

# 資訊檢索與文字探勘導論

## 作業三

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)、numpy、pandas 及 math 套件

使用 pip install nltk、pip install numpy、pip install pandas

(2) preprocessor.py 放入 hw1 裡面文字前處理的 function，stopwords.txt 裡面放入要被讀進來的 stopwords，training\_index 裡面放入每一個 class 的訓練集文章的編號，IRTM 資料夾裡放入所有文章的 txt 檔

(3) 執行 204 到 211 行即可得到結果

```
204 training_index_txt, training_index_all = get_training_index()
205 df_list, tf_list = calculate_df_tf(training_index_txt)
206 df_all, tf_all = calculate_df_tf_all(df_list, tf_list)
207 total_training_doc_num = calculate_num_doc(training_index_txt)
208 features = chi_square_feature_selection(df_list, df_all, total_training_doc_num)
209 conditional_prob, prior = train_NB_clf(tf_list, tf_all, features)
210 result_class = test_NB_clf(training_index_all, prior, features, conditional_prob)
211 save_result(result_class)
```

### 4. 程式處理邏輯

(0) import 需要套件

```
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 import pandas as pd
7 import math
```

## -----前處理-----

### (1) 建立兩個 list

training\_index\_txt (list of list) : 外層의 list 代表每一個 class，內層의 list 放入所有 training document 的編號並且在結尾加上.txt

training\_index\_all (list) : 將所有的 training document 的編號放入一個 list 內

```
12 def get_training_index():
13     """
14     training_index_txt : list of list, each element consists of document number in each class with .txt at the end
15     training_index_all : list, each element consists of document number in each class
16     """
17     training_index_txt = []
18     training_index_all = []
19     f = open('training_index.txt')
20     lines = f.readlines()
21     for line in lines:
22         line = line.split(' ')
23         line = line[1:16]
24         temp = []
25         for docNum in line:
26             training_index_all.append(int(docNum))
27             temp.append(docNum + '.txt')
28         training_index_txt.append(temp)
29     return training_index_txt, training_index_all
```

(2) 計算每個 class 內每個 term 的 term frequency 和 document frequency 並且存入 df\_list 及 tf\_list 內。兩個 list 中的每一個元素都是一個 dictionary(總共 13 個)，放入每個 term 的 df 及 tf

```
31 #for each class, calculate document frequency and term frequency for each term and record them by using dictionary
32 def calculate_df_tf(training_index_txt):
33     """
34     df_list : 13 dictionary in a list, each dictionary contains word and its document frequency
35     tf_list : 13 dictionary in a list, each dictionary contains word and its term frequency
36     """
37     df_list = []
38     tf_list = []
39     for index_by_class in training_index_txt:
40         temp_list_df = [] #contains dictionary with keys = words, values = 1
41         temp_dict_tf = dict()
42         temp_dict_df = dict()
43         #iterate each document in class
44         for index in index_by_class:
45             raw_text = get_text(index)
46             text = text_preprocessing(raw_text)
47
48             #count df by combining dictionaries
49             temp_list_df.append(dict((word, 1) for word in set(text)))
50
51             #iterate each word in document
52             for word in text:
53                 if word not in temp_dict_tf.keys():
54                     temp_dict_tf[word] = 1
55                 else:
56                     temp_dict_tf[word] += 1
57
58             for my_dict in temp_list_df:
59                 for key, value in my_dict.items():
60                     temp_dict_df.setdefault(key, 0)
61                     temp_dict_df[key] += value
62
63             tf_list.append(temp_dict_tf)
64             df_list.append(temp_dict_df)
65
66     #do not record word:''
67     for my_dict in tf_list:
68         my_dict.pop('')
69
70     for my_dict in df_list:
71         my_dict.pop('')
72
73     return df_list, tf_list
```

(3) 對所有 term 計算所有 class 合併的 df 及 tf 留著備用

```
82 #calculate df and tf for each term for all training documents
83 ▼ def calculate_df_tf_all(df_list, tf_list):
84     df_all = {}
85     ▼ for my_dict in df_list:
86         ▼ for key, value in my_dict.items():
87             df_all.setdefault(key, 0)
88             df_all[key] += value
89
90     tf_all = {}
91     ▼ for my_dict in tf_list:
92         ▼ for key, value in my_dict.items():
93             tf_all.setdefault(key, 0)
94             tf_all[key] += value
95
96     return df_all, tf_all
```

(4) 計算每一個 class document 的數量及所有 class 加起來 document 的數量留著備用

```
106 #calculate number of training document
107 ▼ def calculate_num_doc(training_index_txt):
108     doc_num = []
109     ▼ for index in training_index_txt:
110         doc_num.append(len(index))
111     total_training_doc_num = sum(doc_num)
112
113     return total_training_doc_num
```

