

# Chapter NLP:II

## II. Corpus Linguistics

- Empirical Research
- Text Corpora
- Corpus Properties
- Data Acquisition
- Data Annotation

# Empirical Research

1. Quantitative research based on numbers and statistics.
2. Studies phenomena and research questions by analyzing data.
3. Derives knowledge from experience rather than from theory or belief.

# Empirical Research

1. Quantitative research based on numbers and statistics.
2. Studies phenomena and research questions by analyzing data.
3. Derives knowledge from experience rather than from theory or belief.

Quantitative versus qualitative research:

- Quantitative.** Characterized by objective measurements.
- Qualitative.** Emphasizes the understanding of human experience.

# Empirical Research

1. Quantitative research based on numbers and statistics.
2. Studies phenomena and research questions by analyzing data.
3. Derives knowledge from experience rather than from theory or belief.

Quantitative versus qualitative research:

- Quantitative.** Characterized by objective measurements.
- Qualitative.** Emphasizes the understanding of human experience.

Descriptive versus inferential statistics:

- Descriptive.** Procedures for summarizing and comprehending a sample or distribution of values. Used to describe phenomena.  
1 2 2 2 → mean  $M = 1.75$
- Inferential.** Procedures that help draw conclusions based on values. Used to generalize inferences beyond a given sample.  
The average number is **significantly greater than 1**.

# Empirical Research

## Research Questions

What is a good research question? [Bartos 1992]

- Asks about the relationship between two or more variables.
- Is testable (i.e., it is possible to collect data to answer the question).
- Is stated clearly and in the form of a question.
- Does not pose an ethical or moral problem for implementation.
- Is specific and restricted in scope.
- Identifies exactly what is to be solved.

# Empirical Research

## Research Questions

What is a good research question? [Bartos 1992]

- Asks about the relationship between two or more variables.
- Is testable (i.e., it is possible to collect data to answer the question).
- Is stated clearly and in the form of a question.
- Does not pose an ethical or moral problem for implementation.
- Is specific and restricted in scope.
- Identifies exactly what is to be solved.

Example of a **poorly formulated** question:

*“What is the effectiveness of parent education when given problem children?”*

# Empirical Research

## Research Questions

What is a good research question? [Bartos 1992]

- Asks about the relationship between two or more variables.
- Is testable (i.e., it is possible to collect data to answer the question).
- Is stated clearly and in the form of a question.
- Does not pose an ethical or moral problem for implementation.
- Is specific and restricted in scope.
- Identifies exactly what is to be solved.

Example of a **poorly formulated** question:

*“What is the effectiveness of parent education when given problem children?”*

Example of a **well-formulated** question:

*“What is the effect of the STEP parenting program on the ability of parents to use natural, logical consequences (as opposed to punishment) with their child who has been diagnosed with bipolar disorder?”*

# **Empirical Research**

## Empirical Research in NLP

- **Corpus linguistics.**

NLP is studied in a corpus-linguistics manner; i.e., approaches are developed and evaluated on collections of text.

- **Evaluation measures.**

An evaluation of the quality of an approach is important, especially of its effectiveness.

- **Experiments.**

The quality of an approach is empirically evaluated on test corpora and compared to alternative approaches.

- **Hypothesis testing.**

Methods which verify whether results of an experiment are meaningful/valid by estimating the odds that the results happened by chance.

# **Empirical Research**

## Empirical Research in NLP

- **Corpus linguistics.**  
NLP is studied in a corpus-linguistics manner; i.e., approaches are developed and evaluated on collections of text.
- **Evaluation measures.**  
An evaluation of the quality of an approach is important, especially of its effectiveness.
- **Experiments.**  
The quality of an approach is empirically evaluated on test corpora and compared to alternative approaches.
- **Hypothesis testing.**  
Methods which verify whether results of an experiment are meaningful/valid by estimating the odds that the results happened by chance.

# **Text Corpora**

## Corpus Linguistics

- The study of language as expressed in principled collections of natural language texts, called text corpora.
- Aims to derive knowledge and rules from real-world text.
- Covers both manual and automatic analysis of text.

# Text Corpora

## Corpus Linguistics

- The study of language as expressed in principled collections of natural language texts, called text corpora.
- Aims to derive knowledge and rules from real-world text.
- Covers both manual and automatic analysis of text.

Three main techniques:

1. **Analysis.** Developing and evaluating methods based on a corpus.
  2. **Annotation.** Coding data with categories to facilitate data-driven research.
  3. **Abstraction.** Mapping of annotated texts to a theory-based model.
- Need for text corpora: Without a corpus, it's hard to develop a strong approach—and impossible to reliably evaluate it.

*“It’s often not the one who has the best algorithm that wins.  
It’s who has the most data.”*

# Text Corpora

## Definition 1 (Text Corpus [Butler 2004])

A text corpus is (an electronically stored) collection of data designed with according to specific corpus design criteria to be maximally representative of (a particular variety of) language or other semiotic systems.

The basic unit for representing text is typically a word (captures meaning).

Examples:

- 200,000 product reviews for sentiment analysis
- 1,000 news articles for part-of-speech tagging



Corpora in NLP:

- NLP approaches are developed and evaluated on text corpora.
- Usually, the corpora contain annotations of the output information type to be inferred.

