

# Power-Law Size Distributions

Last updated: 2020/09/12, 13:39:25 EDT

Principles of Complex Systems, Vol. 1 | @pocsvox  
CSYS/MATH 300, Fall, 2020

Prof. Peter Sheridan Dodds | @peterdodds

Computational Story Lab | Vermont Complex Systems Center  
Vermont Advanced Computing Core | University of Vermont



Licensed under the [Creative Commons Attribution-NonCommercial-ShareAlike 3.0 License](https://creativecommons.org/licenses/by-nc-sa/3.0/).

## Outline

Our Intuition

Definition

Examples

Wild vs. Mild

CCDFs

Zipf's law

Zipf  $\leftrightarrow$  CCDF

References

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



1 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



2 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

Two of the many things we struggle with cognitively:

1. Probability.

- Ex. [The Monty Hall Problem](#).
- Ex. [Daughter/Son born on Tuesday](#). (see next two slides; Wikipedia entry [here](#).)

2. Logarithmic scales.

On counting and logarithms:



- Listen to Radiolab's 2009 piece: ["Numbers."](#)
- Later: [Benford's Law](#).

Also to be enjoyed: the magnificence of the [Dunning-Kruger effect](#)



3 of 64

## Homo probabilisticus?

The set up:

- A parent has two children.

Simple probability question:

- What is the probability that both children are girls?

The next set up:

- A parent has two children.
- We know one of them is a girl.

The next probabilistic poser:

- What is the probability that both children are girls?

Try this one:

- A parent has two children.
- We know one of them is a girl born on a Tuesday.

Simple question #3:

- What is the probability that both children are girls?

Last:

- A parent has two children.
- We know one of them is a girl born on December 31.

And ...

- What is the probability that both children are girls?

Let's test our collective intuition:



Money  
 $\equiv$   
Belief

Two questions about wealth distribution in the United States:

- Please estimate the percentage of all wealth owned by individuals when grouped into quintiles.
- Please estimate what you believe each quintile should own, ideally.
- Extremes: 100, 0, 0, 0 and 20, 20, 20, 20, 20



6 of 64

Wealth distribution in the United States: [13]

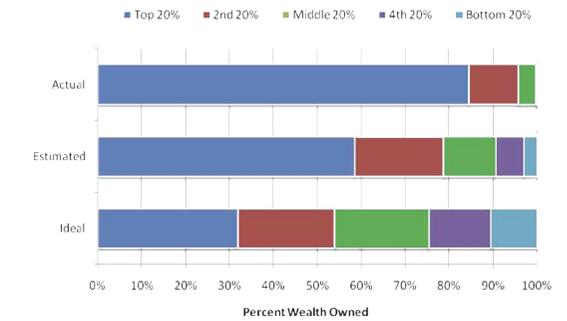


Fig. 2. The actual United States wealth distribution plotted against the estimated and ideal distributions across all respondents. Because of their small percentage share of total wealth, both the "4th 20%" value (0.2%) and the "Bottom 20%" value (0.1%) are not visible in the "Actual" distribution.

"Building a better America—One wealth quintile at a time" Norton and Ariely, 2011. [13]

Wealth distribution in the United States: [13]

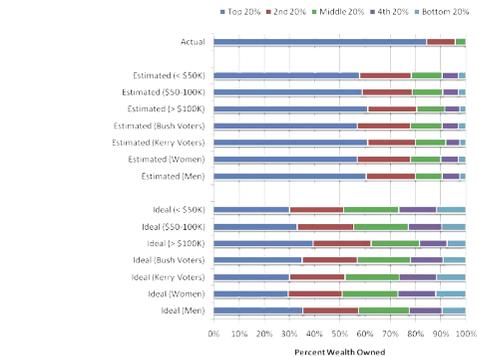


Fig. 3. The actual United States wealth distribution plotted against the estimated and ideal distributions of respondents of different income levels, political affiliations, and genders. Because of their small percentage share of total wealth, both the "4th 20%" value (0.2%) and the "Bottom 20%" value (0.1%) are not visible in the "Actual" distribution.

A highly watched video based on this research is

The sizes of many systems' elements appear to obey an inverse power-law size distribution:

$$P(\text{size} = x) \sim cx^{-\gamma}$$

where  $0 < x_{\min} < x < x_{\max}$  and  $\gamma > 1$ .

- $x_{\min}$  = lower cutoff,  $x_{\max}$  = upper cutoff

- Negative linear relationship in log-log space:

$$\log_{10} P(x) = \log_{10} c - \gamma \log_{10} x$$

- We use base 10 because we are **good people**.

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



4 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



5 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



6 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



7 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



9 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



12 of 64

## Size distributions:

Usually, only the tail of the distribution obeys a power law:

$$P(x) \sim c x^{-\gamma} \text{ for } x \text{ large.}$$

Still use term 'power-law size distribution.'

Other terms:

- Fat-tailed distributions.
- Heavy-tailed distributions.

Beware:

Inverse power laws aren't the only ones: lognormals, Weibull distributions, ...

## Size distributions:

Many systems have discrete sizes  $k$ :

- Word frequency
- Node degree in networks: # friends, # hyperlinks, etc.
- # citations for articles, court decisions, etc.

$$P(k) \sim c k^{-\gamma}$$

where  $k_{\min} \leq k \leq k_{\max}$

- Obvious fail for  $k = 0$ .
- Again, typically a description of distribution's tail.

