

# HW3 Report

R07725021 資管碩 1 洪靖雯

## 1.執行環境

- Sublime
- Cmder

## 2.程式語言

- Python 3.6.6

## 3.執行方式

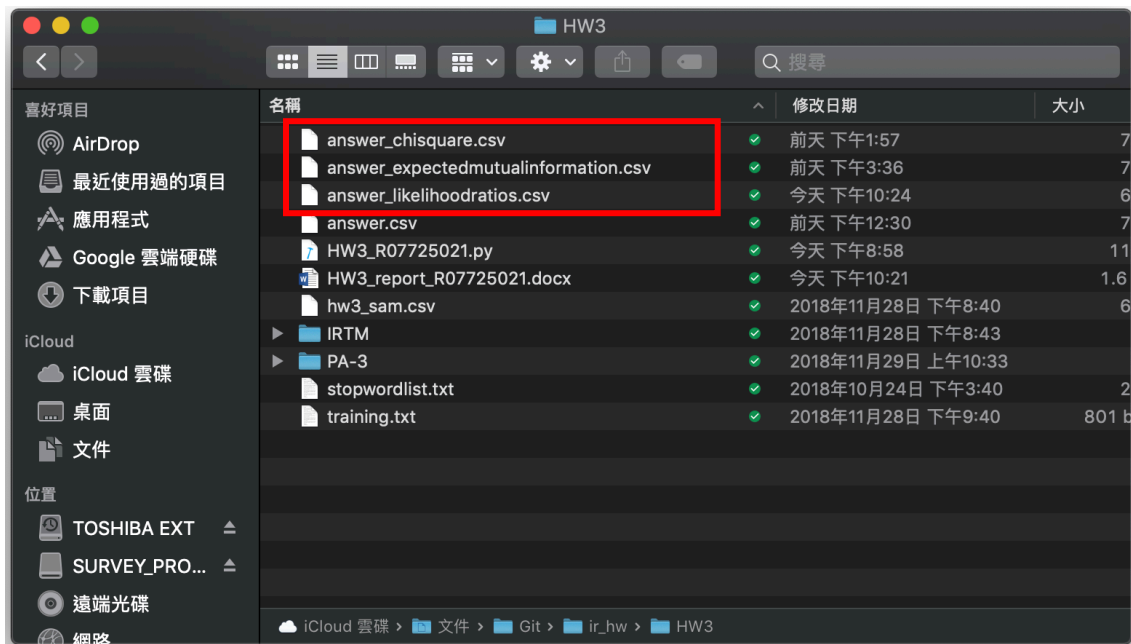
- 套件: nltk
  - 有另外用 conda 4.5.8 架虛擬環境安裝，直接用在 system 裡應該也可行
- 編譯方式: 使用 `python HW3_R07725021.py` 即可