# Text Corpora

## On Representativeness

- “*extent to which a sample includes the full range of variability in a population*”  
[Biber 1993]

Here: Sample is our corpus, population is all of the language variety.

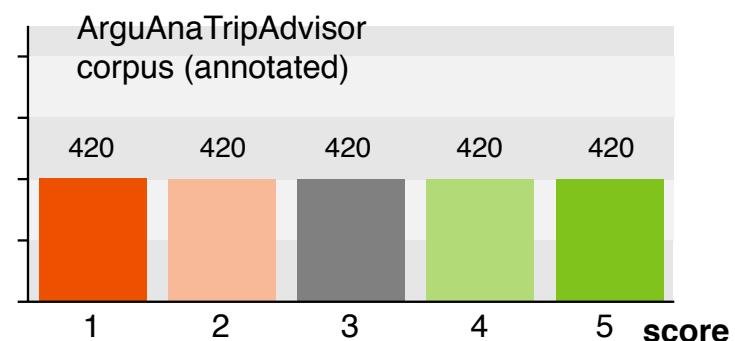
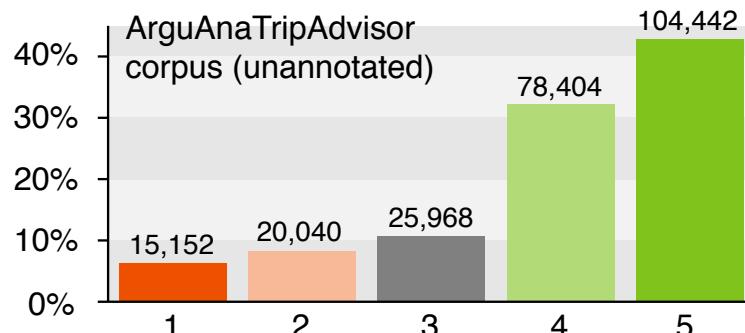
- “*A corpus is thought to be representative of the language variety it is supposed to represent if the findings based on its contents can be generalized to the said language variety.*” [Leech 1991]

Question: If we find certain features in the corpus, are we likely to find the same features in further data of that type?

- But—what is representative to the users of language?  
*“According to claims, the most likely document that an ordinary English citizen will cast his or her eyes over is The Sun newspaper”* [Sinclair 2005]
- Keyword: reception versus production
- Corpus representativeness is important for generalization, since the corpus governs what can be learned about a given domain.

# Text Corpora

## Representative Data versus Balanced Data



- A corpus is representative for some output information type  $C$ , if it includes the full range of variability of texts with respect to  $C$ .
- The distribution of texts over the values of  $C$  should be representative for the real distribution.
- Balance with respect to a feature means that no value/level of the feature dominates; equally distributed with respect to a feature (e.g. genre, category of linguistic phenomena).
- A balanced distribution, where all values are evenly represented, may be favorable (particularly for machine learning).

# Text Corpora

## Text as Data

**Bits:** A sequence of bits that symbolize text when decoded into glyphs [cf [WT:II-166 ff.](#)]

**String:** concatenation of glyphs (alphabet elements)

- “Hello world！”, “”, “00010111100010101”, “To be or not to be...”
- essential, elementary data type in computer linguistics
- common operations: e.g.
  - concatenation: “Hello” + “World!” + “!” → “Hello World!”
  - splitting: split(“Hello World！”, “ ”) → {“Hello”, “World!”}
  - case conversion: uppercase(“Hello”) → “HELLO”
  - substring: substr(“Hello”, start = 0, length = 4) → “Hell”

**Document:** compound data type

- (collection of) strings (e.g. title, body) [+ Metadata]

**Corpus:** collection of documents

# Text Corpora

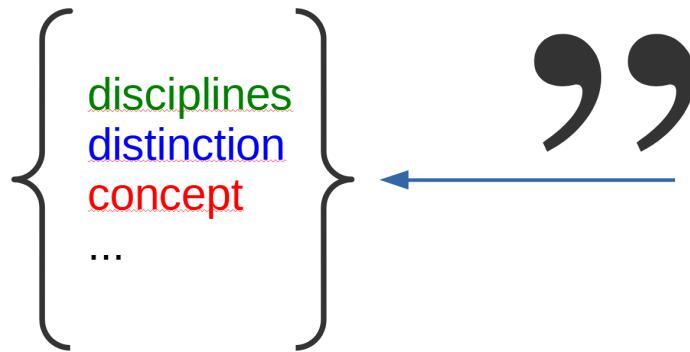
## Text as Data (continued)

### Type: (cp. class)

- (abstract) string representing a meaningful concept, e.g. words

### Token: (cp. object)

- (concrete) string as instance of a meaningful concept



In disciplines such as knowledge representation and philosophy, the type–token distinction is a distinction that separates a concept from the objects which are particular instances of the concept.“

(Wikipedia → Type–token distinction)

### Vocabulary:

- complete set of all types occurring in a [document | collection]

