

```

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 from math import log
14
15 import xml.etree.ElementTree as ET
16
17 import json
18
19 import os
20 import glob
21 import string
22 import re
23 import itertools, nltk, string
24
25 #####
26
27 stop_words = set(stopwords.words('english'))
28
29 if not os.path.exists('Dataset_original'):
30     os.makedirs('Dataset_original')
31 originalDS = 'Dataset_original/'
32
33 if not os.path.exists('Dataset_as_txt'):
34     os.makedirs('Dataset_as_txt')
35 asTextDS = 'Dataset_as_txt/'
36
37 if not os.path.exists('Dataset_preprocessed'):
38     os.makedirs('Dataset_preprocessed')
39 preprocDS = 'Dataset_preprocessed/'
40
41 if not os.path.exists('Keyphrases_golden_set'):
42     os.makedirs('Keyphrases_golden_set')
43 goldenSet = 'Keyphrases_golden_set/'
44
45 if not os.path.exists('Keyphrases_experimental'):
46     os.makedirs('Keyphrases_experimental')
47 exprSet = 'Keyphrases_experimental/'
48
49 if not os.path.exists('Metrics'):
50     os.makedirs('Metrics')
51 metricsF = 'Metrics/'
52
53 #####
54 #
55 # XML/TREE OPERATIONS
56 #
57 #####
58
59 def get_root_from_xml(file):
60     return ET.parse(file).getroot()

```

```

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

```

```

121     for key, value in d.items():
122         open(goldenSet + key + ".txt", 'w', encoding='utf-8').close()
123         with open(goldenSet + key + ".txt", 'a+', encoding='utf-8') as gkps:
124             for i in value:
125                 if len([word for word in i[0].split()]) in [1, 2, 3]:
126                     gkps.write(i[0] + "\n")
127
128 #####
129 def preprocess_file(docname, outfile_name):
130     with open(docname, 'r', encoding='utf-8') as file:
131         outfile = open(outfile_name, 'w', encoding='utf-8')
132
133         for line in file:
134             print(preprocess_sentence(line), end='', file=outfile)
135         outfile.close()
136
137     return outfile_name
138
139
140 def preprocess_list(inlist, outfile_name):
141     # only run this once for 20newsgroup to save the preprocessed DS in memory
142     i = 1
143     for article in inlist:
144         sentences = sent_tokenize(article)
145         outfile = open(outfile_name + str(i) + '.txt', 'w',
146                       encoding='utf-8')
147         for sentence in sentences:
148             preprocessed_sentence = preprocess_sentence(sentence)
149             print(preprocessed_sentence, end='', file=outfile)
150         outfile.close()
151         i += 1
152     return (outfile_name, i)
153
154
155 def preprocess_sentence(sentence):
156     processed = sentence.lower()
157     processed = re.sub(r'\d+', '', processed)
158     tokenizer = RegexpTokenizer(r'\w+')
159     tokens = tokenizer.tokenize(processed)
160     #filtered_words = [w for w in tokens if not w in stop_words]
161     filtered_words = [w for w in tokens]
162     if sentence[-1] == '\n':
163         return ' '.join(filtered_words) + ' '
164     return ' '.join(filtered_words)
165
166 #####
167 def select_candidates(text):
168     grammar = r'KT: {( <JJ>* <NN.*>+ <IN>)? <JJ>* <NN.*>+ }'
169
170     stop_words = set(nltk.corpus.stopwords.words('english'))
171
172     chunker = nltk.chunk.regexp.RegexpParser(grammar)
173     tagged_sents = nltk.pos_tag_sents(nltk.word_tokenize(sent) for sent in
nltk.sent_tokenize(text))
174
175     all_chunks =
list(itertools.chain.from_iterable(nltk.chunk.tree2conlltags(chunker.parse(tagged_sen
t))
176                                     for tagged_sent in tagged_sents))
177