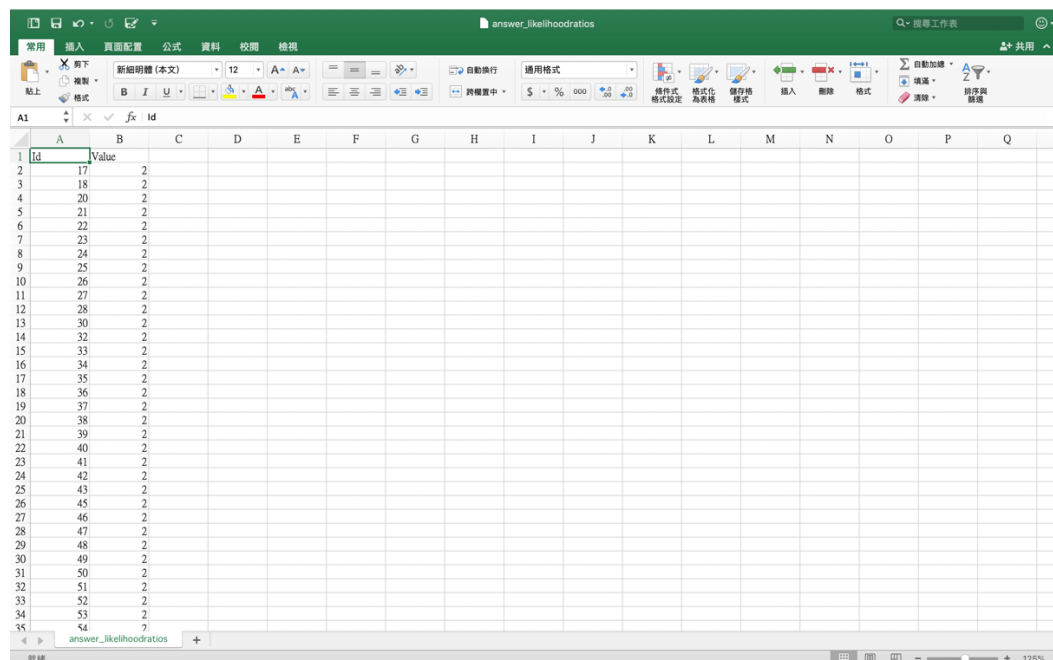
```
HW3 — -bash — 80x24
Last login: Tue Dec  4 20:44:15 on ttys000
[glpc1n247:~ hosi$ cd Documents/Git/ir_hw/
[glpc1n247:ir_hw hosi$ ls
HW1          HW2          HW3          README.md
[glpc1n247:ir_hw hosi$ cd HW3
[glpc1n247:HW3 hosi$ python HW3_R07725021.py
[glpc1n247:HW3 hosi$
```

## ➤ 輸出結果

- i. answer\_xxx.csv (總共有三個檔案，分別是用 chi-square, likelihood ratios 和 expected mutual information(EMI) feature selection 做出的結果，預測出來效果略有差異)



- answer\_likelihood.csv (目前效果最好)，共 900 個 testing data



## 4.作業處理邏輯說明

### i. HW3\_R07725021\_1.py

- Import 必要的套件

```
HW3_R07725021.py x
1 from nltk.tokenize import word_tokenize
2 from nltk.stem import PorterStemmer
3 from collections import defaultdict
4 import math
5 import csv
6
```

- 先將 1095 個 documents 做 tokenize(運用之前 HW1 的方法), 運用 dictionary 的結構來儲存, key 為 docid, value 為該文章的 terms。

```
7 def tokenize(docid): #term
8     file = open("IRTM/"+docid+".txt","r")
9     data = file.read()
10    file.close()
11
12    #tokenize(用nltk的tokenizer)
13    tokendata = word_tokenize(data)
14    regex = [c for c in tokendata if c.isalpha()] #去除non-alphanumeric
15
16    #lowercasing
17    lowerString = [elements.lower() for elements in regex]
18
19    #stemming
20    stemmer = PorterStemmer()
21    stemString = [stemmer.stem(lowerstring) for lowerstring in lowerString]
22
23    #stopword removal
24    file2 = open("stopwordlist.txt","r")
25    stopwordList = file2.read()
26    file2.close()
27    removeStopword = [word for word in stemString if word not in stopwordList]
28
29    return removeStopword #回傳tokenize後該文章的term (字串)
30
```

```
168 dict_doc = {} #儲存tokenize後的1095文章的term
169 for i in range(1,1096): #將1095個doc做tokenize
170     dict_doc[str(i)] = tokenize(str(i)) #key為docid,value為該doc的所有term (i和docid轉成str)
171
```

- 讀取 training.txt 得到各 class 及對應的 training doc, 並將 training data 做 tokenize! 一樣用 dictionary 結構儲存, dict\_class 儲存 traing.txt, key 為 classid, value 為 15 個 docid; dict\_train\_doc 儲存 tokenize 後的 training data, key 為 docid, value 為 terms。

```

172 #讀取training.txt，得到各class及對應的training docid
173 dict_class = {} # key為classid,value為docid_list
174 with open('training.txt','r') as f:
175     for line in f:
176         (classID,docID_list) = line.split(' ',1)
177         docID_set = docID_list.split() #去除docid中最後面\n
178         dict_class[classID] = docID_set
179
180 #tokenize training document
181 dict_train_doc = {} #key為docid,value為terms
182 for key,value in dict_class.items(): # key為classid,value為docid_list
183     for docid in value:
184         tokens = tokenize(docid)
185         dict_train_doc[docid] = tokens
186

