

資訊檢索與文字探勘導論

作業二

會計四
B06702064
林聖硯

1. 執行環境

Anaconda Spyder

2. 程式語言

Python 3.6.8

3. 執行方式

- (1) 需要 import nltk(preprocessor.py 檔需要)及 numpy 套件
使用 pip install nltk、pip install numpy
- (2) preprocessor.py 放入 hw1 裡面文字前處理的 function
- (3) 請先在建一個名為「result」的資料夾(存放 Task2 tfidf 的結果)
(放在 zip 檔內的 Doc1 是我手動移動出來的)
- (4) 打開 pa2.py 並執行 217-241 行即可得到 cosine similarity 的結果

*備註

- (1) 請依序執行(後面的 Task 會需要用到前面生成的東西)
- (2) 整個執行過程大約需要 6 分鐘(我一直用 for + list 抱歉 QQ)
- (3) txt 檔及 stopwords 檔已經寫死在程式裡
- (4) 若要計算不同文章的 cosine similarity 請更改 num_doc1 及 num_doc2

傳入的參數，直接輸入想要的兩篇文章號碼即可

```
216 #####Task 1: get document frequency#####
217 number_txt = generate_txt_list()
218 raw_text = get_raw_text(number_txt)
219 bag_of_words = get_BOW(raw_text)
220 df = get_df(bag_of_words, raw_text)
221 words_index = get_word_index(bag_of_words)
222 save_df(words_index, bag_of_words, df)
223 #195秒完成
224
225 #####Task 2: Transfer each document into a tf-idf unit vector#####
226 num_doc = get_doc_length(raw_text)
227 len_BOW = get_BOW_length(bag_of_words)
228 idf_array = get_idf_array(df)
229 tf_matrix = get_tf_matrix(raw_text, bag_of_words)
230 tfidf_matrix = get_tfidf_matrix(tf_matrix, idf_array, len_BOW)
231 unit_tfidf_matrix = get_unit_tfidf_matrix(num_doc, tfidf_matrix)
232 save_unit_tfidf_matrix(num_doc, unit_tfidf_matrix, number_txt)
233 #140秒完成
234
235 #####Task 3: compute cosine similarity#####
236 num_doc1 = 1
237 num_doc2 = 2
238 len_BOW = get_BOW_length(bag_of_words)
239 index_1, index_2, unit_tfidf_1, unit_tfidf_2 = retrieve_index_and_tfidf(num_doc1, num_doc2)
240 cosine_similarity = compute_cosine_similarity(index_1, index_2, unit_tfidf_1, unit_tfidf_2, len_BOW)
241 print(f"Cosine similarity of document{num_doc1} and document{num_doc2} is", round(cosine_similarity, 5))
```

結果呈現

```
In [32]: print(f"Cosine similarity of document{num_doc1} and document{num_doc2} is",
round(cosine_similarity, 5))
Cosine similarity of document1 and document2 is 0.18093
```

4. 作業處理邏輯說明

(0)import 套件

導入 numpy 及 Preprocessor 檔案

```
1 #import necessary modules from hw1
2 from Preprocessor import text_preprocessing, get_text
3
4 #import necessary module
5 import numpy as np
6
```

(Task 1-1)讀檔

generate_txt_list:製造出['1.txt', '2.txt',...]的 list

get_raw_text:讀取每個文本並且經過 preprocessing 之後存到 list

```
8 #####Task 1: caculate document frequency#####
9 def generate_txt_list():
10     '''
11     generate a list containing ['1.txt', '2.txt', ...]
12     '''
13     number_txt = []
14     for i in range(1095):
15         number_txt.append(str(i+1) + ".txt")
16     return number_txt
17
18 def get_raw_text(number_txt):
19     '''
20     retrieve 1095 documents from txt files
21     output : list of list ([news[preprocessed words, prerpocessed words, ...]])
22     '''
23     raw_text = []
24     for txt in number_txt:
25         text = get_text(txt)
26         text_final = text_preprocessing(text)
27         raw_text.append(text_final)
28     return raw_text
```