## Word frequency:

Brown Corpus (~ 10<sup>6</sup> words):

rank	word	% q	rank	word	% q
1.	the	6.8872	1945.	apply	0.0055
2.	of	3.5839	1946.	vital	0.0055
3.	and	2.8401	1947.	September	0.0055
4.	to	2.5744	1948.	review	0.0055
5.	a	2.2996	1949.	wage	0.0055
6.	in	2.1010	1950.	motor	0.0055
7.	that	1.0428	1951.	fifteen	0.0055
8.	is	0.9943	1952.	regarded	0.0055
9.	was	0.9661	1953.	draw	0.0055
10.	he	0.9392	1954.	wheel	0.0055
11.	for	0.9340	1955.	organized	0.0055
12.	it	0.8623	1956.	vision	0.0055
13.	with	0.7176	1957.	wild	0.0055
14.	as	0.7137	1958.	Palmer	0.0055
15.	his	0.6886	1959.	intensity	0.0055

## Jonathan Harris's Wordcount:

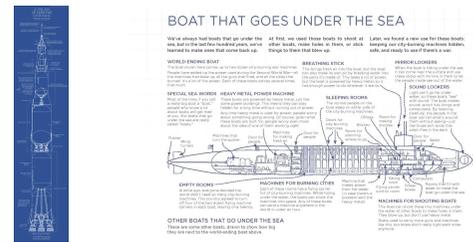
A word frequency distribution explorer:

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

13 of 64

## "Thing Explainer: Complicated Stuff in Simple Words"

by Randall Munroe (2015).<sup>[11]</sup>



Up goer five

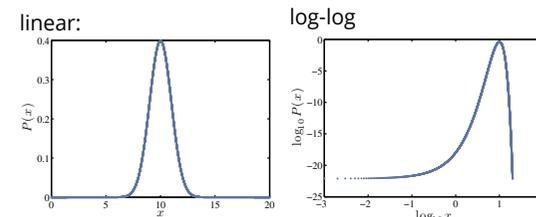
Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

14 of 64

## The statistics of surprise—words:

First—a Gaussian example:

$$P(x)dx = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} dx$$



mean  $\mu = 10$ , variance  $\sigma^2 = 1$ .

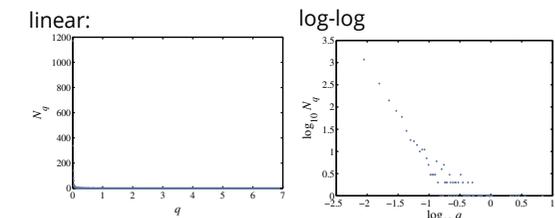
Activity: Sketch  $P(x) \sim x^{-1}$  for  $x = 1$  to  $x = 10^7$ .

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

15 of 64

## The statistics of surprise—words:

Raw 'probability' (binned) for Brown Corpus:



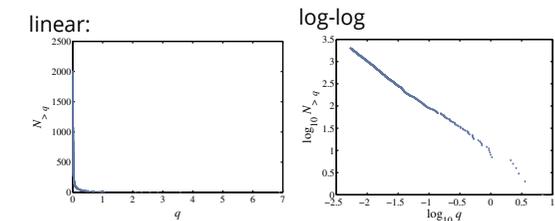
- $q_w$  = normalized frequency of occurrence of word  $w$  (%).
- $N_q$  = number of distinct words that have a normalized frequency of occurrence  $q$ .
- e.g.  $q_{\text{the}} \approx 6.9\%$ ,  $N_{q_{\text{the}}} = 1$ .

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

16 of 64

## The statistics of surprise—words:

Complementary Cumulative Probability Distribution  $N_{>q}$ :



Also known as the 'Exceedance Probability.'

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

17 of 64

## My, what big words you have ...



Test capitalizes on word frequency following a heavily skewed frequency distribution with a decaying power-law tail.

This Man Can Pronounce Every Word in the Dictionary (story here)

Best of Dr. Bailly

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

18 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

19 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

20 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

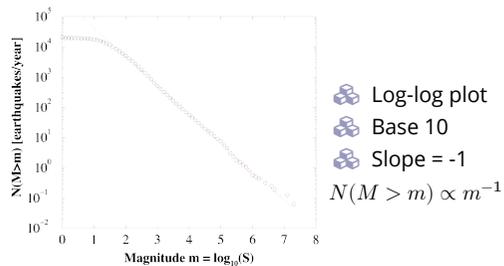
Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

21 of 64

# The statistics of surprise:

## Gutenberg-Richter law



From both the very awkwardly similar Christensen et al. and Bak et al.:  
 "Unified scaling law for earthquakes"<sup>[4, 1]</sup>



22 of 64

PoCS, Vol. 1  
 @pocsvox  
 Power-Law Size Distributions

Our Intuition  
 Definition  
**Examples**  
 Wild vs. Mild  
 CCDFs  
 Zipf's law  
 Zipf ⇔ CCDF  
 References



23 of 64

PoCS, Vol. 1  
 @pocsvox  
 Power-Law Size Distributions

Our Intuition  
 Definition  
**Examples**  
 Wild vs. Mild  
 CCDFs  
 Zipf's law  
 Zipf ⇔ CCDF  
 References



