

```

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 from sklearn.preprocessing import StandardScaler
7 from sklearn.linear_model import Perceptron
8
9 import numpy as np
10 import nltk
11 from nltk.corpus import stopwords
12 from nltk.tokenize import word_tokenize
13 from nltk.tokenize import RegexpTokenizer
14 from nltk.tokenize import sent_tokenize
15
16 from math import log
17
18 import xml.etree.ElementTree as ET
19
20 import json
21
22 import os
23 import glob
24 import string
25 import re
26 import itertools, nltk, string
27
28 #####
29
30 stop_words = set(stopwords.words('english'))
31
32 if not os.path.exists('Keyphrases_experimental'):
33     os.makedirs('Keyphrases_experimental')
34 exprSet = 'Keyphrases_experimental/'
35
36 if not os.path.exists('Metrics'):
37     os.makedirs('Metrics')
38 metricsF = 'Metrics/'
39
40 if not os.path.exists('Dataset_preprocessed'):
41     os.makedirs('Dataset_preprocessed')
42 preprocFullDS = 'Dataset_preprocessed/'
43
44 goldenSetTrain = 'Goldenset_train/'
45 goldenSetTest = 'Goldenset_test/'
46 trainSet = 'Train_Set_Preprocessed/'
47 testSet = 'Test_Set_Preprocessed/'
48
49 #####
50 def ngrams(docname, low, high):
51     with open(docname, 'r', encoding='utf-8') as document:
52         result = []
53
54         c_vec = CountVectorizer(ngram_range=(low, high))
55         ngrams = c_vec.fit_transform(document)
56         vocab = c_vec.vocabulary_
57         count_values = ngrams.toarray().sum(axis=0)
58
59         for ngram in sorted([[count_values[i], k] for (k, i) in vocab.items()],
reverse=True):

```

```

60         if ngram[1] not in stop_words:
61             result.append(ngram)
62
63     return result
64
65 def idf(dataset, low, high):
66     # assuming dataset is already a list (like train.data)
67
68     vectorizer = TfidfVectorizer(strip_accents='unicode', ngram_range=(low, high))
69
70     vectorizer.fit_transform(dataset)
71
72     idf = vectorizer.idf_
73
74     return dict(zip(vectorizer.get_feature_names(), idf))
75
76 #####
77 def rates(gs, exp):
78     tp = 0
79     fp = 0
80     fn = 0
81     tp_at_5 = 0
82
83     for i in exp:
84         if i in gs:
85             tp += 1
86         if i not in gs:
87             fp += 1
88
89     for i in gs:
90         if i not in exp:
91             fn += 1
92
93     for i in exp[:5]:
94         if i in gs:
95             tp_at_5 += 1
96
97     return tp, fp, fn, tp_at_5
98
99 def precision(tp, fp):
100     return 0 if tp + fp == 0 else tp/(tp+fp)
101
102 def recall(tp, fn):
103     return 0 if tp + fn == 0 else tp/(tp+fn)
104
105 def f_one(r, p):
106     return 0 if r + p == 0 else 2 * r * p / (r + p)
107
108 def precision_at_n(tp_at_n, n):
109     return tp_at_n/n
110
111 def mean_average_precision(golden, experimental):
112     precisions = []
113     for i in range(1, len(experimental)+1):
114         tp, fp, _, _ = rates(golden, experimental[:i])
115         precisions.append(precision(tp, fp))
116
117     meavp = 0
118     if len(precisions) != 0:
119         meavp = sum(precisions)/len(precisions)

```

```

120     return mavp
121
122 def metrics(golden, experimental):
123     with open(golden, 'r', encoding = 'utf-8') as gs_file, open(experimental, 'r',
124         encoding = 'utf-8') as exp_file:
125         gs = []
126         exp = []
127         for line in gs_file:
128             gs.append(line)
129         for line in exp_file:
130             exp.append(line)
131
132         tp, fp, fn, p_at_5 = rates(gs, exp)
133
134         prec = precision(tp, fp)
135
136         rec = recall(tp, fn)
137
138         p_at_5 = p_at_5/5
139
140         f1 = f_one(rec, prec)
141
142         mav = mean_average_precision(gs, exp)
143
144         file = metricsF + os.path.basename(golden)
145         s1 = "Precision: " + str(prec) + "\n" + "Recall: " + str(rec) + "\n" + "F1: "
146         + str(f1) + "\n" + "P@5: " + str(p_at_5) + "\n" + "MAV: " + str(mav)
147
148         open(file, 'w', encoding = 'utf-8').close()
149         with open(file, 'w', encoding = 'utf-8') as f:
150             f.write(s1)
151
152     return prec, rec, f1, p_at_5, mav
153 #####
154 class Dataset:
155     # models a dataset
156     def __init__(self, directory):
157         self.directory = directory
158         self.asList = list()
159         self.idf = None
160
161     def from_files_to_list(self):
162         for filename in glob.glob(self.directory + '*.txt'):
163             with open(filename, 'r', encoding='utf-8') as document: # utf-8 is
164                 breaking at (spanish n)
165                 string = document.read()
166                 self.asList = self.asList + [string]
167
168     def compute_idf(self, low, high):
169         self.idf = idf(self.asList, low, high)
170
171
172 class Document:
173     # models a processed doc
174     def __init__(self, filename):
175         self.file = filename
176         self.ngrams = None # candidates

