

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

Лабораторная работа №5
по курсу «Автоматическая обработка текстов»
«Лемматизация»

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

Проверила:
Лукашевич Н. В.

Москва 2018

1 Цель работы

Обработать текст и построить частотный список получившихся лемм.

2 Текст программы

```
1 import re
2 import pymorphy2
3 from collections import Counter
4
5 morph = pymorphy2.MorphAnalyzer()
6
7 file_name = "a_dance_with_dragons.txt"
8 N_WORDS_TO_PRINT = 50
9
10 def read_words(filename: str) -> list:
11     contents = open(filename, 'r').read().casefold()
12     clean = re.sub(r'^a-z0-9', '', contents) # add cyrillic
13     splitted = re.compile(r'\s+').split(clean)
14     return splitted[:-1] # remove last empty word
15
16 def main():
17     word_list = read_words(file_name)
18
19     normalized = list(map(lambda word: morph.parse(word)[0].normal_form,
20                           word_list))
21     sorted_by_rank = Counter(normalized).most_common()
22
23     with open('out.txt', 'w+') as out:
24         for word, count in sorted_by_rank[:N_WORDS_TO_PRINT]:
25             print(f'{word}\t\t{count}', file=out)
26
27 if __name__ == '__main__':
28     main()
```

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 ходе работы текст был обработан и лемматизирован с помощью библиотеки `pymorph2`. Результаты представлены в таблице 1.

Из документации к `Pymorphy2`: «*Pymorphy2 использует алгоритм нахождения нормальной формы, который работает наиболее быстро (берется первая форма в лексеме) - поэтому все причастия нормализуются в инфинитивы.*»

4 Выводы

В ходе работы была произведена обработка и лемматизация литературного произведения. Общее количество уникальных слов сократилось, а частота встречи лемм гораздо выше частоты встречи аналогичных слов.

Таблица 1: Наиболее часто встречаемые леммы

лемма	частота
и	10198
он	7911
в	6640
не	5337
на	4774
быть	3990
с	3654
что	3420
они	3158
это	2546
но	2418
она	2358
как	2094
из	2002
а	1777
весь	1772
за	1645
ты	1618
свой	1577
мы	1537
вы	1357
её	1356
тот	1352
когда	1298
мочь	1291
то	1269
от	1246
один	1169
бы	1111
мой	1110
чтобы	1066
лорд	1049
если	1031
человек	1021
так	973
сказать	956
этот	928
джон	871
рука	835
ещё	830
тириона	696