

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

Лабораторная работа №8
по курсу «Информационный поиск»

«Ранжирование документов по запросу.
Языковая модель»

Выполнил:
студент группы ИУ9-21М
Беляев А. В.

Проверила:
Лукашевич Н. В.

Москва 2019

1 Цель работы

Для фактов из Л.Р.1 и соответствующих статей, необходимо среди этих статей найти предложения, содержащие факты. Поиск осуществляется с помощью языковой модели

2 Ход работы

В Л.Р.1 были предложены следующие факты:

- В рецензии на компьютерную игру критик пожаловался на то, что смерть заставляет начинать уровень заново.
- Под стенами осажденной шведами русской крепости немцы побили шотландцев за пиво.
- Есть версия, что Джек Потрошитель был женщиной.

В этих фактах содержались ссылки на следующие статьи:

- Earth Shaker (игра)
- Осада Везенберга 1574 года
- Пиво
- Джек Потрошитель
- Мэри Пирси

Перед началом ранжирования документы и запросы были предобработаны: были убраны разделители и статьи были разбиты на предложения с помощью NLTK, убрана пунктуация, убраны слишком короткие предложения (менее трех слов), произведена нормализация предложений.

Функции чтения файлов, а также «очистки» и нормализации предложений были взяты из Л.Р. 3.

3 Текст программы

```
1 from lab3 import read_sentences, clean_up_sentence, normalize_sentence
2
3 ARTICLES = [
4     'art1-game.txt',
5     'art2-siege.txt',
6     'art3-beer.txt',
7     'art4-jack.txt',
8     'art5-jack-rule-63.txt'
9 ]
10
```

```

11 FACTS = [
12     'В рецензии на компьютерную игру критик пожаловался на то, что
    ↳ смерть заставляет начинать уровень заново.',
13     'Под стенами осажденной шведами русской крепости немцы побили
    ↳ шотландцев за пиво.',
14     'Есть версия, что Джек Потрошитель был женщиной.'
15 ]
16
17 COLLECTION = []
18
19 RESULTS_TO_OUTPUT = 5
20
21 LAMBDA = 0.5
22
23
24 def lang_model(query: str, docs: list) -> list:
25     def collection_model (term: str) -> float:
26         return COLLECTION.count(term) / len(COLLECTION)
27
28     def document_model(term: str, document: list) -> float:
29         return document.count(term) / len(document) #  $tf(t,d) / len(d)$ 
30
31     similarity = {}
32     for doc in docs:
33         p = 1
34         for term in query.split(' '):
35             p *= (1 - LAMBDA) * collection_model(term) + LAMBDA *
36                 ↳ document_model(term, doc)
37             assert 0 != p # make sure smoothing works
38             similarity[doc] = p
39
40     sorted_by_weight = sorted(similarity.items(), key=lambda doc:
41                               ↳ doc[1], reverse=True)
42     return sorted_by_weight[:RESULTS_TO_OUTPUT]
43
44 def main():
45     sentences = []
46     for article_filename in ARTICLES:
47         sentences.extend(read_sentences(article_filename))
48
49     # create collection from words of all sentences
50     for sentence in sentences:
51         COLLECTION.extend(sentence.split(' '))

```

```

52     # remove empty word that could appear by mistake :)
53     if '' in COLLECTION:
54         COLLECTION.remove('')
55
56     clean_queries = [clean_up_sentence(s) for s in FACTS]
57     queries = [normalize_sentence(s) for s in clean_queries]
58
59     # smoothing: add words from query to collection to avoid 'p(term)==0'
60     for q in queries:
61         COLLECTION.extend(q.split(' '))
62
63     for q in queries:
64         matched_docs = lang_model(q, sentences)
65         print(q)
66         for match in matched_docs:
67             print(f'\t{match[1]} {match[0]}')
68
69
70 if __name__ == '__main__':
71     main()

```

4 Результаты работы

Результаты ранжирования представлены в таблице 1. В таблице представлены очищенные нормализованные данные, с которыми работал алгоритм, а также итоговые веса документов по отношению к запросу.

Основная проблема - слова в запросе и документах были сопоставлены «как есть». Если бы допускалось использование синонимов, результаты были бы значительно лучше.

Также заметно тяготение алгоритма к более коротким предложениям. В целом же результат ранжирования получился схожим с векторной моделью и TFIDF. Точно так же, как и в других моделях, документы были выданы верно в 2х из 3х случаев, с незначительными отличиями.

В качестве отличительной черты стоит отметить простоту алгоритма (в том числе при нехватке данных (сглаживание)) по сравнению с двумя предыдущими.

5 Выводы

В ходе работы были изучена языковая модель, позволяющая ранжировать документы по степени соответствия запросу. Наивная реализация этого метода далека от идеала. Тем не менее, результат получился схожим с предыдущими реализациями.

Таблица 1: Языковая модель

	рецензия компьютерный игра критик пожаловаться смерть заставлять начинать уровень заново
8.34e-33	вскрывать брюшной полость джек потрошитель начинать уже смерть жертва
7.75e-33	потеря жизнь уровень запускаться заново жизнь последний игра заканчиваться
8.24e-34	метр выпуск книга компьютерный игра быть указать игра хорошесть графика отличный звуковой сопровождение
7.74e-34	заклучение критик сообщить впечатлеть
7.46e-34	автор книга компьютерный мир посчитать игра отличный график очень неплохой музыка
	стен осажда нной швед русской крепость немец побить шотландец пиво
7.87e-33	неоднократный попытка швед совершить подкоп взорвать стена вовремя пресекается защитник крепость
4.35e-34	результат бойня погибнуть немец шотландец
1.78e-34	шведский финский солдат присоединиться шотландский преимущественно пехота немецкий преимущественно конница артиллерия наёмник
1.40e-34	почти девятимесячный осада ревелеть оплот шведский владычество прибалтика увенчаться успех войско иван грозный удалться взять несколько меньший крепость частность вейсенштейн
1.28e-34	шотландец бежать немец русский гарнизон везенберг быть поздний доставить москва
	версия джек потрошитель быть женщиной
4.32e-12	основа быть взять женский версия убийца
2.19e-12	придерживаться версия пять жертва
1.71e-12	один версия имя джек потрошитель скрываться душевнобольной польский еврей эмигрант аарон косминский
1.36e-12	ставить этот версия сомнение однозначный доказательство тот жертва быть задушить существовать
1.22e-12	утверждать джек потрошитель быть льюис кэрролл