inning	battingteam	outs	balls	strikes	pitches	batter	batt
1	ANA201104080	1	0	0	1	2	FBSX	davir003	R
2	ANA201104080	1	0	1	0	0	X	nix-j001	R
3	ANA201104080	1	0	2	1	2	CBCS	bautj002	R
4	ANA201104080	1	1	0	3	1	СВВВВ	iztum001	L
5	ANA201104080	1	1	0	2	2	BCSBS	kendh001	R
6	ANA201104080	1	1	1	2	1	CBB1>S	abreb001	L
7	ANA201104080	1	1	2	3	2	CBB1>S.FBFB	abreb001	L
8	ANA201104080	1	1	2	0	2	CCX	huntt001	R
9	ANA201104080	1	1	2	1	0	BX	wellv001	R
10	ANA201104080	2	0	0	2	2	CBBFX	linda001	L
11	ANA201104080	2	0	1	1	2	BFCX	hilla001	R
12	ANA201104080	2	0	2	2	2	FBCBFX	rivej001	R
13	ANA201104080	2	1	0	3	2	CBBSBB	calla001	L
14	ANA201104080	2	1	0	0	0	X	trumm001	R

Why baseball?

- Fans know statistics (e.g. batting average, wins above replacement...)
- Data-driven managing
- Great opportunity for education

Why big data?

- Lots of data online
- Lots of interesting questions to ask

Why SciServer?

- Online analysis with Python (or R)
- Easy to share data and scripts

Data size and completeness

The first game in the dataset

May 4, 1871:

Cleveland Forest Cities at Fort Wayne Kekiongas

Retrosheet Expan	ded Box	Sc	ore	•					
Game of Thursday	, 5/4/1	871		C 1	leve	elan	d at	Ft.	Wayr
Cleveland									
Ft. Wayne	010 01	.0 0	99-	2					
CLEVELAND	AB	R	Н	BI	ВВ	S0	P0	Α	
D.White, c							 9		
G.Kimball, 2b									
C.Pabor, 1f									
A.Allison, cf	4	0	1	0	0	2	2	0	
E.White, rf A.Pratt, p	3	0	0	0	0	3	1	0	
A.Pratt, p	2	0	0	0	1	0	1	3	
E.Sutton, 3b	3	0	1	0	0	0	0	1	
J.Carleton, 1b	3	0	0	0	0	1	6	0	
J.Bass, ss	3	0	0	0	0	0	1	4	
Totals	30	0	4	0	1	6	27	9	
BATTING	D 11-41								
2B: D.White (off RBI, scoring pos A.Pratt 0-1.				har	1 2	out	s: G	.Kim	ball
BASERUNNING									
SB: A.Allison (2	nd base	of	f E	.Ma	athe	ews/	B.Le	nnon).
CS: A.Allison (2	nd base	by	В.	Mat	he	vs/B	.Len	non)	
Team LOB: 4									
FIELDING									
PR: D.White 3.									

F.Sellman, 3b	4	0	0	1	a	0	2	1
•	4							
J.Foran, 1b	4							
•	3							
B.Lennon, c	4							
T.Carey, 2b	3							
E.Mincher, lf								
J.McDermott, cf								
	3							
Totals	31	2	4	2	1	0	27	 3
2B: B.Lennon (off 2-out RBI: J.McDer RBI, scoring posit T.Carey 0-1; E.M BASERUNNING Team LOB: 3 FIELDING	mott. ion, inche	les er 0	s t	ı			s: F umbl	

