

The Inspection Paradox is Everywhere

Allen Downey
Olin College



Go to greenteapress.com/ip and follow instructions.

Experiment

Suppose you want to know average family size in the U.S.

And you have a convenience sample.

greenteapress.com/ip



Could ask, "How many children do you have?"

But the respondents are young.

And we want complete family sizes.

greenteapress.com/ip

Idea!

Go up a generation.

"How many children does
your mother have?"

greenteapress.com/ip



Survey

Go to

greenteapress.com/ip

and do the survey.

Survey of Family Size

An attempt to estimate the average family size in the U.S. by asking adults how big their families are.

* Required

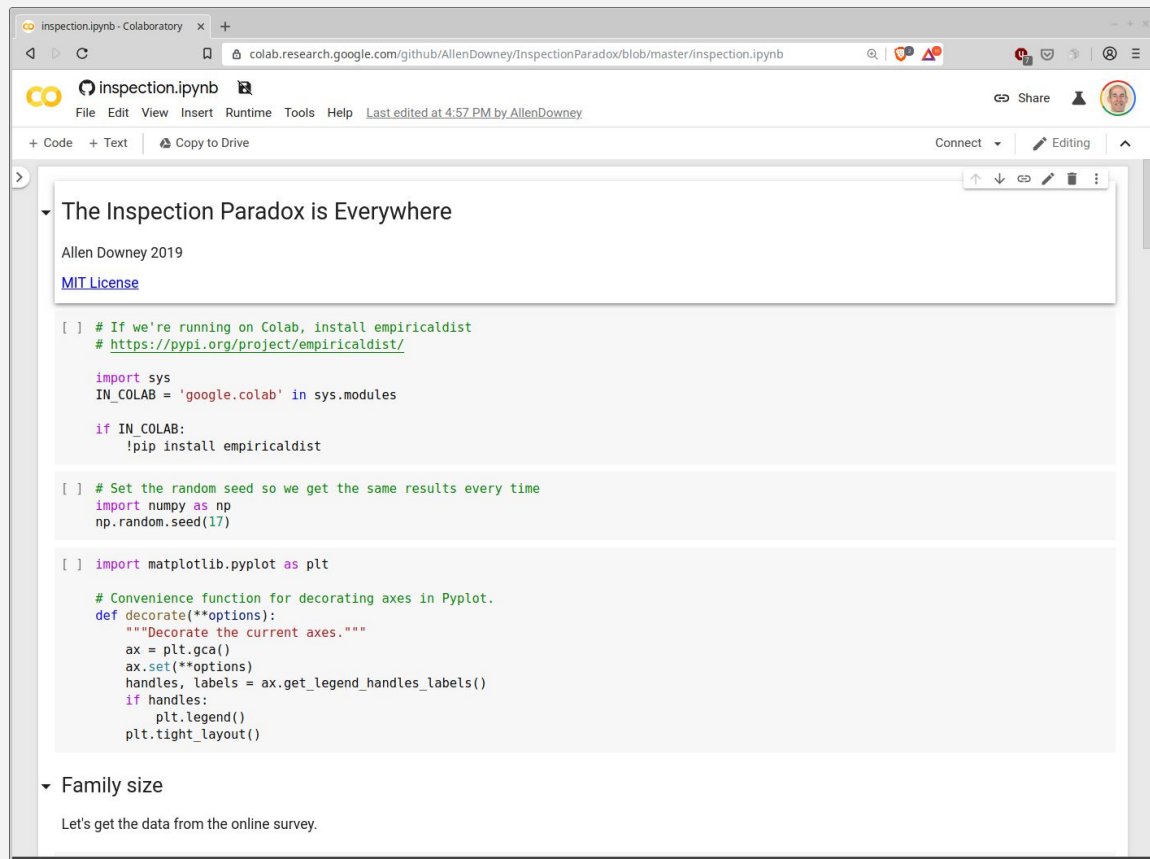
How many children has your biological mother born that were alive at birth? *

Your answer

Submit

When you are done,
follow the link to the
Jupyter notebook.

greenteapress.com/ip



The screenshot shows a web browser window displaying a Google Colaboratory Jupyter notebook. The browser's address bar shows the URL: `colab.research.google.com/github/AllenDowney/InspectionParadox/blob/master/inspection.ipynb`. The notebook interface includes a top menu bar with options like 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', and 'Help'. Below the menu is a toolbar with icons for '+ Code', '+ Text', and 'Copy to Drive'. The notebook content is organized into sections, with the first section titled 'The Inspection Paradox is Everywhere'. This section contains a text block with the author 'Allen Downey 2019' and a link to the 'MIT License'. Below the text are three code blocks. The first code block contains installation instructions for 'empiricaldist' on Colab. The second code block sets a random seed for reproducibility. The third code block imports 'matplotlib.pyplot' and defines a 'decorate' function for customizing plot axes. The notebook also shows a second section titled 'Family size' with the introductory text 'Let's get the data from the online survey.'

```
[ ] # If we're running on Colab, install empiricaldist
# https://pypi.org/project/empiricaldist/

import sys
IN_COLAB = 'google.colab' in sys.modules

if IN_COLAB:
    !pip install empiricaldist

[ ] # Set the random seed so we get the same results every time
import numpy as np
np.random.seed(17)

[ ] import matplotlib.pyplot as plt

# Convenience function for decorating axes in Pyplot.
def decorate(**options):
    """Decorate the current axes."""
    ax = plt.gca()
    ax.set(**options)
    handles, labels = ax.get_legend_handles_labels()
    if handles:
        plt.legend()
    plt.tight_layout()
```

Let's see some results.

Why might the PyData audience come from big families?

- Education
- Affluence
- Race and ethnicity
- etc.

Lots of possible sampling bias.

And one more thing...

Families with no children are not represented.

Families with no children are not represented.

Families with many children are over-represented.

Families with no children are not represented.

Families with many children are over-represented.

In general,

families with x children are over-represented by a factor of x .

Length-biased sampling

Sampling process where members of the population are sampled in proportion to size, length, duration, etc.

Inspection paradox

Subtly different sampling processes yield
surprisingly different results.

Inspection paradox

- Common error, but not well known.
- Once you know about it, you see it everywhere.
- Often problematic, but sometimes useful for experimental design.

Average class size

Ask teachers how big their classes are. Average = 31

Ask students how big their classes are. Average = 56

Who's lying?

Both right

They are averages across different populations.



[Home](#)
[Fast Facts](#)
[Students](#)
[Instruction and Student Life](#)
[Faculty and Staff](#)
[Diversity](#)
[Finance](#)
[Facilities](#)
[Research](#)

[Strategic Plan](#)
[Peer University Comparisons](#)
[Additional Facts and Figures](#)
[Regional Campuses](#)
[System-wide](#)
[Definitions](#)

Home > Instruction and Student Life

Distribution of Undergraduate¹ Classes by Course Level and Class Size

(for Fall 2012)

Download a [PDF](#) of this page ([Adobe Acrobat Reader](#) Required).

Course Level	Class Size								Total
	1	2-9	10-19	20-29	30-39	40-49	50-99	100+	
000-199	38	164	659	917	241	70	99	123	2,311
200-299	82	108	370	486	307	84	109	134	1,680
Lower Level	120	272	1,029	1,403	548	154	208	257	3,991
Percent of Lower Level Total	3.0%	6.8%	25.8%	35.2%	13.7%	3.9%	5.2%	6.4%	100.0%
300-399	4	148	387	314	115	96	186	53	1,303
400-499	14	132	256	190	83	67	64	17	823
Upper Level	18	280	643	504	198	163	250	70	2,126
Percent of Upper Level Total	0.8%	13.2%	30.2%	23.7%	9.3%	7.7%	11.8%	3.3%	100.0%
500-599	0	79	102	67	43	29	23	2	345
600-699	0	4	14	5	7	8	6	4	48
800-899	0	0	0	0	0	0	0	0	0
Dual Level	0	83	116	72	50	37	29	6	393
Percent of Dual Level Total	0.0%	21.1%	29.5%	18.3%	12.7%	9.4%	7.4%	1.5%	100.0%
Total All Classes	138	635	1,788	1,979	796	354	487	333	6,510
Percent of All Classes	2.1%	9.8%	27.5%	30.4%	12.2%	5.4%	7.5%	5.1%	100.0%

¹"Undergraduate" Classes refers to organized classes with one or more undergraduate students enrolled.

Average across students

138 classes with 1 student = 138 students

333 classes with 100+ students = 33,300+

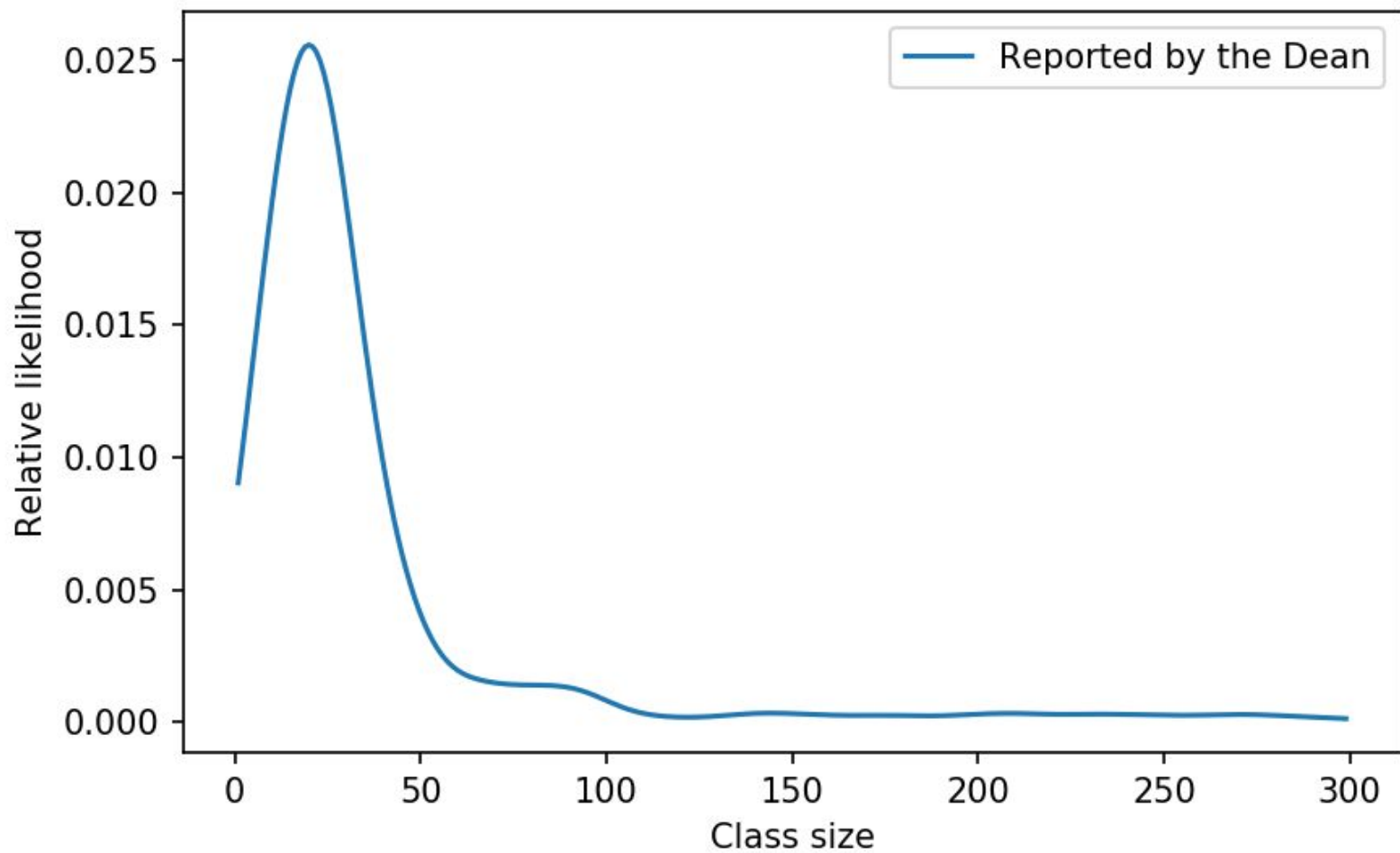
Large classes get oversampled.

