

```

1  #!/usr/bin/python
2  # -*- coding: utf-8 -*-
3  from sklearn.datasets import fetch_20newsgroups
4  from sklearn.feature_extraction.text import CountVectorizer
5  from sklearn.feature_extraction.text import TfidfVectorizer
6
7  import nltk
8  from nltk.corpus import stopwords
9  from nltk.tokenize import word_tokenize
10 from nltk.tokenize import RegexpTokenizer
11 from nltk.tokenize import sent_tokenize
12
13 import xml.etree.ElementTree as ET
14
15 import json
16
17 import os
18 import glob
19 import string
20 import re
21
22 #####
23
24 stop_words = set(stopwords.words('english'))
25
26 if not os.path.exists('Dataset_original'):
27     os.makedirs('Dataset_original')
28 originalDS = 'Dataset_original/'
29
30 if not os.path.exists('Dataset_as_txt'):
31     os.makedirs('Dataset_as_txt')
32 asTextDS = 'Dataset_as_txt/'
33
34 if not os.path.exists('Dataset_preprocessed'):
35     os.makedirs('Dataset_preprocessed')
36 preprocDS = 'Dataset_preprocessed/'
37
38 if not os.path.exists('Keyphrases_golden_set'):
39     os.makedirs('Keyphrases_golden_set')
40 goldenSet = 'Keyphrases_golden_set/'
41
42 if not os.path.exists('Keyphrases_experimental'):
43     os.makedirs('Keyphrases_experimental')
44 exprSet = 'Keyphrases_experimental/'
45
46 if not os.path.exists('Metrics'):
47     os.makedirs('Metrics')
48 metricsF = 'Metrics/'
49
50 #####
51 #
52 # XML/TREE OPERATIONS
53 #
54 #####
55
56 def get_root_from_xml(file):
57     return ET.parse(file).getroot()
58
59
60 def from_root_get_document(root):

```

```

61     return list(root)[0]
62
63
64 def from_document_get_sentences(doc):
65     return list(doc)[0]
66
67
68 def from_sentences_get_nth_sentence(sentences, n):
69     return list(sentences)[n]
70
71
72 def from_sentence_get_tokens(sentence):
73     return list(sentence)[0]
74
75
76 def from_tokens_get_full_sentence(tokens, n):
77     sentence = []
78     for i in tokens:
79         sentence.append(i[0].text)
80     return sentence
81
82
83 #####
84 def from_xml_to_txt():
85     for filename in glob.glob(originalDS + '*.xml'):
86         with open(filename, 'r', encoding='utf-8') as document:
87             noPath = os.path.basename(filename)
88             asTxt = os.path.splitext(noPath)[0] + '.txt'
89
90             open(asTextDS + asTxt, 'w', encoding='utf-8').close()
91             with open(asTextDS + asTxt, 'a+', encoding='utf-8') as out:
92                 root = get_root_from_xml(document)
93                 doc = from_root_get_document(root)
94                 sentences = from_document_get_sentences(doc)
95                 for i in range(len(sentences)):
96                     sentence = \
97                         from_sentences_get_nth_sentence(sentences, i)
98                     tokens = from_sentence_get_tokens(sentence)
99                     for j in range(len(tokens)):
100                         full_sentence = \
101                             from_tokens_get_full_sentence(tokens, j)
102                     for u in full_sentence:
103                         if u != '.':
104                             out.write(u + ' ')
105                         else:
106                             out.write(u + '\n')
107
108 def preproc_all_txt():
109     for filename in glob.glob(asTextDS + '*.txt'):
110         noPath = os.path.basename(filename)
111         out = preprocDS + noPath
112
113         preprocess_file(filename, out)
114
115 def parse_json():
116     with open("test.reader.json", 'r', encoding='utf-8') as fjs:
117         d = json.load(fjs)
118         for key, value in d.items():
119             open(goldenSet + key + ".txt", 'w', encoding='utf-8').close()
120             with open(goldenSet + key + ".txt", 'a+', encoding='utf-8') as gkps:

```

```

121         for i in value:
122             if len([word for word in i[0].split()]) in [1, 2, 3]:
123                 gkps.write(i[0] + "\n")
124
125 #####
126 def preprocess_file(docname, outfile_name):
127     with open(docname, 'r', encoding='utf-8') as file:
128         outfile = open(outfile_name, 'w', encoding='utf-8')
129
130         for line in file:
131             print(preprocess_sentence(line), end='', file=outfile)
132         outfile.close()
133
134     return outfile_name
135
136
137 def preprocess_list(inlist, outfile_name):
138     # only run this once for 20news group to save the preprocessed DS in memory
139     i = 1
140     for article in inlist:
141         sentences = sent_tokenize(article)
142         outfile = open(outfile_name + str(i) + '.txt', 'w',
143                        encoding='utf-8')
144         for sentence in sentences:
145             preprocessed_sentence = preprocess_sentence(sentence)
146             print(preprocessed_sentence, end='', file=outfile)
147         outfile.close()
148         i += 1
149     return (outfile_name, i)
150
151
152 def preprocess_sentence(sentence):
153     processed = sentence.lower()
154     processed = re.sub(r'\d+', '', processed)
155     tokenizer = RegexpTokenizer(r'\w+')
156     tokens = tokenizer.tokenize(processed)
157     #filtered_words = [w for w in tokens if not w in stop_words]
158     filtered_words = [w for w in tokens]
159     if sentence[-1] == '\n':
160         return ' '.join(filtered_words) + ' '
161     return ' '.join(filtered_words)
162
163 #####
164 def ngrams(docname, low, high):
165     with open(docname, 'r', encoding='utf-8') as document:
166         result = []
167
168         c_vec = CountVectorizer(ngram_range=(low, high))
169         ngrams = c_vec.fit_transform(document)
170         vocab = c_vec.vocabulary_
171         count_values = ngrams.toarray().sum(axis=0)
172
173         for ngram in sorted([[count_values[i], k] for (k, i) in vocab.items()],
174                             reverse=True):
175             result.append(ngram)
176
177     return result
178
179 def idf(dataset, low, high):

```

```

180     # assuming dataset is already a list (like train.data)
181
182     vectorizer = TfidfVectorizer(strip_accents='unicode', ngram_range=(low, high))
183
184     vectorizer.fit_transform(dataset)
185
186     idf = vectorizer.idf_
187
188     return dict(zip(vectorizer.get_feature_names(), idf))
189
190 #####
191
192 def rates(gs, exp):
193     tp = 0
194     fp = 0
195     fn = 0
196     tp_at_5 = 0
197
198     for i in exp:
199         if i in gs:
200             tp += 1
201         if i not in gs:
202             fp += 1
203
204     for i in gs:
205         if i not in exp:
206             fn += 1
207
208     for i in exp[:5]:
209         if i in gs:
210             tp_at_5 += 1
211
212     return tp, fp, fn, tp_at_5
213
214 def precision(tp, fp):
215     return 0 if tp + fp == 0 else tp/(tp+fp)
216
217 def recall(tp, fn):
218     return 0 if tp + fn == 0 else tp/(tp+fn)
219
220 def f_one(r, p):
221     return 0 if r + p == 0 else 2 * r * p / (r + p)
222
223 def precision_at_n(tp_at_n, n):
224     return tp_at_n/n
225
226 def mean_average_precision(golden, experimental):
227     precisions = []
228     for i in range(1, len(experimental)+1):
229         tp, fp, _, _ = rates(golden, experimental[:i])
230         precisions.append(precision(tp, fp))
231     meavp = sum(precisions)/len(precisions)
232     return meavp
233
234 def metrics(golden, experimental):
235     with open(golden, 'r', encoding='utf-8') as gs_file, open(experimental, 'r',
236                                                                encoding='utf-8') as exp_file:
237         gs = []
238         exp = []
239         for line in gs_file:

```

```

239         gs.append(line)
240     for line in exp_file:
241         exp.append(line)
242
243     tp, fp, fn, p_at_5 = rates(gs, exp)
244
245     prec = precision(tp, fp)
246
247     rec = recall(tp, fn)
248
249     p_at_5 = p_at_5/5
250
251     f1 = f_one(rec, prec)
252
253     mav = mean_average_precision(gs, exp)
254
255     file = metricsF + os.path.basename(golden)
256     s1 = "Precision: " + str(prec) + "\n" + "Recall: " + str(rec) + "\n" + "F1: "
+ str(f1) + "\n" + "P@5: " + str(p_at_5) + "\n" + "MAV: " + str(mav)
257
258     open(file, 'w', encoding = 'utf-8').close()
259     with open(file, 'w', encoding = 'utf-8') as f:
260         f.write(s1)
261
262
263     return prec, rec, f1, p_at_5, mav
264
265
266
267 #####
268 class Dataset:
269     # models a dataset
270     def __init__(self, directory):
271         self.directory = directory
272         self.asList = list()
273         self.idf = None
274
275     def from_files_to_list(self):
276         for filename in glob.glob(self.directory + '*.txt'):
277             with open(filename, 'r', encoding='utf-8') as document: # utf-8 is
breaking at (spanish n)
278                 string = document.read()
279                 self.asList = self.asList + [string]
280
281     def compute_idf(self, low, high):
282         self.idf = idf(self.asList, low, high)
283
284
285
286 class Document:
287     # models a processed doc
288     def __init__(self, filename):
289         self.file = filename
290         self.ngrams = None # candidates
291         self.keyphrases = list()
292
293     def compute_ngrams(self, low, high):
294         self.ngrams = ngrams(self.file, low, high)
295
296

```

```

297     def get_keyphrases(self, ds_idf):
298         for phrase in self.ngrams:
299             tf = phrase[0]
300             idf = ds_idf[phrase[1]]
301             tfidf_with_len = tf * idf * len(phrase[1].split())
302             self.keyphrases.append((phrase[1], tfidf_with_len))
303
304     # Class helpers
305     def write_keyphrases(self, n):
306         keyphrases = sorted(self.keyphrases, key=lambda x: x[1], reverse=True)[:n]
307         open(exprSet + os.path.basename(self.file), 'w', encoding = 'utf-8').close()
308         with open(exprSet + os.path.basename(self.file), 'a+', encoding = 'utf-8') as
out:
309             for i in keyphrases:
310                 out.write(i[0] + '\n')
311
312
313     def get_preprocessed_text(self):
314         with open(self.file, 'r', encoding='utf-8') as document:
315             string = document.read()
316             return string
317
318 #####
319
320 def main():
321     #test.reader.json and dataset_original needed, others are created
322     #from_xml_to_txt() #only #1 time
323     #preproc_all_txt() #only #1 time
324     #parse_json() #only #1 time
325
326     ds = Dataset(preprocDS)
327     ds.from_files_to_list()
328     ds.compute_idf(1, 3)
329     for file in glob.glob(preprocDS + '*.txt'):
330         doc = Document(file)
331         doc.compute_ngrams(1, 3)
332         doc.get_keyphrases(ds.idf)
333         doc.write_keyphrases(10)
334
335     precision_sum = 0
336     recall_sum = 0
337     f1_sum = 0
338     p_at_5_sum = 0
339     mav_sum = 0
340     total_files = 0
341
342     for file in glob.glob(goldenSet + '*.txt'):
343         total_files += 1
344         p, r, f1, p_at_5, mav = metrics(file, exprSet + os.path.basename(file))
345         precision_sum += p
346         recall_sum += r
347         f1_sum += f1
348         p_at_5_sum += p_at_5
349         mav_sum += mav
350
351
352     print("Average Precision: " + str(precision_sum/total_files))
353     print("Average Recall: " + str(recall_sum/total_files))
354     print("Average F1: " + str(f1_sum/total_files))
355     print("Average P@5: " + str(p_at_5_sum/total_files))

```

```
356     print("Average MAV: " + str(mav_sum/total_files))
357
358
359 if __name__ == '__main__':
360     main()
```