

Programming Assignment 3

執行環境: Visual Studio Code

程式語言: Python 3.11.5

作業架構:

R12725048

- |— report.pdf
- |— training.txt
- |— pa3.py
- |— data/

執行方式

- 使用VS code跑pa3.py檔
- 需要下載的套件有:
 - `pip install stop_words`: 刪除不太帶有資訊的單詞所需
 - `pip install nltk`: 使用Porter' s algorithm.所需
- 直接按全部執行即可

處理邏輯

- Step 1 : import 所需套件

```
1  ✓ import math
2  from nltk.stem import PorterStemmer
3  import re
4  from stop_words import get_stop_words
5  import csv
```

處理邏輯

- Step 2 : 根據HW2 一樣先依據上次的news dataset建立class TF_TFIDF()，並建立其 object : tf_idf
 - 此步驟會做資料的前處理、建好dictionary、tf-idf table
 - 皆是上次作業的內容，因此簡報就不多花篇幅敘述

```
211     # step1:藉由HW2先根據所有data建立dictionary&tfidf table
212     file_path = "./IRTM"
213     num_docs = 1095
214     tf_idf = TF_IDF(file_path, num_docs)
```

處理邏輯

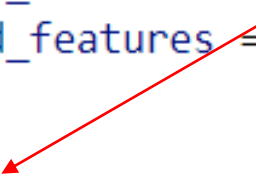
- Step 3 : 讀取training.txt，並將內容分別儲存成以下形式
- training_dataset = [1,2,3,...700,730,...,1019] → 將所有要拿來訓練的文章id儲存起來並照id大小排序
- training_list = {1:[11,19,29,113,...],2:[1,2,3,...],...} → 用dict儲存每個class各別需要訓練的doc有哪些
- classes = [1,2,3,...13] → 用list紀錄所有的類別

```
217 training_dataset = [] #用list來儲存所有要拿來training的doc id
218 training_list = {} #用dict分別儲存每個class中有那些training的doc id
219 classes = [] #用list來儲存class有哪些
220 f = open('training.txt')
221 for line in f.readlines():
222     input = line.split(' ')
223     if '\n' in input:
224         input.remove('\n')
225     input = [int(i) for i in input]
226     training_list[input[0]] = input[1:]
227     training_dataset.extend(input[1:])
228     classes.append(input[0])
229 f.close()
230 training_dataset.sort()
```

處理邏輯

- 採取Log likelihood ratio的方式來挑features
- 而又因為作業為Multiple Classifiers，所以使用select the top k/n features for each n classifiers的方式來挑選500個features
- Step 4(1): 重要features挑選
 - 要把最後的features union起來，因為不同類別可能挑到相同的feature

```
232 # step2: 根據training dataset先做重要term的篩選
233 selected_features = FeaturesSelection(training_list, classes, int(500/13))
234 selected_features = set(selected_features)
```



```
198 def FeaturesSelection(training_list:dict, classes:list, k:int) -> list:
199     vocabulary = []
200     score_list = {} #用來儲存每個class中的每個term的llr score e.g.score_list={'apple':0.025,'banana':3.521,...}
201     #針對各個class各挑k個重要features最後合併起來
202     for c in classes:
203         V = ExtractVocabulary(training_list[c])
204         score_list = likelihood_ratio(c, V)
205         #將score_list依照values(llr score)排序
206         sorted_score = dict(sorted(score_list.items(), key=lambda item: item[1], reverse=True))
207         #每個class挑最高分的k個features合併起來
208         vocabulary.extend(list(sorted_score.keys())[:k])
209
210     return vocabulary
```

處理邏輯


- Step 4(2):
likelihood ratio

```
171 def likelihood_ratio(c:int, V:list) -> dict:
172     score = {}
173     for term in V:
174         n11 = 0
175         n01 = 0
176         n10 = 0
177         n00 = 0
178         for doc in training_dataset:
179             #若此doc為此類別(on topic)
180             if doc in training_list[c]:
181                 #若term為present
182                 if term in ExtractTokensFromDoc(doc):
183                     n11+=1
184                 #若term為absent
185                 else:
186                     n10+=1
187             #若此doc為不屬於此類別(off topic)
188             else:
189                 #若term為present
190                 if term in ExtractTokensFromDoc(doc):
191                     n01+=1
192                 #若term為absent
193                 else:
194                     n00+=1
195         pt = (n11+n01)/len(training_dataset)
196         p1 = n11/(n11+n10)
197         p2 = n01/(n01+n00)
198         H1_likelihood = (math.pow(pt,n11)) * (math.pow((1-pt),n10)) * (math.pow(pt,n01)) * (math.pow((1-pt),n00))
199         H2_likelihood = (math.pow(p1,n11)) * (math.pow((1-p1),n10)) * (math.pow(p2,n01)) * (math.pow((1-p2),n00))
200         llr = (-2)*(math.log(H1_likelihood)-math.log(H2_likelihood))
201         score[term] = llr
202     return score
```

