

```
[1]:
Created on Mon Nov 18 05:34 2022

@author: Brandon Botter - bbb0103

1) Use the following codes to load the assignment12.txt
    which contains file names. How many file names are in it? (10 points)

file = open("Assignment12.txt", 'r')
text1 = file.read()
file.close()

2) Identify the pattern of the file names, and find out how many file names
    are in the pattern (20 points)

3) Find out file names who don't match with the pattern you designed. (20 points)

4) Use the following codes to read the text from
    "arxiv_annotated13_1.txt"Download arxiv_annotated13_1.txt"
in
file = open("arxiv_annotated13_1.txt", 'r')
text = file.read()
file.close()

Normalize the words and find out their counts.

#Imports
import os
import nltk
import re
import pandas as pd

In [2]:
#1) Use the following codes to load the assignment12.txt
# which contains file names. How many file names are in it? (10 points)

#set the read path
readpath = "%s\Degree\PennState\DAAN_86\Week12\Homeworkk"
#change the directory
os.chdir(readpath)
file_name = "Assignment12.txt"

#open the file
file = open(file_name, 'r')
#read the text as a string
text1 = file.read()
#close the file
file.close()

#use regex to split text1 by any number of blank spaces
file_list = re.split('\s+', text1)

#number of filenames in text1
num_filenames = len(file_list)

print("Splitting the files by spaces, there are " + str(num_filenames) +
      " files in " + file_name)

#Splitting the files by spaces, there are 90 files in Assignment12.txt

In [3]:
#2) Identify the pattern of the file names, and find out how many file names
#match the pattern (20 points)

#The commented section here works as well as the function below

#cleare the pattern
pattern = r'([a-zA-Z]+-[a-zA-Z]+[0-9]+-[0-9]+[0-9]+\.[a-zA-Z]+)'
#compile the pattern
regex = re.compile(pattern, flags = re.IGNORECASE)
#get the list of cases that match the pattern
reg_list = regex.findall(text1)
#find the length of the list of matching cases
reg_len = len(reg_list)

print("Using regex, there are " + str(reg_len) + " files in " + file_name)
print("This was due to typos in some of the file names.")

#I tried to use the regex-generator for the pattern
#but it wouldn't work correctly. I used it as a start and then
#modified it.

#the pattern from regex-generator shown here:
#r"([a-zA-Z]+_?([a-zA-Z]+[0-9]+[0-9]+[0-9]+\.[a-zA-Z]+)"

#compile the regex object
regex = re.compile(r'([a-zA-Z]+-[a-zA-Z]+[0-9]+-[0-9]+[0-9]+\.[a-zA-Z]+)',
                  flags=re.IGNORECASE)
#find all cases that match the pattern and return them
return regex.findall(input_text)

#the list of cases that match the pattern
reg_list = use_regex(text1)
#find the length of the list of matching cases
reg_len = len(reg_list)

print("Using regex, there are " + str(reg_len) + " files in " + file_name)

In [4]:
#3) Find out file names who don't match with the pattern you designed. (20 points)

#use sets to effectively look for the missing pieces

#bring the lists into sets
set_a = set(file_list)
set_b = set(reg_list)

#put it back into a list so all of my lists of files are var type list
typeo_list = list(set_a - set_b)

print("\nthe files which were spelled incorrectly are:\n")
print(typeo_list)

The files which were spelled incorrectly are:

['1m_annotated17_2.txt', '15m_annotated120_1.txt', 'p1os_annotated1_6_2.txt', 'p1os_annotated5_1375_1.txt',
'jdm_annotated32_2.txt', 'p1os_annotated1233_2.txt']

In [5]:
#4) Use the following codes to read the text from
    "arxiv_annoted13_1.txt"

#open and read the file
file_name2 = "arxiv_annoted13_1.txt"
file = open(file_name2, 'r')
text2 = file.read()
file.close()

#Normalize the words and find out their counts.