```

```

178     candidates = [' '.join(word for word, pos, chunk in group).lower()
179                   for key, group in itertools.groupby(all_chunks, lambda
wordposchunk: wordposchunk[2] != '0') if key]
180
181     result = [cand for cand in candidates
182               if cand not in stop_words and len(cand.split()) < 4]
183
184     result = list(dict.fromkeys(result))
185
186     #if len(result) < 10:
187         #grammar = r'KT: {<DT>? <JJ>* (<NN>|<NP>|<PRN>)+}'
188
189         #chunker = nltk.chunk.regexp.RegexpParser(grammar)
190
191         #all_chunks =
list(itertools.chain.from_iterable(nltk.chunk.tree2conlltags(chunker.parse(tagged_sen
t))
192                                     #for tagged_sent in
tagged_sents))
193         #candidates = [' '.join(word for word, pos, chunk in group).lower()
194                       #for key, group in itertools.groupby(all_chunks, lambda
wordposchunk: wordposchunk[2] != '0') if key]
195
196         #result = [cand for cand in candidates
197                   #if cand not in stop_words and len(cand.split()) < 4 and not all(char
in punct for char in cand)]
198
199         return result
200
201 def ngrams(docname, low, high):
202     with open(docname, 'r', encoding='utf-8') as document:
203         result = []
204
205         c_vec = CountVectorizer(ngram_range=(low, high))
206         ngrams = c_vec.fit_transform(document)
207         vocab = c_vec.vocabulary_
208         count_values = ngrams.toarray().sum(axis=0)
209
210         for ngram in sorted([(count_values[i], k) for (k, i) in vocab.items()],
reverse=True):
211             result.append(ngram)
212
213     return result
214
215
216 def n_t(list1, list2):
217     c=[]
218     for i in range(0,len(list1)):
219         n=0
220         for j in range(0,len(list2)):
221             if (list1[i] in list2[j]) == True:
222                 n=n+1
223         c.append((n,list1[i]))
224     return c
225
226 def frequency(list1, doc1):
227     c=[]
228     for i in range(0,len(list1)):
229         c.append(doc1.count(list1[i]))
230     return c

```

```

231
232 def avgdl(list1):
233     s=0
234     for i in range(0,len(list1)):
235         s = s+len(list1[i])
236     total = s/len(list1)
237     return total
238
239 def bm25(candidates, dataset, doc):
240     N = len(dataset)
241     D = len(doc)
242     k1 = 1.2
243     b = 0.75
244     avgdl1 = avgdl(dataset)
245     frequency1 = frequency(candidates, doc)
246     n_t1 = n_t(candidates, dataset)
247     final = dict()
248     for i in range(0,len(candidates)):
249         score = log((N - n_t1[i][0] + 0.5)/(n_t1[i][0] + 0.5))*((frequency1[i]*
250 (k1+1))/((frequency1[i]+(k1*(1-b+(b*(D/avgdl1)))))))
251         final[candidates[i]] = score
252     return final
253 #####
254 def rates(gs, exp):
255     tp = 0
256     fp = 0
257     fn = 0
258     tp_at_5 = 0
259
260     for i in exp:
261         if i in gs:
262             tp += 1
263         if i not in gs:
264             fp += 1
265
266     for i in gs:
267         if i not in exp:
268             fn += 1
269
270     for i in exp[:5]:
271         if i in gs:
272             tp_at_5 += 1
273
274     return tp, fp, fn, tp_at_5
275
276 def precision(tp, fp):
277     return 0 if tp + fp == 0 else tp/(tp+fp)
278
279 def recall(tp, fn):
280     return 0 if tp + fn == 0 else tp/(tp+fn)
281
282 def f_one(r, p):
283     return 0 if r + p == 0 else 2 * r * p / (r + p)
284
285 def precision_at_n(tp_at_n, n):
286     return tp_at_n/n
287
288 def mean_average_precision(golden, experimental):
289     precisions = []

