

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Курсовая работа по курсу «Информационный поиск»

Студент: И. Д. Недосеков
Преподаватель: А. А. Кухтичев
Группа: М8О-406Б
Дата:
Оценка:
Подпись:

Москва, 2022

Курсовая работа

Необходимо разработать систему исправления опечаток в поисковом запросе.

1 Описание

Для исправления опечаток будем использовать биграммы. N-грамма — последовательность из n элементов. С семантической точки зрения, это может быть последовательность звуков, слогов, слов или букв. Для нас подойдет биграмма из букв. Индекс биграмм будем строить по готовому корпусу. Будем считать что слово нужно исправлять, если слова нет в словаре. Для слов с опечаткой будем находить все слова в которых содержатся некоторые биграммы из слова с опечаткой. Для слова и кандидата будем считать коэффициент Жаккара. Если коэффициент больше 0.8, то заменяем слова на кандидата. Иначе ищем слово с коэффициентом большим 0.5 и делаем запрос с этим словом и без него (поисковую выдачу смешиваем).

2 Исходный код

```
1 from helpers.mongodb_connector import mongo_client
2 from helpers.bigrams import get_bigrams_of_word
3 from config import Config
4
5
6
7 class Dimension:
8
9 def __init__(self, word: str) -> None:
10     self.word = word
11     self.words_bigrams = set(get_bigrams_of_word(word))
12
13 def distance_to(self, candidate: str) -> float:
14     if abs(len(self.word) - len(candidate)) > 3:
15         return 0
16     words_bigrams = self.words_bigrams
17     candidates_bigrams = set(get_bigrams_of_word(candidate))
18     distance = abs(len(words_bigrams) - len(words_bigrams & candidates_bigrams))
19     return abs(1 - distance / len(words_bigrams | candidates_bigrams))
20
21
22 async def get_fuzzy_words(word: str) -> list[str]:
23     bigrams = get_bigrams_of_word(word)
24     word_dimension = Dimension(word)
25     fuzzy_words_rating: set[tuple[str, float]] = set()
26     async for bigram_with_words in mongo_client.test.bigram_words.find({'bigram':
27         ↳ {'$in': bigrams}}):
28         fuzzy_words_rating.update(
29             (candidate, word_dimension.distance_to(candidate))
30             for candidate in bigram_with_words['words']
31             if word_dimension.distance_to(candidate) >=
32                 ↳ Config.minimum_distance_for_fuzzy_words
33         )
34
35     return [
36         (word, rank)
37         for word, rank in reversed(sorted(fuzzy_words_rating, key=lambda row: row[1]))
38     ]
39
40 async def is_dictionary_word(word: str) -> bool:
41     return bool(
```

```

40     await mongo_client.test.word_bigrams.count_documents({'word': word})
41 )
42
43 def get_transformed_queries(enrich_query: EnrichQuery) -> list[list[str]]:
44     queries: list[list[str]] = [[]]
45
46     for word in enrich_query.query:
47         if word not in enrich_query.fuzzy_mapping:
48             for query in queries:
49                 query.append(word)
50
51         elif enrich_query.fuzzy_mapping[word][0][1] >
52             ↪ Config.minimum_distance_for_fuzzy_words_for_replace:
53             for query in queries:
54                 query.append(enrich_query.fuzzy_mapping[word][0][0])
55
56         elif enrich_query.fuzzy_mapping[word][0][1] >
57             ↪ Config.minimum_distance_for_fuzzy_words_for_combine_request:
58             transformed = deepcopy(queries)
59             for query in transformed:
60                 query.append(enrich_query.fuzzy_mapping[word][0][0])
61
62             queries.extend(transformed)
63
64     return queries

```

```

----- построение индекса биграмм -----
1  import asyncio
2  import asyncstdlib
3  from loguru import logger
4  from helpers.mongodb_connector import mongo_client
5
6
7  async def build_bigram_words_dictionary():
8      total_documents = await mongo_client.test.word_bigrams.count_documents({})
9      percentage = 0
10     async for i, word_with_bigrams in
11         ↪ asyncstdlib.enumerate(mongo_client.test.word_bigrams.find()):
12         if i / total_documents * 100 > percentage:
13             logger.debug('Process {:.0%} documents', i / total_documents)
14             percentage += 10

```

```
14
15     word = word_with_bigrams['word']
16     await asyncio.gather(
17         *(
18             mongo_client.test.bigram_words.find_one_and_update(
19                 {'bigram': bigram},
20                 {'$push': {'words': word}},
21                 upsert=True,
22
23             )
24         for bigram in word_with_bigrams['bigrams']
25     )
26 )
```

3 Выводы

Выполнив первую курсовую работу по курсу «Информационный поиск», познакомился с основной теорией по исправлению поискового запроса, коэффициентом Жаккара. Для коротких слов исправление через биграммы плохо подходит, так как лишком большое число биграмм относительно общего количества в слове поврежденно (для "малоко"повреждено 2 из 5 биграмм). Для длинных слов такой метод вполне применим.

Список литературы

- [1] Маннинг, Рагхаван, Шютце *Введение в информационный поиск* — Издательский дом «Вильямс», 2011. Перевод с английского: доктор физ.-мат. наук Д. А. Ключина — 528 с. (ISBN 978-5-8459-1623-4 (рус.))