#first we'll remove special characters
#taken from "Text Analytics with Python, 2nd ed."
def remove_special_characters(text, remove_digits = False):
    pattern = r'([a-zA-Z0-9\s])' if not remove_digits else r'([a-zA-Z0-9\s])'
    text = re.sub(pattern, '', text)
    return text

text2 = remove_special_characters(text2)

In [6]:
#break the large string into a list of words
#use regex to split text1 by any number of blank spaces
words = re.split('\s+', text2)
#try a nltk tokenizer for possible different results
words_token = nltk.word_tokenize(text2)
#try a split function for possible different results
words_split = text2.split(' ')

#it will go with the word tokenizer as it removed the empty string

In [7]:
#find some word counts
#number of words in the list
num_words = len(words)
#number of unique words
unique_words = list(set(words_token))
#number of unique words
num_unique_words = len(unique_words)

print("\nfrom the Arxiv file, there are by the NLTK word tokenizer:")
print(" " + str(num_words) + " words," "word count")
print(" " + str(num_unique_words) + " unique words.\n")

From the Arxiv file, there are by the NLTK word tokenizer:
844 words.
333 unique words.

In [8]:
#use can also do this with a FreqDist
from nltk.probability import FreqDist
from nltk.tokenize import word_tokenize

#find the unique words from the word list
text2 = FreqDist(words_token)
#the FreqDist function likewise reports
print("The FreqDist function likewise reports " +
      str(len(dist)) + " unique words.")

print("\nSome of these words may come from the same root stem/word.")
print("\nI would look at both root stems and root words here.")

The FreqDist function likewise reports 333 Unique words.

Some of these words may come from the same root stem/word.
We'll look at both root stems and root words here.

In [9]:
#Root Stemming
#Build the Porter Stemmer object
porter = nltk.PorterStemmer()
#Get the root stems for each word
word_stems = [porter.stem(word) for word in words_token]

#Get the distribution of the stem/words
stem_dist = FreqDist(word_stems)

#Download for functionality if needed
nltk.download('punkt')

#lemmatization
#Build the WordNetLemmatizer object
WNlemma = nltk.WordNetLemmatizer()
#Get the lemmatization of the words
word_lemmas = [WNlemma.lemmatize(word) for word in words_token]
#lemmas distribution
lemma_dist = FreqDist(word_lemmas)

[nltk_data] Downloading package omw-1.4 to
[nltk_data] C:\Users\bbbot\anaconda3\lib\nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!

In [13]:
#display the counts
#display the counts for all of the different word grabs and counts

#dist
#stem_dist
#lemma_dist

#Get tuples of the distributions
tup_words = tuple(dist.items())
tup_stem = tuple(stem_dist.items())
tup_lemma = tuple(lemma_dist.items())

#Create data frames and join them together...
#there is probably a better way to do this using a pivot but here we are.

twidf = pd.DataFrame(data = tup_words, columns = ["words", "word_count"])
tsidf = pd.DataFrame(data = tup_stem, columns = ["stems", "stem_count"])

#Join the tuple stem data frame to the tuple word data frame on the left
tsidf keeps all values and will join the frames together
df = twidf.join(tsidf, how='left', sort = True)

tsidf = pd.DataFrame(data = tup_lemmas, columns = ["lemmas", "lemma_count"])

#Join the lemma data frame to the word/stem data frame on the left
tsidf keeps all values and will join the frames together
df = df.join(tsidf, how='left', sort = True)

print("\nthe word, stem, and lemma distribution counts:\n")
print("You'll notice that eventually the stem and lemma counts will " +
      "not match up with the words column as the tokenizers interpret " +
      "differently. In the case of lemma, it missed the word 'a'. " +
      "It also had trouble with strange breaks such as eg, which " +
      "in the text file was presented as 'e g' and thus taken as two words.\n")

print("\nI attempted to make use of Dipyran's sequence of text to compare " +
      "normalization results but there seems to be an issue with a " +
      "dependency I was unable to resolve. ")

#Print all of the words, counts, stems, stem_counts, lemmas, and lemma_counts
print(df.to_string())

You'll notice that, eventually the stem and lemma columns will not match up with the words column as the tokenizers interpret differently. In the case of lemma, it missed the word 'a'. It also had trouble with strange breaks such as 'eg' which in the text file was presented as 'e g' and thus taken as two words.