```

- 除去 training data,剩下的為 testing data,也做 tokenize 用 dict 儲存，key 為 docid, value 為 terms。

```

187 #去除train doc，剩下1095-195=900個test doc
188 dict_test_doc = dict_train_doc.copy() #複製1095個doc
189 for classid,docid_list in dict_class.items():
190     for docid in docid_list:
191         del dict_test_doc[docid] #刪除在train data的doc，留下來都當test data
192

```

- 將 training data 做 feature selection，計算每個 term 的數值，用 dict 儲存，key 為 term, value 為 feature selection 值

```

193 #feature selection
194 #計算每個train doc的term的feature selection(chi-square, likelihood ratios, EMI)
195 dict_feature_selection = {} #key為term,value為feature selection值
196 for terms_list in dict_train_doc.values(): #key為docid,value為terms，取得每一個term
197     for term in terms_list:
198         feature_selection = contingencyTable(term,dict_class,dict_train_doc)
199         if(dict_feature_selection.get(term)): #有重複的term，跳過
200             continue
201         dict_feature_selection[term] = feature_selection
202

```

#### ◆ 先算出 contingency table

有些值像是全部文章數為 1095 篇, 13 個 class 每個 class 全部文章為 15 篇, 非該 class 的文章為  $195 - 15 = 180$  篇; 其他皆以 list 來儲存個別 class 的每個 term 的 present 和 absent 儲存，以及 on class 和 not on class 也儲存。後半則以選擇 feature selection 方法來做呼叫。

```

77
78 def contingencyTable(term,dict_class,dict_train_doc): #計算每個term在13個class各別的chi-square
79 #計算個chi squares需要的值 (c-p,c-a,notc-p,notc-a,c-all,p-all,doc_all)
80 doc_all = 195
81 c_all = 15
82 notc_all = 180
83
84 #計算在每一個class各出現幾個doc(df) [class-present]
85 c_p = []
86 for classid,docid_list in dict_class.items(): #key為classid,value為docid和terms
87     count_df = 0
88     for docid in docid_list: #取每一個docid
89         if term in dict_train_doc[docid]:#看terms list有沒有含這個term,有就代表這個class中
90             count_df += 1
91     c_p.append(count_df) #這個term在這個class的df值
92
93 c_a = [] #[class-absent]
94 notc_p = [] #[not class-present]
95 notc_a = [] #[not class-absent]
96
97 for i in range(0,13): #第幾個class
98     c_a.append(15-c_p[i]) #c_all
99
100     count_notc_p = 0 #計算各個class中not c-present的值
101     for j in range(0,13): #加總not c-present總合
102         if(j == i):
103             continue;
104         count_notc_p += c_p[j]
105     notc_p.append(count_notc_p)
106     notc_a.append(notc_all - notc_p[i])
107
108 #計算各別class的p_all,a_all
109 p_all = []
110 a_all = []
111 for i in range(0,13):
112     p_all.append(c_p[i]+notc_p[i])
113     a_all.append(c_a[i]+notc_a[i])
114
115 #chi-square: 計算該term在全部class的4個chi square值,四個chi-square值加總,c_all,notc_all一樣
116 #得到這個term在13個class中最大的chi-square值
117 # chi_square_sum = countChiSquare(c_p,c_a,notc_p,notc_a,c_all,notc_all,p_all,a_all,doc_all)
118
119 #likelihood-ratios:計算該term在全部class的likelihood-ratios,得到最大的
120 likelihood_ratios_sum = countLikelihoodRatios(c_p,c_a,notc_p,notc_a,doc_all)
121
122 #expected-mutual-information(EMI):計算該term在全部class的4個EMI值,加總後回傳在13個class中最
123 # expected_mutual_information_sum = countExpectedMutualInformation(c_p,c_a,notc_p,notc_a,doc_all)
124
125 # return chi_square_sum
126 return likelihood_ratios_sum
127 # return expected_mutual_information_sum
128

