

```

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 os
14 import glob
15 import string
16 import re
17
18 #####
19 #some global vars such as file nanes and directories
20
21 preprocessedFile = "preprocessed.txt"
22
23
24 keyphrasesFile = "keyphrases.txt"
25
26 if not os.path.exists('Dataset'):
27     os.makedirs('Dataset')
28 preprocDS = 'Dataset/'
29
30 outFileDSBase = "Dataset/DS_preproc_"
31
32 inFile = "original.txt"
33
34 stop_words = set(stopwords.words('english'))
35
36 train = fetch_20newsgroups(subset='train', remove=('headers', 'footers',
37 'quotes')).data
38
39 test = fetch_20newsgroups(subset='test', remove=('headers', 'footers',
40 'quotes')).data
41
42 full_ds = train + test
43
44 #####
45 def preprocess_file(docname, outfile_name):
46     #given a file docname, preprocesses and saves it at file outfile_name
47     with open(docname, 'r', encoding = 'utf-8') as file:
48         outfile = open(outfile_name, 'w', encoding = 'utf-8')
49
50         for line in file:
51             print(preprocess_sentence(line), end='', file=outfile)
52         outfile.close()
53
54     return outfile_name
55
56 def preprocess_list(inList, outfile_name):
57     #given a dataset like 20 newsgroup, preprocess it and save it on files.
58     i = 1
59     for article in inList:

```

```

59     sentences = sent_tokenize(article)
60     outfile = open(outfile_name + str(i) + ".txt", "w", encoding = 'utf-8')
61     for sentence in sentences:
62         preprocessed_sentence = preprocess_sentence(sentence)
63         print(preprocessed_sentence, end='', file=outfile)
64     outfile.close()
65     i += 1
66     return outfile_name, i
67
68 def preprocess_sentence(sentence):
69     #preprocess a sentece: lowercase, remove numbers, stopwords, punctuation.
70     processed = sentence.lower()
71     processed = re.sub(r'\d+', "", processed)
72     tokenizer = RegexpTokenizer(r'\w+')
73     tokens = tokenizer.tokenize(processed)
74     filtered_words = [w for w in tokens if not w in stop_words]
75     if(sentence[-1] == '\n'): return " ".join(filtered_words) + " "
76     return " ".join(filtered_words)
77
78
79 #####
80 def ngrams(docname, low, high):
81     #given a file, get ngrams from range low to high.
82     with open(docname, 'r', encoding = 'utf-8') as document:
83         result = []
84
85         c_vec = CountVectorizer(ngram_range=(low, high))
86         ngrams = c_vec.fit_transform(document)
87         vocab = c_vec.vocabulary_
88         count_values = ngrams.toarray().sum(axis=0)
89
90         for ngram in sorted([[count_values[i],k] for k,i in vocab.items()],
91 reverse=True):
92             result.append(ngram)
93
94     return result
95
96 def idf(dataset, low, high, candidates):
97     #given a dataset of files, compute the idf score for words in list candidates,
98     #ngrams ranging from low to high.
99     vectorizer = TfidfVectorizer(strip_accents='unicode', ngram_range=(low, high),
100 vocabulary = candidates)
101
102     vectorizer.fit_transform(dataset)
103
104     idf = vectorizer.idf_
105
106     return dict(zip(vectorizer.get_feature_names(), idf))
107
108 #####
109 class Dataset:
110     # models a dataset. assuming files already preprocessed
111     def __init__(self, directory):
112         self.directory = directory #directory of the files
113         self.asList = list()
114         self.idf = None
115
116     def from_files_to_list(self):
117         for filename in glob.glob(self.directory + '*.txt'):
118             with open(filename, 'r', encoding='utf-8') as document:

```

```

117         string = document.read()
118         self.asList = self.asList + [string]
119
120     def compute_idf(self, preprocessed_file, low, high, candidates):
121         self.idf = idf(self.asList + [preprocessed_file], low, high, candidates)
122
123     def get_idf(self):
124         return self.idf
125
126
127 class Document:
128     #models a document. assuming the document is already preprocessed.
129     def __init__(self, filename):
130         self.file = filename
131         self.ngrams = None
132         self.keyphrases = list()
133
134     def compute_ngrams(self, low, high):
135         self.ngrams = ngrams(self.file, low, high)
136
137     def find_keyphrases(self, ds_idf):
138         for phrase in self.ngrams:
139             tf = phrase[0]
140             idf = ds_idf[phrase[1]]
141             tfidf_with_len = tf * idf * len(phrase[1].split())
142             self.keyphrases.append((phrase[1], tfidf_with_len))
143
144     def get_preprocessed_text(self):
145         with open(self.file, 'r', encoding='utf-8') as document:
146             string = document.read()
147         return string
148
149     def get_ngrams_without_tf(self):
150         ngrams = list()
151         for i in self.ngrams:
152             ngrams.append(i[1])
153         return ngrams
154
155     def print_keyphrases(self, n):
156         kp = sorted(self.keyphrases, key = lambda x: x[1], reverse = True)[:n]
157
158         open(keyphrasesFile, 'w', encoding='utf-8').close()
159         with open(keyphrasesFile, 'a+', encoding='utf-8') as file:
160             for phrase in kp:
161                 file.write(phrase[0] + '\n')
162         print(kp)
163
164
165 #####
166 def main():
167     #preprocess_list(full_ds, outFileDSBase) #run this only once
168     #preprocess_file(inFile, preprocessedFile) #run this only once
169
170     ds = Dataset(preprocDS)
171     ds.from_files_to_list()
172
173     doc = Document(preprocessedFile)
174     doc.compute_ngrams(1, 3)
175     ds.compute_idf(doc.get_preprocessed_text(), 1, 3, doc.get_ngrams_without_tf())
176

```

```

177     doc.find_keyphrases(ds.get_idf())
178     doc.print_keyphrases(5)
179
180 if __name__ == "__main__":
181     main()
182
183

```