I attempted to make use of Dipyran's sequence of text to compare normalization results but there seems to be an issue with a dependency I was unable to resolve.

words word_count stems stem_count lemmas lemma_count
0 abstract 1 abstract 1.0 abstract 1.0
1 MSC 20 mic 20.0 MSC 20.0
2 although 1 although 1.0 although 1.0
3 the 44 the 44.0 the 44.0
4 internet 15 internet 15.0 internet 15.0
5 as 28 as 28.0 a 47.0
6 level 6 level 6.0 level 6.0
7 topology 8 topology 8.0 topology 1.0
8 has 4 ha 4.0 ha 4.0
9 been 1 been 1.0 been 1.0
10 extensively 1 extens 1.0 extensively 1.0
11 studied 1 studi 2.0 studied 1.0
12 over 1 over 1.0 over 1.0
13 past 1 past 1.0 past 1.0
14 little 2 few 2.0 few 2.0
15 years 1 year 1.0 year 1.0
16 little 1 littl 1.0 little 1.0
17 is 8 as 8.0 is 8.0
18 known 1 known 1.0 known 1.0
19 about 1 about 1.0 about 1.0
20 details 1 detail 1.0 detail 1.0
21 of 34 of 34.0 of 34.0
22 taxonomy 6 taxonomi 6.0 taxonomy 6.0
23 an 7 an 7.0 an 7.0
24 node 1 node 1.0 node 1.0
25 can 4 can 4.0 can 4.0
26 represent 1 repres 3.0 represent 1.0
27 a 19 a 19.0 a 19.0
28 wide 2 wide 2.0 variety 1.0
29 variety 1 variety 1.0 organization 2.0
30 organization 1 organ 1.0
31 e 2 e 2.0 g 2.0
32 large 2 g 2.0 large 3.0
33 large 3 larg 3.0 larg 3.0
35 or 3 or 3.0 small 4.0
36 small 1 small 1.0 small 1.0
37 private 2 privat 2.0 business 3.0
38 business 3 busi 3.0 university 4.0
39 university 1 univers 4.0 with 1.0
40 with 9 with 9.0 vastly 12.0
41 vastly 1 vastli 1.0 different 2.0
42 differ 1 differ 1.0 different 1.0
43 network 8 network 10.0 characteristic 12.0
44 characteristics 2 characterist 2.0 external 1.0
45 external 1 exten 1.0 heterogenous 2.0
46 connectivity 2 connect 3.0 pattern 3.0
47 patterns 3 patern 3.0 growth 4.0
48 growth 1 growth 1.0 tendency 1.0
49 tendencies 1 tendenc 1.0 and 24.0
50 and 28 and 24.0 other 6.0
51 other 1 other 1.0 property 1.0
52 properties 2 properti 2.0 that 12.0
53 that 12 that 12.0 we 20.0
54 we 20 we 20.0 we 20.0
55 hardly 1 hardli 1.0 neglect 1.0
56 neglect 1 neglect 1.0 while 1.0
57 while 1 while 1.0 working 1.0
58 working 1 work 3.0 on 9.0
59 on 9 on 9.0 voracious 1.0
60 voracious 1 veraci 1.0 representation 1.0
61 representations 1 represent 1.0 list 19.0
62 in 19 in 19.0 simulation 1.0
63 simulation 1 simul 1.0 priorior 2.0
64 environments 2 environ 2.0 AIMX 3.0
65 AIMX 3 aimx 3.0 this 5.0
66 this 1 this 1.0 paper 3.0
67 paper 1 paper 1.0 introduce 2.0
68 introduce 2 introduc 3.0 radically 3.0
69 radically 1 radiat 1.0 router 3.0
70 new 3 new 3.0 approach 3.0
71 approach 3 approach 3.0 based 3.0
72 based 1 bas 1.0 machine 3.0
73 machine 3 machin 3.0 learning 3.0
74 learning 3 learn 3.0 technique 2.0
75 technique 1 techniq 1.0
76 to 18 to 18.0 map 1.0
77 map 1 map 2.0 all 1.0
78 all 1 all 1.0 last 12.0
79 asse 12 ase 12.0 into 2.0
80 into 2 into 2.0 success 1.0
81 natural 1 natur 1.0 naturally 1.0
82 ONNX 12 onnx 12.0 successfully 1.0
83 successfully 1 success 1.0 classify 4.0
84 classify 1 classifi 1.0 pressure 3.0
85 KUMBER 15 number 10.0 percent 1.0