```

## ◆ Chi-square feature selection

```

30
31 def countChiSquare(c_p,c_a,notc_p,notc_a,c_all,notc_all,p_all,a_all,doc_all): #計算該term在
32     chi_square_list = []
33     chi_square_sum = 0 #總和
34     for i in range(0,13): #計算每一個class的chi-square
35         list_p = [p_all[i],a_all[i],p_all[i],a_all[i]]
36         list_c = [c_all,c_all,notc_all,notc_all]
37         list_observed = [c_p[i],c_a[i],notc_p[i],notc_a[i]]
38         for j in range(0,4): #四個chi-square加總
39             observed_frequency = list_observed[j]
40             expected_count = doc_all*(list_p[j]/doc_all)*(list_c[j]/doc_all)
41             chi_square = ((observed_frequency- expected_count)**2)/expected_count
42             chi_square_sum += chi_square
43         chi_square_list.append(chi_square_sum)
44
45     return max(chi_square_list) #回傳最高的chi-square值
46

```

## ◆ Likelihood ratios feature selection

```

def countLikelihoodRatios(c_p,c_a,notc_p,notc_a,doc_all): #計算likelihood ratios
    likelihood_ratios_list = []
    likelihood_ratios_sum = 0
    for i in range(0,13):
        upcount = (((c_p[i]+notc_p[i])/doc_all)**c_p[i]) * ((1-((c_p[i]+notc_p[i])/doc_all)
        downcount = (((c_p[i]/(c_p[i]+c_a[i]))**c_p[i]) * ((1-(c_p[i]/(c_p[i]+c_a[i])))**c_
        likelihood_ratios_sum = (-2) * math.log(upcount/downcount)
        likelihood_ratios_list.append(likelihood_ratios_sum)
    return max(likelihood_ratios_list)

```

## ◆ Expected mutual information

```

57 def countExpectedMutualInformation(c_p,c_a,notc_p,notc_a,c_all,notc_all,p_all,a_all,doc_al
58 expected_mutual_information_list = []
59 expected_mutual_information_sum = 0
60 for i in range(0,13):
61     list_p = [p_all[i],a_all[i],p_all[i],a_all[i]]
62     list_c = [c_all,c_all,notc_all,notc_all]
63     list_observed = [c_p[i],c_a[i],notc_p[i],notc_a[i]]
64     for j in range(0,4):
65         # print("j:"+str(j)+"\n")
66         upcount = list_observed[j]/doc_all
67         # print("upcount:"+str(upcount)+"\n")
68         downcount = (list_p[j]/doc_all) * (list_c[j]/doc_all)
69         # print("downcount:"+ str(downcount)+"\n")
70         if(upcount == 0) : #不知為何有list_observed(c_p)為0的情況，會導致log0情況，所以直接加0
71             expected_mutual_information_sum += 0
72         else:
73             expected_mutual_information_sum += (list_observed[j]/doc_all) * math.log(
74             expected_mutual_information_list.append(expected_mutual_information_sum)
75
76 return max(expected_mutual_information_list)
77

```

以上述任一種 feature selection 方法，計算出每個 terms 在各 13 個 class 的數值，選擇最大的數值回傳，以做後續的比較。

- 計算出 term 的數值後，排序 dict，取前 500 個數值高 terms 數來做 vocabulary。(目前發現 terms 取 150 個預測準確度最高)

```

203 #取前500個chi-square大的term
204 feature_selection_list = []
205 count = 1
206 for key,value in sorted(dict_feature_selection.items(), key = lambda x:x[1],reverse=True):
207     # print("%s %s\n" % (key,value))
208     if(count > 150): #500,450,430,400,350,300,200,150,125,100;目前150跑出來最好
209         break
210     feature_selection_list.append(key)
211     count += 1
212

