

# СЛОВА — ЛЕКСИЧЕСКАЯ СТАТИСТИКА

## КВАНТИТАТИВНЫЙ АНАЛИЗ ТЕКСТА

---

Кирилл Александрович Маслинский

07.02.2022 / 01

Институт русской литературы (Пушкинский Дом) РАН

# Предисловие



СѢВАСНАНАКАКОПОДОБА  
ИТЬУЛѢКОУБИТИ:



ѢВѢСТЬУЛѢКОУИДѢТИ  
ПАЧЕВЪСЕГОЖИТИА  
ДАНЕПРИЛЕЖИТЬИ  
ИЖЕНЪЗЪЛО. ИЖЕБЕ  
СИЖЕЗДѢРЖАННІЕ. ОУДО  
БРЕННІЕ. НОР. ОВОУ. ГЛАСОУ  
ОУДНЛЕННІЕ. ОУДОБРЕННІЕ.  
ЖДЕННІЕ. ИПИТИ. БЕЗГО  
БОРА. СХОУДѢРЖАННІЕ. ДЛѢ  
ПРѢДЪСТАРИЦИ. ДЛѢТАНИ  
И. ПРѢДЪЛОУДРИ. ДЛИ. ПО

WHEN IRISH MONKS BEGAN SEPARATING WORDS IN  
MANUSCRIPTS BY SPACES IN THE SEVENTH CENTURY,  
LITTLE COULD THEY KNOW THAT THEY WERE  
PERFORMING A CENTRAL TASK OF COMPUTATIONAL  
TEXT ANALYSIS.

*ANDREW PIPER, ENUMERATIONS: DATA AND  
LITERARY STUDY (2018) P. 42*

# ЛЕКСИЧЕСКАЯ СТАТИСТИКА

---

## КАК НАЧАТЬ СЧИТАТЬ СЛОВА

---

- сколько разных слов в сообщении в мессенджере?
- на странице научной статьи?
- сколько новых слов на каждой следующей странице,
- сколько новых слов во второй половине книги (по сравнению с первой)?
- сколько новых слов в книге по сравнению с другими книгами?
- когда наконец мы перестанем встречать новые слова?

## КАК НАЧАТЬ СЧИТАТЬ СЛОВА

---

- сколько разных слов в сообщении в мессенджере?
- на странице научной статьи?
- сколько новых слов на каждой следующей странице,
- сколько новых слов во второй половине книги (по сравнению с первой)?
- сколько новых слов в книге по сравнению с другими книгами?
- когда наконец мы перестанем встречать новые слова?

## КАК НАЧАТЬ СЧИТАТЬ СЛОВА

---

- сколько разных слов в сообщении в мессенджере?
- на странице научной статьи?
- сколько новых слов на каждой следующей странице,
- сколько новых слов во второй половине книги (по сравнению с первой)?
- сколько новых слов в книге по сравнению с другими книгами?
- когда наконец мы перестанем встречать новые слова?



## КАК НАЧАТЬ СЧИТАТЬ СЛОВА

---

- сколько разных слов в сообщении в мессенджере?
- на странице научной статьи?
- сколько новых слов на каждой следующей странице,
- сколько новых слов во второй половине книги (по сравнению с первой)?
- сколько новых слов в книге по сравнению с другими книгами?
- когда наконец мы перестанем встречать новые слова?

## КАК НАЧАТЬ СЧИТАТЬ СЛОВА

---

- сколько разных слов в сообщении в мессенджере?
- на странице научной статьи?
- сколько новых слов на каждой следующей странице,
- сколько новых слов во второй половине книги (по сравнению с первой)?
- сколько новых слов в книге по сравнению с другими книгами?
- когда наконец мы перестанем встречать новые слова?