86 percent 3 percent 3.0 expected 1.0
87 expected 1 expect 2.0 accuracy 2.0
88 accuracy 2 accuraci 2.0 release 1.0
89 release 1 releas 1.0 community 1.0
90 community 1 commun 1.0 management 1.0
91 dataset 2 dataset 2.0 augmented 1.0
92 augmented 1 augment 3.0 information 2.0
93 information 1 inform 2.0 part 3.0
94 set 3 set 3.0 attribute 1.0
95 attributes 1 attrib 1.0 used 3.0
96 attribute 1 attrib 1.0 believe 1.0
97 believe 1 believ 1.0 will 3.0
98 will 3 will 3.0 serve 1.0
99 serve 1 serv 1.0 invaluable 1.0
100 invaluable 1 inval 1.0 adio 1.0
101 addition 1 addit 1.0 further 2.0
102 further 1 furthe 1.0 understanding 3.0
103 understanding 3 understand 3.0 structure 3.0
104 structure 3 structur 4.0 evolution 5.0
105 evolution 1 evoluit 1.0 introduction 1.0
106 introduction 1 introduct 1.0 rapid 1.0
107 rapid 1 rapid 1.0 expansion 1.0
108 expansion 1 expans 1.0
109 last 1 last 1.0 two 1.0
110 two 1 two 1.0 decade 1.0
111 decade 1 deced 1.0 produced 1.0
112 produced 1 produc 1.0 scale 1.0
113 scale 1 scale 1.0 system 2.0
114 system 1 system 2.0 dial 1.0
115 thousands 2 thousand 2.0 diverse 1.0
116 diverse 1 divers 2.0 independently 1.0
117 independently 1 independ 1.0 collectively 1.0
118 managed 1 manag 1.0 management 1.0
119 networks 2 collect 1.0 provide 1.0
120 collectively 1 provid 1.0 globally 1.0
121 provide 1 global 2.0 across 1.0
122 global 1 across 1.0 spectrum 1.0
123 spectrum 1 spectru 1.0 from 6.0
124 spectral 1 geopolit 1.0 number 1.0
125 geopolitical 1 from 6.0 number 1.0
126 from 1 routab 1.0 routable 1.0
127 number 1 identifi 4.0 routable 1.0
128 globally 1 increas 1.0 identifier 1.0
129 routable 1 less 1.0 done 1.0
130 identifiers 1 than 2.0 le 1.0
131 increased 1 more 1.0 than 2.0
132 less 1 exact 1.0 than 1.0
133 than 2 extact 1.0 exerting 1.0
134 more 1 presur 1.0 significant 1.0
135 match 1 interdomain 1.0 pressure 1.0
136 significant 1 rout 2.0 interdomain 1.0
137 pressure 1 well 2.0 routing 2.0
138 interdomain 1 function 1.0 representation 2.0
139 routing 2 part 1.0 functional 1.0
140 well 2 impress 1.0 structural 1.0
141 functional 1 result 1.0 part 1.0
142 structural 1 heterogen 1.0 impressive 1.0
143 parts 1 highli 1.0 resulted 1.0
144 complex 1 complex 1.0 understanding 1.0
145 resulted 1 challeng 1.0 highly 1.0
146 heterogenous 1 accur 3.0 complex 1.0
147 highly 1 realist 1.0 intamix 1.0
148 complex 1 model 5.0 accurate 2.0
149 challenges 1 infrastruc 2.0 realistic 2.0
150 accurate 1 particula 1.0 heterogenous 2.0
151 realistic 2 intermix 1.0 infrastructure 2.0
152 modeling 4 own 1.0 particular 1.0
153 operation 1 operat 1.0 router 1.0
154 particular 1 by 4.0 owned 1.0
155 intermix 1 mani 2.0 operated 1.0
156 operation 1 backbon 1.0 service 1.0
157 operated 1 region 1.0 many 2.0
158 by 4 access 1.0 backbone 1.0
159 backbone 1 compani 4.0 priorior 2.0
160 backbone 1 statist 2.0 regional 1.0
161 providers 4 faith 1.0 access 1.0
162 character 1 character 2.0 element 1.0
163 access 1 type 11.0 statistical 1.0
164 universities 1 critic 1.0 community 1.0
165 companies 1 path 1.0 publicli 1.0
