# Text statistics

September 12, 2023

(Chapter 5.1 IIR)

# Overview

# Outline

# The Reuters collection

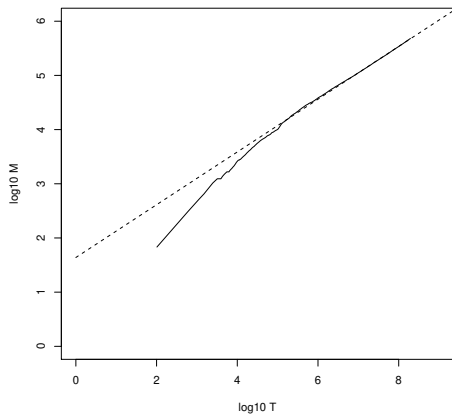| symbol | statistic | value |
|--------|-----------|-------|
| $N$ | documents | 800,000 |
| $L$ | avg. # word tokens per document | 200 |
| $M$ | word types | 400,000 |
|  | avg. # bytes per word token (incl. spaces/punct.) | 6 |
|  | avg. # bytes per word token (without spaces/punct.) | 4.5 |
|  | avg. # bytes per word type | 7.5 |
| $T$ | non-positional postings | 100,000,000 |

# How big is the term vocabulary?

- That is, how many distinct words are there?
- Can we assume there is an upper bound?
- Not really: At least $70^{20} \approx 10^{37}$ different words of length 20.
- The vocabulary will keep growing with collection size.
- Heaps' law:

$$M = kT^b \tag{1}$$

- $M$ is the size of the vocabulary, $T$ is the number of tokens in the collection.
- Typical values for the parameters $k$ and $b$ are: $30 \leq k \leq 100$ and $b \approx 0.5$.
- Heaps' law is linear in log-log space.
  - It is the simplest possible relationship between collection size and vocabulary size in log-log space.
  - Empirical law

# Heaps' law for Reuters



Vocabulary size $M$ as a function of collection size $T$ (number of tokens)

$$M = kT^b \qquad (2)$$

the best least squares fit for Reuters-RCV1.
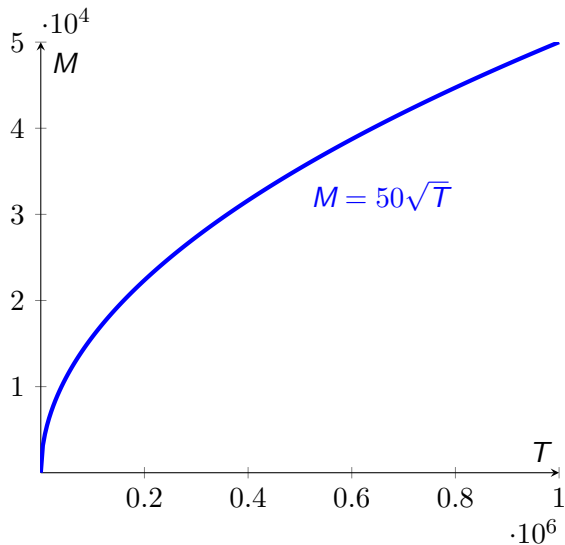
$$\log_{10} M = 0.49 * \log_{10} T + 1.64 \quad (3)$$

In other words,

$$M = 10^{1.64} T^{0.49} \qquad (4)$$

$$k = 10^{1.64} \approx 44 \qquad (5)$$

$$b = 0.49 \qquad (6)$$

# Heaps' law without loglog scale



$M = 50\sqrt{T}$

# Empirical fit for Reuters

- Example: for the first 1,000,020 tokens Heaps' law predicts 38,323 terms:

$$44 \times 1{,}000{,}020^{0.49} \approx 38{,}323$$

- The actual number is 38,365 terms, very close to the prediction.
- Empirical observation: fit is good in general.

# Exercise

1. What is the effect of including spelling errors vs. automatically correcting spelling errors on Heaps' law?

2. Compute vocabulary size $M$
   - Looking at a collection of web pages, you find that there are 3000 different terms in the first 10,000 tokens and 30,000 different terms in the first 1,000,000 tokens.
   - Assume a search engine indexes a total of 20,000,000,000 ($2 \times 10^{10}$) pages, containing 200 tokens on average per page
   - What is the size of the vocabulary of the indexed collection as predicted by Heaps' law?