```

```

290     for i in range(1,len(experimental)+1):
291         tp, fp, _, _ = rates(golden, experimental[:i])
292         precisions.append(precision(tp, fp))
293     meavp = sum(precisions)/len(precisions)
294     return meavp
295
296 def metrics(golden, experimental):
297     with open(golden, 'r', encoding = 'utf-8') as gs_file, open(experimental, 'r',
298 encoding = 'utf-8') as exp_file:
299         gs = []
300         exp = []
301         for line in gs_file:
302             gs.append(line)
303         for line in exp_file:
304             exp.append(line)
305
306         tp, fp, fn, p_at_5 = rates(gs, exp)
307
308         prec = precision(tp, fp)
309
310         rec = recall(tp, fn)
311
312         p_at_5 = p_at_5/5
313
314         f1 = f_one(rec, prec)
315
316         mav = mean_average_precision(gs, exp)
317
318         file = metricsF + os.path.basename(golden)
319         s1 = "Precision: " + str(prec) + "\n" + "Recall: " + str(rec) + "\n" + "F1: "
320 + str(f1) + "\n" + "P@5: " + str(p_at_5) + "\n" + "MAV: " + str(mav)
321
322         open(file, 'w', encoding = 'utf-8').close()
323         with open(file, 'w', encoding = 'utf-8') as f:
324             f.write(s1)
325
326         return prec, rec, f1, p_at_5, mav
327
328 #####
329 #####
330 class Dataset:
331     # models a dataset
332     def __init__(self, directory):
333         self.directory = directory
334         self.asList = list()
335         self.bm25 = None
336
337     def from_files_to_list(self):
338         for filename in glob.glob(self.directory + '*.txt'):
339             with open(filename, 'r', encoding='utf-8') as document: # utf-8 is
340                 breaking at (spanish n)
341                 string = document.read()
342                 self.asList = self.asList + [string]
343
344 class Document:
345     # models a processed doc
346     def __init__(self, filename):
347         self.file = filename

```

```

347     self.ngrams = None
348     self.bm25 = None
349     self.grammar_candidates = None
350     self.keyphrases_idf = list()
351     self.keyphrases_bm25 = list()
352
353     def compute_ngrams(self, low, high):
354         self.ngrams = ngrams(self.file, low, high)
355
356     def get_grammar_candidates(self):
357         self.grammar_candidates = select_candidates(self.get_preprocessed_text())
358
359     def get_keyphrases_bm25(self):
360         #for phrase in self.ngrams:
361         for phrase in self.grammar_candidates:
362             score = self.bm25[phrase]
363             #score_len = 0.5*score + 0.5*len(phrase.split())
364             self.keyphrases_bm25.append((phrase, score))
365
366     def compute_bm25(self, dataset):
367         self.bm25 = bm25(self.grammar_candidates, dataset,
368 self.get_preprocessed_text())
369
370     def write_keyphrases_bm25(self, n):
371         keyphrases = sorted(self.keyphrases_bm25, key=lambda x: x[1], reverse=True)
372         [:n]
373         with open(exprSet + os.path.basename(self.file), 'a+', encoding = 'utf-8') as
out:
374             for i in keyphrases:
375                 out.write(i[0] + '\n')
376
377     def get_ngrams_without_tf(self):
378         ngrams = list()
379         for i in self.ngrams:
380             ngrams.append(i[1])
381         return ngrams
382
383     def get_preprocessed_text(self):
384         with open(self.file, 'r', encoding='utf-8') as document:
385             string = document.read()
386             return string
387
388 #####
389
390 def main():
391     #test.reader.json and dataset_original needed, others are created
392     #from_xml_to_txt() #only #1 time
393     #preproc_all_txt() #only #1 time
394     #parse_json() #only #1 time
395
396     ds = Dataset(preprocDS)
397     ds.from_files_to_list()
398
399     for file in glob.glob(preprocDS + '*.txt'):
400         doc = Document(file)
401         doc.compute_ngrams(1, 3)
402         doc.get_grammar_candidates()
403         doc.compute_bm25(ds.asList)

```

```

404         doc.get_keyphrases_bm25()
405         doc.write_keyphrases_bm25(10)
406
407         precision_sum = 0
408         recall_sum = 0
409         f1_sum = 0
410         p_at_5_sum = 0
411         mav_sum = 0
412         total_files = 0
413
414         for file in glob.glob(goldenSet + '*.txt'):
415             total_files += 1
416             p, r, f1, p_at_5, mav = metrics(file, exprSet + os.path.basename(file))
417             precision_sum += p
418             recall_sum += r
419             f1_sum += f1
420             p_at_5_sum += p_at_5
421             mav_sum += mav
422
423
424         print("Average Precision: " + str(precision_sum/total_files))
425         print("Average Recall: " + str(recall_sum/total_files))
426         print("Average F1: " + str(f1_sum/total_files))
427         print("Average P@5: " + str(p_at_5_sum/total_files))
428         print("Average MAV: " + str(mav_sum/total_files))
429
430
431 if __name__ == '__main__':
432     main()

```