24 of 64

## "On a class of skew distribution functions"

Herbert A. Simon,  
 Biometrika, **42**, 425–440, 1955. <sup>[15]</sup>

## "Power laws, Pareto distributions and Zipf's law"

M. E. J. Newman,  
 Contemporary Physics, **46**, 323–351, 2005. <sup>[12]</sup>

## "Power-law distributions in empirical data"

Clauset, Shalizi, and Newman,  
 SIAM Review, **51**, 661–703, 2009. <sup>[5]</sup>

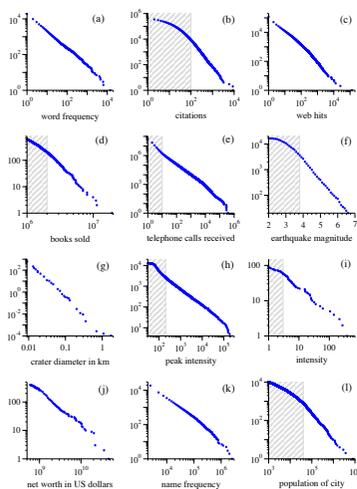


FIG. 1. Cumulative distribution of "rich frequency" of words. The distributions were computed as described in Appendix A. Data in the shaded region were evaluated from the calculations of the exponents in Table 1. Source references for the data are given in the text. (a) Numbers of occurrences of words in the novel *Harry Potter and the Chamber of Secrets*. (b) Numbers of hits on web sites by 60,000 users of the America Online Internet service for the day of 1 December 1997. (c) Numbers of hits on web sites by 60,000 users of the America Online Internet service for the day of 1 December 1997. (d) Numbers of books sold in the US for a random date. (e) Magnitude of earthquakes in California between January 1955 and May 1992. Magnitude is proportional to the logarithm of the maximum amplitude of the earthquake, and hence the distribution obeys a power law. (f) Telephone calls received in the US for a random date. (g) Diameter of craters in the US for a random date. (h) Peak, gamma-ray intensity of solar flares in counts per second, measured from Earth or that between February 1980 and November 1989. (i) Intensity of waves from 1916 to 1980, measured as backscatters per 1000 of the population of the US for a random date. (j) Peak, gamma-ray intensity of solar flares in counts per second, measured from Earth or that between February 1980 and November 1989. (k) Intensity of waves from 1916 to 1980, measured as backscatters per 1000 of the population of the US for a random date. (l) Population of US cities in the year 1990. (m) Frequency of occurrence of family names in the US in the year 1990.

## Size distributions:

### Some examples:

- Earthquake magnitude (Gutenberg-Richter law)<sup>[9, 1]</sup>  $P(M) \propto M^{-2}$
  - # war deaths:<sup>[14]</sup>  $P(d) \propto d^{-1.8}$
  - Sizes of forest fires<sup>[8]</sup>
  - Sizes of cities:<sup>[15]</sup>  $P(n) \propto n^{-2.1}$
  - # links to and from websites<sup>[2]</sup>
- Note: Exponents range in error



24 of 64

## Size distributions:

### More examples:

- # citations to papers:<sup>[6, 7]</sup>  $P(k) \propto k^{-3}$ .
- Individual wealth (maybe):  $P(W) \propto W^{-2}$ .
- Distributions of tree trunk diameters:  $P(d) \propto d^{-2}$ .
- The gravitational force at a random point in the universe:<sup>[10]</sup>  $P(F) \propto F^{-5/2}$ . (See the Holtmark distribution<sup>[7]</sup> and stable distributions<sup>[8]</sup>.)
- Diameter of moon craters:<sup>[12]</sup>  $P(d) \propto d^{-3}$ .
- Word frequency:<sup>[15]</sup> e.g.,  $P(k) \propto k^{-2.2}$  (variable).
- # religious adherents in cults:<sup>[5]</sup>  $P(k) \propto k^{-1.8 \pm 0.1}$ .
- # sightings of birds per species (North American Breeding Bird Survey for 2003):<sup>[5]</sup>  $P(k) \propto k^{-2.1 \pm 0.1}$ .
- # species per genus:<sup>[17, 15, 5]</sup>  $P(k) \propto k^{-2.4 \pm 0.2}$ .



25 of 64

PoCS, Vol. 1  
 @pocsvox  
 Power-Law Size Distributions

Our Intuition  
 Definition  
**Examples**  
 Wild vs. Mild  
 CCDFs  
 Zipf's law  
 Zipf ⇔ CCDF  
 References



26 of 64

PoCS, Vol. 1  
 @pocsvox  
 Power-Law Size Distributions

Our Intuition  
 Definition  
**Examples**  
 Wild vs. Mild  
 CCDFs  
 Zipf's law  
 Zipf ⇔ CCDF  
 References



27 of 64

## Table 3 from Clauset, Shalizi, and Newman<sup>[5]</sup>:

Basic parameters of the data sets described in section 6, along with their power-law fits and the corresponding p-values (statistically significant values are denoted in bold).