# Outline

# Zipf's law

- Now we have characterized the growth of the vocabulary in collections.
- We also want to know how many frequent vs. infrequent terms we should expect in a collection.
- In natural language, there are a few very frequent terms and very many very rare terms.
- Zipf's law: The $i^{\text{th}}$ most frequent term has frequency $\mathrm{cf}_i$ proportional to $1/i$.
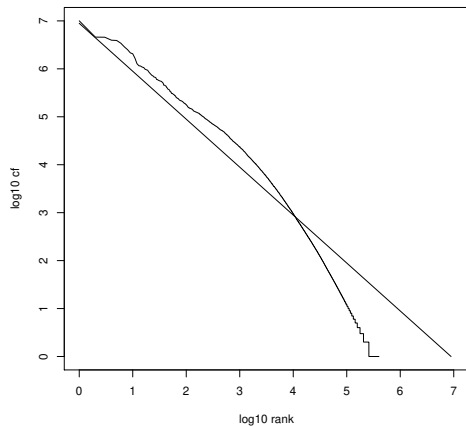
$$\mathrm{cf}_i \propto \frac{1}{i} \tag{7}$$

- $\mathrm{cf}_i$ is collection frequency: the number of occurrences of the term $t_i$ in the collection.
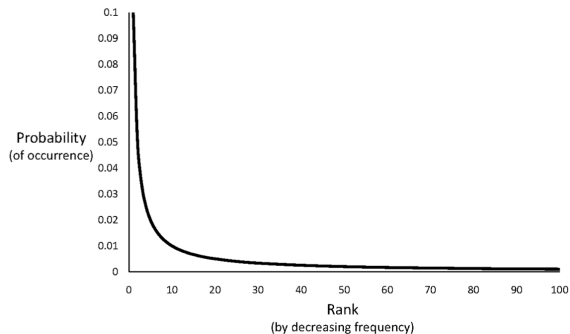
# Zipf's law

$$\mathrm{cf}_i \propto \frac{1}{i} \tag{8}$$

- If the most frequent term (*the*) occurs $\mathrm{cf}_1$ times, then the second most frequent term (*of*) has half as many occurrences $\mathrm{cf}_2 = \frac{1}{2}\mathrm{cf}_1$ …
- …and the third most frequent term (*and*) has a third as many occurrences $\mathrm{cf}_3 = \frac{1}{3}\mathrm{cf}_1$ etc.
- Equivalent: $\mathrm{cf}_i = ci^k$ and $\log \mathrm{cf}_i = \log c + k \log i$ (for $k = -1$)
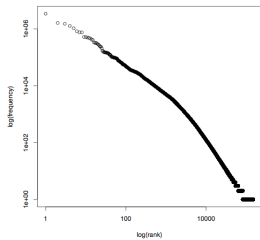- Example of a power law

# Zipf's law for Reuters



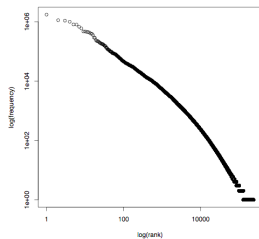Fit is not great. What is important is the key insight: Few frequent terms, many rare terms.
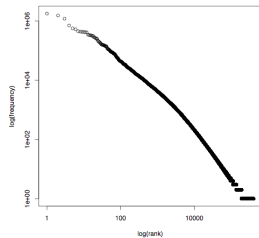
# Zipf's law if not log-loged
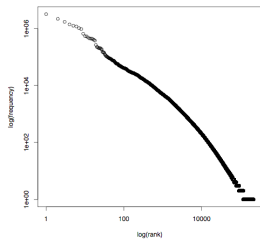
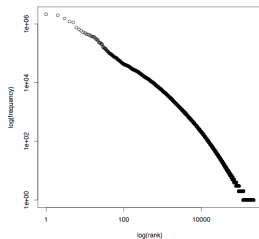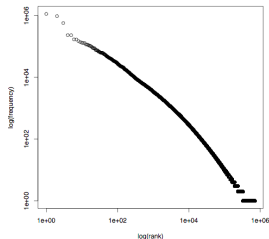# Zipf's law in different languages (parliament documents)