Event data format

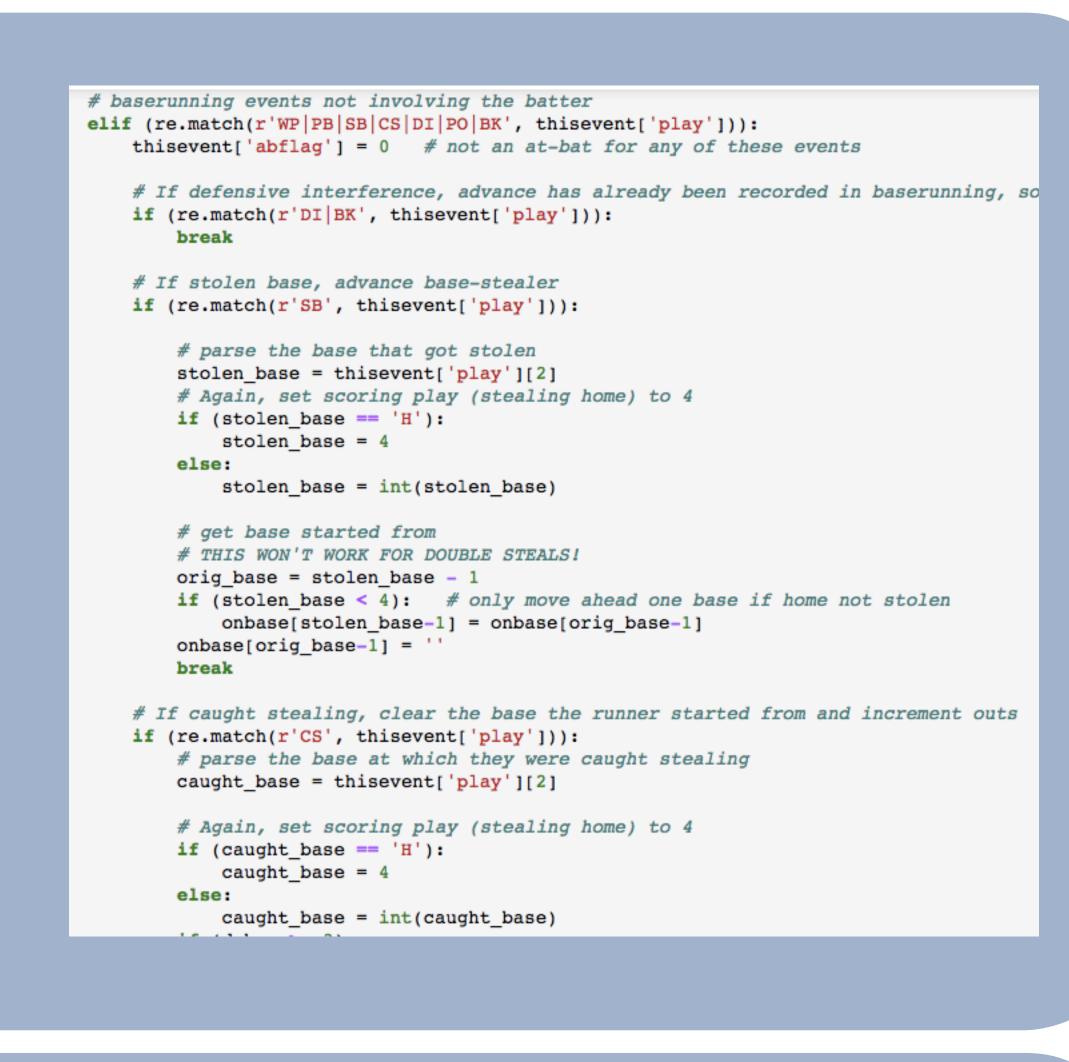


	innning	batting_team	play	visitor_score	home_score
0	7	TBA	43	0	1
1	7	TBA	63	0	1
2	7	TBA	3	0	1
3	8	TBA	K	0	1
4	8	TBA	K	0	1
5	8	TBA	K	0	1
6	9	TBA	K	0	1
7	9	TBA	63	0	1
8	9	TBA	K	0	1

Fun events in 2016

Date	Batting	Fielding	Inn.	Score	Batter	Scorecard says
June 17	Arizona Diamond- backs	Philadelphia Phillies	5	3-2	Wellington Castillo	CS3(1545E3)
July 6	Milwaukee Brewers	Washington Nationals	3	3-4	Chris Carter	K+SBH;SB2
October 1	Colorado Rockies	Milwaukee Brewers	9	3-2	Stephen Cardullo	5!6!3