Quantity	n	$\langle x \rangle$	$\sigma$	$P_{max}$	$x_{min}$	$\hat{\alpha}$	$n_{tail}$	$P$
count of word use	18,855	11.14	1.84(3)	11,086	7 ± 2	1.95(2)	2938 ± 987	<b>0.49</b>
protein interaction degree	1846	2.34	3.05	56	5 ± 2	3.1(3)	204 ± 263	<b>0.31</b>
metabolic degree	1641	5.68	17.81	468	4 ± 1	2.8(1)	748 ± 136	<b>0.00</b>
Internet degree	22,688	5.63	37.83	2583	21 ± 9	2.12(9)	770 ± 1124	<b>0.29</b>
telephone calls received	51,360,423	3.88	179.09	375,746	120 ± 49	2.09(1)	102,292 ± 210,147	<b>0.63</b>
intensity of wars	115	15.70	49.97	382	2.1 ± 3.5	1.7(2)	70 ± 14	<b>0.20</b>
terrorist attack severity	9101	4.35	31.58	2749	12 ± 4	2.4(2)	547 ± 1663	<b>0.68</b>
HTTP size (kilobytes)	239,386	7.36	57.94	10,971	36.35 ± 22.74	2.48(5)	673 ± 2232	<b>0.00</b>
species per genus	509	5.59	6.94	56	4 ± 2	2.4(2)	233 ± 138	<b>0.10</b>
bird species sightings	591	3384.36	10,952.34	138,705	6679 ± 2463	2.1(2)	66 ± 41	<b>0.55</b>
blackouts ( $\times 10^3$ )	211	253.87	610.31	7500	239 ± 90	2.3(3)	59 ± 35	<b>0.62</b>
sales of books ( $\times 10^3$ )	633	1986.67	1396.60	19,077	2400 ± 430	3.7(3)	139 ± 115	<b>0.66</b>
population of cities ( $\times 10^3$ )	19,447	9.00	77.83	8006	52.46 ± 11.88	2.37(8)	580 ± 177	<b>0.76</b>
email address books size	4581	12.45	21.49	333	57 ± 21	3.5(6)	196 ± 449	<b>0.16</b>
forest fire size (acres)	203,785	0.90	20.99	4121	6324 ± 3487	2.2(9)	521 ± 6801	<b>0.05</b>
solar flare intensity	12,773	689.41	6520.59	231,300	323 ± 89	1.79(2)	1711 ± 384	<b>1.00</b>
quake intensity ( $\times 10^3$ )	19,902	24.54	563.83	63,096	0.794 ± 80.198	1.64(4)	11,697 ± 2159	<b>0.00</b>
religious followers ( $\times 10^3$ )	103	27.26	136.64	1050	3.85 ± 1.09	1.8(1)	39 ± 26	<b>0.42</b>
freq. of surnames ( $\times 10^3$ )	2753	50.59	113.99	2502	111.92 ± 40.67	2.5(2)	239 ± 215	<b>0.20</b>
net worth (mil. USD)	400	2388.69	4167.35	46,000	990 ± 364	2.3(1)	302 ± 77	<b>0.00</b>
citations to papers	412,229	16.17	44.02	8904	160 ± 35	3.16(6)	345 ± 1859	<b>0.20</b>
papers authored	401,445	7.21	16.52	1416	133 ± 13	4.3(1)	988 ± 377	<b>0.90</b>
hits to web sites	119,724	9.83	392.52	129,641	2 ± 13	1.81(8)	50,981 ± 16,898	<b>0.00</b>
links to web sites	241,428,853	9.15	106,871.65	1,199,466	3684 ± 151	2.336(9)	29,986 ± 1560	<b>0.00</b>

We'll explore various exponent measurement techniques in assignments.

## power-law size distributions

### Gaussians versus power-law size distributions:

- Mediocristan versus Extremistan
  - Mild versus Wild (Mandelbrot)
  - Example: Height versus wealth.
- 
- See "The Black Swan" by Nassim Taleb.<sup>[16]</sup>
  - Terrible if successful framing: Black swans are not that surprising ...



27 of 64

PoCS, Vol. 1  
 @pocsvox  
 Power-Law Size Distributions

Our Intuition  
 Definition  
**Examples**  
 Wild vs. Mild  
 CCDFs  
 Zipf's law  
 Zipf ⇔ CCDF  
 References



28 of 64

PoCS, Vol. 1  
 @pocsvox  
 Power-Law Size Distributions

Our Intuition  
 Definition  
**Examples**  
 Wild vs. Mild  
 CCDFs  
 Zipf's law  
 Zipf ⇔ CCDF  
 References



28 of 64

PoCS, Vol. 1  
 @pocsvox  
 Power-Law Size Distributions

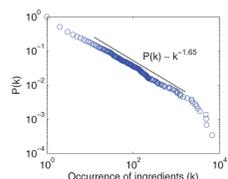
Our Intuition  
 Definition  
**Examples**  
 Wild vs. Mild  
 CCDFs  
 Zipf's law  
 Zipf ⇔ CCDF  
 References



30 of 64

## "Geography and similarity of regional cuisines in China"

Zhu et al.,  
 PLoS ONE, **8**, e79161, 2013. <sup>[18]</sup>



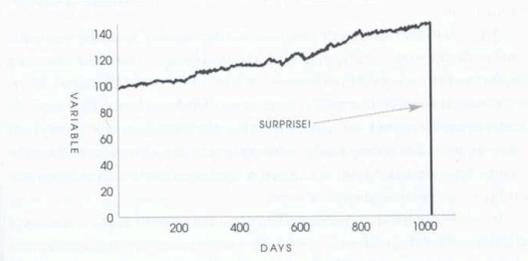
- Fraction of ingredients that appear in at least  $k$  recipes.
- Oops in notation:  $P(k)$  is the Complementary Cumulative Distribution  $P_{\geq}(k)$



24 of 64

# Turkeys ...

FIGURE 1: ONE THOUSAND AND ONE DAYS OF HISTORY



A turkey before and after Thanksgiving. The history of a process over a thousand days tells you nothing about what is to happen next. This naïve projection of the future from the past can be applied to anything.

