



Xi'an Jiaotong-Liverpool University

西交利物浦大学

Lab 6: Feature Generation

Author **Jeongyeong Park**

Module **INT104**

Teacher **Xiaobo Jin**

Date **27th/April/202**

Introduction

Feature generation is the method of generating new features from one or more existing features. In this lab, the features are generated from the text collections. All files in subdirectories of the dataset should be read and texts in these files should be split into words. In order to deal with files more efficiently in this lab, the name of them would better save into 2-dimensional list. Since files might contain some stopwords which appear frequently but do not convey much meaning, they should be removed. For TF-IDF, unique words which appear even once in the files should be obtained. Finally, as going through all files in the dataset, words in each file are compared to obtained unique words and TF-IDF will be computed. When words are dealt in this lab, they should be the key of dictionary for each file so they can be tracked easily.

Design & Implementation

1. Design

1) Design for reading the text files and split the document text into words

First, names of subdirectories and the files in each subdirectories should be recorded in lists(subdirec_List and textfile_List respectively) so they can be used for reading files. Subdirec_List will be a linear list in order to save the name of each subdirectory, and textfile_List will be a 2-dimensional list in order to save all files in each subdirectory in each list. In order to read the files and save splitted words from the document text, each list for each subdirectory is created. Then, as we go through each file, the document texts are read and splitted into words and the splitting separator is non-alphabet letters. Therefore, when this process is done, non-alphabet letters are replaced into ' '(one space), all words in the document text are changed into lower-case. Finally, all words are separated with the default separator which is any whitespace since all non-alphabet letters changed into ' '. The following line will be used :

```
re.sub(r'^a-zA-Z',' ',word.lower()).split()
```

2) Design for removing stopwords from the text collection

In order to remove stopwords that are appearing frequently but having no information, a text file containing stopwords should be read first. After reading this text file, each stopwords will be saved in a list. Since words in every file are splitted with a splitting separator which is non-alphabet, stopwords also should be splitted with the same splitting separator. The following line is used :

```
re.sub(r'^a-zA-Z',' ',stopword)
```

Then, the stopwords list will be converted to set in order to remove duplicate words. Words in each word list of each file will be compared to these stopwords in the stopwords list and if a word is the same as the stopwords in the stopwords list, it will be removed. In python list, if there are two identical elements in the list, and that element should be removed, only the first one will be removed from the list. Thus, following line is used :

```
list1=[word for word in list1 if word not in list2]
```

to make sure all stopwords in the words lists of files are removed.

3) Design for performing word stemming to remove the word suffix

For stemming, following line is used :

```
list=[stemmer.stem(plural) for plural in list]
```

4) Design for computing the frequency of each word in each document

For this step, unique words appearing in files located in the dataset should be obtained first. By using set, unique words appearing in files in each subdirectory are obtained first. Set will automatically remove duplicate words. Then, these sets are combined together in order to get final unique words set with the following line :

```
final_set=set1.union(set2.union(set3.union(set4.union(set5))))
```

After obtaining unique words, this set is converted to list to have index. A dictionary is created for saving the frequency of each word in each document. In order to keep tracking filename and the frequency of each word together, nested dictionary is created such as following example :

```
dictionary={'file_name1':{'word1' : 2, 'word2':0.....}, 'file_name2':{'word1' : 0, 'word2' : 2.....},.....}
```

For counting words, first, the inner dictionaries are created for all files in the dataset. Then, unique words will be set as keys for inner dictionaries and value will be set for 0. Example of key line :

```
inner_dictionary=dict.fromkeys(unique_word_list,0)
```

Then, as going through the list of words of the file, every time the word in the list of unique words is detected in the list of words of the file, one is added to the value of having a key as that word. Example of key lines :

```
for word in word_list_of_file1:
```

```
    dictionary[file1][word]+=1
```

This process is done for each subdirectory.

5) Design for computing document frequency

In order to obtain document frequency, a dictionary for it is needed and set unique words in list as keys and set value for it as 0. The following line will be used :

```
dictionary=dict.fromkeys(unique_words_list,0)
```

Then, as going through a dictionary of frequency of a word, if a value of a particular word is not 0 in that file, add one to the value having key as that word. Example of key lines :

```
for word in unique_word_list:
```

```
    if (dictionary_of_frequency[file1][word]!=0):
```

```
        dictionary_of_document_fre[word]+=1
```