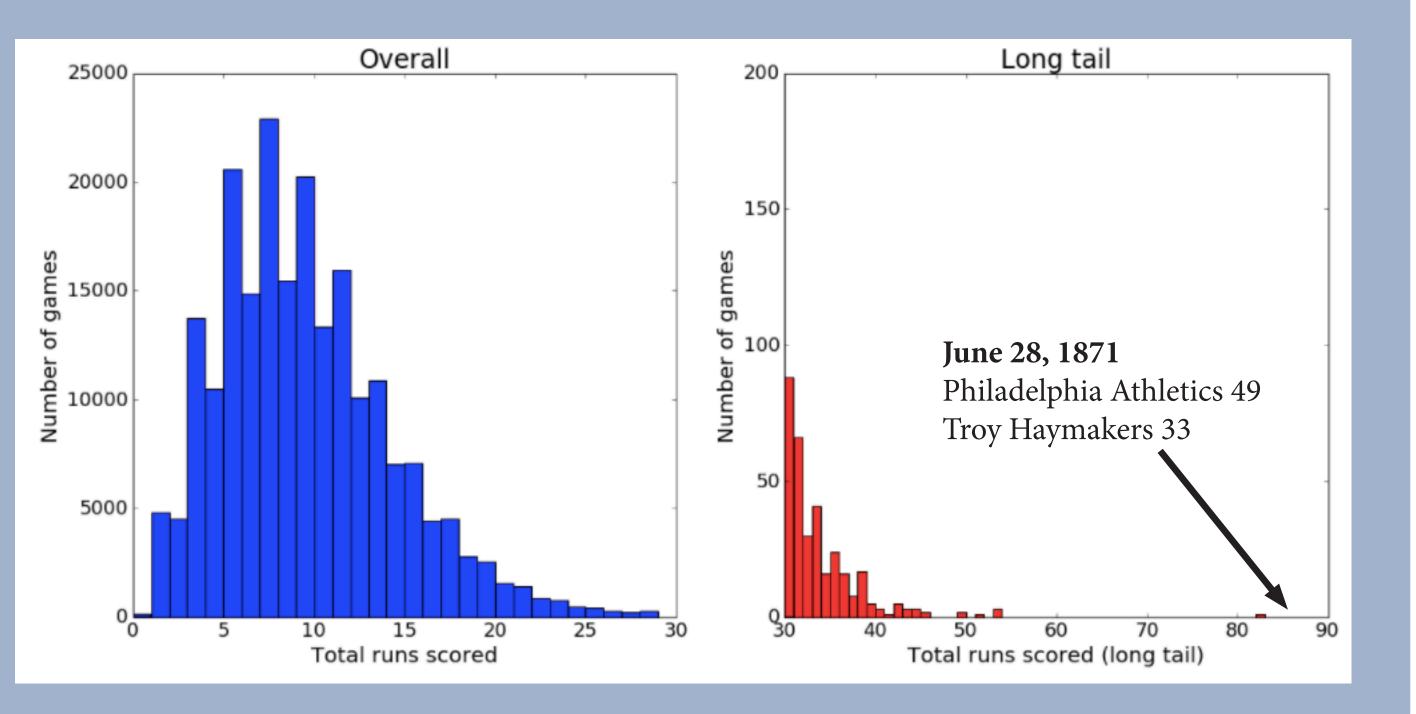
script example



Data

- Source: retrosheet.org
- Online volunteers transcribe historic baseball
- Games: box scores, lineups, game stats; one row per game
- Events: outcomes of at-bats or similar (e.g. stolen base); one row per event
- Pitches: outcome of individual pitches (e.g. ball, strike); one row per pitch

Histogram of total runs per game



Future: Research Questions

- Ideas for other research questions?
- Write your email and I'll contact you

What else could we study?

Not a baseball fan?

Take me out to the movies...

"Bechdel test": a simple metric for representation of women in movies.

Does this movie have...