# Text Corpora

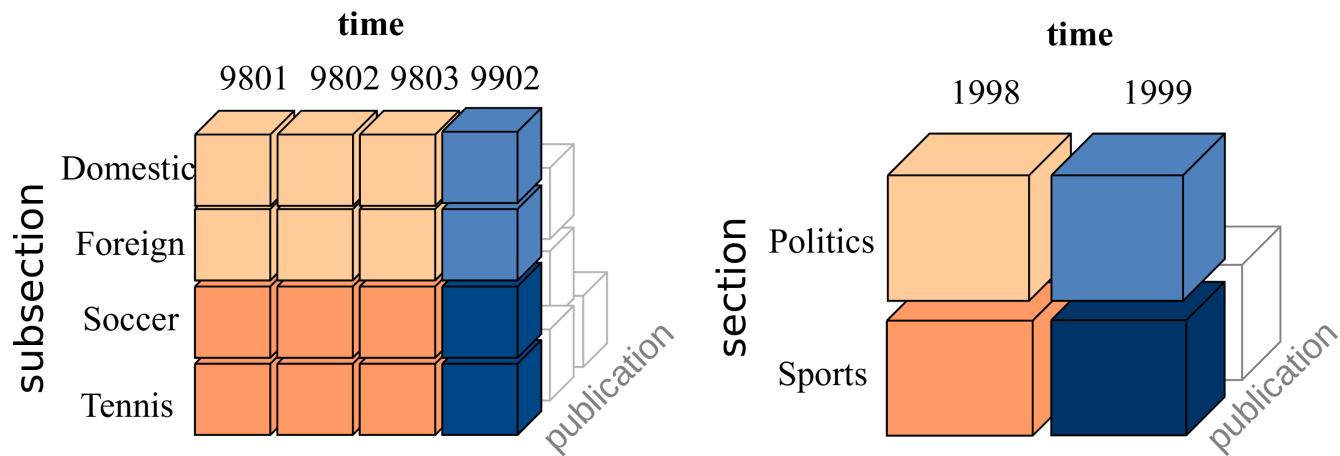
## Metadata

Metadata = text external context / covariate

Metadata = data facet

- Subselections of sources
- Aggregation / differentiation of results

context → contrast → meaning



# Text Corpora

## Research in Language Use

Concordance: (alphabetical) list of principal words (or phrases) used in a book (nowadays: corpus) listing every instance of each with immediate context

The screenshot shows the Sketch Engine Concordance interface. The search query "CQL \"in\" \"the\"? []? context" has returned 706,992 results. The results are displayed in a table with columns: Details, Left context, KWIC (16), and Right context (17). The KWIC column highlights the search term in red. The interface includes a sidebar with various icons and a toolbar with numbered buttons (2-20) and other functions like KWIC dropdown, search, and download.

Details	Left context	KWIC (16)	Right context (17)
391 ① earlychildhoodmagazine... ice violence against children	ice violence against children	in humanitarian contexts	, thereby improving the physic
392 ① nsta.org	isks and activities that occur	in the social contexts	of day-to-day living, whether o
393 ① ancientdragon.org	universal truth can only exist	in the context	of some particular situation. <
394 ① edtalks.org	<s> He discusses open-ness	in the social context	, the technical area, and educat
395 ① theoldgeek.nz	ord immoral has no meaning	in this context	. </s><s> We are stuck saying
396 ① dangcongsan.vn	in the EU market, particularly	in the (19) text	of the strengthening euro. </s
397 ① fifthestate.org	writer Paul Goodman insisted	in the context	of 1960s movements, there m
398 ① bsa.govt.nz	ster therefore concluded that	in the context	of a news item reporting on a
399 ① wisc.edu	he consequences of tracking	in contexts	beyond the US and the UK, wh
400 ① dukeandduchessofcam... have to picture wildlife crime		in the context	of the overall damage that's b

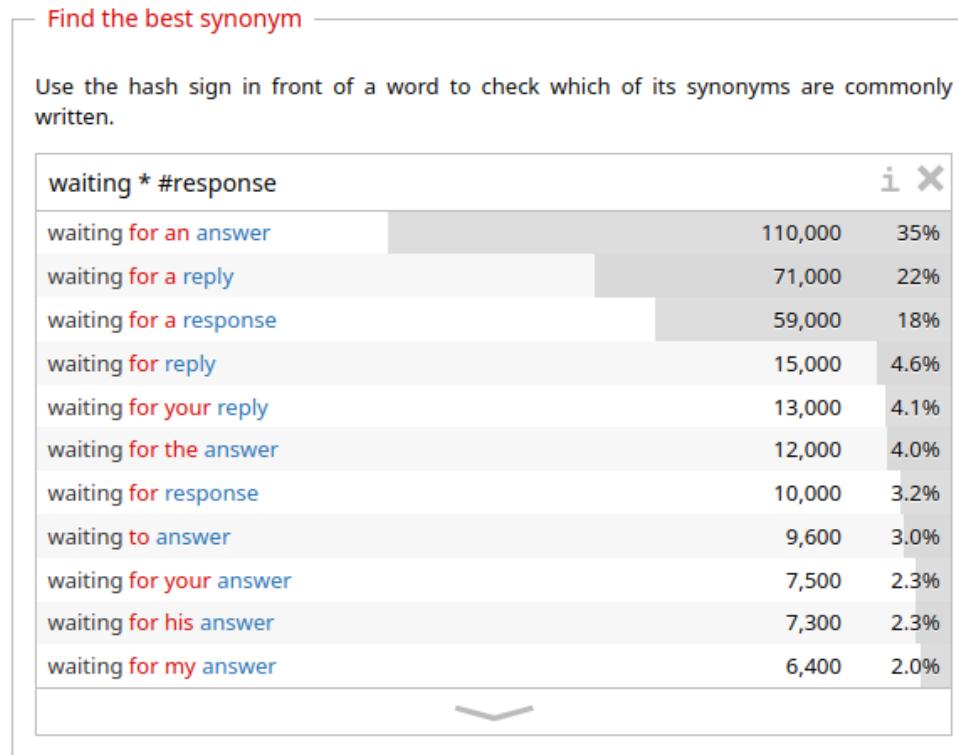
Rows per page: 10 391–400 of 706,992 | < < 40 > >|