(Task 1-2)

get_BOW:創造出 bag of words 的 list (裡面的 token 順序就固定下來)

get_df:計算 BOW 裡面每個 token 在幾個文本出現過

```
30 def get_BOW(raw_text):
31     '''
32     generate bag of words model
33     '''
34     bag_of_words = []
35     for text in raw_text:
36         bag_of_words.extend(text)
37     bag_of_words = dict.fromkeys(bag_of_words)
38     bag_of_words = sorted(bag_of_words)
39     bag_of_words = bag_of_words[1:]
40     return bag_of_words
```

```

42 def get_df(bag_of_words, raw_text):
43     '''
44     generate document frequency in a list
45     the order of the words is unchanged
46     '''
47     df = []
48     for word_find in bag_of_words:
49         count = 0
50         for news in raw_text:
51             for word_news in news:
52                 if (word_find == word_news):
53                     count += 1
54                     break
55             df.append(count)
56     return df

```

(Task 1-3)

get_word_index: 替 BOW 裡面的每個 token 創造出對應的 term_index(從 1 開始)

save_df: 將 term_index, term, document frequency 存到 txt 裡

```

58 def get_word_index(bag_of_words):
59     '''
60     generate index for each word in BOW
61     '''
62     words_index = [i+1 for i in range(len(bag_of_words))]
63     return words_index
64
65 def save_df(words_index, bag_of_words, df):
66     '''
67     save document frequency of each word into dictionary.txt
68     '''
69     with open('dictionary.txt', 'w') as f:
70         f.write('{:<15} {:<15} {}'.format('t_index', 'term', 'df'))
71         f.write('\n')
72         for i in range(len(words_index)):
73             f.write('{:<15} {:<15} {}'.format(words_index[i], bag_of_words[i], df[i]))
74             f.write('\n')

```

(Task 2-1)

get_doc_length: 計算有幾篇 document，先存著備用

get_BOW_length: 計算 BOW 裡面有幾個 token，存著備用

get_idf_array: 使用 df 計算 idf，並且存到 numpy array 方便後面計算

```

76 #####Task 2: Transfer each document into a tf-idf unit vector#####
77 def get_doc_length(raw_text):
78     '''
79     return number of total documents
80     '''
81     num_doc = len(raw_text)
82     return num_doc
83
84 def get_BOW_length(bag_of_words):
85     '''
86     return length of BOW
87     '''
88     len_BOW = len(bag_of_words)
89     return len_BOW
90
91 def get_idf_array(df):
92     '''
93     turn df into idf numpy array
94     '''
95     df_array = np.array(df)
96     num_doc = len(raw_text)
97     idf_array = np.log10(num_doc/df_array)
98     return idf_array

```

(Task 2-2)

get_tf_matrix: 生成 tf 矩陣

get_tfidf_matrix: 使用 idf vector 和 tf matrix 生成 tfidf 矩陣

```
100 def get_tf_matrix(raw_text, bag_of_words):
101     """
102     building tf matrix (14269, 1095)
103     row : term, column : document
104     """
105     tf = []
106     #1095, 14299
107     for news in raw_text:
108         temp = []
109         for word_find in bag_of_words:
110             count = 0
111             for word_news in news:
112                 if (word_find == word_news):
113                     count += 1
114             temp.append(count)
115         tf.append(temp)
116     #tf-idf
117     tf_matrix = np.array(tf)
118     tf_matrix = tf_matrix.transpose()
119     return tf_matrix
120
121 def get_tfidf_matrix(tf_matrix, idf_array, len_BOW):
122     """
123     combine tf matrix and idf vector into tfidf matrix
124     return (14269, 1095) (term, document)
125     """
126     tfidf_list = []
127     for i in range(len_BOW):
128         tfidf = tf_matrix[i] * idf_array[i]
129         tfidf_list.append(tfidf)
130     tfidf_matrix = np.array(tfidf_list)
131     return tfidf_matrix
```

