

Московский авиационный институт
(национальный исследовательский университет)

Факультет информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Курсовая работа по курсу «Информационный поиск»

Студент: А. К. Киреев
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. Индекс биграмм будем строить по готовому корпусу, по которому уже отработал nltk. Будем считать что слово нужно исправлять, если слова нет в словаре (словарь формируем из слов корпуса). 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. Выбираем из полученной выборки сто наиболее похожих слов (с наибольшим коэффициентом Жаккара), то есть делаем грубую оценку, а потом для них и исходного слова считаем расстояние Дамерау-Левенштейна и выбираем три слова с наименьшим расстоянием. Такая схема была выбрана, так как индекс биграмм быстро работает, но хуже умеет исправлять ошибки, чем алгоритм Дамерау-Левенштейна.

2 Исходный код

```
bigram_index.py
1  from typing import List, Dict, Set
2  from logging import Logger
3
4  from db import get_words_by_bigrams, check_if_exists
5
6  from search_helper.common import get_bigrams
7  from search_helper.metrics import jaccard_coef, damerau_levenshtein_distance
8
9
10 class BigramIndex:
11     def __init__(self, enriched_request: List[str],
12                 count_bound: int = 3, distance_bound: float = 3) -> None:
13         self.req = enriched_request
14         self.count_bound = count_bound
15         self.distance_bound = distance_bound
16         self.search_dict: Dict = dict()
17
18     async def build(self, logger: Logger) -> None:
19         for word in self.req:
20             if await check_if_exists(word):
21                 logger.debug('Word "%s" exists', word)
22                 self.search_dict[word] = [word]
23                 continue
24
25             self.search_dict[word] = []
26             bigrams: Set[str] = get_bigrams(word)
27             index: Dict = await get_words_by_bigrams(bigrams)
28
29             words: Set[str] = set()
30             for _, others in index.items():
31                 for other in others:
32                     words.add(other)
33
34             coefs: List[List[float]] = []
35             for other in words:
36                 coefs.append([jaccard_coef(word, other), other])
37             coefs.sort(reverse=True)
38
39             most_similar: List = [w for _, w in coefs[:100]]
40             distances: List[List[float]] = []
41             for other in most_similar:
```

```

42     distances.append([damerau_levenshtein_distance(word, other), other])
43     distances.sort()
44
45     for d, supposed in distances[:self.count_bound]:
46         logger.debug('Supposed word "%s" with LD-distance %s', supposed, d)
47         if d <= self.distance_bound:
48             self.search_dict[word].append(supposed)
49         else:
50             self.search_dict[word].append(word)
51
52     def get_supposed(self, word: str) -> Set[str]:
53         return self.search_dict[word]
54
55     def get_search_dict(self) -> Dict:
56         return self.search_dict
57

```

```

----- metrics.py -----
1  from typing import Set, Dict
2
3  from search_helper.common import get_bigrams
4
5
6  def damerau_levenshtein_distance(lhs: str, rhs: str) -> float:
7      d: Dict = dict()
8      for i in range(-1, len(lhs) + 1):
9          d[(i, -1)] = i + 1
10     for j in range(-1, len(rhs) + 1):
11         d[(-1, j)] = j + 1
12
13     for i in range(len(lhs)):
14         for j in range(len(rhs)):
15             if lhs[i] == rhs[j]:
16                 cost = 0
17             else:
18                 cost = 0.9
19             d[(i, j)] = min(
20                 d[(i - 1, j)] + 1,
21                 d[(i, j - 1)] + 1,
22                 d[(i - 1, j - 1)] + cost,
23             )
24     if i and j and lhs[i] == rhs[j - 1] and lhs[i - 1] == rhs[j]:
25         d[(i, j)] = min(d[(i, j)], d[(i - 2, j - 2)] + 1)

```

```

26     return d[len(lhs) - 1, len(rhs) - 1]
27
28
29
30 def jaccard_coef(lhs: str, rhs: str):
31     lhs_bigrams: Set[str] = get_bigrams(lhs)
32     rhs_bigrams: Set[str] = get_bigrams(rhs)
33
34     coef: float = len(lhs_bigrams & rhs_bigrams) / len(lhs_bigrams | rhs_bigrams)
35     return coef

```

```

----- save_bigrams.py -----
1  import asyncio
2
3  from typing import List
4
5  from motor.motor_asyncio import AsyncIOMotorClient
6
7  client = AsyncIOMotorClient()
8
9
10 async def save_bigrams() -> None:
11     async for record in client.IR.WordsStorage.find():
12         word: str = record['word']
13         bigrams: List[str] = record['bigrams']
14
15         await asyncio.gather(
16             *(
17                 client.IR.BigramStorage.find_one_and_update(
18                     {'bigram': bigram},
19                     {'$push': {'words': word}},
20                     upsert=True
21                 )
22                 for bigram in bigrams
23             )
24         )

```

3 Выводы

Выполнив курсовую работу по курсу «Информационный поиск», познакомился с основной теорией по исправлению поискового запроса, коэффициентом Жаккара, вспомнил алгоритм подсчета расстояния Дамерау-Левенштейна. Следует заметить, что для коротких слов исправление через биграммы плохо подходит, так как лишком большое число биграмм относительно общего количества в слове поврежденно, поэтому после выбора ста наиболее подходящих слов я использую алгоритм Дамерау-Левенштейна.

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.))