

```

#####
#   REQUIRES MATPLOTLIB AND NUMPY   #
#####
1 import random; random.seed(123)
2 import codecs
3 import string
4 from nltk.stem.porter import PorterStemmer
5 from nltk.probability import FreqDist
6 import gensim
7 import matplotlib
8 import matplotlib.pyplot as plt
9 import numpy as np
10
11
12
13
14 ##### Task 1 #####
15 # Open the file
16 f = codecs.open("pg3300.txt", "r", "utf-8")
17
18 f = f.readlines()
19 text = ""
20 for line in f:
21     text += line
22
23 # Split the text into paragraphs.
24 text = text.split("\n\n")
25
26 # Remove lines that contains Gutenberg
27 without_gutenberg = [par for par in text if "gutenberg" not in par.lower()]
28
29 # Remove all empty lines:
30 filtered_list = [par for par in without_gutenberg if len(par) > 0]
31
32 # Splits each paragraph into words
33 words = [par.split() for par in filtered_list]
34
35 # Initialize stemmer and freqDist
36 stemmer = PorterStemmer()
37 freqDist = FreqDist()
38
39 # Stem words and remove punctuation. Also update freqDist.
40 for i, par in enumerate(words):
41     for j, word in enumerate(par):
42         words[i][j] = stemmer.stem(word.strip(string.punctuation + "\n\r\t").lower())
43         freqDist[word] += 1
44     words[i] = list(filter(None, words[i])) # Filter out empty strings
45
46 # Filter out empty lists
47 words = list(filter(None, words))
48
49
50 # Start plot
51 freqDist = sorted(freqDist.items(), key=lambda x: x[1], reverse=True)[:15]
52
53 labels = [x[0] for x in freqDist]
54 values = [x[1] for x in freqDist]
55
56 # This code for plotting was taken from https://matplotlib.org/3.1.1/gallery/lines_bars_and_markers/barchart.html#sphx-glr-gallery-lines-bars-and-markers-barchart-py
57 # With minor changes.
58 x = np.arange(len(labels))
59 width = 0.40
60
61 fig, ax = plt.subplots()
62 rects1 = ax.bar(x - width/2, values, width, label='Freq')
63 ax.set_ylabel('Freq')
64 ax.set_title('Top 15 most frequent words')
65
66 ax.set_xticks(x)
67 ax.set_xticklabels(labels)
68 ax.legend()
69
70 def autolabel(rects):
71     """Attach a text label above each bar in *rects*, displaying its height. """
72     for rect in rects:
73         height = rect.get_height()
74         xy=(rect.get_x() + rect.get_width() / 2, height),
75         xytext=(0, 3),
76         textcoords="offset points",
77         ha='center', va='bottom')
78
79 autolabel(rects1)
80 fig.tight_layout()
81
82 plt.show()
83 # End plot
84
85 #####

```

```

86 #
87 # TO MAKE THE FOLLOWING CODE RUN #
88 # CLOSE THE PLOT WINDOW! #
89 #
90 #####
91 ##### Task 2
92 # Create dictionary
93 dictionary = gensim.corpora.Dictionary(words)
94
95 # Open stopwords and make a list on the form ["word1", "word2", ... "wordN"]
96 stopwords = open('stopwords.txt', 'r').read().split(",")
97
98 # Stem the stopwords
99 stopwords = [stemmer.stem(w) for w in stopwords]
100
101 # Create list with stopwords ids
102 stopword_ids = []
103 for w in stopwords:
104     try:
105         stopword_ids.append(dictionary.token2id[w])
106     except:
107         continue
108
109 # Filter out stopwords
110 dictionary.filter_tokens(stopword_ids)
111
112 # Create a bag of words
113 bag_of_words = list()
114 for par in words:
115     bag_of_words.append(dictionary.doc2bow(par))
116
117 tfidf_model = gensim.models.TfidfModel(bag_of_words)
118
119 corpus_tfidf = tfidf_model[bag_of_words]
120
121 # Create matrix similarity
122 matrixsim = gensim.similarities.MatrixSimilarity(bag_of_words)
123
124 # Do the same for LSI:
125 lsi = gensim.models.LsiModel(corpus_tfidf, id2word=dictionary, num_topics=100)
126 lsi_corpus = lsi[corpus_tfidf]
127
128
129 # Print three first topics.
130 print('#####')
131
132 print('Three first lsi topics')
133 print(lsi.show_topics(3))
134 print('#####\n\n\n')
135
136 ##### Task 4 #####
137 # Function for preprocessing of query
138 def preprocessing(q):
139     return [stemmer.stem(w.strip(string.punctuation).lower()) for w in q.split()]
140
141 # Preprocess the query and convert to bag of words
142 q = preprocessing("What is the function of money?")
143 q = dictionary.doc2bow(q)
144
145
146 # Report weights
147 q_tfidf = tfidf_model[q]
148 index = gensim.similarities.MatrixSimilarity(corpus_tfidf)
149
150 print('#####')
151 print('Reporting weights')
152 for pair in q_tfidf:
153     weight = pair[1]
154     word = dictionary.get(pair[0])
155     print(f'Word: {word}, weight: {weight}')
156 print('#####\n\n\n')
157
158 def report_relevant_paragraphs(q):
159     docs2similarity = enumerate(index[q])
160     sorted_docs = sorted(docs2similarity, key=lambda kv: -kv[1])[0:3] #Sort the docs.
161     relevant_paragraphs = list()
162     for pair in sorted_docs:
163         relevant_paragraphs.append(pair[0]) # Append the paragraph index to relevant paragraphs
164
165     for par in relevant_paragraphs:
166         print(f'[Paragraph: {par+1}]')
167         lines = filtered_list[par].splitlines(6)[0:6] # Split the paragraphs into lines and get the first 5
168         print(" ".join(lines) + '\n') # Print the lines.
169
170
171 print('#####')
172 print('Relevant paragraphs for query "What is the function of money?")
173 report_relevant_paragraphs(q_tfidf) # Print the relevant docs (first 5 lines).
174 print('#####\n\n\n')

```

```
175 print('#####')
176
177 print('#####')
178 print("Top 3 topics for \"What is the function of money?\")
179 lsi_corpus_q = lsi[q_tfidf]
180 sorted_lsi = (sorted(lsi_corpus_q, key=lambda kv: -abs(kv[1]))[:3] )
181 all_test_topics = (lsi.show_topics())
182 for i in sorted_lsi:
183     print("[Topic ", i[0], "]")
184     print(all_test_topics[i[0]][1])
185 print('#####')
```