處理邏輯

- Step 5:建立class NaiveBayes()，並建立其object: NB，將NB的attributes&methods初始化

```
239 # step3: 用Naive Bayes進行training
240 NB = NaiveBayes(tf_idf, classes, training_dataset, training_list)
241 prior, condprob = NB.TrainMultinomialNB(selected_features)
```



```
82 class NaiveBayes():
83     def __init__(self, tf_idf:TF_IDF, classes:list, training_dataset:list, training_list:dict):
84         self.tf_idf = tf_idf
85         self.classes = classes
86         self.training_dataset = training_dataset
87         self.training_list = training_list
88
```


處理邏輯

- Step 6: 進行training

```
239 # step3: 用Naive Bayes進行training
240 NB = NaiveBayes(tf_idf, classes, training_dataset, training_list)
241 prior, condprob = NB.TrainMultinomialNB(selected_features)
```

```
89 def TrainMultinomialNB(self, V:list) -> list | dict :
90     #宣告prior陣列存放各類別的P(c)值
91     prior = [0 for i in range(0, len(self.classes)+1)]
92     condprob = {}
93     N = len(self.training_dataset)
94
95     for c in self.classes:
96         Nc = CountDocsInClass(c)
97         prior[c] = Nc/N
98         text_c = ConcatenateTextOfAllDocsInClass(self.training_list[c])
99
100         num_of_term_in_class_c = sum(list(text_c.values()))
101
102         for term in V:
103
104             T_ct = CountTokensOfTerm(text_c, term)
105             if term not in condprob:
106                 condprob[term] = [1/(num_of_term_in_class_c+len(V)) for i in range(0, len(self.classes)+1)]
107
108             condprob[term][c] = (T_ct+1)/(num_of_term_in_class_c+len(V))
109
110     return prior, condprob
```

紀錄每個term在每個class的分數，並做smoothing處理

處理邏輯

```
249 # step4: 將剩下的資料集進行apply
250 mapping_class_of_doc = [0 for i in range(1, num_docs+2)]
251 for doc in range(1, num_docs+1):
252     mapping_class_of_doc[doc] = NB.ApplyMultinomialNB(selected_features, classes, prior, condprob, doc)
```

- Step 7: 將所有dataset拿去驗證

```
111
112 def ApplyMultinomialNB(self, V:list, categories:list, prior:list, condprob:dict, d:int) -> int:
113     word_in_d = ExtractTokensFromDoc(d)
114     score = [0 for i in range(0, len(categories)+1)]
115     max_score = -100000000
116     mapping = 0
117     for c in categories:
118         score[c] = math.log(prior[c])
119         for term in word_in_d:
120             #只計算selected_features裡的字，其他不用算分數
121             if term in V:
122                 score[c] += math.log(condprob[term][c]) * tf_idf.tf[d][term]
123             #紀錄最高分的類別是哪個
124             if score[c] > max_score:
125                 max_score = score[c]
126                 mapping = c
127     return mapping
```

處理邏輯

- Step 8: 將result寫到csv檔

```
256 # step5: write result to 'HW3_csv.csv'
257 header = ['Id','Value']
258 data = []
259 for i in range(1,num_docs+1):
260     if i not in training_dataset:
261         data.append([i,mapping_class_of_doc[i]])
262
263
264 with open('HW3_csv.csv', 'w', encoding='UTF8', newline='') as f:
265     writer = csv.writer(f)
266
267     # write the header
268     writer.writerow(header)
269
270     # write multiple rows
271     writer.writerows(data)
```

Kaggle 分數

8

r12725048



0.99111

14

13h



Your Best Entry!

Your most recent submission scored 0.99111, which is an improvement of your previous score of 0.99000. Great job!

Tweet this