Class size x gets oversampled by x .

I used the data in their table to generate an unbiased sample of class size.

And kernel density estimation (KDE) to plot the distribution.

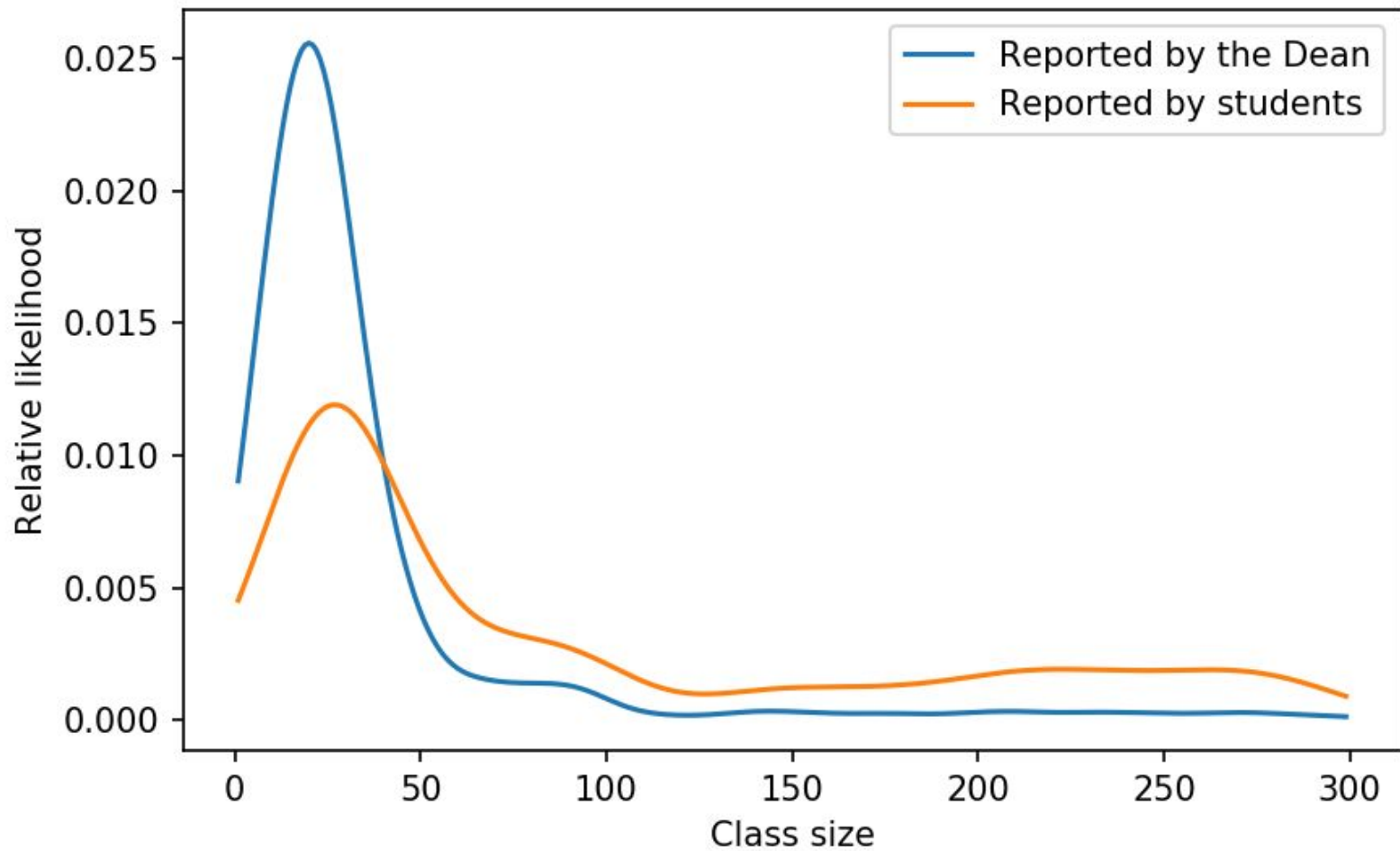
Distribution of class sizes



```
def resample_weighted(sample, weights):  
    """Generate a biased sample.  
  
    sample: NumPy array  
    weights: NumPy array  
  
    returns: NumPy array  
    """  
    n = len(sample)  
    p = weights / np.sum(weights)  
    return np.random.choice(sample, n, p=p)
```

```
biased = resample_weighted(unbiased, unbiased)
```

Distribution of class sizes



If you are not careful,
this kind of biased sampling is a problem.

If you are clever, you can use it.

Suppose your school doesn't publish class sizes...

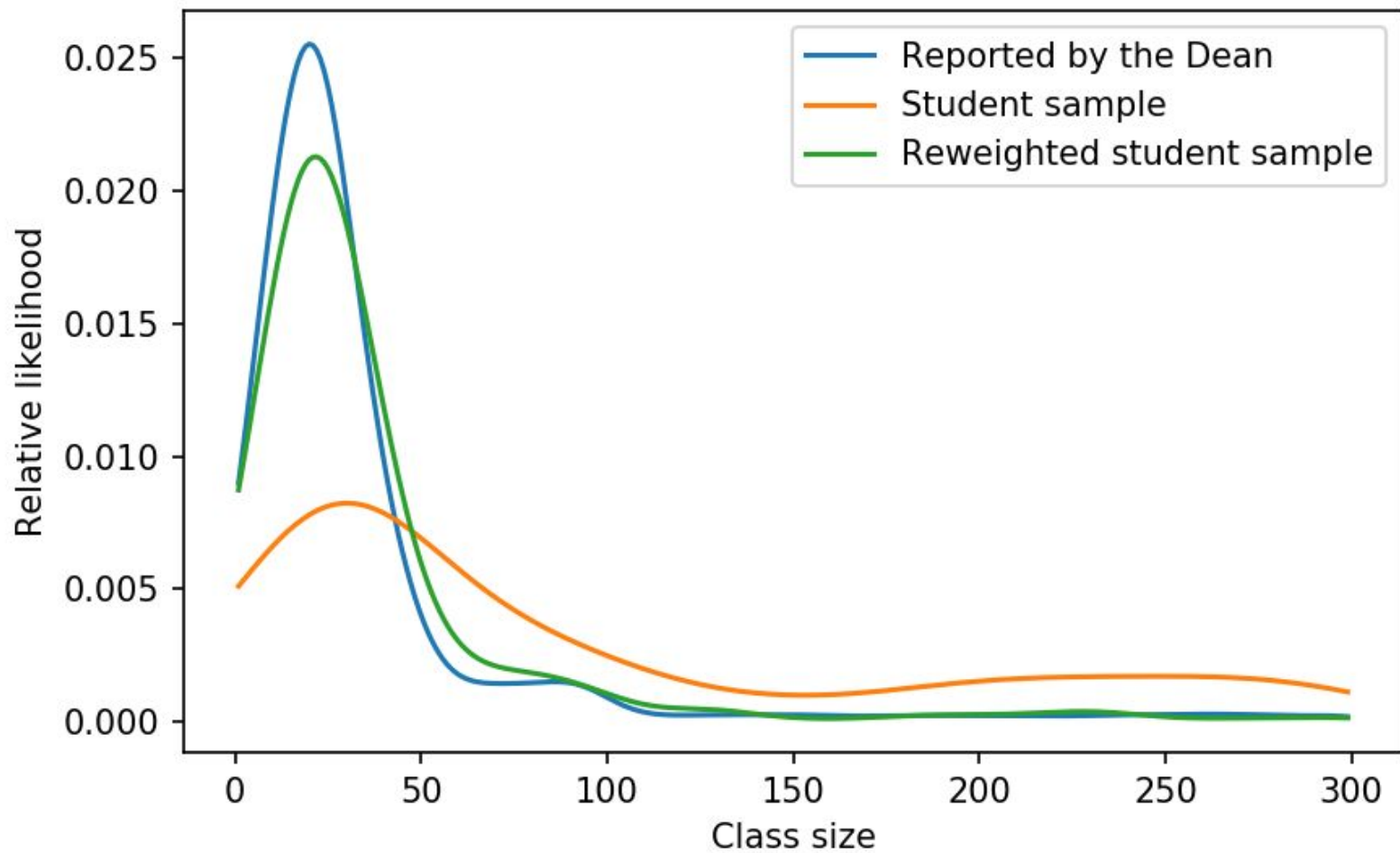
Experimental design

1. Sample students.
2. Resample with inverse weights.

```
sample = np.random.choice(biased, 500)
```

```
reweighted = resample_weighted(sample, 1/sample)
```

Distribution of class sizes



It's everywhere

Airlines:

“We are losing money because too many planes are nearly empty.”

Passengers:

“Flying is miserable because the planes are always full.”