[[www.sketchengine.eu](http://www.sketchengine.eu)]

# Text Corpora

## Research in Language Use (continued)

Compare usages of a word, analyse keywords, analyse frequencies, find phrases, idioms, etc.



[netspeak.org]

# Corpus Properties

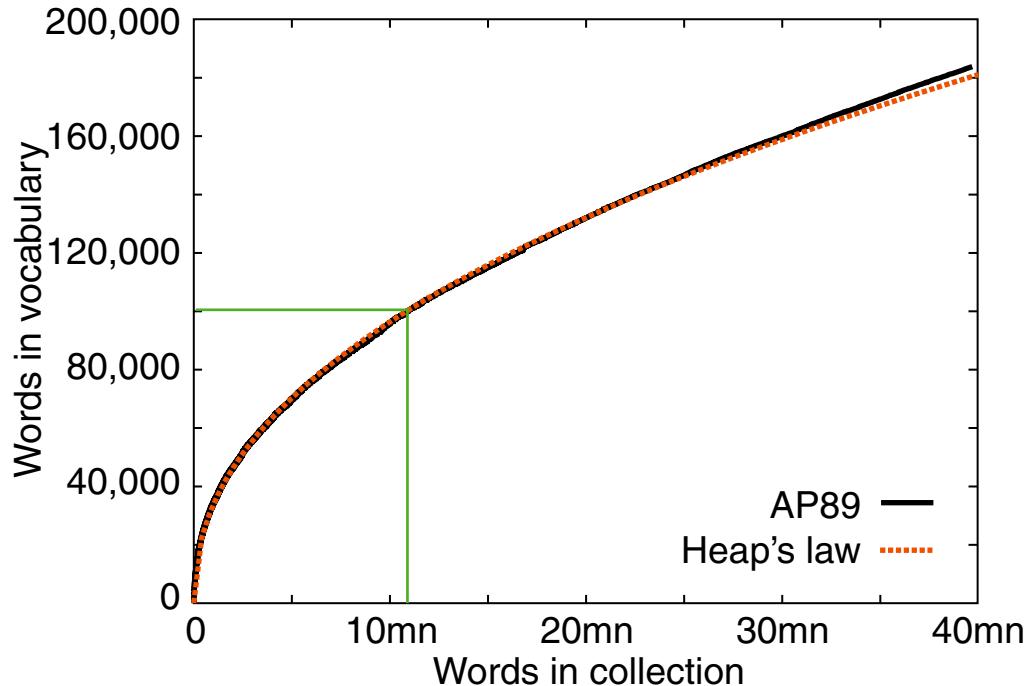
## Vocabulary Growth: Heaps' Law

The vocabulary  $V$  of a collection of documents grows with the collection.

Vocabulary growth can be modeled with Heaps' Law:

$$|V| = k \cdot n^\beta,$$

where  $n$  is the number of non-unique words, and  $k$  and  $\beta$  are collection parameters.



- Corpus: AP89
- $k = 62.95$ ,  $\beta = 0.455$
- At 10,879,522 words:  
100,151 predicted,  
100,024 actual.
- At < 1,000 words:  
poor predictions