166 statistical 1 toward 1.0 fully 11.0
167 faithfully 1 for 7.0 critial 1.0
168 characterizes 1 it 1.0 toward 1.0
169 types 10 knowledge 1.0 toward 1.0
170 critical 1 mandatori 1.0 for 7.0
171 critical 1 synthet 1.0 toward 1.0
172 toward 1 construct 2.0 knowledge 1.0
173 for 7 measur 1.0 mandatory 1.0
174 intra 1 intra 1.0 related 1.0
175 knowledge 1 inter 1.0 synthetically 1.0
176 mandatory 1 router 4.0 constructed 1.0
177 enhanced 1 exampl 1.0 composed 1.0
178 synthetically 1 dual 2.0 intra 1.0
179 constructed 1 home 3.0 inter 1.0
180 measured 1 be 1.0 dual 1.0
181 topologies 3 drastic 1.0 example 3.0
182 intra 1 like 1.0 expect 1.0
183 inter 1 contain 1.0 dual 2.0
184 router 3 dozen 1.0 homed 2.0
185 example 3 intern 1.0 be 1.0
186 edit 1 host 1.0 drastically 1.0
187 dual 2 element 1.0 likely 1.0
188 homed 2 switch 1.0 contain 1.0
189 host 1 server 1.0 done 1.0
190 drastically 1 firewall 1.0 internal 1.0
191 firewall 1 hand 2.0 host 1.0
192 dual 1 lower 1.0 element 1.0
193 contain 1 probabl 1.0 switch 1.0
194 dozens 1 have 1.0 server 1.0
195 intermix 1 simpl 1.0 firewall 1.0
196 routers 1 simpl 1.0 build 2.0
197 hosts 1 sinc 2.0 most 1.0
198 element 1 relat 1.0 probably 1.0
199 switches 1 such 1.0 have 1.0
200 servers 1 among 1.0 single 1.0
201 firewall 1 often 1.0 single 1.0
202 hand 2 appropri 1.0 since 2.0
203 most 1 if 1.0 there 1.0
204 probably 1 compos 1.0 dual 1.0
205 have 1 moreov 1.0 diversity 1.0
206 since 1 annot 1.0 among 1.0
207 simple 1 their 1.0 accurately 1.0
209 there 1 exhibit 1.0 augment 1.0
210 do 1 servic 1.0 appropriate 1.0
211 diversity 1 grow 2.0 if 1.0
212 among 1 attract 1.0 characterize 1.0
213 innerts 1 custom 1.0 management 1.0
214 accurately 1 engag 1.0 moreover 1.0
215 augment 1 agreement 1.0 annotating 1.0
216 appropriate 1 through 1.0 pre-requisite 1.0
217 if 1 one 1.0 pre-requisite 1.0
218 characterize 1 do 1.0 exhibit 1.0
219 composing 1 significantli 1.0 growth 2.0
220 moreover 1 time 1.0 engaging 1.0
221 annotating 1 thru 1.0 attracting 1.0
222 environment 1 categor 1.0 customer 1.0
223 prerequisite 1 necessari 2.0 engaging 1.0
224 exhibit 1 develop 2.0 agreement 1.0
225 environment 1 also 1.0 connect 1.0
226 grow 2 ip 2.0 connect 1.0
227 attracting 1 address 2.0 through 1.0
228 customers 1 usef 1.0 do 1.0
229 engaging 1 traffic 1.0 do 1.0
230 agreements 1 analysi 1.0 significantly 1.0
231 network 1 often 1.0 significantly 1.0
232 connect 1 requir 1.0 thus 1.0
233 through 1 distinguish 1.0 categorizing 1.0
234 expect 1 one 1.0 between 1.0
235 do 1 packet 1.0 identify 1.0
236 significantly 1 come 1.0 develop 2.0
237 time 1 given 1.0 dual 1.0
238 thus 1 possibl 1.0 also 1.0
239 categorizing 1 realiz 1.0 mapping 1.0
240 the necessary 1 goal 1.0 related 1.0
241 identify 2 check 1.0 address 2.0
242 develop 2 origin 1.0 user 2.0
243 model 1 high 1.0 traffic 1.0
244 also 1 which 1.0 analysis 1.0
245 mapping 1 simpl 1.0 study 1.0
246 high 1 empir 1.0 required 1.0
247 addresses 1 algorith 1.0 requested 1.0
248 users 2 observ 1.0 distinguish 1.0
