

Кирилл Маслинский  
«Встреча Байеса и Ципфа, или Как измерить  
различия в частотности слов»

27 января 2024 года





# Вероятность появления слова постоянна в рамках корпуса следовательно

Различия в частотности слова между корпусами мы можем трактовать как различия в вероятности слова в разных корпусах.



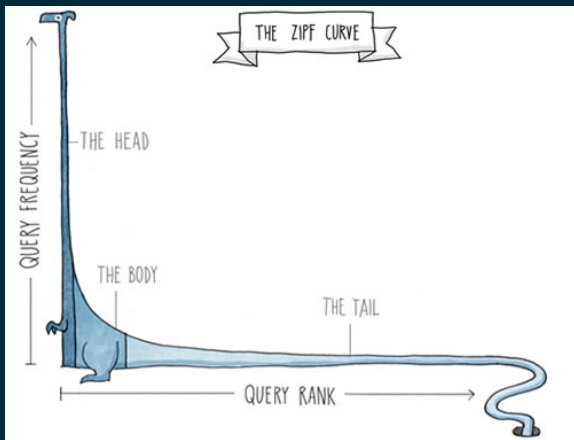
Закон Ципфа Предсказывает частотность слова по его рангу в частотном списке:

$$f(w) \sim \frac{1}{r(w)} \quad (1)$$

$f(w)$  — частотность слова  $w$

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

## Персонаж: Джордж Ципф



## ЭКСПЕРИМЕНТ С ЦАРЯМИ И ЦВЕТОЧКАМИ

# ОТНОШЕНИЕ ПРАВДОПОДОБИЯ

---



## Accurate Methods for the Statistics of Surprise and Coincidence

Ted Dunning\*  
New Mexico State University

*Much work has been done on the statistical analysis of text. In some cases reported in the literature, inappropriate statistical methods have been used, and statistical significance of results have not been addressed. In particular, asymptotic normality assumptions have often been used unjustifiably, leading to flawed results.*

*This assumption of normal distribution limits the ability to analyze rare events. Unfortunately rare events do make up a large fraction of real text.*

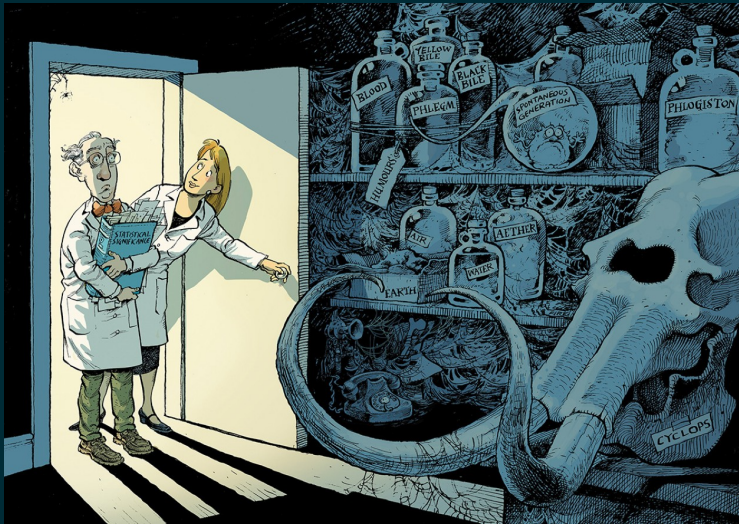


Цитата с сайта: <https://ucrel.lancs.ac.uk/llwizard.html>

The higher the G2 value, the more significant is the difference between two frequency scores. For these tables, a G2 of 3.8 or higher is significant at the level of  $p < 0.05$  and a G2 of 6.6 or higher is significant at  $p < 0.01$ .

- 95th percentile; 5% level;  $p < 0.05$ ; critical value = 3.84
- 99th percentile; 1% level;  $p < 0.01$ ; critical value = 6.63
- 99.9th percentile; 0.1% level;  $p < 0.001$ ; critical value = 10.83
- 99.99th percentile; 0.01% level;  $p < 0.0001$ ; critical value = 15.13

# Уходящая эпоха СТАТИСТИЧЕСКОЙ ЗНАЧИМОСТИ



Amrhein et al. "Scientists rise up against statistical significance" (2019)

## SIMPLE MATHS

---

## Simple maths (by Adam Kilgarriff)

«это слово встречается в этом корпусе вдвое чаще, чем в том»

- Самый простой подход
  - Нормализовать частотности
    - употреблений на тысячу или употреблений на миллион (IPM)
  - Вычислить отношение нормализованных частотностей
  - Отсортировать список слов по значению отношения



## ПРОБЛЕМА 1: НЕЛЬЗЯ ДЕЛИТЬ НА 0

слово	fc	rc	отношение
редкость	10	0	?
помешивать	100	0	?
вкуснотища	1000	0	?

Стандартное решение: прибавить 1:

слово	fc	rc	отношение
редкость	11	1	11
помешивать	101	1	101
вкуснотища	1001	1	1001



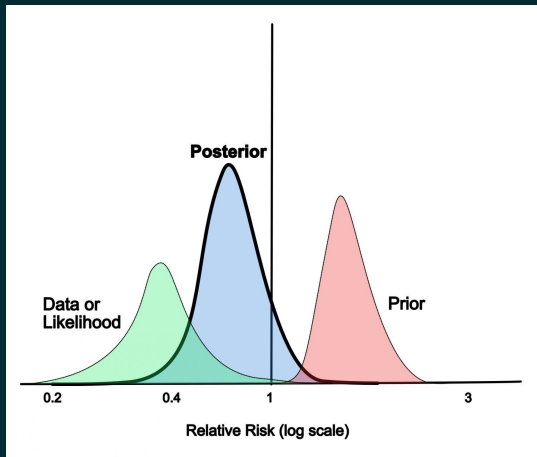


Теорема Байеса

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (2)$$

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}} \quad (3)$$

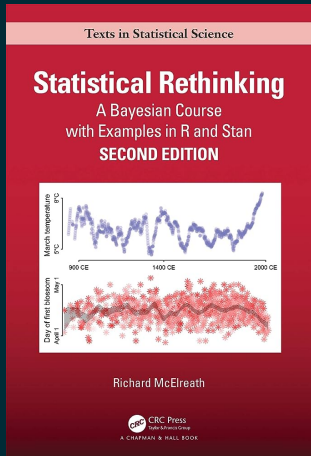




# Текстуальная гомогенность (пренебрегаем структурой текста)

следовательно

Можем рассматривать частотность слова в корпусе как результат пуассоновского процесса (с параметром  $\lambda$ , отражающим вероятность слова в корпусе).



AND free lecture series on  
Youtube







[https://github.com/  
maslinych/bayes-zipf](https://github.com/maslinych/bayes-zipf)