

# Statistical Natural Language Processing

P. Rychlý

NLP Centre, FI MU, Brno

September 23, 2022

- 1 Word lists
- 2 Collocations
- 3 Language Modeling
- 4 N-grams
- 5 Evaluation of Language Models

# Statistical Natural Language Processing

- statistics provides a summary (of a text)
- highlights important or interesting facts
- can be used to model data
- foundation of estimating probabilities
- fundamental statistics: size (+ domain, range)

# Statistical Natural Language Processing

- statistics provides a summary (of a text)
- highlights important or interesting facts
- can be used to model data
- foundation of estimating probabilities
- fundamental statistics: size (+ domain, range)

	lines	words	bytes
Book 1	3,715	37,703	223,415
Book 2	1,601	16,859	91,031

# Word list

- list of all words from a text
- list of most frequent words
- words, lemmas, senses, tags, domains, years ...

Book 1	Book 2
the, and, of, to, you, his, in, said, that, I, will, him, your, he, a, my, was, with, s, for, me, He, is, , , it, them, be, The, all, , have, from, , on, her, , , , are, their, were, they, which, , t, up, , had, there	the, I, to, a, of, is, that, , you, he, and, said, was, , in, it, not, me, my, have, And, are, one, for, But, his, be, The, It, at, all, with, on, will, as, very, had, this, him, He, from, they, , so, them, no, You, do, would, like

# Word list

- list of all words from a text
- list of most frequent words
- words, lemmas, senses, tags, domains, years ...

Book 1	Book 2
the, and, of, to, you, his, in, said, that, I, will, him, your, he, a, my, was, with, s, for, me, He, is, <b>father</b> , , it, them, be, The, all, <b>land</b> , have, from, , on, her, , <b>son</b> , , are, their, were, they, which, <b>sons</b> , t, up, , had, there	the, I, to, a, of, is, that, <b>lit- tle</b> , you, he, and, said, was, , in, it, not, me, my, have, And, are, one, for, But, his, be, The, It, at, all, with, on, will, as, very, had, this, him, He, from, they, <b>planet</b> , so, them, no, You, do, would, like

## Word list

- list of all words from a text
- list of most frequent words
- words, lemmas, senses, tags, domains, years ...

Book 1	Book 2
the, and, of, to, you, his, in, said, that, I, will, him, your, he, a, my, was, with, s, for, me, He, is, <b>father</b> , <b>God</b> , it, them, be, The, all, <b>land</b> , have, from, <b>Jacob</b> , on, her, <b>Yahweh</b> , <b>son</b> , <b>Joseph</b> , are, their, were, they, which, <b>sons</b> , t, up, <b>Abraham</b> , had, there	the, I, to, a, of, is, that, <b>lit- tle</b> , you, he, and, said, was, <b>prince</b> , in, it, not, me, my, have, And, are, one, for, But, his, be, The, It, at, all, with, on, will, as, very, had, this, him, He, from, they, <b>planet</b> , so, them, no, You, do, would, like