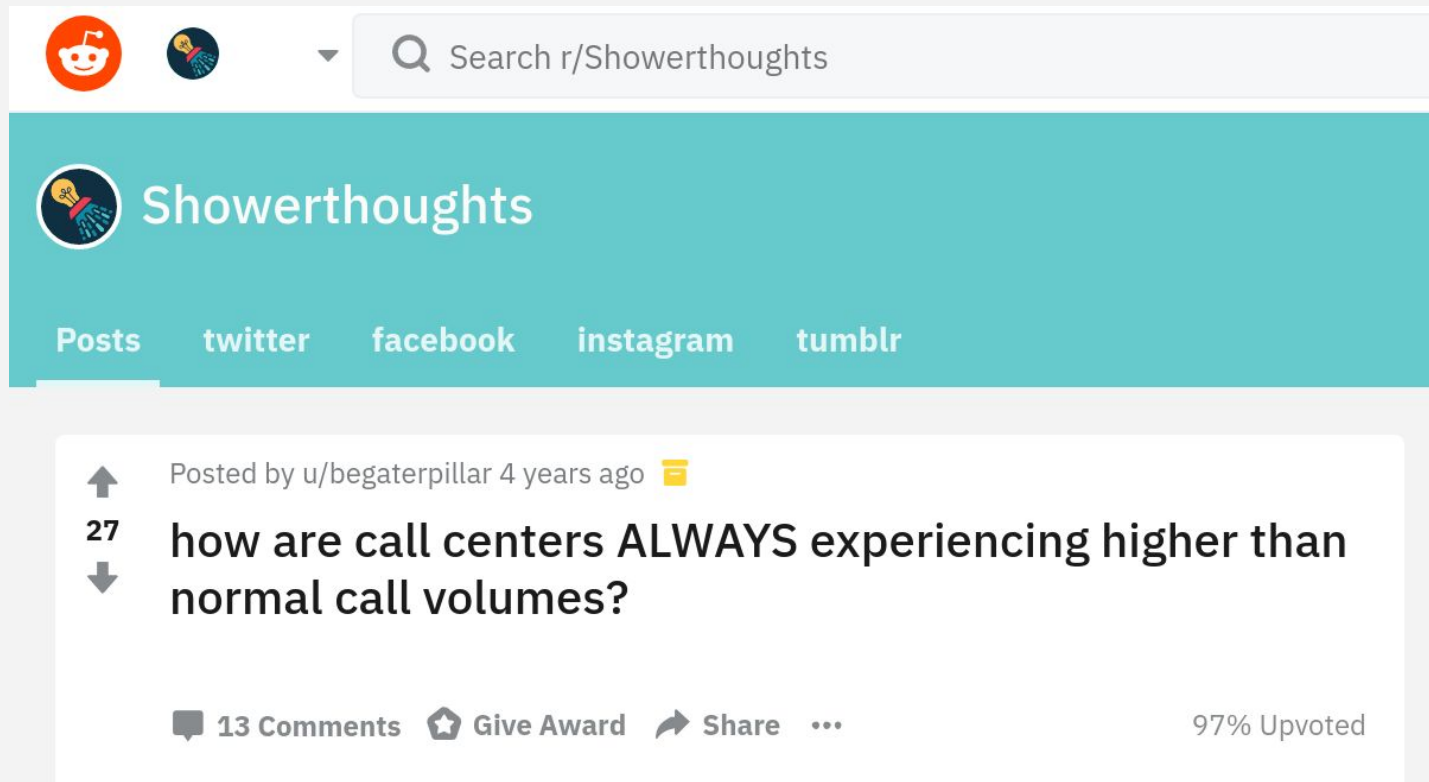
Both could be true

Few people enjoy a near-empty flight.

Lots of people suffer on a full one.

A plane with x passengers is oversampled by x .

It's everywhere



BUSINESS

Why You Can't Get a Taxi

And how an upstart company may change that

MEGAN MCARDLE MAY 2012 ISSUE

WHERE I LIVE in Washington, D.C., about a mile and a half north of the Capitol, you can sometimes get a taxi in two minutes flat. And sometimes, after spending 20 minutes wistfully waving two fingers in the air while the traffic hurtles past, you have to give up and trudge to the train.

Waiting for the bus in the rain

Suppose buses run every 20 minutes on average.

You expect to wait 10 minutes, on average.

Right?

Nope.

If there's any variation,
there are long intervals and short ones.

You are more likely to arrive during a long one.

Nope.

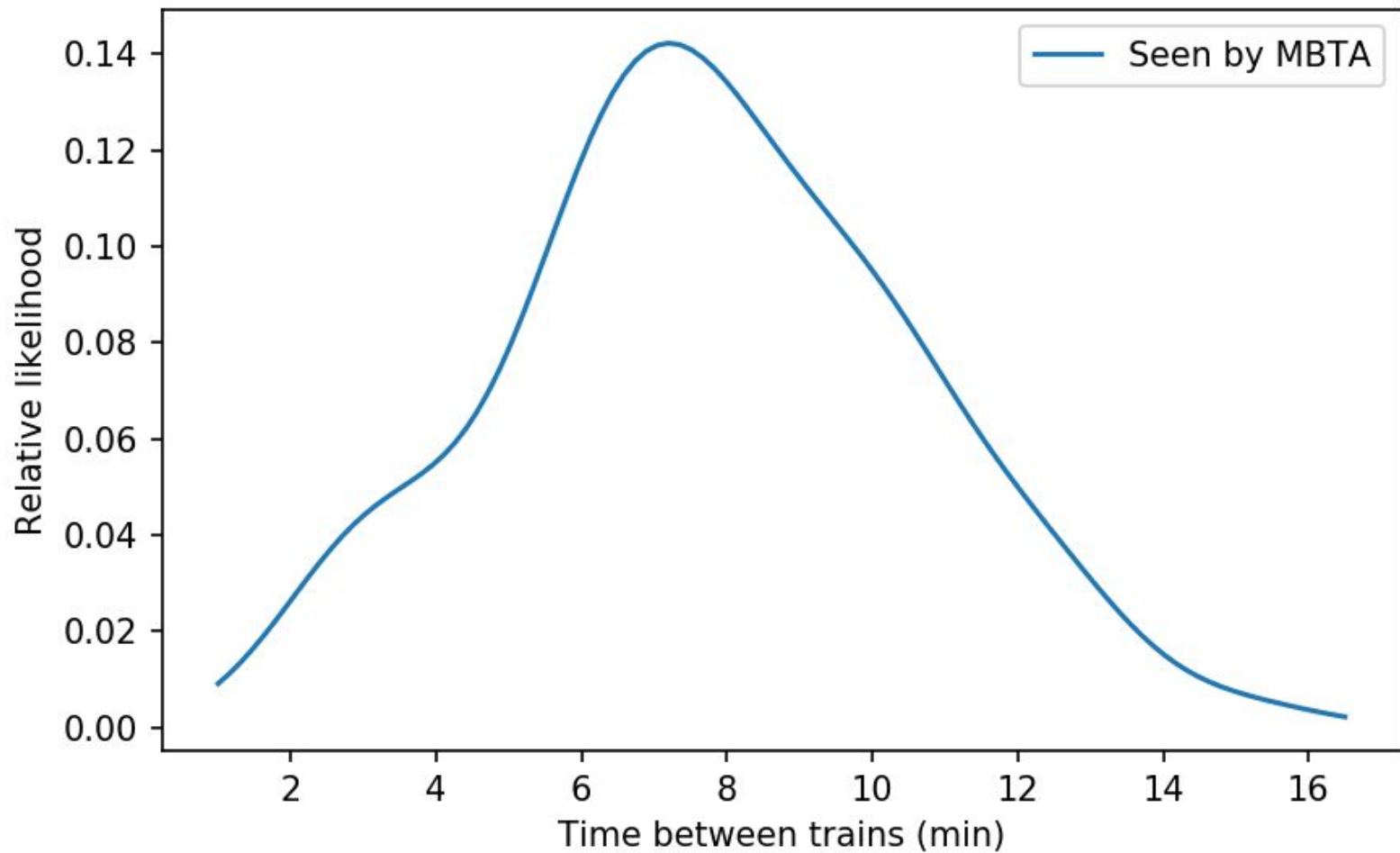
If there's any variation,
there are long intervals and short ones.

You are more likely to arrive during a long one.

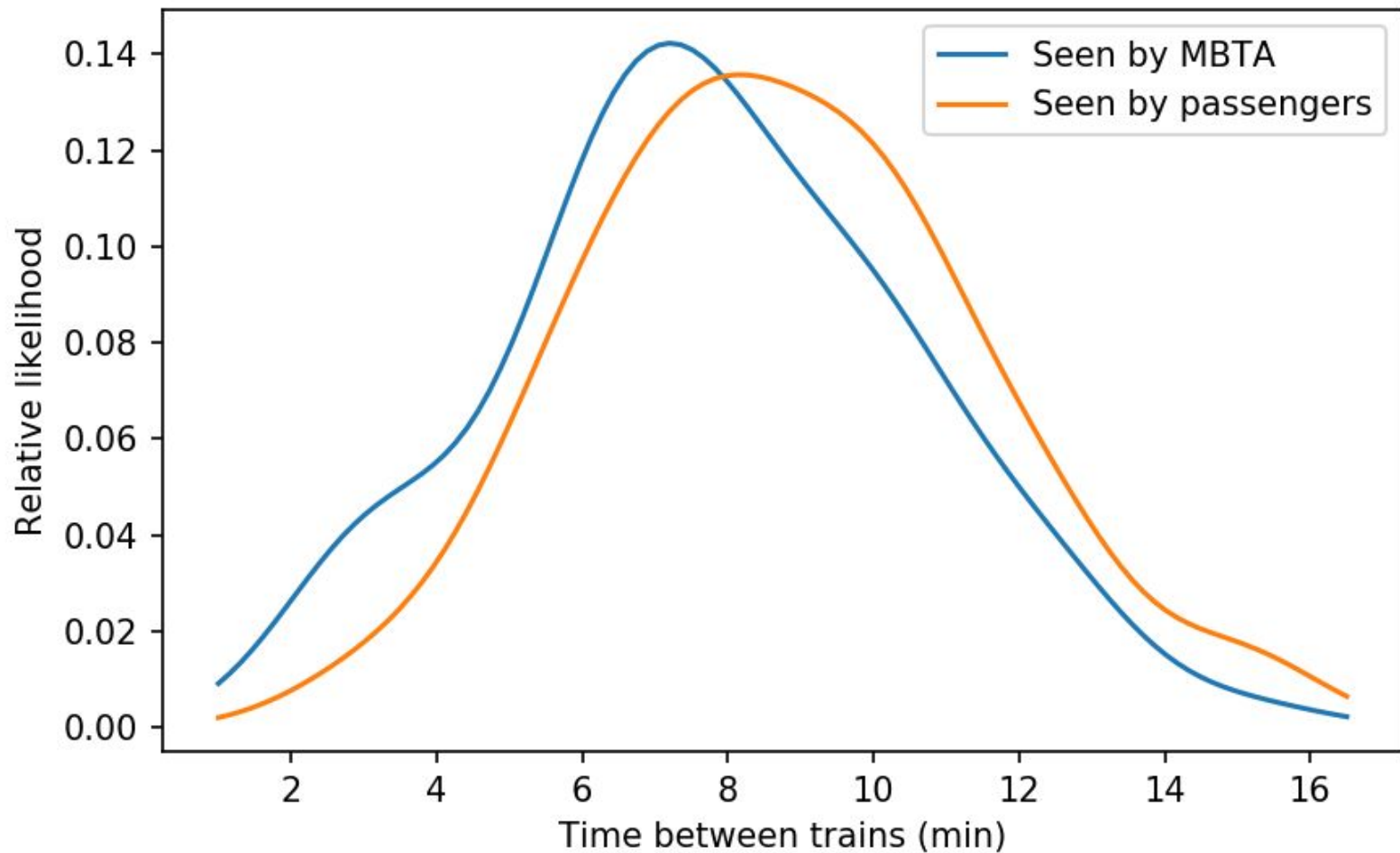
Intervals of duration t are oversampled by a factor of t .



Distribution of time between trains



Distribution of time between trains



Average reported by MBTA: 7.8 minutes.

Average observed by passengers: 8.8 minutes.




In this example the difference is moderate
because the variance is moderate.

It can be much bigger.

Let me ask you a question...

Are you popular?

Hint: no.