From "The Black Swan" [16]

# Taleb's table [16]

## Mediocristan/Extremistan

- Most typical member is mediocre/Most typical is either giant or tiny
- Winners get a small segment/Winner take almost all effects
- When you observe for a while, you know what's going on/It takes a very long time to figure out what's going on
- Prediction is easy/Prediction is hard
- History crawls/History makes jumps
- Tyranny of the collective/Tyranny of the rare and accidental

# Size distributions:



Power-law size distributions are sometimes called Pareto distributions after Italian scholar Vilfredo Pareto.

- Pareto noted wealth in Italy was distributed unevenly (80-20 rule; misleading).
- Term used especially by practitioners of the Dismal Science.

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

# Devilish power-law size distribution details:

## Exhibit A:

Given  $P(x) = cx^{-\gamma}$  with  $0 < x_{\min} < x < x_{\max}$ , the mean is ( $\gamma \neq 2$ ):

$$\langle x \rangle = \frac{c}{2-\gamma} (x_{\max}^{2-\gamma} - x_{\min}^{2-\gamma})$$

- Mean 'blows up' with upper cutoff if  $\gamma < 2$ .
- Mean depends on lower cutoff if  $\gamma > 2$ .
- $\gamma < 2$ : Typical sample is large.
- $\gamma > 2$ : Typical sample is small.

Insert question from assignment 2

# And in general ...

## Moments:

- All moments depend only on cutoffs.
- No internal scale that dominates/matters.
- Compare to a Gaussian, exponential, etc.

## For many real size distributions: $2 < \gamma < 3$

- mean is finite (depends on lower cutoff)
- $\sigma^2$  = variance is 'infinite' (depends on upper cutoff)
- Width of distribution is 'infinite'
- If  $\gamma > 3$ , distribution is less terrifying and may be easily confused with other kinds of distributions.

Insert question from assignment 3

# Moments

## Standard deviation is a mathematical convenience:

- Variance is nice analytically ...
- Another measure of distribution width:

$$\text{Mean average deviation (MAD)} = \langle |x - \langle x \rangle| \rangle$$

- For a pure power law with  $2 < \gamma < 3$ :

$$\langle |x - \langle x \rangle| \rangle \text{ is finite.}$$

- But MAD is mildly unpleasant analytically ...
- We still speak of infinite 'width' if  $\gamma < 3$ .

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

# How sample sizes grow ...

Given  $P(x) \sim cx^{-\gamma}$ :

- We can show that after  $n$  samples, we expect the largest sample to be<sup>1</sup>

$$x_1 \approx c' n^{1/(\gamma-1)}$$

- Sampling from a finite-variance distribution gives a much slower growth with  $n$ .
- e.g., for  $P(x) = \lambda e^{-\lambda x}$ , we find

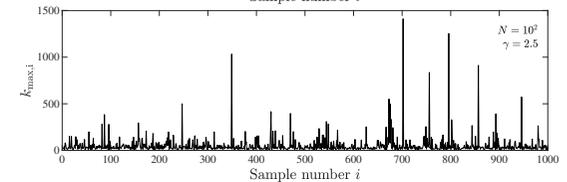
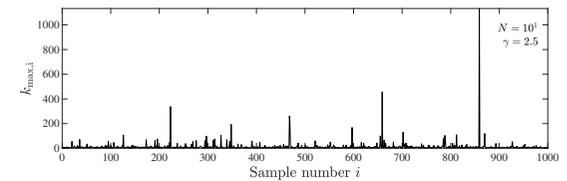
$$x_1 \approx \frac{1}{\lambda} \ln n.$$

Insert question from assignment 4

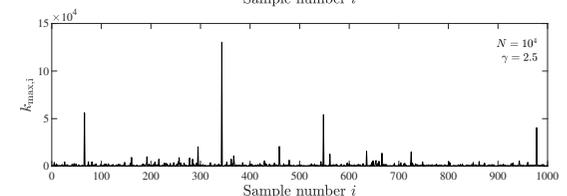
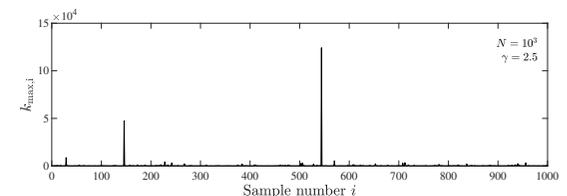
Insert question from assignment 6

<sup>1</sup>Later, we see that the largest sample grows as  $n^\rho$  where  $\rho$  is the Zipf exponent

$\gamma = 5/2$ , maxima of  $N$  samples,  $n = 1000$  sets of samples:



$\gamma = 5/2$ , maxima of  $N$  samples,  $n = 1000$  sets of samples:

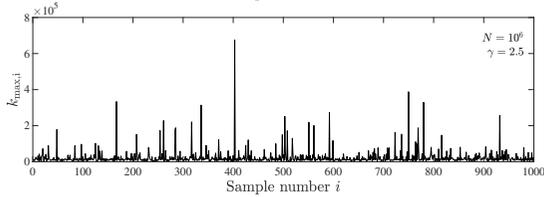
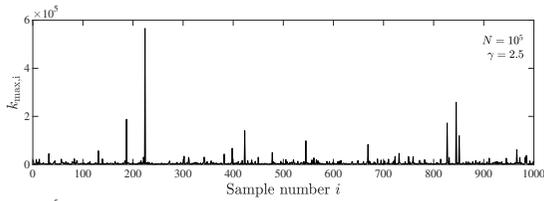


Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf ⇔ CCDF  
References

$\gamma = 5/2$ , maxima of  $N$  samples,  $n = 1000$  sets of samples:

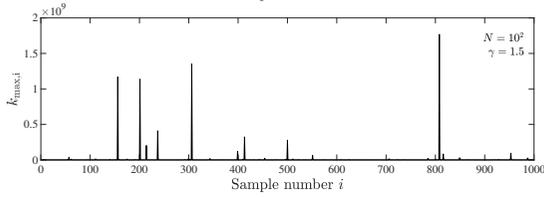
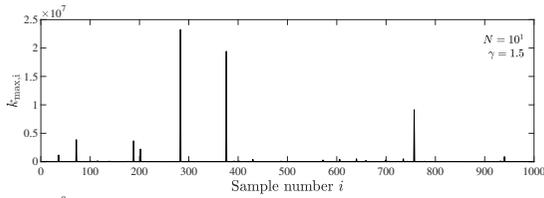


Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



40 of 64

$\gamma = 3/2$ , maxima of  $N$  samples,  $n = 1000$  sets of samples:

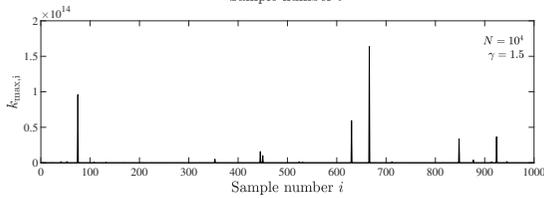
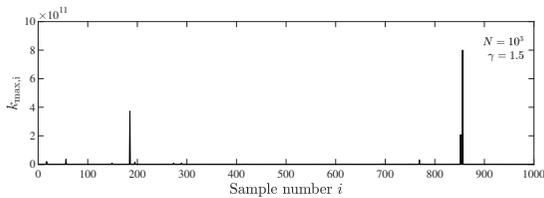


Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



41 of 64

$\gamma = 3/2$ , maxima of  $N$  samples,  $n = 1000$  sets of samples:

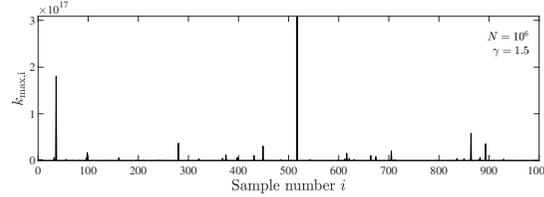
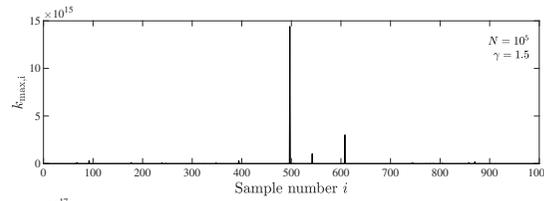


Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



42 of 64

$\gamma = 3/2$ , maxima of  $N$  samples,  $n = 1000$  sets of samples:

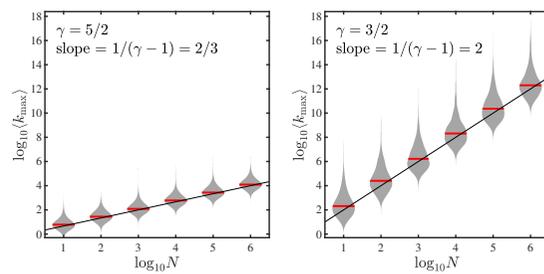


Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



43 of 64

Scaling of expected largest value as a function of sample size  $N$ :



Fit for  $\gamma = 5/2$ :  $k_{\max} \sim N^{0.660 \pm 0.066}$  (sublinear)  
Fit for  $\gamma = 3/2$ :  $k_{\max} \sim N^{2.063 \pm 0.215}$  (superlinear)

295% confidence interval

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



44 of 64

Complementary Cumulative Distribution Function:  
CCDF:

$$P_{\geq}(x) = P(x' \geq x) = 1 - P(x' < x)$$

$$= \int_{x'=x}^{\infty} P(x') dx'$$

$$\propto \int_{x'=x}^{\infty} (x')^{-\gamma} dx'$$

$$= \frac{1}{-\gamma + 1} (x')^{-\gamma + 1} \Big|_{x'=x}^{\infty}$$

$$\propto x^{-\gamma + 1}$$

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



45 of 64

Complementary Cumulative Distribution Function:

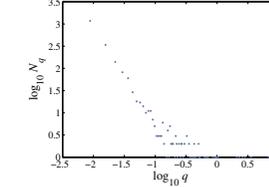
CCDF:



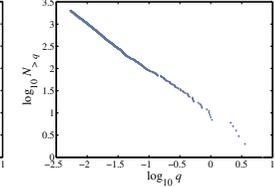
$$P_{\geq}(x) \propto x^{-\gamma + 1}$$

- Use when tail of  $P$  follows a power law.
- Increases exponent by one.
- Useful in cleaning up data.