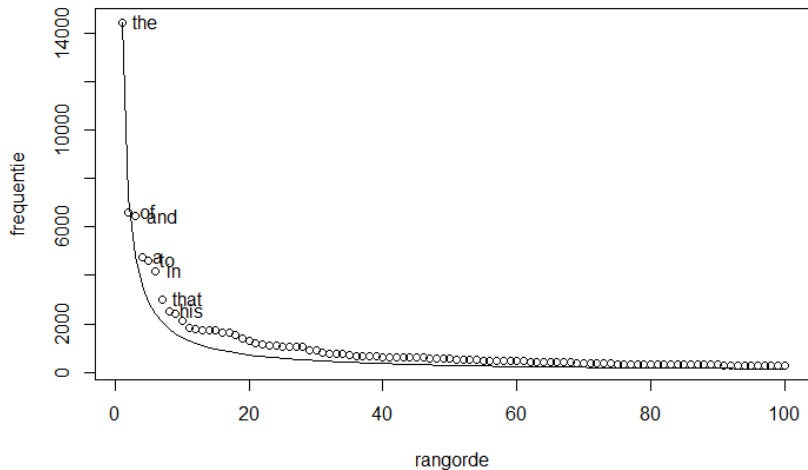
# Frequency

- number of occurrences (raw frequency)
- relative frequency (hits per million)
- document frequency (number of documents with a hit)
- reduced frequency (ARF, ALDf)  
 $1 < \textit{reduced} < \textit{raw}$
- normalization for comparison
- hapax legomena (= 1 hit)

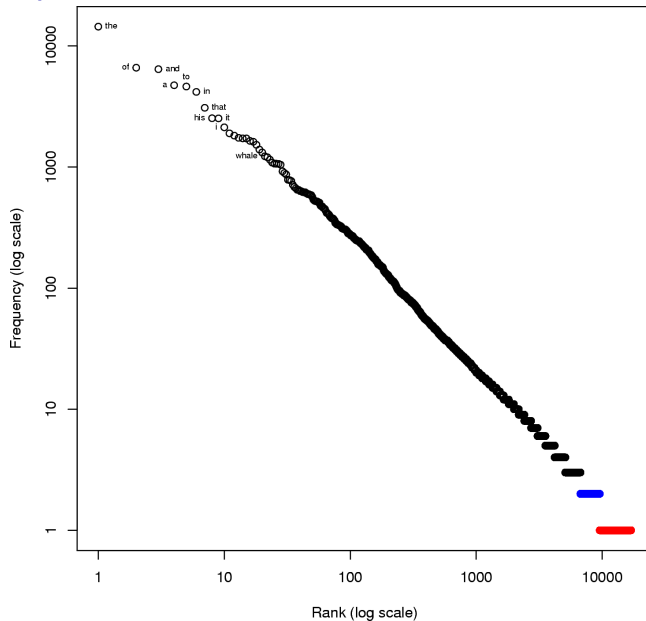


# Zipf's Law

- rank-frequency plot
- $\text{rank} \times \text{frequency} = \text{constant}$



# Zipf's Law



# Keywords

- select only *important* words from a word list
- compare to reference text (norm)
- simple math score:

$$\text{score} = \frac{\text{freq}_{\text{focus}} + N}{\text{freq}_{\text{reference}} + N}$$

<b>Genesis</b>	<b>Little Prince</b>
son God father Jacob Yahweh Joseph Abraham wife behold daughter	prince planet flower little fox never too drawing reply star

# Collocations

- meaning of words is defined by the context
- collocations a *salient* words in the context
- usually not the most frequent
- filtering by part of speech, grammatical relation
- compare to *reference* = context for other words
- many statistics (usually single use only) based on frequencies
- MI-score, t-score,  $\chi^2$ , ...
- logDice – scalable