English

italian

germany

spanish

portuguese

finnish

# Zipf's law in scientific writing



relativity by einstein



on the origin of species by charles darwin

# Zipf's law for kid's books



alice                                                peterpan

| Rank | Word | Frequency |
|---|---|---|
| 1000 | concern | 5,100 |
| 1001 | spoke | 5,100 |
| 1002 | summit | 5,100 |
| 1003 | bring | 5,099 |
| 1004 | star | 5,099 |
| 1005 | immediate | 5,099 |
| 1006 | chemical | 5,099 |
| 1007 | african | 5,098 |

| Word | Freq. | r | $P_r(\%)$ | $r.P_r$ | Word | Freq | r | $P_r(\%)$ | $r.P_r$ |
|------|-------|---|-----------|---------|------|------|---|-----------|---------|
| the | 2,420,778 | 1 | 6.49 | 0.065 | has | 136,007 | 26 | 0.37 | 0.095 |
| of | 1,045,733 | 2 | 2.80 | 0.056 | are | 130,322 | 27 | 0.35 | 0.094 |
| to | 968,882 | 3 | 2.60 | 0.078 | not | 127,493 | 28 | 0.34 | 0.096 |
| a | 892,429 | 4 | 2.39 | 0.096 | who | 116,364 | 29 | 0.31 | 0.090 |
| and | 865,644 | 5 | 2.32 | 0.120 | they | 111,024 | 30 | 0.30 | 0.089 |
| in | 847,825 | 6 | 2.27 | 0.140 | its | 111,021 | 31 | 0.30 | 0.092 |
| said | 504,593 | 7 | 1.35 | 0.095 | had | 103,943 | 32 | 0.28 | 0.089 |
| for | 363,865 | 8 | 0.98 | 0.078 | will | 102,949 | 33 | 0.28 | 0.091 |
| that | 347,072 | 9 | 0.93 | 0.084 | would | 99,503 | 34 | 0.27 | 0.091 |
| was | 293,027 | 10 | 0.79 | 0.079 | about | 92,983 | 35 | 0.25 | 0.087 |
| on | 291,947 | 11 | 0.78 | 0.086 | i | 92,005 | 36 | 0.25 | 0.089 |
| he | 250,919 | 12 | 0.67 | 0.081 | been | 88,786 | 37 | 0.24 | 0.088 |
| is | 245,843 | 13 | 0.65 | 0.086 | this | 87,286 | 38 | 0.23 | 0.089 |
| with | 223,846 | 14 | 0.60 | 0.084 | their | 84,638 | 39 | 0.23 | 0.089 |
| at | 210,064 | 15 | 0.56 | 0.085 | new | 83,449 | 40 | 0.22 | 0.090 |
| by | 209,586 | 16 | 0.56 | 0.090 | or | 81,796 | 41 | 0.22 | 0.090 |
| it | 195,621 | 17 | 0.52 | 0.089 | which | 80,385 | 42 | 0.22 | 0.091 |
| from | 189,451 | 18 | 0.51 | 0.091 | we | 80,245 | 43 | 0.22 | 0.093 |
| as | 181,714 | 19 | 0.49 | 0.093 | more | 76,388 | 44 | 0.21 | 0.090 |
| be | 157,300 | 20 | 0.42 | 0.084 | after | 75,165 | 45 | 0.20 | 0.091 |
| were | 153,913 | 21 | 0.41 | 0.087 | us | 72,045 | 46 | 0.19 | 0.089 |
| an | 152,576 | 22 | 0.41 | 0.090 | percent | 71,956 | 47 | 0.19 | 0.091 |
| have | 149,749 | 23 | 0.40 | 0.092 | up | 71,082 | 48 | 0.19 | 0.092 |
| his | 142,285 | 24 | 0.38 | 0.092 | one | 70,266 | 49 | 0.19 | 0.092 |
| but | 140,880 | 25 | 0.38 | 0.094 | people | 68,988 | 50 | 0.19 | 0.093 |

# Zipf's law is a type of power law

- degree distribution of
  - the web,
  - online social networks,
  - citation networks
  - software networks
- wealth distribution
- there are variants of the power law, such as Mandelbrot law.