[Subscribe](#)[News & Features](#)[Topics](#)[Blogs](#)[Videos & Podcasts](#)[Education](#)[Citizen Science](#)[More Science](#) » [February 2011](#) » [Advances](#) 18 ::  Email ::  Print

Why You're Probably Less Popular Than Your Friends

Where averages and individual perspectives diverge

By [John Allen Paulos](#) | Jan 18, 2011

Are your friends more popular than you are? There doesn't seem to be any obvious reason to suppose this is true, but it probably is. We are all



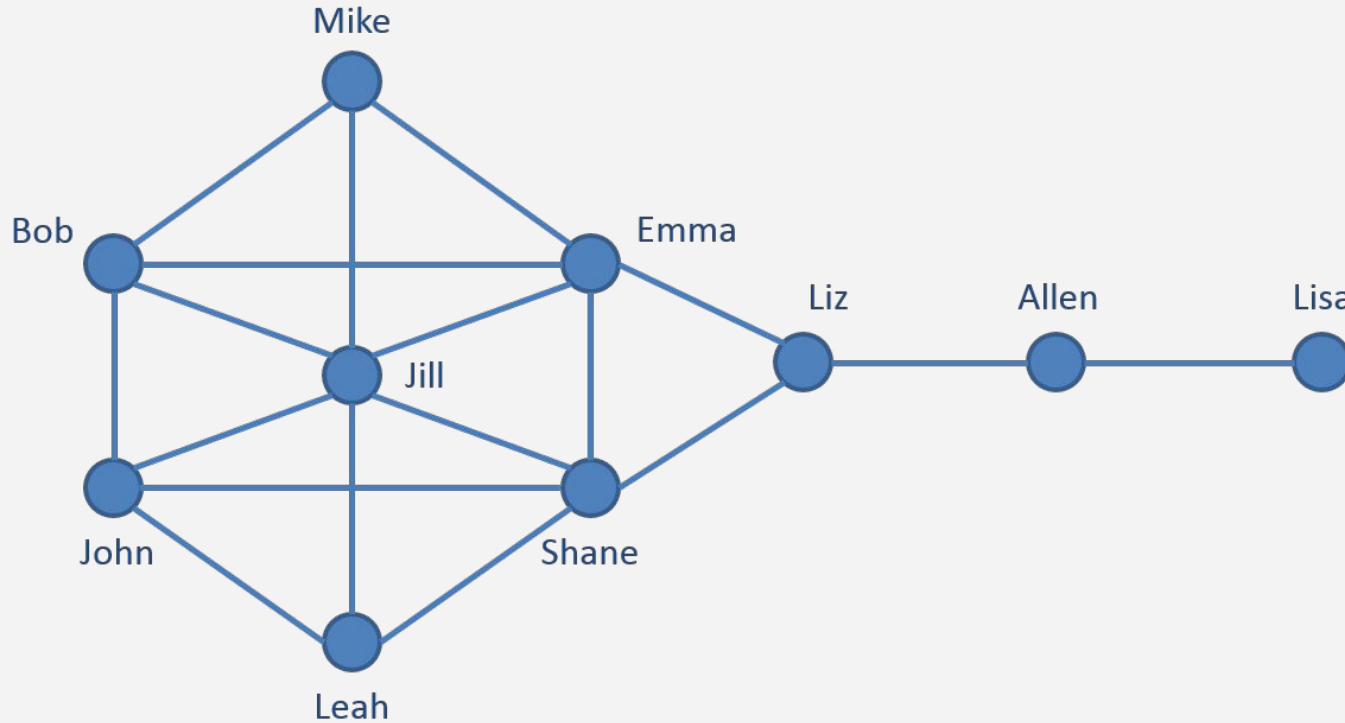
Friendship paradox

How many friends do you have on Facebook?

Suppose I pick one of your friends at random.

Chances are ~80% that they have more friends than you.

Think of a graph



One more time

1. Choose a random node.
2. Choose a random edge and follow it.
3. Query that node.

A node with degree x is oversampled by x .



- SNAP for C++ ▶
- SNAP for Python ▶
- SNAP Datasets ▶
- What's new
- People
- Papers
- Citing SNAP
- Links
- About
- Contact us

Social circles: Facebook

Dataset information

This dataset consists of 'circles' (or 'friends lists') from Facebook. Facebook data was collected from survey participants using this [Facebook app](#). The dataset includes node features (profiles), circles, and ego networks.

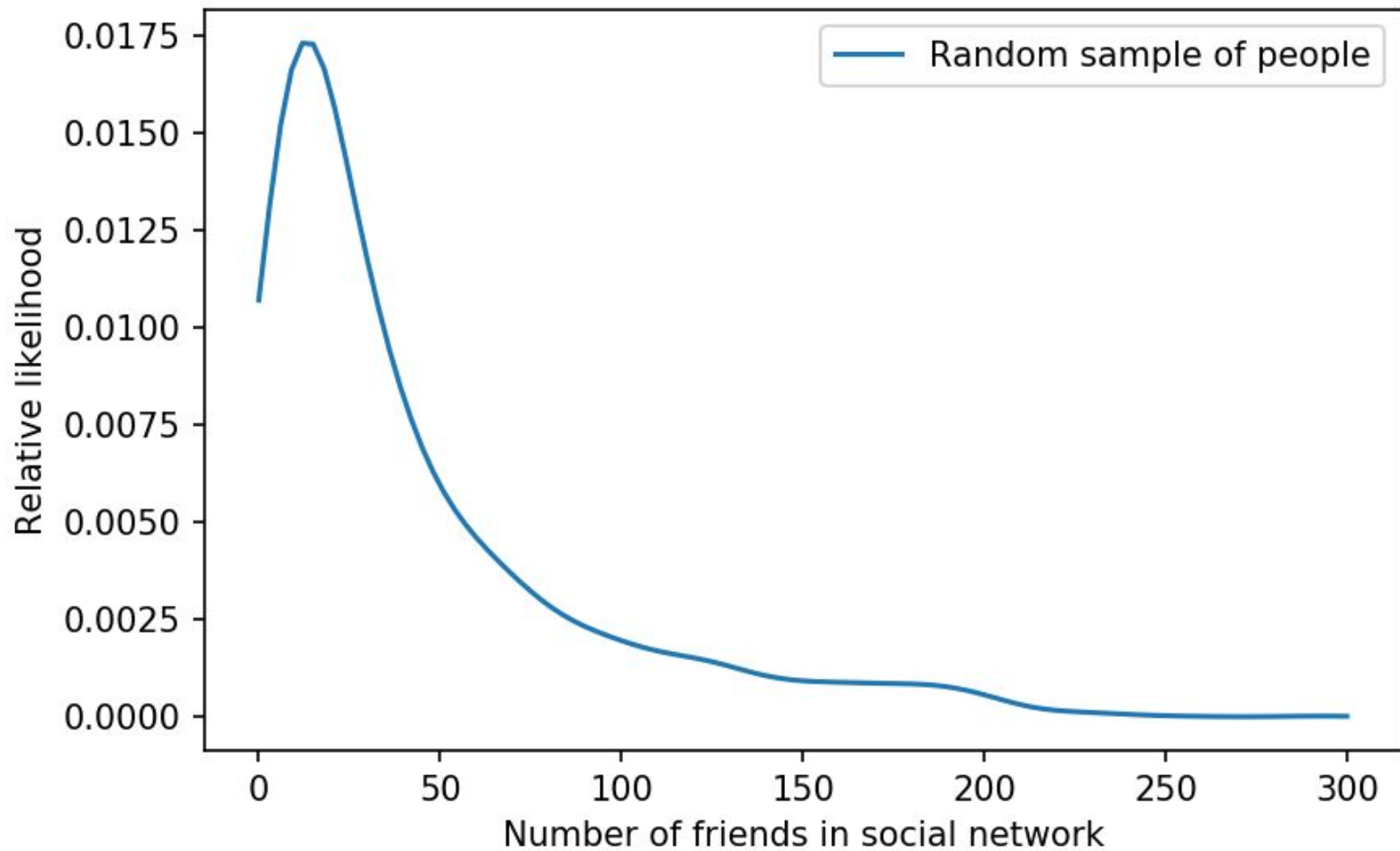
Facebook data has been anonymized by replacing the Facebook-internal ids for each user with a new value. Also, while feature vectors from this dataset have been provided, the interpretation of those features has been obscured. For instance, where the original dataset may have contained a feature "political=Democratic Party", the new data would simply contain "political=anonymized feature 1". Thus, using the anonymized data it is possible to determine whether two users have the same political affiliations, but not what their individual political affiliations represent.

Data is also available from [Google+](#) and [Twitter](#).

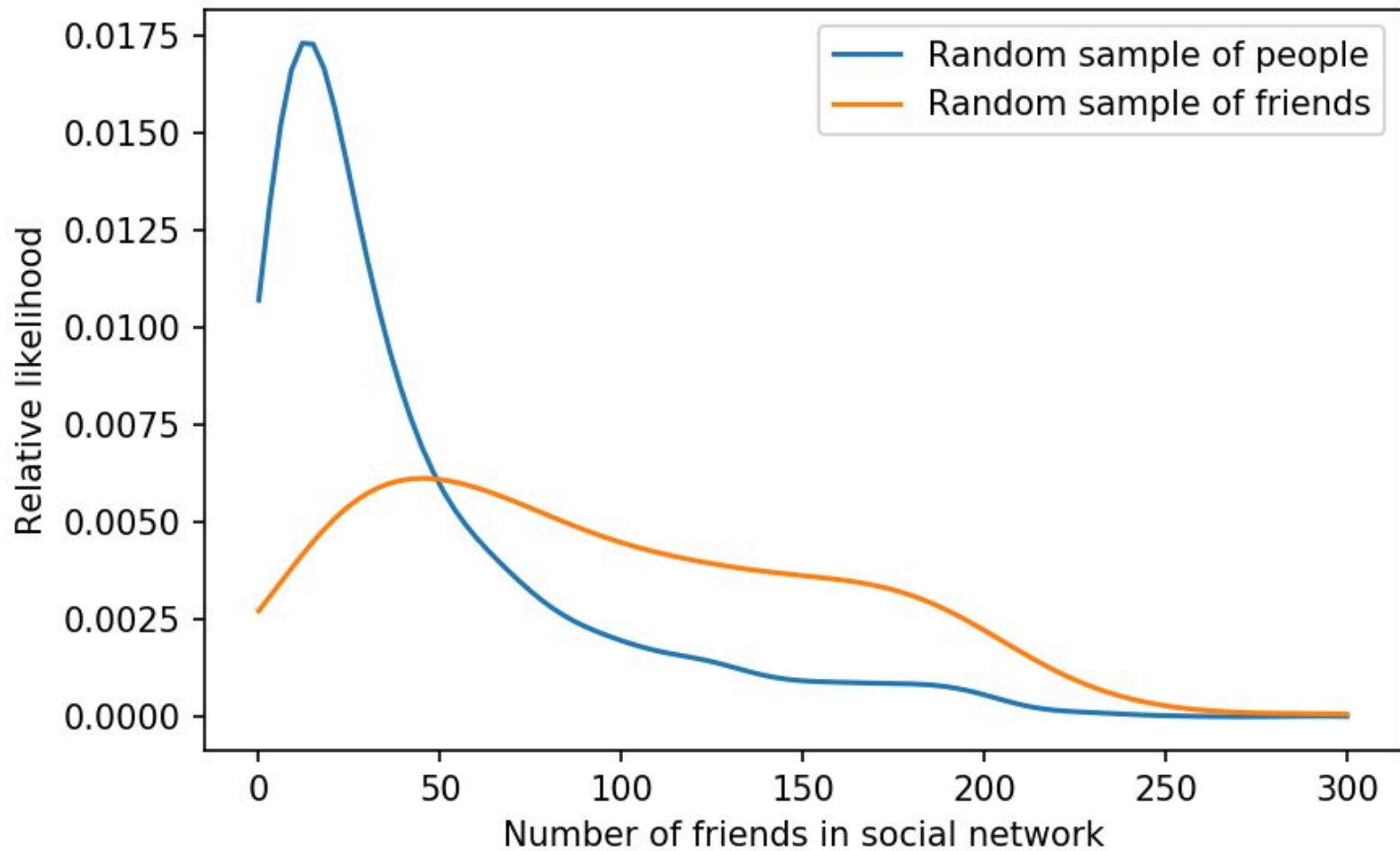
Dataset statistics

Nodes	4039
Edges	88234
Nodes in largest WCC	4039 (1.000)
Edges in largest WCC	88234 (1.000)
Nodes in largest SCC	4039 (1.000)
Edges in largest SCC	88234 (1.000)
Average clustering coefficient	0.6055
Number of triangles	1612010
Fraction of closed triangles	0.2647
Diameter (longest shortest path)	8
90-percentile effective diameter	4.7