- 1. Two or more female characters who...
- 2. talk to each other...
- 3. about something other than a man?

Processing

 Python scripts in SciServer
 Data type Games 1901-2016 1874-1900 213,307 1974-2016 1921-1973 ~10 million 2009-2016 1930-2008 ~40 million
 SciServer

Compute:

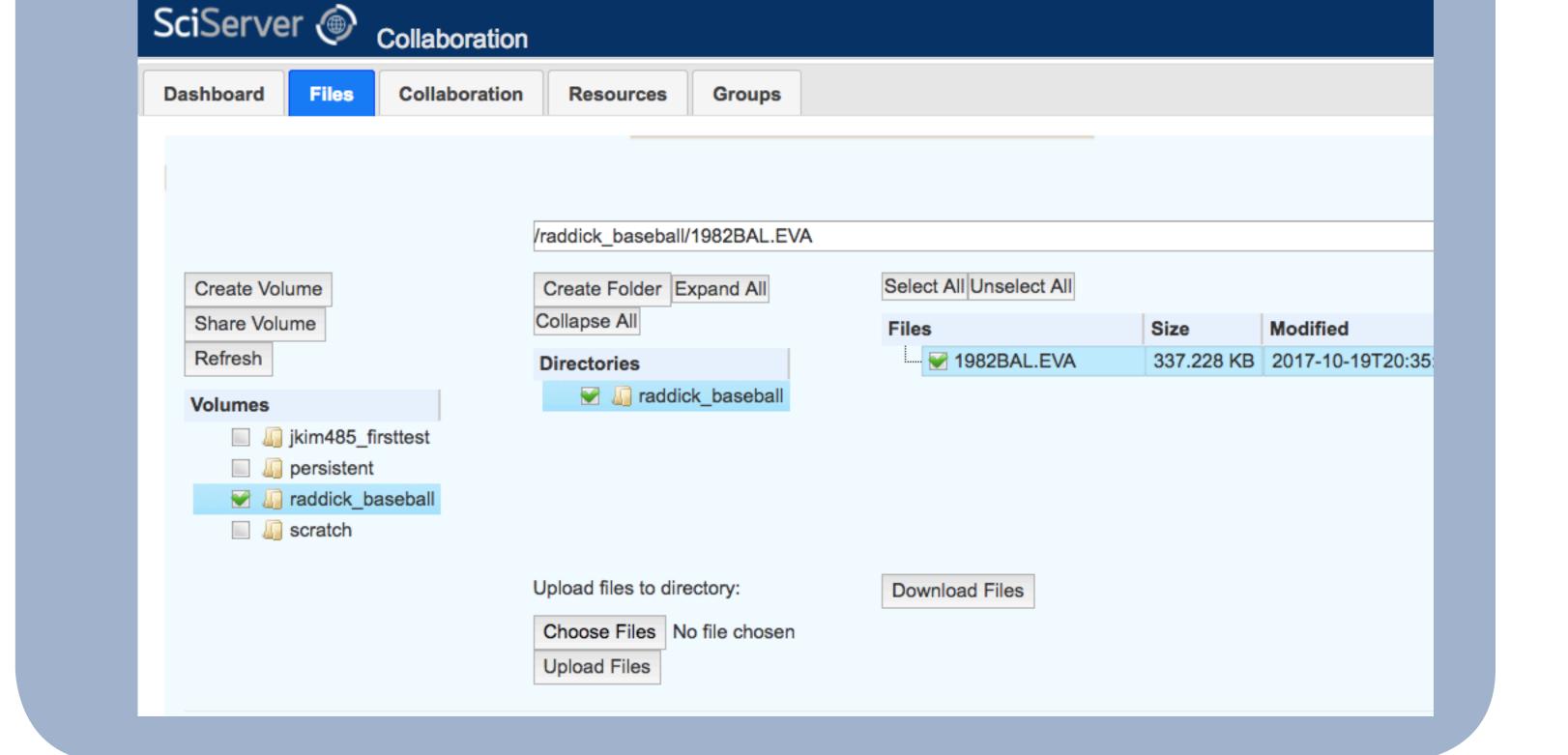
- 1. Copy data from retrosheet.org
- 2. Load game data as JSON
- 3. Parse game and event data into CSV
- 4. Calculate game progress (outs, score)

My scripts are inefficient, please help!