6) Design for computing document-by-word matrix

In order to obtain document-by-word matrix, a_{ik} which is the weight of word k in document i. For a_{ik} , $a_{ik} = f_{ik} * \log(\frac{N}{n_k})$ formula is used. In order to calculate, N which is the number of documents in the dataset should be found. Then, an array in order to save a_{ik} is created. Additionally, in order to save each a_{ik} for each word, N numbers of dictionaries will be created inside this array and starting from the first file of first subdirectory dictionary will be indexed 0 from N-1. By using dictionary for the frequency of word in document and dictionary for the document frequency, computing a_{ik} . Example of key lines:

```
for word in unique_word_list:
```

```
    array_for_weight[index][word]=dict_frequency[file][word]*math.log(N/dict_n[word])
```

After that, list is created to save value for final matrix. Inside this linear list, N number of dictionaries will be created to get value A_{ik} for each word in each file. A_{ik} is computed with formula :

$$A_{ik} = \frac{a_{ik}}{\sqrt{\sum_{j=1}^D a_{ij}^2}} \quad D=\text{the number of the unique words in the document}$$

Words in same document share same denominator in this formula. Therefore, for each file, calculating denominator first, and then use it to calculate to calculate each A_{ik} value. After calculating every A_{ik} in the same file, denominator should be set 0 in order to calculate new denominator for words in next file. Example of key lines :

```
denominator=0
```

```
for word in unique_word_list:
```

```
denominator=math.sqrt(float(array_weight[index][word]))+denominator
```

```
for word in unique_word_list:
```

```
    document_by_word[index][word]=array_weight[index][word]/denominator
```

```
denominator=0
```

After computing A_{ik} for all files in the dataset, convert this list within dictionaries into DataFrame. Since the $A_{N \times D}$ matrix should be computed, convert this DataFrame into numpy array. Finally, save this numpy array into npz file.

2. Implementation

1) Code for reading the text files and split the document text into words

```
1 import os
2 import numpy as np
3 from nltk.stem.porter import *
4 import pandas as pd
5 import math
6
7
8 #Listing name of subdirectories in dataset
9 subdirec_List=os.listdir('/dataset')
10
11 #Listing files in each subdirectories(using jagged list)
12 textfile_List=[]
13 for i in range(len(subdirec_List)):
14     textfile_List.append([])
15     textfile_List[i]=os.listdir('/dataset/'+subdirec_List[i])
```

Figure 1. Code for reading the name of files

```

17 wd_alt = []
18 wd_comp = []
19 wd_rec = []
20 wd_soc = []
21 wd_talk = []
22 #Open and read files in each subdirectories
23 for i in range(len(subdirec_List)):
24     for j in textfile_List[i]:
25         if(i==0):
26             f=open('/dataset/'+subdirec_List[i]+'/' + j, 'r', encoding='Latin1')
27             r = f.read()
28             wd_alt.append(re.sub(r'[^a-zA-Z]', ' ', r.lower()).split())
29         elif(i==1):
30             f = open('/dataset/' + subdirec_List[i] + '/' + j, 'r', encoding='Latin1')
31             r = f.read()
32             wd_comp.append(re.sub(r'[^a-zA-Z]', ' ', r.lower()).split())
33         elif(i==2):
34             f = open('/dataset/' + subdirec_List[i] + '/' + j, 'r', encoding='Latin1')
35             r = f.read()
36             wd_rec.append(re.sub(r'[^a-zA-Z]', ' ', r.lower()).split())
37         elif(i==3):
38             f = open('/dataset/' + subdirec_List[i] + '/' + j, 'r', encoding='Latin1')
39             r = f.read()
40             wd_soc.append(re.sub(r'[^a-zA-Z]', ' ', r.lower()).split())
41         else:
42             f = open('/dataset/' + subdirec_List[i] + '/' + j, 'r', encoding='Latin1')
43             r = f.read()
44             wd_talk.append(re.sub(r'[^a-zA-Z]', ' ', r.lower()).split())

```

Figure 2. Code for splitting the document text into words

2) Code for removing stopwords from the text collection

```

46 #save stopwords in the list
47 sw = open('/stopwords.txt', 'r', encoding='Latin1')
48 stopword=sw.read().splitlines()
49 for j in range(len(stopword)):
50     stopword[j]=re.sub(r'[^a-zA-Z]', ' ', stopword[j])
51
52 #removing stopwords
53 #for alt.atheism
54 for i in range(len(wd_alt)):
55     wd_alt[i]=[word for word in wd_alt[i] if word not in stopword]
56 #for comp.graphics
57 for i in range(len(wd_comp)):
58     wd_comp[i]=[word for word in wd_comp[i] if word not in stopword]
59 #for rec.motorcycles
60 for i in range(len(wd_rec)):
61     wd_rec[i]=[word for word in wd_rec[i] if word not in stopword]
62 #for soc.religion.christian
63 for i in range(len(wd_soc)):
64     wd_soc[i]=[word for word in wd_soc[i] if word not in stopword]
65 #for talk.politics.misc
66 for i in range(len(wd_talk)):
67     wd_talk[i]=[word for word in wd_talk[i] if word not in stopword]
68
69