# Corpus Properties

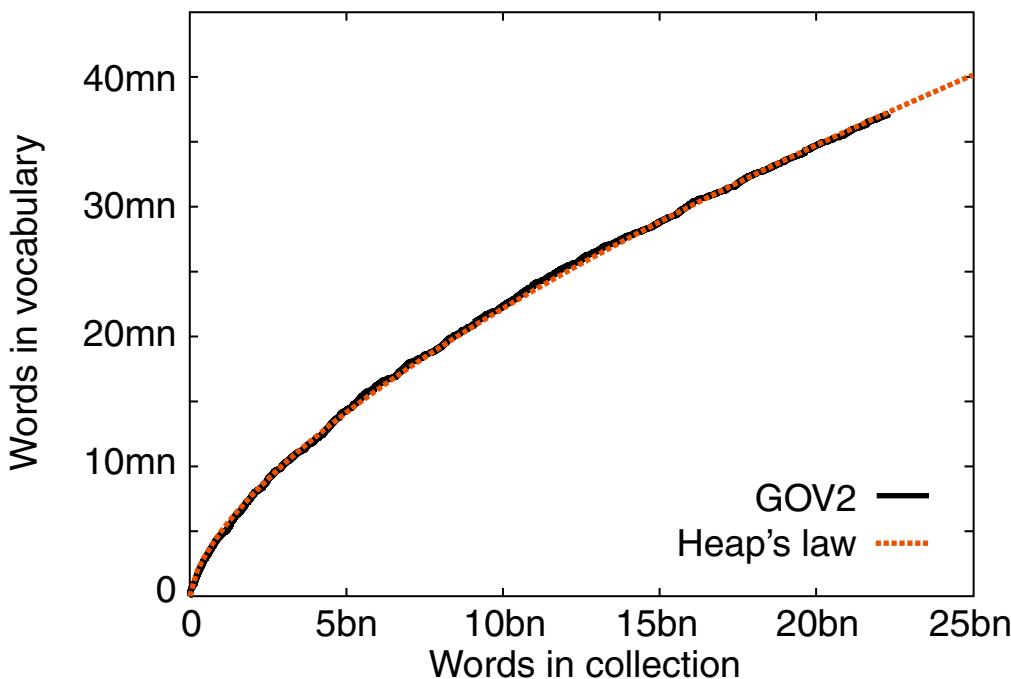
## Vocabulary Growth: Heaps' Law

The vocabulary  $V$  of a collection of documents grows with the collection.

Vocabulary growth can be modeled with Heaps' Law:

$$|V| = k \cdot n^\beta,$$

where  $n$  is the number of non-unique words, and  $k$  and  $\beta$  are collection parameters.



- Corpus: GOV2
- $k = 7.34, \beta = 0.648$
- Vocabulary continuously grows in large collections
- New words include spelling errors, invented words, code, other languages, email addresses, etc.

# Corpus Properties

## Term Frequency: Zipf's Law

- The distribution of word frequencies is very *skewed*: Few words occur very frequently, many words hardly ever.
- For example, the two most common English words (*the, of*) make up about 10% of all word occurrences in text documents. In large text samples, about 50% of the unique words occur only once.



George Kingsley Zipf, an American linguist, was among the first to study the underlying statistical relationship between the frequency of a word and its rank in terms of its frequency, formulating what is known today as Zipf's law.

For natural language, the "Principle of Least Effort" applies.

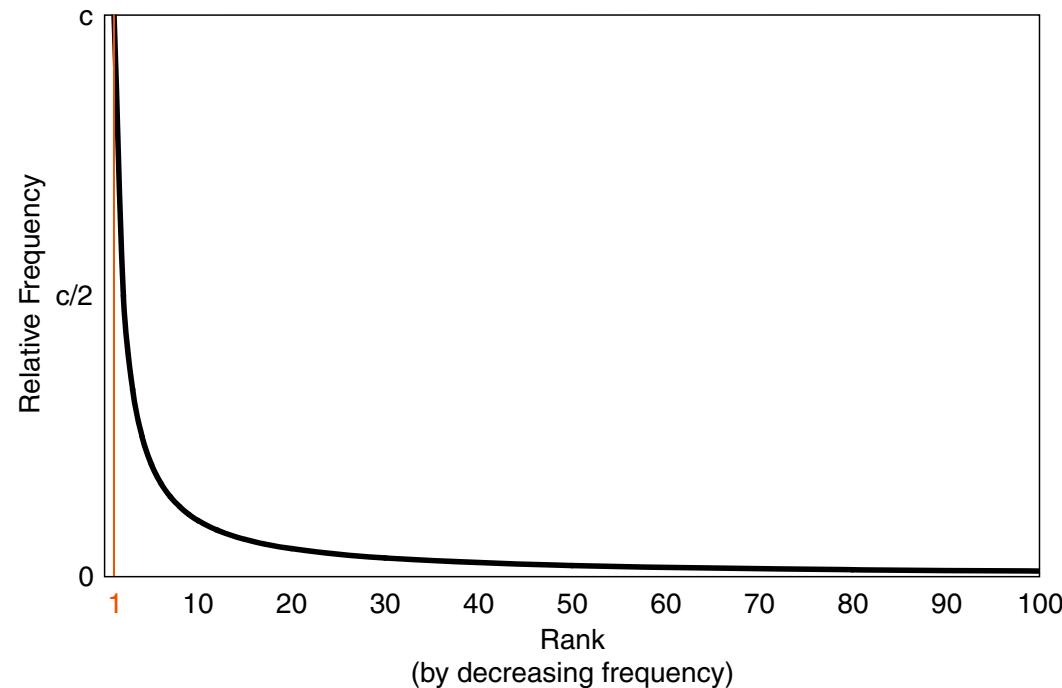
# Corpus Properties

## Term Frequency: Zipf's Law (continued)

The relative frequency  $P(w)$  of a word  $w$  in a sufficiently large text (collection) inversely correlates with its frequency **rank**  $r(w)$  in a power law:

$$P(w) = \frac{c}{(r(w))^a} \quad \Leftrightarrow \quad P(w) \cdot r(w)^a = c,$$

where  $c$  is a constant and the exponent  $a$  is language-dependent; often  $a \approx 1$ .



# Corpus Properties

## Term Frequency: Zipf's Law (continued)

Example: Top 50 most frequent words from AP89. Have a guess at *c*?

