

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ ИМЕНИ Н. Э. БАУМАНА
Факультет информатики и систем управления
Кафедра теоретической информатики и компьютерных технологий

Лабораторная работа №1
по курсу «Автоматическая обработка текстов»
«Закон Ципфа»

Выполнил:
студент группы ИУ9-11М
Беляев А. В.

Проверила:
Лукашевич Н. В.

Москва 2018

1 Цель работы

Проверить выполнение закона Ципфа для художественного произведения

2 Текст программы

```
1 import re
2 import pymorphy2
3 import matplotlib
4 import numpy as np
5 from collections import Counter
6
7 matplotlib.use('TkAgg') # macos fix
8
9 morph = pymorphy2.MorphAnalyzer()
10
11
12 def draw_plot(x_values, y_values):
13     from matplotlib import pyplot as plt
14     plt.plot(np.array(x_values), np.array(y_values))
15     plt.show()
16
17
18 def main():
19     with open("a_dance_with_dragons.txt", encoding="utf-8", mode="r") as f:
20         contents = f.read().lower()
21
22         clean_text = re.sub("[^a-z0-9]", " ", contents)
23         splitted = re.compile("\s+").split(clean_text)
24         normalized = list(map(lambda word: morph.parse(word)[0].normal_form
25                               , splitted))
26         sorted_by_rank = Counter(normalized).most_common()
27
28         with open('out.txt', 'w+') as out:
29             for word, count in sorted_by_rank[:1000]:
30                 print(f'{word}\t{count}', file=out)
31
32         frequencies_list = list(map(lambda word_count_tuple:
33                                     word_count_tuple[1], sorted_by_rank))
34         draw_plot(list(range(1, len(frequencies_list) + 1)),
35                   frequencies_list)
```

Листинг 1: Исходный код программы

3 Результаты тестирования

В качестве художественного произведения был взят роман «Танец с драконами» из цикла «Песнь Льда и Пламени» Джорджа Р. Р. Мартина. В результате тестирования был построен график распределения частоты использования слов. Как видно из графика на Рис. 1, закон Ципфа выполняется и частота встречаемости слова обратно пропорциональна рангу слова.

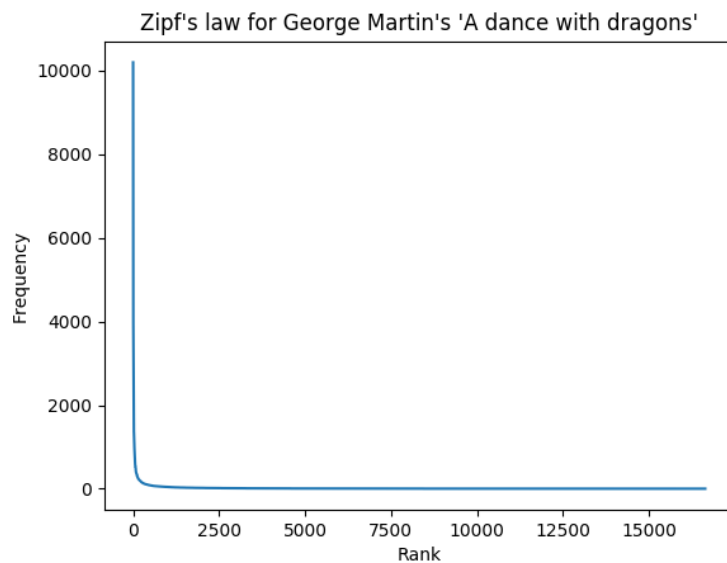


Рис. 1:

Таблица 1: Частоты некоторых слов с рангом от 30 до 50

слово	частота
один	1169
лорд	1049
человек	1021
так	973
сказать	956
джон	871
рука	835
тирион	696
кто	680
только	677
большой	606
королева	593

Частоты некоторых слов в ранге от 30 до 50 представлены в Таблице 1

4 Выводы

В ходе работы было проверено выполнение закона Ципфа для романа «Танец с драконами». Тестирование показало, что закон выполняется.