## КАК НАЧАТЬ СЧИТАТЬ СЛОВА

---

- сколько разных слов в сообщении в мессенджере?
- на странице научной статьи?
- сколько новых слов на каждой следующей странице,
- сколько новых слов во второй половине книги (по сравнению с первой)?
- сколько новых слов в книге по сравнению с другими книгами?
- когда наконец мы перестанем встречать новые слова?

## Лексическая статистика

### Закон Ципфа

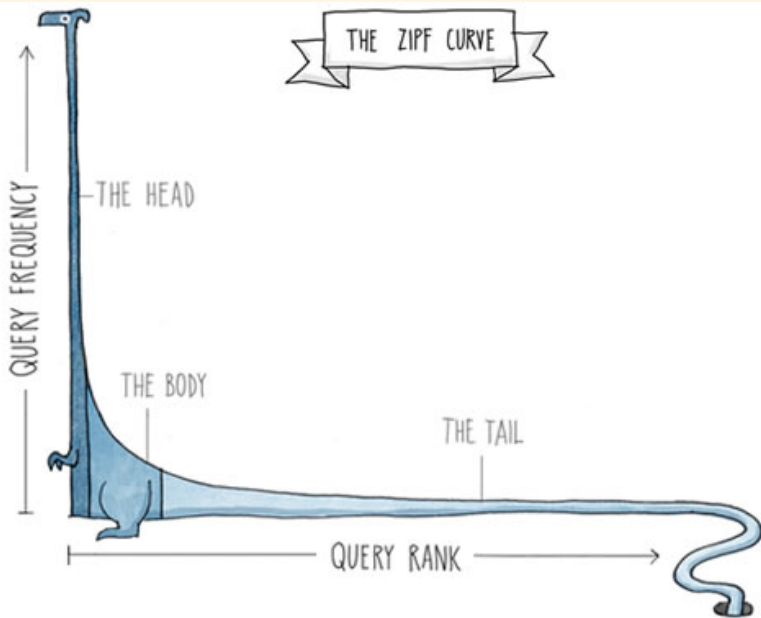
Частотность и состав лексикона

Размер и скорость роста словаря

Размер лексикона/Лексическое разнообразие

Практические следствия закона Ципфа

# THE ZIPF CURVE



Предсказывает частотность слова по его рангу в частотном списке:

$$f(w) = \frac{C}{r(w)^a} \quad (1)$$

$f(w)$  — частотность слова  $w$

$r(w)$  — ранг слова  $w$  в частотном списке

$C$  — константа

$a$  — константа, близкая к 1.

## ПРЕДСКАЗАНИЯ ЗАКОНА ЦИПФА

При  $a = 1$ ,  $C = 60000$  закон Ципфа предсказывает:

$$f(w) = \frac{C}{r(w)}$$

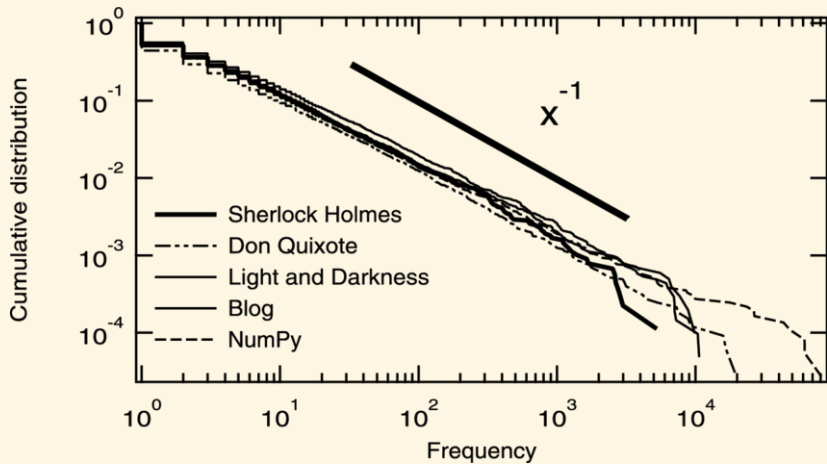
- самое частотное слово встретится  $f(w) = C/1 = 60000$  раз
- второе по частотности слово  $C/2 = 30000$  раз
- третье по частотности слово  $C/3 = 20000$  раз
- сотое  $C/100 = 600$  раз
- сто первое  $C/101 = 594,06$  раз (около 99% частотности сотого)
- и длинный хвост из 80000 слов с частотностью между 1,5 и 0,5.