$r$	$w$	frequency	$P \cdot 100$	$P \cdot r$
1	the	2,420,778	6.09	0.061
2	of	1,045,733	2.63	0.053
3	to	968,882	2.44	0.073
4	a	892,429	2.25	0.090
5	and	865,644	2.18	0.109
6	in	847,825	2.13	0.128
7	said	504,593	1.27	0.089
8	for	363,865	0.92	0.073
9	that	347,072	0.87	0.079
10	was	293,027	0.74	0.074
11	on	291,947	0.73	0.081
12	he	250,919	0.63	0.076
13	is	245,843	0.62	0.080
14	with	223,846	0.56	0.079
15	at	210,064	0.53	0.079
16	by	209,586	0.53	0.084
17	it	195,621	0.49	0.084
18	from	189,451	0.48	0.086
19	as	181,714	0.46	0.087
20	be	157,300	0.40	0.079
21	were	153,913	0.39	0.081
22	an	152,576	0.38	0.084
23	have	149,749	0.38	0.087
24	his	142,285	0.36	0.086
25	but	140,880	0.35	0.089

$r$	$w$	frequency	$P \cdot 100$	$P \cdot r$
26	has	136,007	0.34	0.089
27	are	130,322	0.33	0.089
28	not	127,493	0.32	0.090
29	who	116,364	0.29	0.085
30	they	111,024	0.28	0.084
31	its	111,021	0.28	0.087
32	had	103,943	0.26	0.084
33	will	102,949	0.26	0.085
34	would	99,503	0.25	0.085
35	about	92,983	0.23	0.082
36	i	92,005	0.23	0.083
37	been	88,786	0.22	0.083
38	this	87,286	0.22	0.083
39	their	84,638	0.21	0.083
40	new	83,449	0.21	0.084
41	or	81,796	0.21	0.084
42	which	80,385	0.20	0.085
43	we	80,245	0.20	0.087
44	more	76,388	0.19	0.085
45	after	75,165	0.19	0.085
46	us	72,045	0.18	0.083
47	percent	71,956	0.18	0.085
48	up	71,082	0.18	0.086
49	one	70,266	0.18	0.087
50	people	68,988	0.17	0.087

# Corpus Properties

## Term Frequency: Zipf's Law (continued)

Example: Top 50 most frequent words from AP89. For English:  $c \approx 0.1$ .

$r$	$w$	frequency	$P \cdot 100$	$P \cdot r$
1	the	2,420,778	6.09	0.061
2	of	1,045,733	2.63	0.053
3	to	968,882	2.44	0.073
4	a	892,429	2.25	0.090
5	and	865,644	2.18	0.109
6	in	847,825	2.13	0.128
7	said	504,593	1.27	0.089
8	for	363,865	0.92	0.073
9	that	347,072	0.87	0.079
10	was	293,027	0.74	0.074
11	on	291,947	0.73	0.081
12	he	250,919	0.63	0.076
13	is	245,843	0.62	0.080
14	with	223,846	0.56	0.079
15	at	210,064	0.53	0.079
16	by	209,586	0.53	0.084
17	it	195,621	0.49	0.084
18	from	189,451	0.48	0.086
19	as	181,714	0.46	0.087
20	be	157,300	0.40	0.079
21	were	153,913	0.39	0.081
22	an	152,576	0.38	0.084
23	have	149,749	0.38	0.087
24	his	142,285	0.36	0.086
25	but	140,880	0.35	0.089

$r$	$w$	frequency	$P \cdot 100$	$P \cdot r$
26	has	136,007	0.34	0.089
27	are	130,322	0.33	0.089
28	not	127,493	0.32	0.090
29	who	116,364	0.29	0.085
30	they	111,024	0.28	0.084
31	its	111,021	0.28	0.087
32	had	103,943	0.26	0.084
33	will	102,949	0.26	0.085
34	would	99,503	0.25	0.085
35	about	92,983	0.23	0.082
36	i	92,005	0.23	0.083
37	been	88,786	0.22	0.083
38	this	87,286	0.22	0.083
39	their	84,638	0.21	0.083
40	new	83,449	0.21	0.084
41	or	81,796	0.21	0.084
42	which	80,385	0.20	0.085
43	we	80,245	0.20	0.087
44	more	76,388	0.19	0.085
45	after	75,165	0.19	0.085
46	us	72,045	0.18	0.083
47	percent	71,956	0.18	0.085
48	up	71,082	0.18	0.086
49	one	70,266	0.18	0.087
50	people	68,988	0.17	0.087

## Remarks:

- Collection statistics for AP89:

---

Total documents	84,678
Total word occurrences	39,749,179
Vocabulary size	198,763
Words occurring > 1000 times	4,169
Words occurring once	70,064

---

# Corpus Properties

## Term Frequency: Zipf's Law (continued)

For relative frequencies,  $c$  can be estimated as follows:

$$1 = \sum_{i=1}^n P(w_i) = \sum_{i=1}^n \frac{c}{r(w_i)} = c \sum_{i=1}^n \frac{1}{r(w_i)} = c \cdot H_t, \quad \rightsquigarrow \quad c = \frac{1}{H_t} \approx \frac{1}{\ln(t)}$$

where  $t$  is the size  $|V|$  of the vocabulary  $V$ , and  $H_n$  is the  $n$ -th harmonic number.