$$\text{logDice} = 14 + \log \frac{f_{AB}}{f_A + f_B}$$

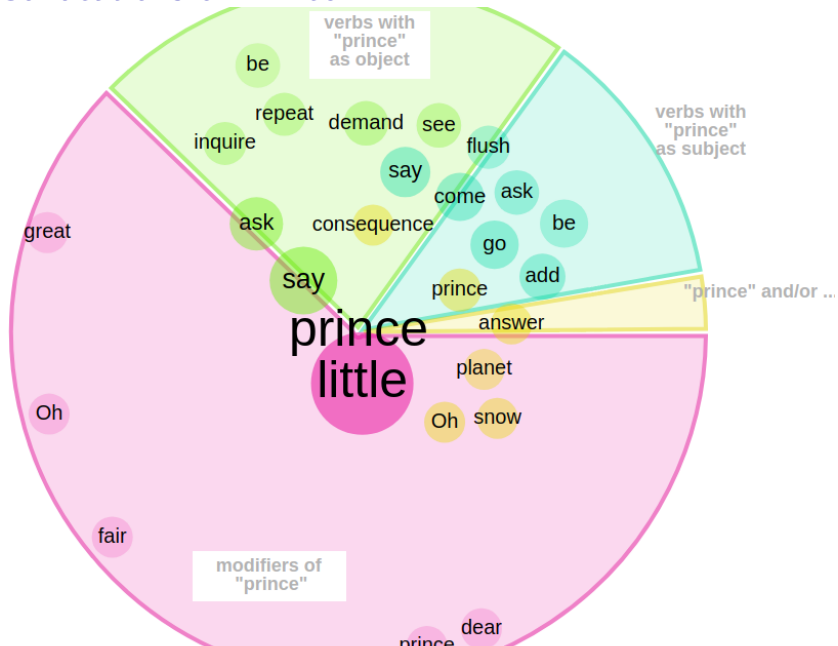
# Collocations of Prince

modifiers of "prince"	
<b>little</b> ...	the little prince
<b>fair</b> ...	fair , little prince
<b>Oh</b> ...	Oh , little prince
<b>dear</b> ...	dear little prince
<b>prince</b> ...	prince , dear little prince
<b>great</b> ...	great prince

verbs with "prince" as object	
<b>say</b> ...	said the little prince
<b>ask</b> ...	asked the little prince
<b>demand</b> ...	demanded the little prince
<b>see</b> ...	when he saw the little prince coming
<b>inquire</b> ...	inquired the little prince
<b>repeat</b> ...	repeated the little prince , who

verbs with "prince" as subject	
<b>say</b> ...	the little prince said to himself
<b>come</b> ...	saw the little prince coming
<b>go</b> ...	And the little prince went away
<b>add</b> ...	the little prince added
<b>ask</b> ...	the little prince asked
<b>flush</b> ...	The little prince flushed

# Collocations of Prince



# Thesaurus

- comparing collocation distributions
- counting same context

son as noun 301x

	Word	Frequency ?
1	brother	161 ...
2	wife	125 ...
3	father	278 ...
4	daughter	108 ...
5	child	80 ...
6	man	187 ...
7	servant	91 ...
8	Esau	78 ...
9	Jacob	184 ...
10	name	85 ...

Abraham as noun 134x

	Word	Frequency ?
1	Isaac	82 ...
2	Jacob	184 ...
3	Joseph	157 ...
4	Noah	41 ...
5	Abram	61 ...
6	Laban	54 ...
7	Esau	78 ...
8	God	234 ...
9	Abimelech	24 ...
10	father	278 ...

## Multi-word units

- meaning of some words is completely different in the context of specific co-occurring word
- *black hole*, is not black and is not a hole
- strong collocations
- uses same statistics with different threshold
- better to compare context distribution instead of only numbers
- terminology – compare to a reference corpus

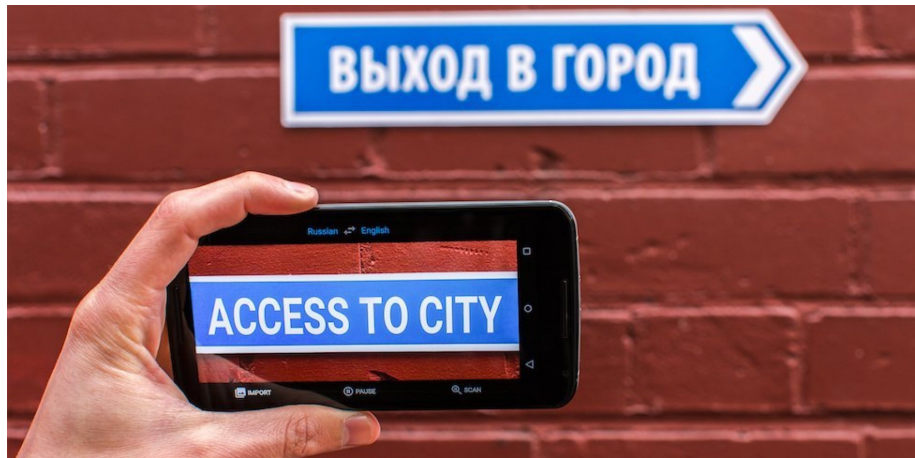


# Language models—what are they good for?

- assigning scores to sequences of words
- predicting words
- generating text

⇒

- statistical machine translation
- automatic speech recognition
- optical character recognition