249 network 1 data 1.0 network 1.0
250 analysis 1 registri 1.0 packet 1.0
251 studies 1 irr 1.0 come 1.0
252 offical 1 citat 1.0 home 1.0
253 required 1 routeview 1.0 given 1.0
254 distinguish 1 intrins 1.0 possible 1.0
255 between 1 then 1.0 realize 1.0
256 packets 1 employ 1.0 goal 1.0
257 come 1 novel 1.0 checking 1.0
258 home 1 build 1.0 origination 1.0
259 given 1 classif 4.0 prefix 1.0
260 possible 1 exploit 1.0 work 2.0
261 realize 1 these 1.0 which 1.0
262 goal 1 six 1.0 work 2.0
263 checking 1 clas 2.0 construct 1.0
264 type 1 reflect 1.0 representational 1.0
265 originates 1 deriv 1.0 algorith 3.0
266 prefix 1 macroscop 1.0 difference 3.0
267 which 1 valid 1.0 empirically 1.0
268 address 1 our 7.0 difference 3.0
269 alies 1 sampl 1.0 use 2.0
270 employ 1 manual 1.0 related 1.0
271 construct 1 demonst 1.0 registry 1.0
272 representative 2 achiev 1.0 irrit 1.0
273 algorithm 1 high 1.0 related 1.0
274 empirically 1 examin 1.0 CITATIONS 1.0
275 observed 1 were 1.0 intrinsic 1.0
276 differences 1 correct 1.0 dual 1.0
277 use 2 final 1.0 employ 1.0
278 network 2 publi 1.0 novel 1.0
279 registries 1 publicli 1.0 build 1.0
280 irr 1 avail 1.0 classification 4.0
281 intrAPIN 2 promot 1.0 exploit 1.0
282 routeview 1 a 1.0 six 1.0
283 intrinsic 1 a 1.0 six 1.0
284 then 1 section 6.0 class 2.0
285 ntk employ 1 start 1.0 reflect 1.0
286 novel 1 brief 1.0 derive 1.0
287 technique 1 discuss 1.0 macroscopic 1.0
288 build 1 relat 1.0 public 1.0
289 classification 3 describ 1.0 validate 2.0
290 exploits 1 specif 1.0 router 7.0
291 these 1 expect 1.0 reas 1.0
292 six 1 them 1.0 using 1.0
293 classes 2 conclud 1.0 sample 1.0
294 text 1 my system 1.0 manually 1.0
295 infrastructures 1 NAN NAN identified 1.0
296 derive 1 NAN NAN validation 1.0
297 macroscopic 1 NAN NAN demonstrates 1.0
298 statistics 1 NAN NAN achieves 1.0
299 validate 2 NAN NAN high 1.0
300 our 1 NAN NAN excess 1.0
301 results 3 NAN NAN were 1.0
302 using 1 NAN NAN correct 2.0
303 complex 1 NAN NAN find 1.0
304 manually 1 NAN NAN make 1.0
305 identified 1 NAN NAN classifier 1.0
306 validation 1 NAN NAN experiment 1.0
307 demonstrates 1 NAN NAN available 1.0
308 achieves 1 NAN NAN promote 1.0
309 high 1 NAN NAN reas 1.0
310 examined 1 NAN NAN section 6.0
311 classifications 1 NAN NAN were 1.0
312 correct 1 NAN NAN brief 1.0
313 finally 1 NAN NAN discussion 1.0
314 make 1 NAN NAN describes 1.0
315 publicly 1 NAN NAN specify 1.0
316 available 1 ntk employ 1.0 manually 1.0
317 promote 1 NAN NAN introduces 1.0
320 research 1 NAN NAN them 1.0
321 tokens 1 NAN NAN conclud 1.0
322 section 6 NAN NAN NAN NAN
323 start 1 NAN NAN NAN NAN
324 brief 1 NAN NAN NAN NAN
325 discussion 1 NAN NAN NAN NAN
326 related 1 NAN NAN NAN NAN
327 describes 1 NAN NAN NAN NAN
328 specify 1 NAN NAN NAN NAN
329 experiments 1 NAN NAN NAN NAN
330 introduces 1 NAN NAN NAN NAN
331 them 1 NAN NAN NAN NAN
332 conclude 1 NAN NAN NAN NAN

In [11]:
#"" Here ends the Homework Submission ""

In [12]:
#Continued here is the non-workable Dipanjan sequence I had dependency issues with.