PDF:



CCDF:



Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



46 of 64

Complementary Cumulative Distribution Function:

Same story for a discrete variable:  $P(k) \sim ck^{-\gamma}$ .



$$P_{\geq}(k) = P(k' \geq k)$$

$$= \sum_{k'=k}^{\infty} P(k')$$

$$\propto k^{-\gamma + 1}$$

Use integrals to approximate sums.

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



47 of 64

Zipfian rank-frequency plots

George Kingsley Zipf:

- Noted various rank distributions have power-law tails, often with exponent -1 (word frequency, city sizes, ...)
- Zipf's 1949 Magnum Opus:

We'll study Zipf's law in depth ...

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References



49 of 64

# Zipfian rank-frequency plots

## Zipf's way:

- Given a collection of entities, rank them by size, largest to smallest.
- $x_r$  = the size of the  $r$ th ranked entity.
- $r = 1$  corresponds to the largest size.
- Example:  $x_1$  could be the frequency of occurrence of the most common word in a text.
- Zipf's observation:

$$x_r \propto r^{-\alpha}$$

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

50 of 64

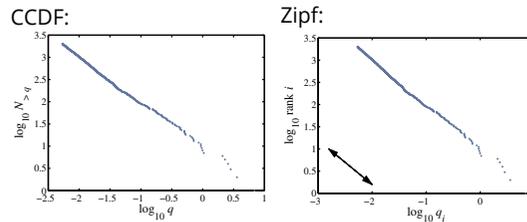
PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

51 of 64

# Size distributions:

## Brown Corpus (1,015,945 words):



- The, of, and, to, a, ... = 'objects'
- 'Size' = word frequency
- Beep: (Important) CCDF and Zipf plots are related
- ...

53 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

54 of 64

## Observe:

- $N_{P_{\geq}(x)}$  = the number of objects with size at least  $x$  where  $N$  = total number of objects.
- If an object has size  $x_r$ , then  $N_{P_{\geq}(x_r)}$  is its rank  $r$ .
- So

$$x_r \propto r^{-\alpha} = (N_{P_{\geq}(x_r)})^{-\alpha}$$

$$\propto x_r^{-(\gamma+1)(-\alpha)} \text{ since } P_{\geq}(x) \sim x^{-\gamma+1}.$$

We therefore have  $1 = (-\gamma + 1)(-\alpha)$  or:

$$\alpha = \frac{1}{\gamma - 1}$$

- A rank distribution exponent of  $\alpha = 1$  corresponds to a size distribution exponent  $\gamma = 2$ .

55 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

52 of 64

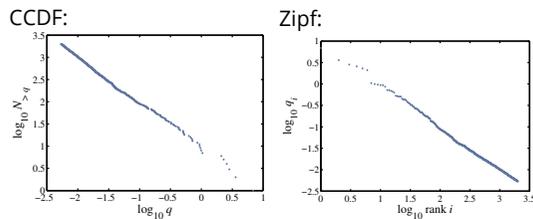


- Examined all games of varying game depth  $d$  in a set of chess databases.
- $n$  = popularity = how many times a specific game path appears in databases.
- $S(n; d)$  = number of depth  $d$  games with popularity  $n$ .
- Show "the frequencies of opening moves are distributed according to a power law with an exponent that increases linearly with the game depth, whereas the pooled distribution of all opening weights follows Zipf's law with universal exponent."
- Propose hierarchical fragmentation model that produces self-similar game trees.

55 of 64

# Size distributions:

## Brown Corpus (1,015,945 words):



- The, of, and, to, a, ... = 'objects'
- 'Size' = word frequency
- Beep: (Important) CCDF and Zipf plots are related
- ...

52 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

53 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

54 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

55 of 64

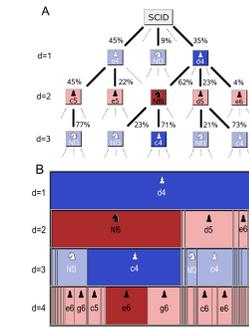
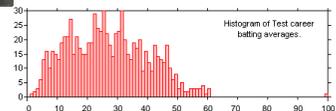


FIG. 1 (color online). (a) Schematic representation of the weighted game tree of chess based on the SCID database and with logarithmic binning. A straight line fit (not shown) yields an exponent of  $\alpha = 2.05$  with a goodness of fit  $R^2 > 0.9992$ . For comparison, the Zipf distribution Eq. (5) with  $\mu = 1$  is indicated as a solid line. Inset: number  $C(n) = \sum_{i=1}^n S(n)$  of openings of depth  $d$  with a given popularity  $n$  for  $d = 16$  and histograms with logarithmic binning for  $d = 4, d = 16$ , and  $d = 22$ . Solid lines are regression lines to the logarithmically binned data ( $R^2 > 0.99$  for  $d < 35$ ). Inset: slope  $\sigma_d$  of the regression line as a function of  $d$  and the analytical estimation Eq. (6) using  $N = 1.4 \times 10^6$  and  $\beta = 0$  (solid line).