```

- 用這些 feature selection 後的 term，來將 training data 和 testing data 做過濾，只剩下有出現在 500 個 feature selection terms 裡的，排除沒看過的字。

```

213 #過濾train data和test data, 只剩這500個term
214 dict_train_doc_filter = {} #key為docid,value為terms
215 dict_test_doc_filter = {} #key為docid,value為terms
216 for docid,terms in dict_train_doc.items(): #將train data過濾
217     filter_terms = [term for term in terms if term in feature_selection_list] #只留下在fea
218     dict_train_doc_filter[docid] = filter_terms
219
220 for docid,terms in dict_test_doc.items(): #將test data過濾
221     filter_terms = [term for term in terms if term in feature_selection_list] #只留下在fea
222     dict_test_doc_filter[docid] = filter_terms
223

```

➤ 接著以 multinomial Naïve Bayes model 來做分類

◆ Training phase, 算出 condprob, prior\_c

```

129 def trainMultinomialNB(dict_class,dict_train_doc_filter,feature_selection_list): #Multinom
130     Nc = 15 # 每個class有15個train doc
131     N = 195 # 13個class, 有13*15 = 195個train doc
132     prior_c = [] #儲存每個class的P(c)
133     prior_c.insert(0,0) #index從1開始
134     condprob = defaultdict(dict) #two dimensional dict
135
136     for classid,docid_list in dict_class.items():
137         prior_c.insert(int(classid),(Nc/N)) #P(c) = Nc / N
138         text_c = 0 #該class所有terms數(textc)
139         # text_term = 0 #每個term在該class的doc中總共出現幾次(Tct)
140         dict_text_term = {}
141         for term in feature_selection_list: #初始化
142             dict_text_term[term] = 0
143
144         for docid in docid_list: #取得每個class的docid list, 得每一個doc
145             text_c += len(dict_train_doc_filter[docid]) #計算這個class總terms數
146             #計算每個term in V(feature selection的500個terms)在這個class中所有doc裡出現幾次
147             for term in dict_train_doc_filter[docid]:
148                 dict_text_term[term] += 1
149
150         for term,frequency in dict_text_term.items():
151             condprob[term][classid] = (frequency+1)/(text_c+len(dict_text_term.keys())) #
152
153     return condprob,prior_c
154

```

◆ Testing phase, 丟入 testing data, 為每一個 test doc 做分類, 得出它在 13 個 class 中數值最高, 便把它分配在這個 class。

```

155 def ApplyMultinomialNB(dict_class,test_data,condprob,prior_c): # multinomial model in test
156     #判斷該test doc屬於哪一個class
157     score = [] #儲存term在每一個class的score,score = logP(c) + logP(X = t | c)
158     score.insert(0,0) #index從1開始
159     for classid,docid_list in dict_class.items():
160         score.insert(int(classid),math.log(prior_c[int(classid)]))
161         for term in test_data: #該test doc的term, 在這個class的分數加總, 就是這個doc屬於這個class
162             score[int(classid)] += math.log(condprob[term][classid])
163     del score[0] #刪除為0的 (index恢復到從0開始)
164     return score.index(max(score))+1 #最大的值, 代表這個class (index為0開始, 所以要加一)
165

```

◆ 最後回傳 test docid 和對應的 classid, 以 dict 儲存, key 為 docid, value 為 classid。



```

224 #分類
225 condprob,prior_c = trainMultinomialNB(dict_class,dict_train_doc_filter,feature_selection_l
226 #test data，算出每一個doc屬於哪一個class
227 dict_answer = {}
228 for docid,terms in sorted(dict_test_doc_filter.items()):
229     doc_class = ApplyMultinomialNB(dict_class,terms,condprob,prior_c) #得到該doc屬於哪個clas
230     dict_answer[int(docid)] = int(doc_class)
231

```

➤ 最後將答案輸出成 csv 檔

```

232 #將答案寫入成csv檔案
233 # with open('answer_chisquare.csv','w',newline='') as csvfile:
234 with open('answer_likelihoodratios.csv','w',newline='') as csvfile:
235 # with open('answer_expectedmutualinformation.csv','w',newline='') as csvfile:
236     #建立csv檔寫入器
237     writer = csv.writer(csvfile)
238
239     #第一行：id,Value
240     writer.writerow(['Id','Value'])
241     #第二行開始輸入答案，格式為docid,classid
242     for docid,classid in sorted(dict_answer.items()):
243         writer.writerow([str(docid),str(classid)])
244

```