Constant  $c$  is language-dependent; e.g., for German  $c = 1/\ln(7.841.459) \approx 0.063$ . [[Wortschatz Leipzig](#)]

Thus, the expected average number of occurrences of a word  $w$  in a document  $d$  of length  $m$  is

$$m \cdot P(w),$$

since  $P(w)$  can be interpreted as a term occurrence probability.

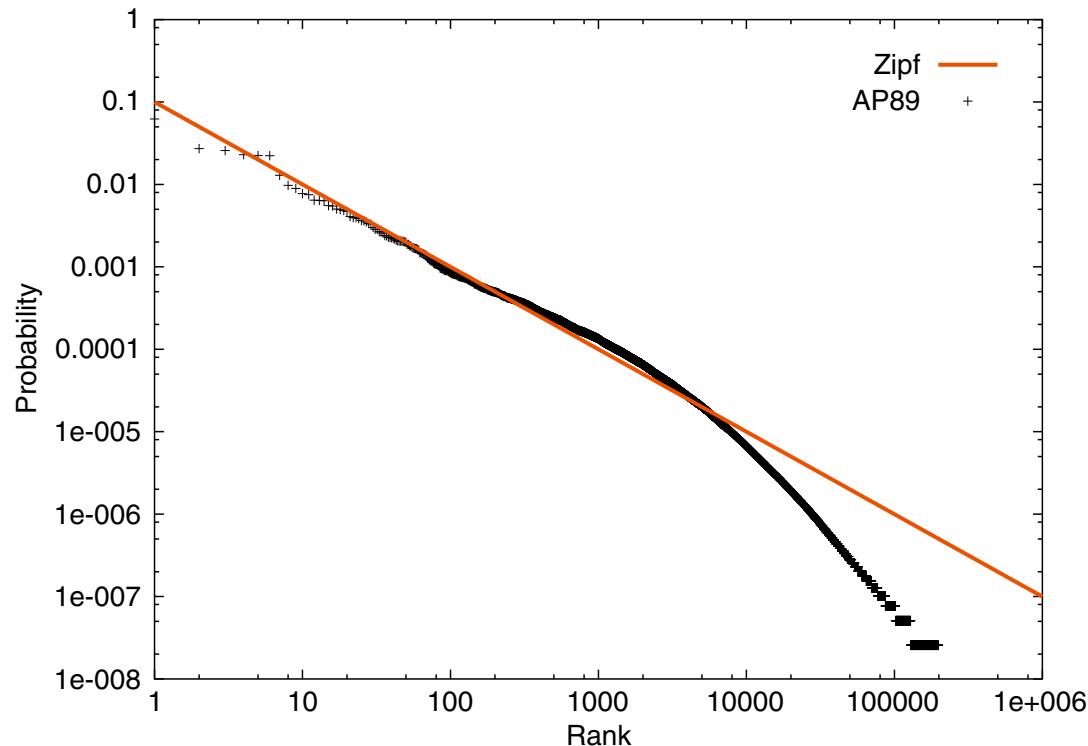
# Corpus Properties

## Term Frequency: Zipf's Law (continued)

By logarithmization a linear form is obtained, yielding a straight line in a plot:

$$\log(P(w)) = \log(c) - a \cdot \log(r(w))$$

Example for AP89:



## Remarks:

- As with all empirical laws, Zipf's law holds only approximately. While mid-range ranks of the frequency distribution fit quite well, this is less so for the lowest ranks and very high ranks (i.e., very infrequent words). The [Zipf-Mandelbrot law](#) is an extension of Zipf's law that provides for a better fit.

$$n \approx \frac{1}{(r(w) + c_1)^{1+c_2}}$$

- Interestingly, this relation cannot only be observed for words and letters in human language texts or music score sheets, but for all kinds of natural symbol sequences (e.g., DNA). It is also true for randomly generated character sequences where one character is assigned the role of a blank space. [\[Li 1992\]](#)
- Independently of Zipf's law, a special case is [Benford's law](#), which governs the frequency distribution of first digits in a number.

# Corpus Properties

## Term Frequency: Zipf's Law (continued)

For the vocabulary,  $t$  (types) is as large as the largest rank of the frequency-sorted list. For words with frequency 1:

$$P(w) = \frac{n_w}{N}, \quad t = r(n_w = 1) = c \times \frac{N}{1} = c \times N \approx e^{1/c}$$

Proportion of word forms that occur only  $n$  time. For  $\mathbf{w}_n$  applies:

$$\mathbf{w}_n = r(n_w) - (r(n_w) + 1) = c \times \frac{N}{n} - c \times \frac{N}{n+1} = \frac{c \times N}{n(n+1)} = \frac{t}{n(n+1)}$$

For  $\mathbf{w}_1$  applies in particular:

$$\mathbf{w}_1 = \frac{t}{2}$$

Half of the vocabulary in a text probably occurs only 1 time.

# Corpus Properties

## Term Frequency: Zipf's Law (continued)

The ratio of words with a given absolute frequency  $n$  can be estimated by

$$\frac{\mathbf{w}_n}{t} = \frac{\frac{t}{n(n+1)}}{t} = \frac{1}{n(n + 1)}$$

Observations:

- Estimations are fairly accurate for small  $x$ .
- Roughly half of all words can be expected to be unique.