# The Don.

## Extreme deviations in test cricket:



- Don Bradman's batting average = 166% next best.
- That's pretty solid.
- Later in the course: Understanding success—the Mona Lisa like Don Bradman?

## A good eye:

- The great Paul Kelly's tribute to the man who was "Something like the tide"

58 of 64

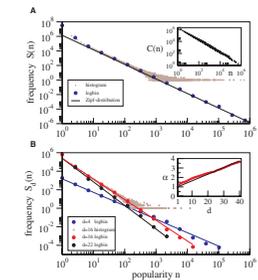


FIG. 2 (color online). (a) Histogram of weight frequencies  $S(n)$  of openings up to  $d = 40$  in the SCID database and with logarithmic binning. A straight line fit (not shown) yields an exponent of  $\alpha = 2.05$  with a goodness of fit  $R^2 > 0.9992$ . For comparison, the Zipf distribution Eq. (5) with  $\mu = 1$  is indicated as a solid line. Inset: number  $C(n) = \sum_{i=1}^n S(n)$  of openings of depth  $d$  with a given popularity  $n$  for  $d = 16$  and histograms with logarithmic binning for  $d = 4, d = 16$ , and  $d = 22$ . Solid lines are regression lines to the logarithmically binned data ( $R^2 > 0.99$  for  $d < 35$ ). Inset: slope  $\sigma_d$  of the regression line as a function of  $d$  and the analytical estimation Eq. (6) using  $N = 1.4 \times 10^6$  and  $\beta = 0$  (solid line).

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

56 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

57 of 64

PoCS, Vol. 1  
@pocsvox  
Power-Law Size Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
References

58 of 64

## References I

- [1] P. Bak, K. Christensen, L. Danon, and T. Scanlon. Unified scaling law for earthquakes. [Phys. Rev. Lett., 88:178501, 2002. pdf](#)
- [2] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. [Science, 286:509–511, 1999. pdf](#)
- [3] B. Blasius and R. Tönjes. Zipf's law in the popularity distribution of chess openings. [Phys. Rev. Lett., 103:218701, 2009. pdf](#)
- [4] K. Christensen, L. Danon, T. Scanlon, and P. Bak. Unified scaling law for earthquakes. [Proc. Natl. Acad. Sci., 99:2509–2513, 2002. pdf](#)

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
[References](#)



60 of 64

## References II

- [5] A. Clauset, C. R. Shalizi, and M. E. J. Newman. Power-law distributions in empirical data. [SIAM Review, 51:661–703, 2009. pdf](#)
- [6] D. J. de Solla Price. Networks of scientific papers. [Science, 149:510–515, 1965. pdf](#)
- [7] D. J. de Solla Price. A general theory of bibliometric and other cumulative advantage processes. [J. Amer. Soc. Inform. Sci., 27:292–306, 1976. pdf](#)
- [8] P. Grassberger. Critical behaviour of the Drossel-Schwabl forest fire model. [New Journal of Physics, 4:17.1–17.15, 2002. pdf](#)

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
[References](#)



61 of 64

## References III

- [9] B. Gutenberg and C. F. Richter. Earthquake magnitude, intensity, energy, and acceleration. [Bull. Seism. Soc. Am., 499:105–145, 1942. pdf](#)
- [10] J. Holtsmark. Über die verbreiterung von spektrallinien. [Ann. Phys., 58:577–, 1919.](#)
- [11] R. Munroe. [Thing Explainer: Complicated Stuff in Simple Words.](#) Houghton Mifflin Harcourt, 2015.
- [12] M. E. J. Newman. Power laws, Pareto distributions and Zipf's law. [Contemporary Physics, 46:323–351, 2005. pdf](#)

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
[References](#)



62 of 64

## References IV

- [13] M. I. Norton and D. Ariely. Building a better America—One wealth quintile at a time. [Perspectives on Psychological Science, 6:9–12, 2011. pdf](#)
- [14] L. F. Richardson. Variation of the frequency of fatal quarrels with magnitude. [J. Amer. Stat. Assoc., 43:523–546, 1949.](#)
- [15] H. A. Simon. On a class of skew distribution functions. [Biometrika, 42:425–440, 1955. pdf](#)
- [16] N. N. Taleb. [The Black Swan.](#) Random House, New York, 2007.

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
[References](#)



63 of 64

## References V

- [17] G. U. Yule. A mathematical theory of evolution, based on the conclusions of Dr J. C. Willis, F.R.S. [Phil. Trans. B, 213:21–87, 1925. pdf](#)
- [18] Y.-X. Zhu, J. Huang, Z.-K. Zhang, Q.-M. Zhang, T. Zhou, and Y.-Y. Ahn. Geography and similarity of regional cuisines in China. [PLoS ONE, 8:e79161, 2013. pdf](#)
- [19] G. K. Zipf. [Human Behaviour and the Principle of Least-Effort.](#) Addison-Wesley, Cambridge, MA, 1949.

PoCS, Vol. 1  
@pocsvox  
Power-Law Size  
Distributions

Our Intuition  
Definition  
Examples  
Wild vs. Mild  
CCDFs  
Zipf's law  
Zipf  $\leftrightarrow$  CCDF  
[References](#)



64 of 64