Distribution of social network size



Distribution of social network size



What a difference an edge makes

Average number of edges “you” have:

42

Average number of edges your friends have:

~100

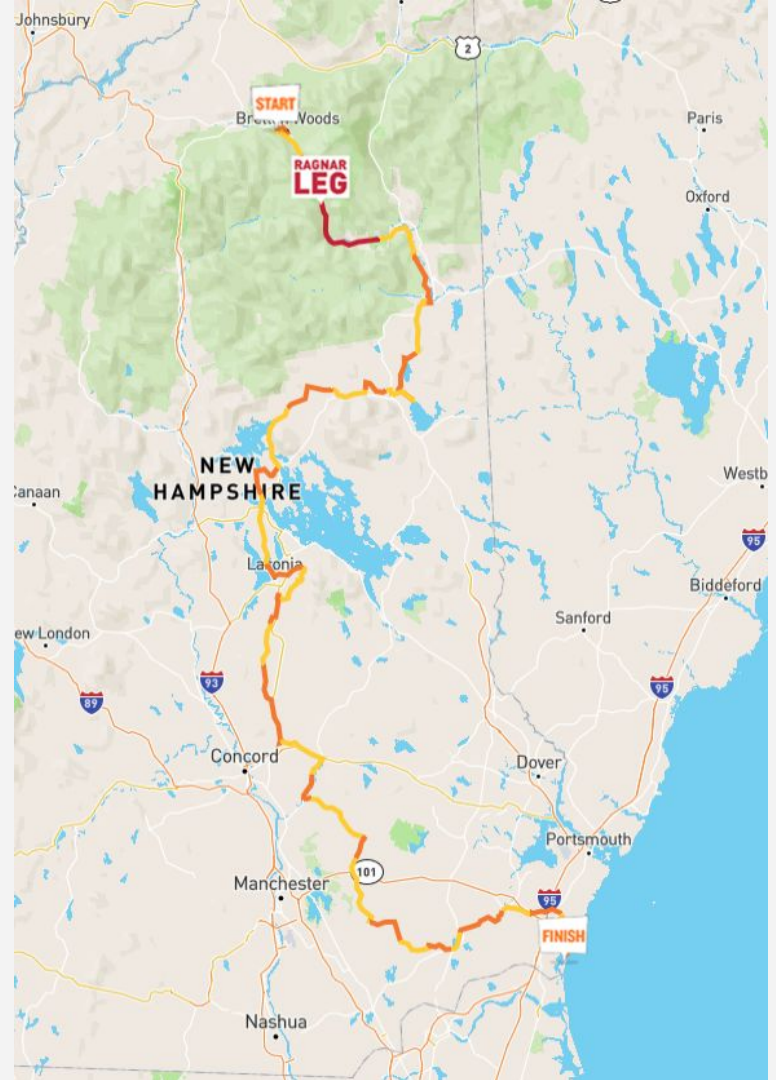
Reach the Beach

200+ mile relay.

When I overtook someone,
I overtook them fast.

When people overtook me,
they blew by me.

Bimodal distribution of runners?



Nope, it's the inspection paradox

Because of the format,
fast and slow runners are spread out.

Almost no relationship between speed and position.

As a runner, who do you see?

Fast runners pass many slow runners,
fewer fast ones.

Slow runners are overtaken by many fast runners,
fewer slow ones.

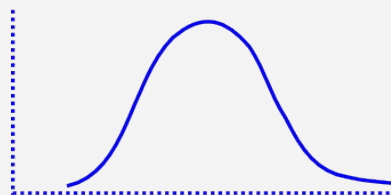
Fast runners pass many slow runners,
fewer fast ones.

Slow runners are overtaken by many fast runners,
fewer slow ones.

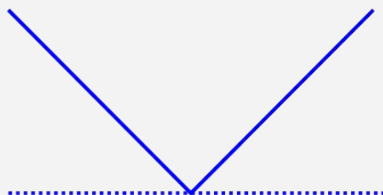
In general, the chance of seeing another runner is
proportional to the **difference in your speeds**.

Prob of observing x is proportional to $\text{abs}(x-v)$

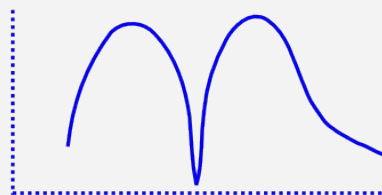
- Actual distribution is single-mode.
- Filter is v-shaped.
- Observed distribution is bimodal.



actual

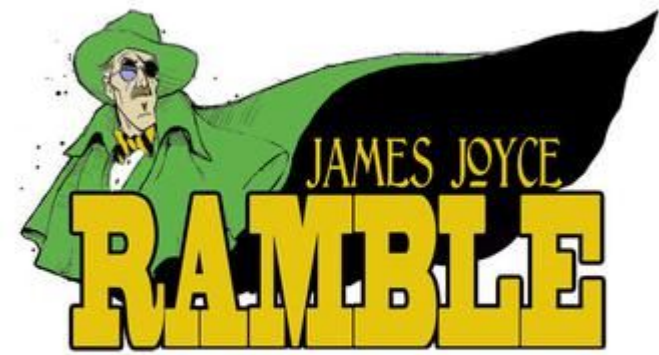




filter



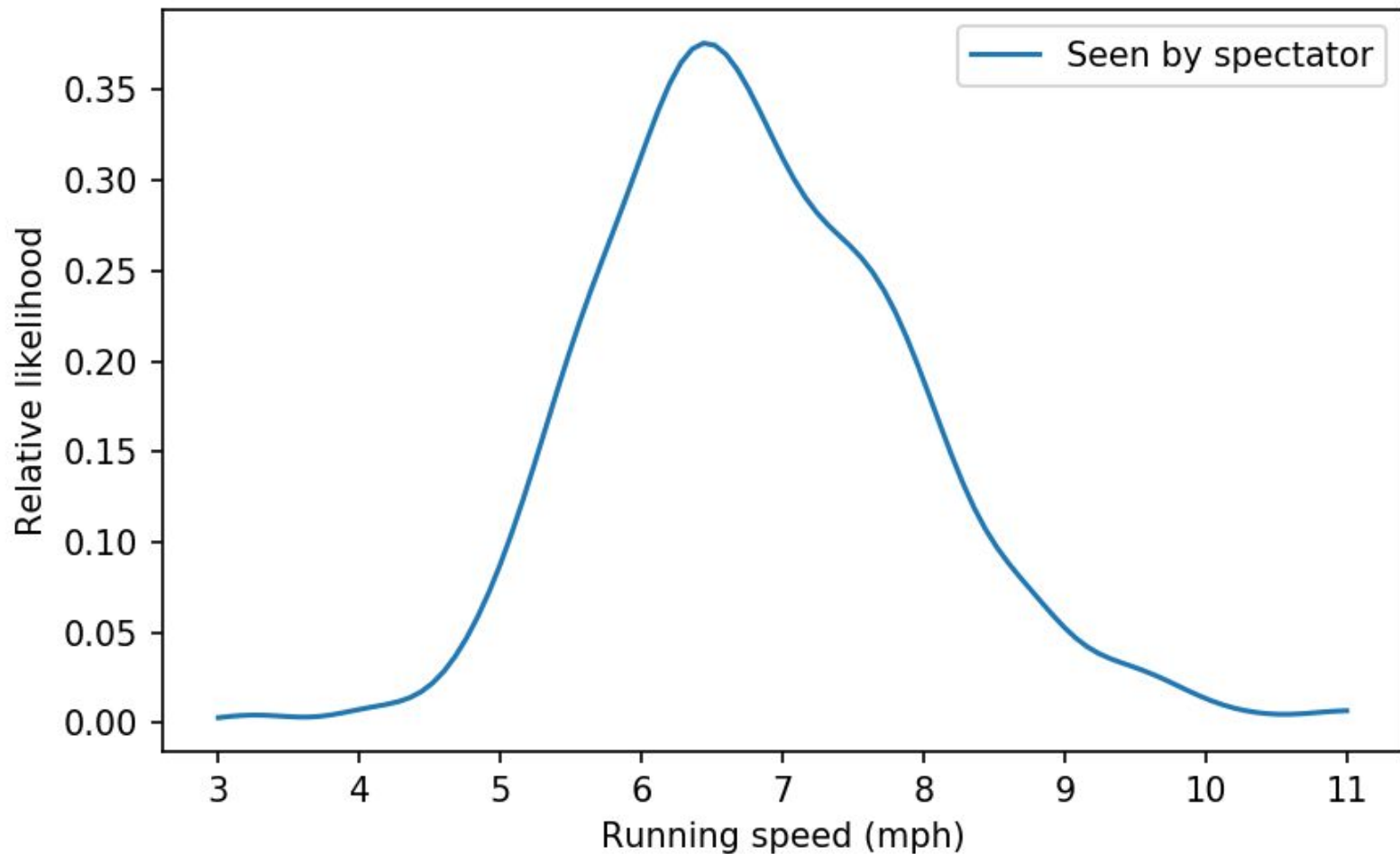
observed

Data from a 10K road race.