Future plans

- Aggregate to get player, team data
- Write tutorials and lesson plans
- Share data and scripts through SciServer
- Look at data from other sports

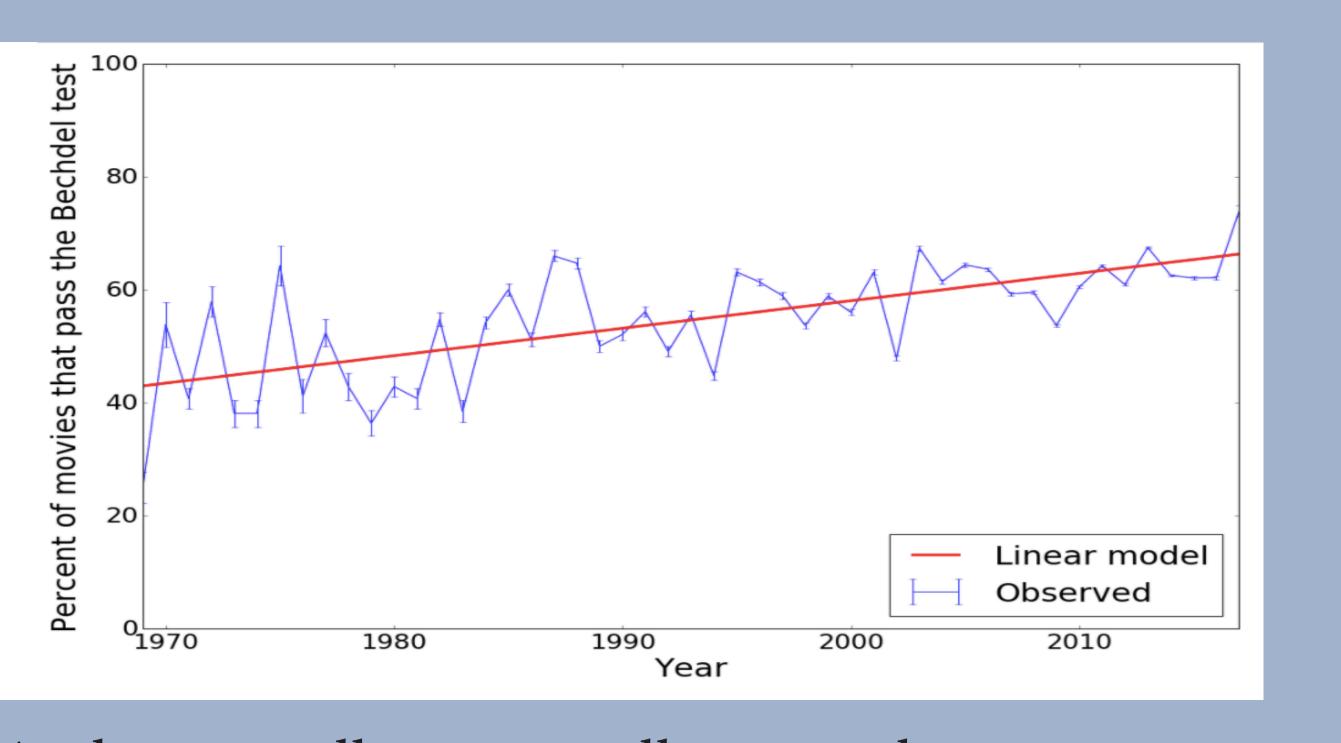
Future: Easy data sharing



Join the team!

- Talk to me about baseball (or other sports)
- Check and improve my scripts
- Find pitchFX data (pitch type, speed, etc.)
- Suggest new research questions
- Share with colleagues and students

Bechdel test results for Hollywood 1969-2017



At this rate, all movies will pass in the year 2087.