$$\log f(w) = \log(C) - a \log r(w) \quad (2)$$

Линейная функция:

$$y = kx + b$$





## ЗАКОН ЦИПФА-МАНДЕЛЬБРОТА (1953)

$$f(w) = \frac{C}{(r(w) + b)^a} \quad (3)$$

При  $C = 60000$ ,  $a = 1$ ,  $b = 1$  предсказанная частотность самого частотного слова:

Закон Ципфа  $\frac{C}{1} = \frac{60000}{1} = 60000$

Закон Ципфа-Мандельброта  $\frac{C}{r+b} = \frac{60000}{(1+1)} = 30000$

1. Психолингвистическое (Ципф):
  - экономия усилий говорящего (меньше разных слов);
  - экономия усилий слушающего (больше разных слов).
2. Теоретико-информационное (Мандельброт):
  - минимизация средней стоимости передачи информации в тексте.
3. Процесс, приводящий к подобному распределению:
  - новые слова с константной вероятностью (Simon 1955);
  - «обезьяна и пишущая машинка» (Miller 1957).

Few Giants — Many dwarfs

## Примеры

- частотности слов;
- размеры городов;
- распределение дохода (закон Парето).

## Лексическая статистика

Закон Ципфа

Частотность и состав лексикона

Размер и скорость роста словаря

Размер лексикона/Лексическое разнообразие

Практические следствия закона Ципфа

## ОТКРЫТЫЕ И ЗАКРЫТЫЕ КЛАССЫ СЛОВ

Словарь языка незамкнут — всё время возникают новые слова.

**Function words, closed-class** Вершину частотного списка занимают служебные части речи (**предлоги, союзы, местоимения**). Все единицы перечислимы, пополняется очень медленно. В тексте выполняют прежде всего грамматическую функцию.

**Content words, open-class** Далее в частотном списке преобладают слова открытых классов (пополняемых), прежде всего **существительные**. В тексте выполняют прежде всего референтную функцию.

## ПРИМЕР: ЧАСТОТНОСТЬ РУССКОЙ ЛЕКСИКИ

Единица измерения частотности:

ipm — вхождений на миллион / instances per million

1	36358.94	и	misc	32600	1.04	мертветь	verb
2	27792.36	в	prep	32601	1.04	сволочной	adj
3	20689.51	не	misc	32602	1.04	втыкаться	verb
4	18942.62	он	pron	32603	1.04	нахлебник	noun
5	16588.14	на	prep	32604	1.04	русоволосый	adj
6	15631.11	я	pron	32605	1.04	автопилот	noun
7	12546.08	что	misc	32606	1.04	иссечение	noun
8	11398.44	тот	adjpron	32607	1.04	бульдожий	adj
9	11223.99	быть	verb	32608	1.04	бренность	noun
10	11150.72	с	prep	32609	1.04	нездоровье	noun
11	9808.61	а	misc	32610	1.04	саргасса	noun
12	8601.72	еще	adjpron	32611	1.04	мелкозерный	adj

## Лексическая статистика

Закон Ципфа

Частотность и состав лексикона

Размер и скорость роста словаря

Размер лексикона/Лексическое разнообразие

Практические следствия закона Ципфа



Чем дальше мы читаем текст, тем реже встречаем новые слова.