Applications:

- Estimation of the number of word forms that occur  $n$  times in the text.
- Estimation of vocabulary size
- Estimation of vocabulary growth as text volume increases
- Analysis of search queries
- Term extraction (for indexing)
- Difference analysis (comparison of documents)

# Sources of Text Data

## Data Sources

### Digitally available texts

- natively digital / born digital
- retro-digitized

### Metadata: “data about data”

- structural metadata
- descriptive metadata

### “Big Data”

- 15,3 Mio .de-Domains (31.12.2012)
- 1.9 Mio articles in F.A.Z. Archive in 1949–2011
- 400 million Twitter tweets per day (2013)

# Sources of Text Data

## Newspapers

archive of political public sphere, societal knowledge or public discourse

### Properties

- representativity (?)
- availability improves

### Difficulties

- licences
- bad OCR

### Example: DIE ZEIT

- <http://www.zeit.de/archiv>
- articles since 1946
- 400.000 articles
- PDF + OCR-ed Text

DIE ZEIT: Jahrgang 1948



Date ← 1948-05-12  
Author(s) ← {"GH", „geh“, „Gerda Heller“}  
Page number ← {1, 1-2}  
Section(s) ← {"Sport", „Leibesübungen“}  
Subsection(s) ← „Handball“  
News agency ← {true|false; „dpa“}

Date  
String[]  
Integer  
String[]  
String  
Boolean

# Sources of Text Data

## Blogs and Forums

Extract of (political) public discourse

## Properties

- expert generated content
- user generated content (comments)

## Properties

- high availability
- lesser license restrictions
- no OCR problems

## Difficulties

- identifying relevant content
- representativity of content?
- Crawling + Web scraping

The screenshot shows the homepage of BlogActiv.eu. At the top, there are language links (English, French, German, etc.) and a search bar. The main header is "BlogActiv.eu" with a small EU flag icon. Below the header is a map of Europe. The left sidebar has a "BROWSE ALL SECTIONS" menu with various categories like Agriculture & Food, Aviation, Climate & Environment, etc. The main content area features a "EDITOR'S CHOICE" section with a thumbnail of a person and the text: "Anti-corruption group report about the dark rooms of the EU". It includes a byline "By Peter Kramer" and a link to "Read & React". Below this is a "LATEST POSTS" section with another article thumbnail and the text: "Why is European energy security not at top of priorities list?". The footer contains a "GURU BLOGGERS" grid of small profile pictures, a call to action "Start your BlogActiv blog today!", and an "ADVERTISEMENT" for "EurActiv jobsite".

```
Date ← 2012-11-12 21:40
Author(s) ← {"E. F."}
Url ← „http://www.blogactiv.eu/blog/31/123“
PolicyField ← „Agriculture“
numberOfComments ← 216
numberOfReadings ← 12002
```

# Sources of Text Data

## Social network

controlled public spheres

### Properties

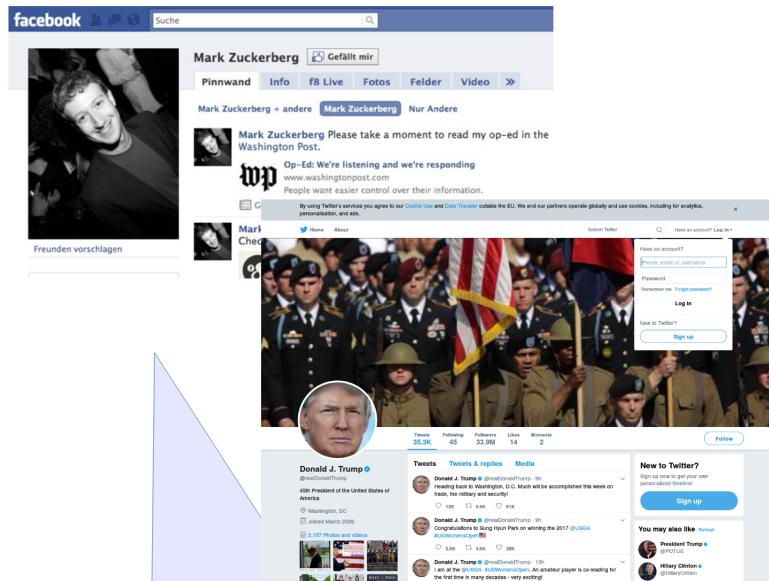
- just in time
- really big data

### Difficulties

- very short text snippets
- typos and special language
- representativity?
- Data acquisition may be complicated

### Data acquisition via APIs

- Twitter sample API (1%)
- Twitter keyword location search
- Facebook API: retrieve user networks and (public) posts, comments, replies from users



```
Type ← {post, comment, reply, tweet}  
Datetime ← 2014-05-12 12:47  
Author ← User_462945  
Reactions ← {like:67, angry:472, sad:12}
```

# Sources of Text Data

## Other Sources

- Emails
- Parliamentary protocols
- Political documents
  - political speeches
  - party manifestos
  - press releases
- Open questions from (online) surveys
- Literature: distant reading of (world) literature
- Scientific publications: lots of science of science studies