#here I now bring in Dipanjan's sequence of text normalization
#it is having an issue with spacy and using fileno for file imports

#script HTML
import re
from bs4 import BeautifulSoup

def strip_html_tags(text):
    soup = BeautifulSoup(text, "html.parser")
    [s.extract() for s in soup.find('iframe', 'script')]
    stripped_text = soup.get_text()
    stripped_text = re.sub(r'(\n|\n|\n|\n)', '\n', stripped_text)
    return stripped_text

#Removing Accented Characters
import unicodedata

def remove_accents_chars(text):
    text = unicodedata.normalize('NFKD', text).encode('ascii', 'ignore').decode('utf-8', 'ignore')
    return text

#Expanding Contractions
from contractions import contractions_dict
import re

def expand_contractions(text, contraction_mapping=contractions_dict):
    contractions_pattern = re.compile('(%s)' % format('\\1', '.join(contraction_mapping.keys())',
    flags=re.IGNORECASE|re.DOTALL))

    def expand_match(contraction):
        match = contraction.group()
        first_char = match[0]
        expanded_contraction = contraction_mapping.get(match)
        if contraction_mapping.get(match):
            expanded_contraction = first_char + expanded_contraction[1:]
        return expanded_contraction

    expanded_text = contractions_pattern.sub(expand_match, text)
    expanded_text = re.sub(r'(\n|\n|\n|\n)', '\n', expanded_text)
    return expanded_text

#Removing Special Characters
#You can find this above

#Case Conversion
def lower(text):
    return text.lower()

#Text Correction
#Correcting Repeating Characters
from nltk.corpus import wordnet
def remove_repeated_characters(tokens):
    def replace_substitution = lambda v:
        if wordnet.synsets(v):
            new_word = repeat_pattern.sub(match_substitution, old_word)
            return replace(new_word if new_word != old_word else new_word)
        correct_tokens = (replace(word) for word in tokens)
    return correct_tokens

#Correcting Spellings
def tokens(text):
    """
    #Get all words from the corpus
    return re.findall('[a-z]+', text.lower())

WORDS = words
WORD_COUNTS = collections.Counter(WORDS)

#Define set of words that are one, two, or three away from our input
def edit1(words):
    """
    #Return all strings that are one edit away
    #from the input word.
    """
    alphabets = 'abcdefghijklmnopqrstuvwxyz'
    return a list of all possible (first, rest) pairs
    #that the input word is one edit off.
    """
    return [(word[:i], word[i:]) for i in range(len(word)+1)]
    pairs = []
    for (a, b) in pairs if b:
        deletes = (a[:i] + b[i+1:] for i in range(len(a)))
        transposes = (a[:i] + b[i+1:] + a[i] + b[i] for (a, b) in pairs if len(b) > 1)
        replaces = (a[:i] + b[i] for (a, b) in pairs if b in alphabet if b != a)
        inserts = (a[:i] + b + a[i+1:] for (a, b) in pairs if b in alphabet)
    return set(deletes + transposes + replaces + inserts)

def edit2(words):
    """
    #Return all strings that are two edits away
    #from the input word.