## -----feature selection-----

### (5) chi-square

對所有 class，去計算裡面每一個字的 chi-square，並且挑出每一個 class 內 chi-square 值最高的 500 個存到 set 內，再將這 13 個 set 放入名為 features 的 list 中。

```
115 #####2. for each class, choose 500 top features (words) by chi-square methods#####
116 def chi_square_feature_selection(df_list, df_all, total_training_doc_num):
117     k = 500
118     features = []
119     for my_dict in df_list:
120         word_list = []
121         score_list = []
122         for word in my_dict.keys():
123             word_list.append(word)
124             '''
125             presentNon : number of document that this word appears in class c
126             absentNon : number of document that this word doesn't appear in class c
127             presentNoff : number of document that this word appears in class other than c
128             absentNoff : number of document that this word doesn't appear in class other than c
129             '''
130             on_topic = 15
131             off_topic = total_training_doc_num - on_topic
132             present = df_all[word]
133             absent = total_training_doc_num - present
134
135             presentNon = my_dict[word]
136             absentNon = on_topic - presentNon
137             presentNoff = present - presentNon
138             absentNoff = absent - absentNon
139
140             E_00 = absent * off_topic / total_training_doc_num
141             E_01 = absent * on_topic / total_training_doc_num
142             E_10 = present * off_topic / total_training_doc_num
143             E_11 = present * on_topic / total_training_doc_num
144
145             chi_00 = ((absentNoff - E_00) ** 2) / E_00
146             chi_01 = ((absentNon - E_01) ** 2) / E_01
147             chi_10 = ((presentNoff - E_10) ** 2) / E_10
148             chi_11 = ((presentNon - E_11) ** 2) / E_11
149
150             score = chi_00 + chi_01 + chi_10 + chi_11
151             score_list.append(score)
152         score_arr = np.asarray(score_list)
153         index = (-score_arr).argsort()[:k]
154         feature_set = set([word_list[i] for i in index])
155         features.append(feature_set)
156     return features
```

## -----NB model-----

### (6) training Multinomial Naïve Bass Classifier

對所有的字，計算在每一個 class 內的條件機率，若有出現在挑出來的 500 個 feature 內才存入最後的字典中。最後，對每一個 class 都蒐集這 500 字的條件機率，並且將 13 個字典存入一個名為 conditional\_prob 的 list。

```
159 def train_NB_clf(tf_list, tf_all, features):
160     conditional_prob = []
161     count = 0
162     prior = []
163     #iterate each class
164     for class_dict in tf_list:
165         conditional_prob_dict = dict()
166         prior.append(15 / 195)
167         word_num = len(tf_all.keys()) ##500 * 13 or 全部不重複的字
168         sum_tf = sum(class_dict.values())
169         for word in tf_all.keys():
170             if word in features[count]:
171                 #add-one smoothing
172                 cond_prob = (class_dict[word] + 1) / (sum_tf + word_num)
173                 conditional_prob_dict[word] = cond_prob
174             conditional_prob.append(conditional_prob_dict)
175             count += 1
176
177     return conditional_prob, prior
```

### (7) predicting class by Multinomial Naïve Bass Classifier

對每個 document 裡面的 term，去對每一個 class 裡面的字典搜尋條件機率，若有找到則取自然對數後加到最後的 score 中。最後每一個 term 都會有 13 個 score，將 class 最高的 score 取出並且存入 result\_class 這個 dictionary 內。(key : document index, value : predicting class)

```
179 def test_NB_clf(training_index_all, prior, features, conditional_prob):
180     result_class = dict()
181     for index in range(1,1096,1):
182         if index not in training_index_all:
183             raw_text = get_text(str(index) + '.txt')
184             text = text_preprocessing(raw_text)
185             score_list = []
186             for i in range(13):
187                 score = 0
188                 score += math.log(prior[i])
189                 for word in text:
190                     if word in features[i]:
191                         score += math.log(conditional_prob[i][word])
192                 score_list.append(score)
193             score_arr = np.asarray(score_list)
194             best_class = (-score_arr).argsort()[0] + 1
195             result_class[index] = best_class[0]
196     return result_class
```

### (8) save result

將結果轉成 dataframe 後存入 'hw3\_result'.csv 並上傳 kaggle

```
198 def save_result(result_class):
199     result = pd.DataFrame(result_class.items(), columns = ['Id', 'Value'])
200     result = result.set_index('Id')
201     result.to_csv('hw3_result.csv')
```