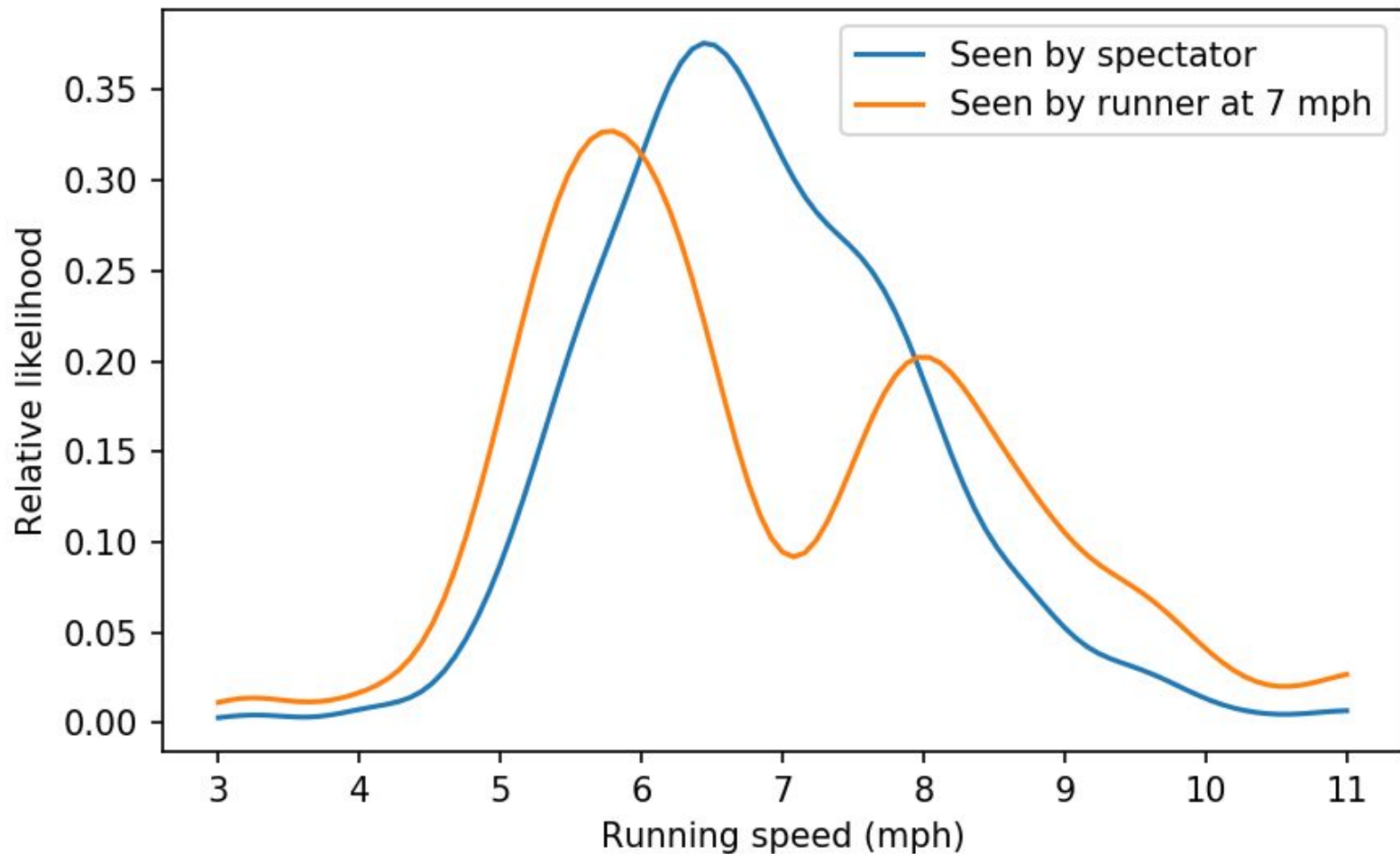


Claim	Steven Rostick M 48 Bib 1382 MA	95	84	14	5:52 min/mi	42:37
Claim	Steve Murray M 46 Bib 1668 VA	96	85	15	6:52 min/mi	42:39
	Allen Downey M 42 Bib 337 Needham, MA	97	86	12	6:53 min/mi	 42:44
Claim	Kate Blake F 32 Bib 107 Dedham, MA	98	12	2	6:54 min/mi	42:48

Distribution of running speed



Distribution of running speed



On the highway

Everybody is going too fast.

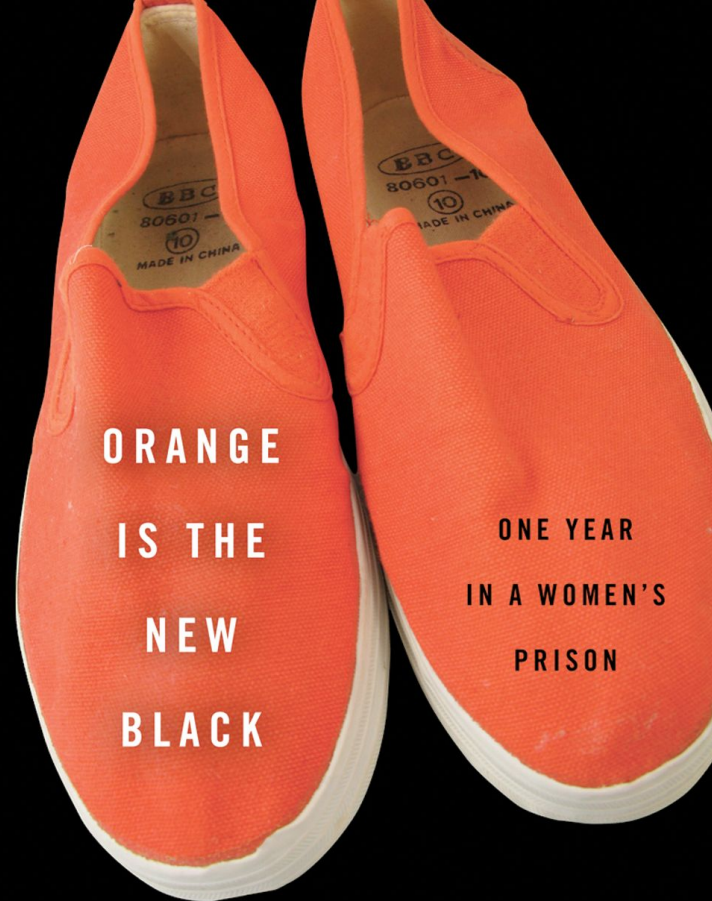
Or too slow.

But hardly any sane, safe drivers like yourself.

Once you see it...

Once you see it...

You see it everywhere.



ORANGE
IS THE
NEW
BLACK

ONE YEAR
IN A WOMEN'S
PRISON

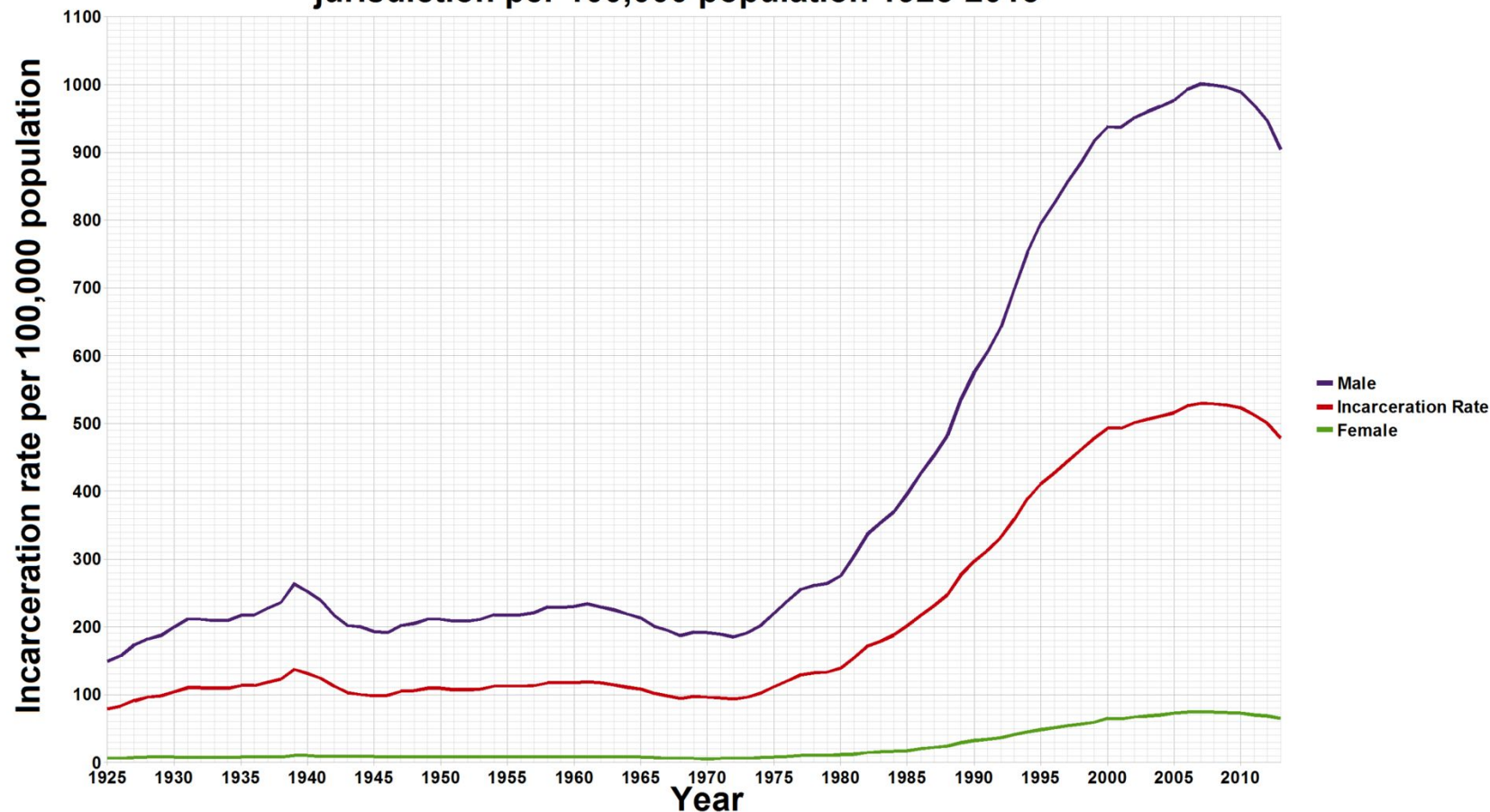
A Memoir

PIPER KERMAN

Long sentences?

Kerman expresses surprise
at the long sentences
her fellow prisoners are serving.

Incarceration rate of inmates incarcerated under state and federal jurisdiction per 100,000 population 1925-2013



But also...

Arrive at a random time.

Choose a random prisoner.

Prisoner with sentence x is oversampled by x .



Statistics

[Inmate Statistics](#)[Population Statistics](#)[Staff Statistics](#)

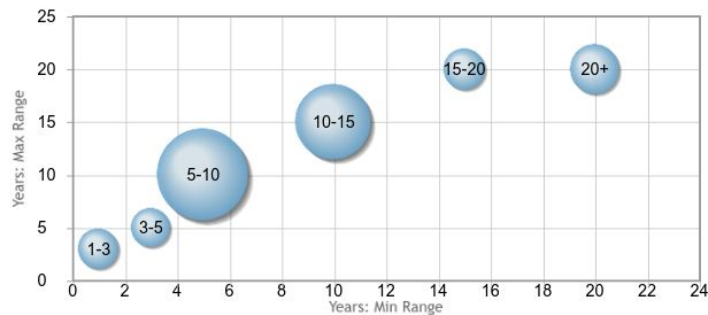
Inmate Statistics

[Age](#)
[Citizenship](#)
[Ethnicity](#)
[Gender](#)
[Offenses](#)
[Prison Safety](#)
[Prison Security Levels](#)
[Programs](#)
[Race](#)
[Release Numbers](#)
[Restricted Housing](#)
[Sentences Imposed](#)

Sentences Imposed

Statistics based on prior month's data -- Last Updated: Saturday, 26 October 2019

Please Note: Data is limited by availability of sentencing information for inmates in BOP custody.