Оценка Гаральда Баайена (Baayen  $G$ ):

$$G = \frac{V(1)}{N} \quad (4)$$

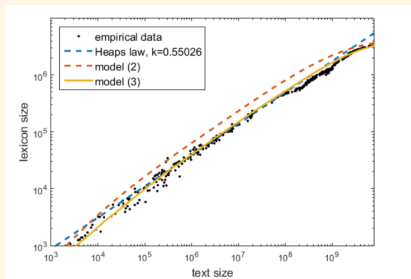
где:

$V(1)$  — количество *hapax legomena* на  $N$  токенов текста

$N$  — количество токенов текста.

Чем дальше мы читаем текст, тем реже встречаем новые слова.

$$V = kN^{\beta} \quad (5)$$



$V$  — размер словаря

$N$  — размер корпуса

$k$  — константа (обычно  
10—100)

$\beta$  — константа  $0 < \beta < 1$   
(обычно 0,4—0,6)

Maximum Likelihood Estimation (MLE) — based on the observed frequency in the corpus:

$$P = \frac{f(w)}{N}$$

где

$P$  — probability;

$f(w)$  — word frequency  $w$ ;

$N$  — corpus size.

Model Based Estimation (LNRE model) — based on the knowledge of the general properties of word distributions:

$$P = \frac{C}{(r(w) + b)^a}$$

$P$  — probability;

$f(w)$  — word rank  $w$  in a frequency list;

$a, b$  — model parameters;

$C$  — normalizing constant.

This is Zipf-Mandelbrot model

## Лексическая статистика

Закон Ципфа

Частотность и состав лексикона

Размер и скорость роста словаря

Размер лексикона/Лексическое разнообразие

Практические следствия закона Ципфа

## КОЭФФИЦИЕНТ ЛЕКСИЧЕСКОГО РАЗНООБРАЗИЯ

Одна из первых и широко используемых мер сложности речи/текста.

$$TTR = \frac{V}{N} \quad (6)$$

$V$  — размер словаря, число разных словоформ/лемм в тексте (types)

$N$  — число словоформ в тексте

Обратная величина: средняя частотность слов в тексте

$$F_{mean} = \frac{N}{V} \quad (7)$$

## КОЭФФИЦИЕНТ ЛЕКСИЧЕСКОГО РАЗНООБРАЗИЯ

Одна из первых и широко используемых мер сложности речи/текста.

$$TTR = \frac{V}{N} \quad (6)$$

$V$  — размер словаря, число разных словоформ/лемм в тексте (types)

$N$  — число словоформ в тексте

Обратная величина: средняя частотность слов в тексте

$$F_{mean} = \frac{N}{V} \quad (7)$$

Классические применения:

- определение авторства
- оценка сложности (детской) речи — развития речи

Проблемы:

- **зависит от длины текста**, длиннее текст — ниже TTR ( $r=0.99$ ).
- зависит от способа выделения types (словоформы/леммы)

Нормализованная версия TTR:

- для сравнения используются фрагменты текста одинаковой длины.



## Лексическая статистика

Закон Ципфа

Частотность и состав лексикона

Размер и скорость роста словаря

Размер лексикона/Лексическое разнообразие

Практические следствия закона Ципфа

1. **Data sparseness** — в сколь угодно большом корпусе:
  - почти все слова встречаются очень редко;
  - небольшая группа частотных слов составляет значительную часть токенов корпуса;
  - LNRE — Large Number of Rare Events.
2. **Рост словаря** — даже очень большие корпуса не содержат всех слов языка:
  - искаженная оценка вероятности слова по частотности в корпусе;
  - нельзя использовать размер словаря для оценки степени лексического разнообразия текста.
3. **Знания о распределении** слов в любом тексте можно использовать для оптимизации и построения моделей.

## ЗАКЛЮЧЕНИЕ

---

- Помни о словах, которые еще не встретились. Делай на них **скидку**.
- Никогда не суди о богатстве словаря автора по количеству разных слов в тексте.
- Откинув небольшое число самых частотных слов, можно резко сократить объем корпуса, сохранив бóльшую часть смысловых слов.