(Task 2-3)

get_unit_tfidf_matrix: normalize tfidf

save_unit_tfidf_matrix: 將 normalized 過後的結果存入 txt 檔

```
133 def get_unit_tfidf_matrix(num_doc, tfidf_matrix):
134     """
135     normalize tfidf matrix with respect to each document
136     return (14269, 1095) (term, document)
137     """
138     unit_tfidf_list = []
139     for i in range(num_doc):
140         tfidf_sum = np.sqrt(np.sum(np.square(tfidf_matrix.transpose()[i])))
141         unit_tfidf_list.append(tfidf_matrix.transpose()[i]/tfidf_sum)
142     unit_tfidf_matrix = np.array(unit_tfidf_list)
143     unit_tfidf_matrix = unit_tfidf_matrix.transpose()
144     return unit_tfidf_matrix
```

```
146 def save_unit_tfidf_matrix(num_doc, unit_tfidf_matrix, number_txt):
147     """
148     save tf-idf file
149     """
150     for i in range(num_doc):
151         doc_tfidf = unit_tfidf_matrix.transpose()[i]
152         #nonzero_index start from zero to 14298
153         nonzero_index = np.nonzero(doc_tfidf)[0]
154         location = './result/doc'+ number_txt[i]
155         with open(location, 'w') as f:
156             f.write('{}'.format(len(nonzero_index)))
157             f.write('\n')
158             f.write('{:<10} {}'.format('t_index', 'tf-idf'))
159             f.write('\n')
160             for index in nonzero_index:
161                 f.write('{:<10} {}'.format(words_index[index], doc_tfidf[index]))
162                 f.write('\n')
```

(Task3-1)

retrieve_index_and_tfidf: 將 unit_tfidf 及其 index 從 txt 檔取出

```
164 #####Task 3: compute cosine similarity#####
165 def retrieve_index_and_tfidf(num_doc1, num_doc2):
166     '''
167     get term indexes and their unit tfidfs from specified documents
168     '''
169     index_1 = []
170     index_2 = []
171     unit_tfidf_1 = []
172     unit_tfidf_2 = []
173     pos1 = './result/doc' + str(num_doc1) + '.txt'
174     pos2 = './result/doc' + str(num_doc2) + '.txt'
175     with open(pos1) as f:
176         count = 0
177         for line in f.readlines():
178             count += 1
179             if (count == 1 or count == 2):
180                 continue
181             else:
182                 index_1.append(int(line.split()[0]))
183                 unit_tfidf_1.append(float(line.split()[1]))
184     with open(pos2) as f:
185         count = 0
186         for line in f.readlines():
187             count += 1
188             if (count == 1 or count == 2):
189                 continue
190             else:
191                 index_2.append(int(line.split()[0]))
192                 unit_tfidf_2.append(float(line.split()[1]))
193     index_1 = np.array(index_1)
194     index_2 = np.array(index_2)
195     unit_tfidf_1 = np.array(unit_tfidf_1)
196     unit_tfidf_2 = np.array(unit_tfidf_2)
197     return index_1, index_2, unit_tfidf_1, unit_tfidf_2
```

(Task3-2)

compute_cosine_similarity: 計算兩篇文章的餘弦相似度，並回傳結果

```
199 def compute_cosine_similarity(index_1, index_2, unit_tfidf_1, unit_tfidf_2, len_BOW):
200     unit_tfidf_array_1 = np.zeros(len_BOW)
201     unit_tfidf_array_2 = np.zeros(len_BOW)
202     count_1 = 0
203     for index in index_1:
204         unit_tfidf_array_1[index-1] = unit_tfidf_1[count_1]
205         count_1 += 1
206
207     count_2 = 0
208     for index in index_2:
209         unit_tfidf_array_2[index-1] = unit_tfidf_2[count_2]
210         count_2 += 1
211     cosine_similarity = np.dot(unit_tfidf_array_1, unit_tfidf_array_2)
212     return cosine_similarity
```