Sentence	# of Inmates	% of Inmates
0 to 1 year*	4,929	2.3 %
> 1 year to < 3 years**	18,726	11.4%
3 years to < 5 years	17,811	10.8%
5 years to < 10 years	41,987	25.6%
10 years to < 15 years	34,917	21.2%
15 years to < 20 years	18,710	11.4%
20 years or more but < Life	22,655	13.8%
Life	4,536	2.8%

* The sentence category "0 to 1 year" includes misdemeanor offenses (0-12 months).

** The sentence category "> 1 to < 3 years" includes the common sentence type: "Twelve months plus 1 day."

Read the fine print...

This is based on a sample of current prisoners.

So the sample is length-biased.

But we can unbiased it!

Prisoners in federal prison typically serve 85% of their nominal sentence.

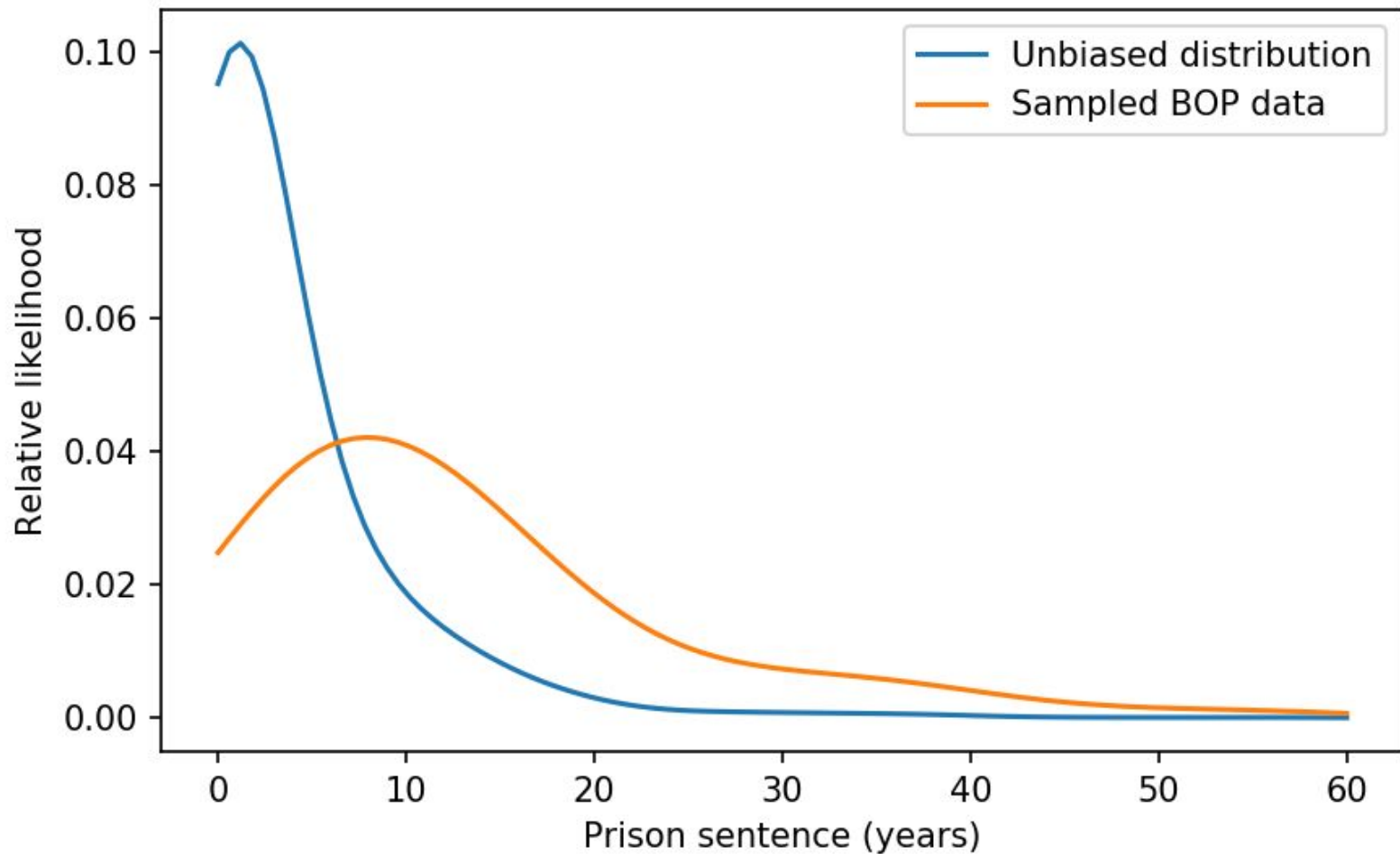
We can take that into account in the weights.

```
weights = 1 / (0.85 * np.array(biased))
```

Here's the unbiased sample.

```
unbiased = resample_weighted(biased, weights)
```

Distribution of federal prison sentences



But it depends on how long you stay

Single observation: biased by x .

Observe for a long period: unbiased.

What if you observe for 13 months?

Interval inspection

If the inspection interval is t
a prisoner with sentence x
is oversampled by $x + t$.

Small t : converges to x .

Large t : converges to constant (unbiased).

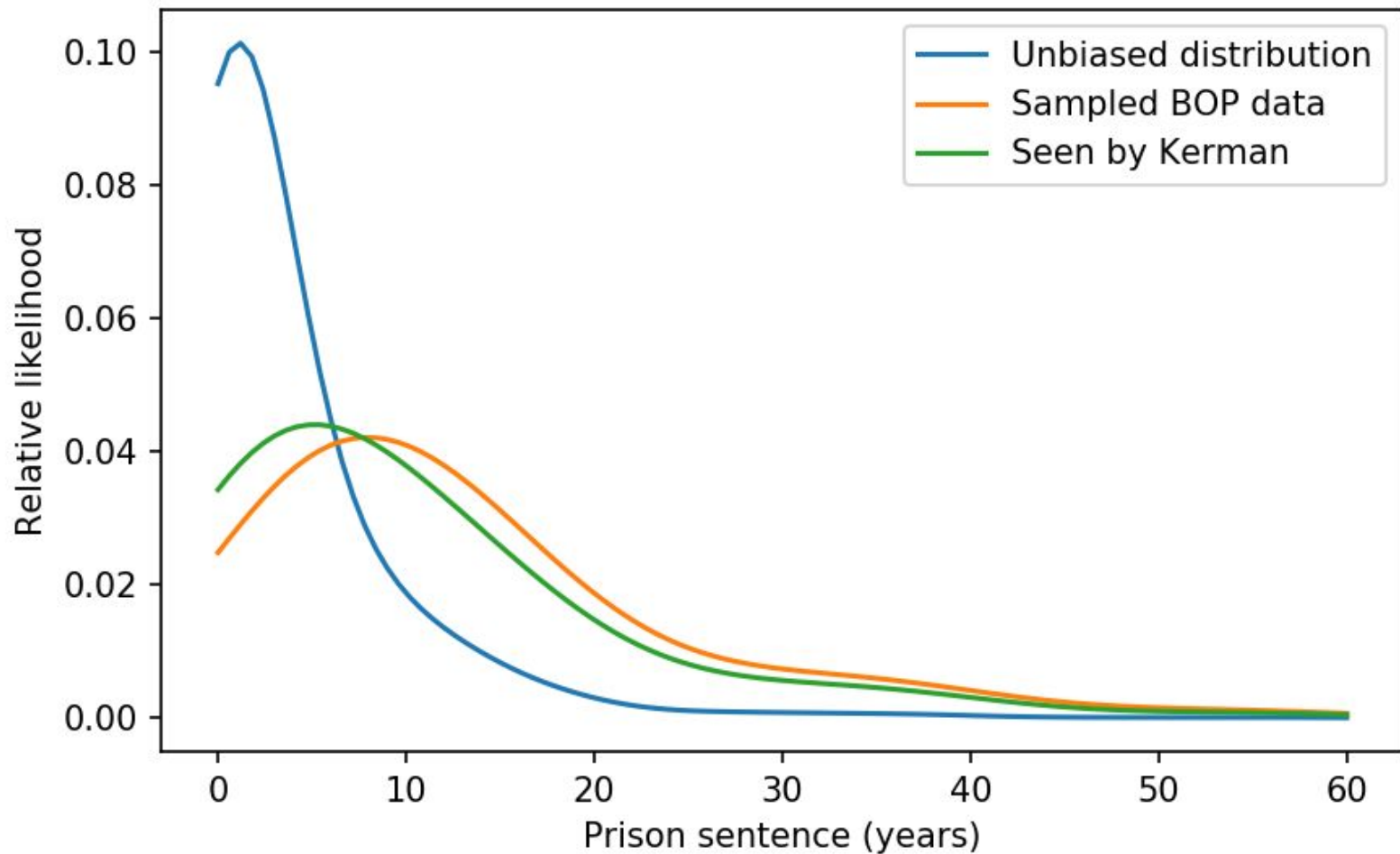
We can also compute the distribution of sentences as seen by someone at the prison for 13 months.

```
x = 0.85 * unbiased  
t = 13 / 12  
  
weights = x + t
```