## Language models – probability of a sentence

- LM is a probability distribution over all possible word sequences.
- What is the probability of utterance of  $s$ ?

### Probability of sentence

$p_{LM}(\text{Catalonia President urges protests})$

$p_{LM}(\text{President Catalonia urges protests})$

$p_{LM}(\text{urges Catalonia protests President})$

...

Ideally, the probability should strongly correlate with fluency and intelligibility of a word sequence.

# N-gram models

- an approximation of long sequences using short n-grams
- a straightforward implementation
- an intuitive approach
- good local fluency

## Randomly generated text

“Jsi nebylo vidět vteřin přestal po schodech se dal do deníku a položili se táhl ji viděl na konci místnosti 101,” řekl důstojník.

## Hungarian

A társaság kötelezettségeiért kapta a középkori temploma az volt, hogy a felhasználók az adottságai, a felhasználó azonosítása az egyesület alapszabályát.

## N-gram models, naïve approach

$$W = w_1, w_2, \dots, w_n$$

$$p(W) = \prod_i p(w_i | w_1 \dots w_{i-1})$$

Markov's assumption

$$p(W) = \prod_i p(w_i | w_{i-2}, w_{i-1})$$

$$p(\textit{this is a sentence}) = p(\textit{this}) \times p(\textit{is} | \textit{this}) \times p(\textit{a} | \textit{this}, \textit{is}) \times p(\textit{sentence} | \textit{is}, \textit{a})$$

$$p(\textit{a} | \textit{this}, \textit{is}) = \frac{|\textit{this is a}|}{|\textit{this is}|}$$

**Sparse data** problem.

# Probabilities, practical issue

- probabilities of words are very small
- multiplying small numbers goes quickly to zero
- limits of floating point numbers:  $10^{-38}$ ,  $10^{-388}$
- using log space:
  - ▶ avoid underflow
  - ▶ adding is faster

$$\log(p_1 \times p_2 \times p_3 \times p_4) = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

## Computing, LM probabilities estimation

Trigram model uses 2 preceding words for probability learning. Using **maximum-likelihood estimation**:

$$p(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\sum_w \text{count}(w_1, w_2, w)}$$

quadrigram: (*lord, of, the, ?*) ( )

<i>w</i>	<b>count</b>	$p(w)$
rings	30,156	0.425
flies	2,977	0.042
well	1,536	0.021
manor	907	0.012
dance	767	0.010
...		

## Large LM – n-gram counts

How many unique n-grams in a corpus?

<b>order</b>	<b>unique</b>	<b>singletons</b>
unigram	86,700	33,447 (38.6%)
bigram	1,948,935	1,132,844 (58.1%)
trigram	8,092,798	6,022,286 (74.4%)
4-gram	15,303,847	13,081,621 (85.5%)
5-gram	19,882,175	18,324,577 (92.2%)

Corpus: Europarl, 30 M tokens.



# Language models smoothing

The problem: an n-gram is missing in the data but is in a *sentence*  $\rightarrow p(\textit{sentence}) = 0$ .

We need to assign non-zero  $p$  for *unseen data*. This must hold:

$$\forall w : p(w) > 0$$

The issue is more pronounced for higher-order models.

Smoothing: an attempt to amend real counts of n-grams to expected counts in any (unseen) data.

Add-one, Add- $\alpha$ , Good-Turing smoothing