    """
    return [e2 for e1 in edit1(word) for e2 in edit1(e1)]

#Define function to return a subset of words from our candidate set of words
#obtained from the edit functions based on if they occur in the vocabulary
def known(words):
    """
    #Return the subset of words that are actually
    #in our WORD_COUNTS dictionary.
    """
    return (w for w in words if w in WORD_COUNTS)

#Select the correct word from a number of candidates
def correct(word):
    """
    #Get the best correct spelling for the input word
    """
    #Priority: 1) edit distance 0, then 1, then 2
    candidates = (known(edit0(word)) if word in WORD_COUNTS else
                  known(edit1(word)) or
                  known(edit2(word)) or
                  known(word))
    return max(candidates, key=WORD_COUNTS.get)

#Now go back and check case and ensure a correct match
def show_the_correct_text(text):
    """
    #Spell-correct word in match.
    #If it doesn't match, return lower/title case.
    """
    word = match.group()
    #For text, upper, lower, title, or just str.
    """
    return (str.upper if text.isupper() else
            str.lower if text.islower() else
            str.title if text.isitle() else
            lambda word: correct(word.lower()))

def correct_text_generic(text):
    """
    #Correct all the words within a text.
    #Returning the corrected text.
    """
    return re.sub('([a-zA-Z]+)', correct_match, text)

#Word correction for spelling can also be done with the textblob library
from textblob import Word
# Word('fiat')
# Word.spellcheck()

#Stemming Words
def simple_tokenize(text):
    ps = nltk.PorterStemmer()
    return ' '.join([ps.stem(word) for word in text.split()])

#Lemmatization
#The book's version is having issues so I'll make my own
import spacy
# Use spacy.load('en_core') if you have downloaded the language model or directly after install spacy
nlp = spacy.load('en_core', parser=True, tagger=True, entity=True)
# Note: My system keeps crashing his crashed yesterday, ours crashes daily!

#Lemmatization
def lemmatize_text(text):
    """
    # Lemmatize text
    # return text = ' '.join([word.lemma_ if word.lemma_ != '-PRON-' else word.text for word in text])
    """
    return text

def lemmatize_text(word):
    """
    # Lemmatize text
    """
    return word

#Remove Stopwords
from nltk.tokenize import ToktokTokenizer
tokenizer = ToktokTokenizer()
stopword_list = nltk.corpus.stopwords.words('english')
def remove_stopwords(text, is_lower_case=False, stopwords=stopword_list):
    if is_lower_case:
        filtered_tokens = [token for token in tokens if token.lower() not in stopwords]
    else:
        filtered_tokens = [token for token in tokens if token not in stopwords]
    filtered_text = ' '.join(filtered_tokens)
    return filtered_text

#All of the Dipanjan functions build into one normalizer
def normalize_corpus(documents, html_stripping=True, contraction_expansion=True,
                    remove_accents=True, char_remove=True, text_lower_case=True,
                    stopword_removal=True, special_char_remove=True,
                    remove_digits=True, remove_punctuation=True):
    normalized_corpus = []
    # Normalize each document in the corpus
    for doc in documents:
        if isinstance(doc, str):
            # Strip HTML
            # Remove accented characters
            if contracted_char_remove:
                doc = remove_contracted_characters(doc)
            # Expand contractions
            if contraction_expansion:
                doc = expand_contractions(doc)
            # Lowercase the text
            if text_lower_case:
                doc = doc.lower()
            # Remove extra newlines
            doc = re.sub(r'(\n|\n|\n|\n)', '\n', doc)
            # Stem text
            if text_stemmer:
                doc = simple_stemmer(doc)
            # Remove special characters and/or digits
            if special_char_remove:
                # Insert spaces between special characters to isolate them
                doc = special_char_pattern.sub(' ', doc)
            # Remove special characters/digits
            doc = re.sub(r'([a-zA-Z0-9\s])', ' ', doc)
            # Remove stopwords
            doc = remove_stopwords(doc, is_lower_case=text_lower_case)
            # Lemmatize corpus
            normalized_corpus.append(doc)
    return normalized_corpus

#And this doesn't work as intended...
text = normalize_corpus(texts)
```



[illegible]