Here's the sample.

```
kerman = resample_weighted(unbiased, weights)
```

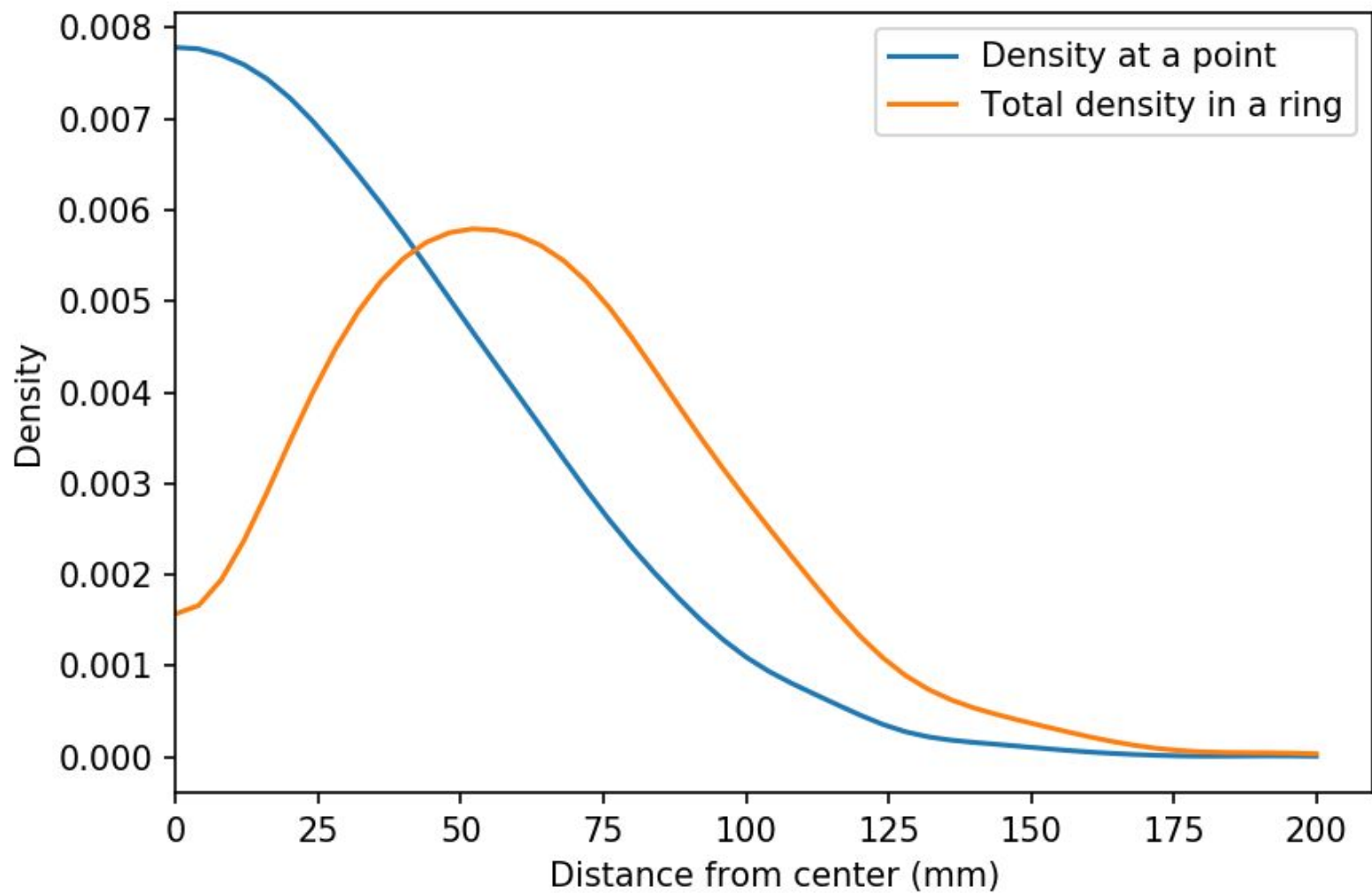
Distribution of federal prison sentences



The dartboard paradox

This happens in multiple dimensions, too.

Even more so.



Details in the notebook.

And on my blog,
Probably Overthinking It.

allendowney.com/blog

PROBABLY OVERTHINKING IT

Data Science, Bayesian Statistics, And Other Ideas

THE DARTBOARD PARADOX

 October 18, 2019  AllenDowney

On November 5, 2019, I will be at PyData NYC to give a talk called [The Inspection Paradox is Everywhere](#). Here's the abstract:

The inspection paradox is a statistical illusion you've probably never

ABOUT ME

Allen Downey is a professor at Olin College and the author of *Think Python*, *Think Bayes*, and other books available from Green Tea Press.

Summary

Length-biased sampling is everywhere.

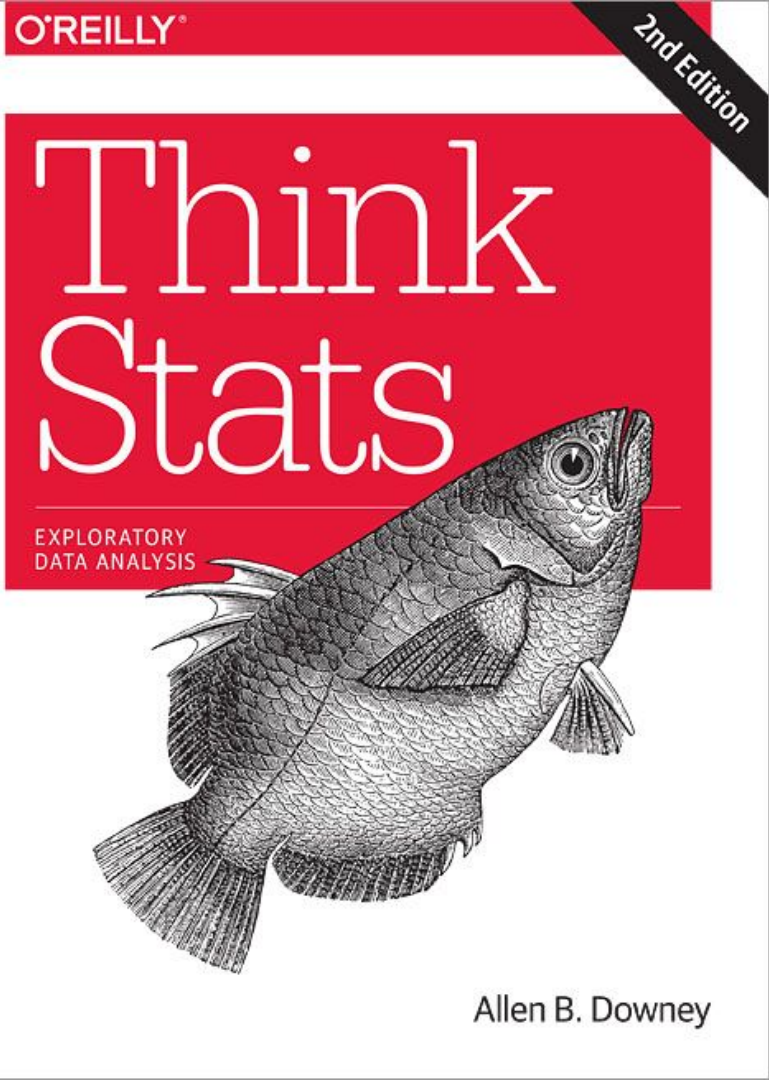
Can be a subtle source of error.

Also an opportunity for clever experimental design.

More reading

Class size example:

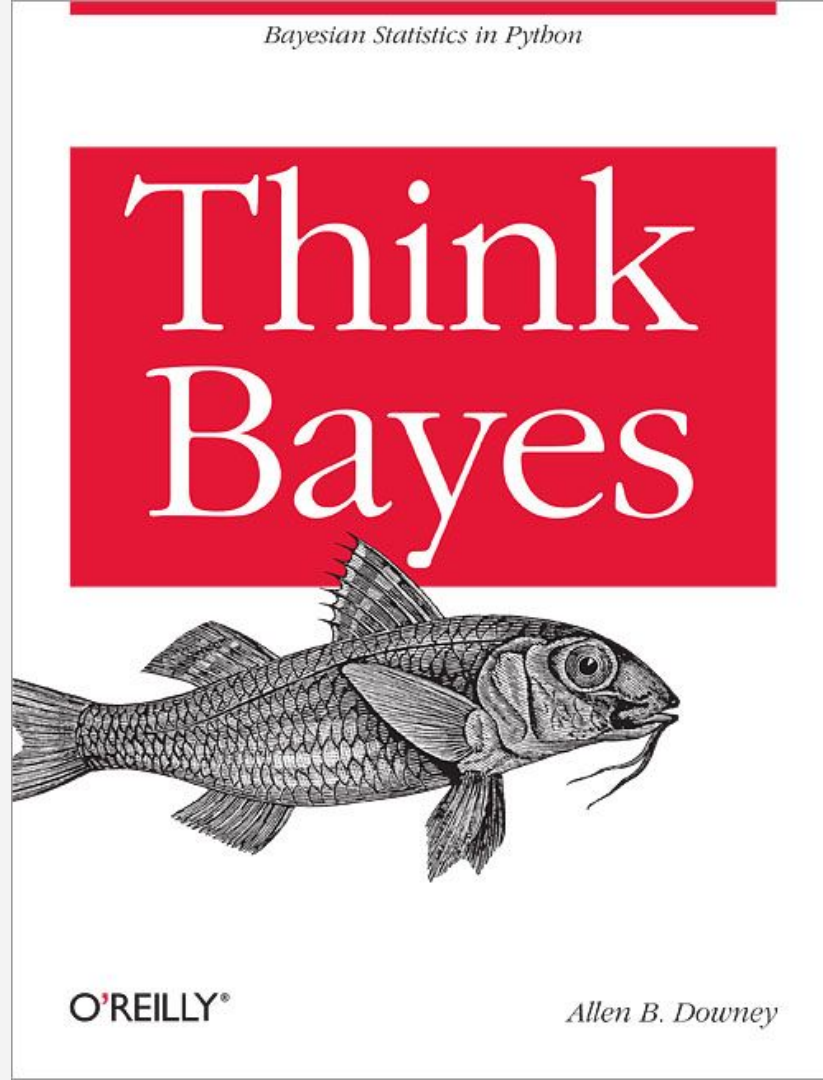
Think Stats, Chapter 3



More reading

Red line example:

Think Bayes, Chapter 8



The Inspection Paradox is Everywhere



Allen Downey

Aug 5 · 8 min read ★



The inspection paradox is a statistical illusion you've probably never heard of. It's a common source of confusion, an occasional cause of error, and an opportunity for clever experimental design. *

And once you know about it, you see it everywhere.

PyI



Part of **PyData** - 153 groups

PyData Boston - Cambridge

Boston, MA

1,181 members · Public group

Organized by **Milos Miljkovic** and 3 others

Share:

[About](#)[Events](#)[Members](#)[Photos](#)[Discussions](#)[More](#)

You're a member

What we're about

PyData is an educational program of NumFOCUS, a 501(c)3 non-profit organization in the United States. PyData provides a forum for the...

[Read more](#)

Past events (12)

[See all](#)

TUE, OCT 8, 6:00 PM

PyData Boston 2.0 Kickoff

HubSpot

Organizers



Milos Miljkovic and 3 others

[Message](#)

Members (1,181)

[See all](#)



Help me plan my sabbatical.

WHAT SHOULD I DO?

📅 September 19, 2019 👤 AllenDowney

I am planning to be on sabbatical from June 2020 to August 2021, so I am thinking about how to spend it. Let me tell you what I *can* do, and you can tell me what I *should* do.

Data Science

I consider myself a data scientist, but that means different things to different people. More specifically, I can contribute in the following areas:

- Data exploration, modeling, and prediction,
- Bayesian statistics and machine learning,
- Scientific computing and optimization,
- Software engineering and reproducible science
- Technical communication, including data visualization.

I have written a series of books related to data science and scientific computing, including *Think Stats*, *Think Bayes*, *Physical Modeling in MATLAB*, and *Modeling and Simulation in Python*.

blog

website

github

downey@allendowney.com/blog

twitter

email



Why Your Friends Have More Friends Than You Do

Scott L. Feld

American Journal of Sociology

Vol. 96, No. 6 (May, 1991), pp. 1464-1477

Published by: [The University of Chicago Press](http://www.jstor.org/stable/2781907)

Stable URL: <http://www.jstor.org/stable/2781907>

Page Count: 14

Article

Thumbnails

References

Viewing page 1464 of pages 1464-1477

Why Your Friends Have More Friends than You Do¹

Scott L. Feld

State University of New York at Stony Brook