```

```

177     self.keyphrases = list()
178     self.CandidatesAsVectors = None
179     self.CandidatesResultVector = None
180     self.grammar_candidates = None
181
182     def compute_ngrams(self, low, high):
183         self.ngrams = ngrams(self.file, low, high)
184
185     def get_keyphrases(self, ds_idf):
186         for phrase in self.ngrams:
187             tf = phrase[0]
188             idf = ds_idf[phrase[1]]
189             tfidf_with_len = tf * idf * len(phrase[1].split())
190             self.keyphrases.append((phrase[1], tfidf_with_len))
191
192     def write_keyphrases(self, n):
193         keyphrases = sorted(self.keyphrases, key=lambda x: x[1], reverse=True)[:n]
194         open(exprSet + os.path.basename(self.file), 'w', encoding = 'utf-8').close()
195         with open(exprSet + os.path.basename(self.file), 'a+', encoding = 'utf-8') as
196 out:
197             for i in keyphrases:
198                 out.write(i[0] + '\n')
199
200     def get_preprocessed_text(self):
201         with open(self.file, 'r', encoding='utf-8') as document:
202             string = document.read()
203             return string
204
205     def predict_perceptron(self, perceptron):
206         for idx, candidate in enumerate(self.CandidatesAsVectors):
207
208             a = perceptron.predict(candidate.reshape(1,-1))
209             if a[0] == 1:
210                 self.keyphrases.append((self.ngrams[idx][1], a))
211
212         open(exprSet + os.path.basename(self.file), 'w', encoding = 'utf-8').close()
213         with open(exprSet + os.path.basename(self.file), 'a+', encoding = 'utf-8') as
214 out:
215             for i in self.keyphrases:
216                 out.write(i[0] + '\n')
217
218     def vectorize_candidates(self, ds_idf, dataset):
219         vecs = []
220         res = []
221         gs = []
222
223         with open(dataset + os.path.basename(self.file), 'r', encoding='utf-8') as
224 golden:
225             for line in golden:
226                 gs.append(line)
227             for i in self.ngrams:
228                 v = np.zeros(3)
229                 tf = i[0]
230                 idf = ds_idf[i[1]]
231                 tfidf_with_len = tf * idf * len(i[1].split())
232                 v[0] = tfidf_with_len
233                 v[1] = len(i[1].split())
234
235                 if i[1] in self.get_preprocessed_text():
236                     v[2] = self.get_preprocessed_text().index(i[1])

```

```

234         else: v[2] = -1
235
236         vecs.append(v)
237
238         res.append(1) if i[1]+"\\n" in gs else res.append(0)
239
240         self.CandidatesAsVectors = vecs
241         self.CandidatesResultVector = res
242
243
244 #####
245
246 def main():
247     ds = Dataset(preprocFullDS)
248     ds.from_files_to_list()
249     ds.compute_idf(1, 3)
250
251     training_vectors = []
252     labels_l = []
253
254     for file in glob.glob(trainSet + '*.txt'):
255         doc = Document(file)
256         doc.compute_ngrams(1, 3)
257         doc.vectorize_candidates(ds.idf, goldenSetTrain)
258         training_vectors += doc.CandidatesAsVectors
259         labels_l += doc.CandidatesResultVector
260
261     labels = np.array(labels_l)
262     arr = np.vstack(training_vectors)
263
264     ppn = Perceptron(max_iter=10000, eta0=0.01, random_state=0)
265     ppn.fit(arr, labels)
266
267     for file in glob.glob(testSet + '*.txt'):
268         doc = Document(file)
269         doc.compute_ngrams(1, 3)
270         doc.vectorize_candidates(ds.idf, goldenSetTest)
271         doc.predict_perceptron(ppn)
272
273
274     precision_sum = 0
275     recall_sum = 0
276     f1_sum = 0
277     p_at_5_sum = 0
278     mav_sum = 0
279     total_files = 0
280
281     for file in glob.glob(goldenSetTest + '*.txt'):
282         total_files += 1
283         p, r, f1, p_at_5, mav = metrics(file, exprSet + os.path.basename(file))
284         precision_sum += p
285         recall_sum += r
286         f1_sum += f1
287         p_at_5_sum += p_at_5
288         mav_sum += mav
289
290
291     print("Average Precision: " + str(precision_sum/total_files))
292     print("Average Recall: " + str(recall_sum/total_files))
293     print("Average F1: " + str(f1_sum/total_files))

```

```

294     print("Average P@5: " + str(p_at_5_sum/total_files))
295     print("Average MAV: " + str(mav_sum/total_files))
296
297 if __name__ == '__main__':
298     main()

```