```

Figure 3. Code for removing stopwords

3) Code for performing word stemming to remove the word suffix

```

70     #stemming
71     stemmer=PorterStemmer()
72     #for alt.atheism
73     for i in range(len(wd_alt)):
74         wd_alt[i]=[stemmer.stem(plural) for plural in wd_alt[i]]
75     #for comp.graphics
76     for i in range(len(wd_comp)):
77         wd_comp[i]=[stemmer.stem(plural) for plural in wd_comp[i]]
78     #for rec.motorcycles
79     for i in range(len(wd_rec)):
80         wd_rec[i] = [stemmer.stem(plural) for plural in wd_rec[i]]
81     #for soc.religion.christian
82     for i in range(len(wd_soc)):
83         wd_soc[i] = [stemmer.stem(plural) for plural in wd_soc[i]]
84     #for talk.politics.misc
85     for i in range(len(wd_talk)):
86         wd_talk[i] = [stemmer.stem(plural) for plural in wd_talk[i]]
87

```

Figure 4. Code for stemming

4) Code for computing the frequency of each word in each document

```

100     #creating set for unique word appear even once in one of files in dataset
101     # for alt.atheism
102     uniwd_alt=set()
103     for i in range(len(wd_alt)):
104         uniwd_alt.update(wd_alt[i])
105     # for comp.graphics
106     uniwd_comp=set()
107     for i in range(len(wd_comp)):
108         uniwd_comp.update(wd_comp[i])
109     # for rec.motorcycles
110     uniwd_rec=set()
111     for i in range(len(wd_rec)):
112         uniwd_rec.update(wd_rec[i])
113     # for soc.religion.christian
114     uniwd_soc=set()
115     for i in range(len(wd_soc)):
116         uniwd_soc.update(wd_soc[i])
117     # for talk.politics.misc
118     uniwd_talk=set()
119     for i in range(len(wd_talk)):
120         uniwd_talk.update(wd_talk[i])
121
122     # final unique word
123     uniwd=set()
124     uniwd=uniwd_alt.union(uniwd_comp.union(uniwd_rec.union(uniwd_soc.union(uniwd_talk))))
125     uniwd=list(uniwd)
126

```

Figure 5. Code for obtaining set of unique word and converting it to list

```

114 #creating dictionary for storing the frequency of word k in document i
115 dict_f={}
116 for i in range(len(textfile_List)):
117     for j in range(len(textfile_List[i])):
118         dict_f[textfile_List[i][j]]=dict.fromkeys(unlwd,0) #Adding unique words as a key of inner dictionary and set 0 for value
119
120 # for files in alt.atheism
121 for i in range(len(textfile_List[0])):
122     for word in wd_alt[i]:
123         dict_f[textfile_List[0][i]][word]+=1
124 #for files in comp.graphics
125 for i in range(len(textfile_List[1])):
126     for word in wd_comp[i]:
127         dict_f[textfile_List[1][i]][word]+=1
128 #for files in rec.motorcycles
129 for i in range(len(textfile_List[2])):
130     for word in wd_rec[i]:
131         dict_f[textfile_List[2][i]][word]+=1
132 #for files in soc.religion.christian
133 for i in range(len(textfile_List[3])):
134     for word in wd_soc[i]:
135         dict_f[textfile_List[3][i]][word]+=1
136 #for files in talk.politics.misc
137 for i in range(len(textfile_List[4])):
138     for word in wd_talk[i]:
139         dict_f[textfile_List[4][i]][word]+=1

```

Figure 6. Code for frequency of word

5) Design for computing document frequency

```

141 #creating dictionary for storing the document frequency
142 dict_n=dict.fromkeys(unlwd,0)
143 for i in range(len(textfile_List)):
144     for j in range(len(textfile_List[i])):
145         for word in unlwd:
146             if (dict_f[textfile_List[i][j]][word]!=0):
147                 dict_n[word]+=1

```

Figure 7. Code for document frequency

6) Design for computing document-by-word matrix

```

149 # finding the number of documents(N) in the dataset
150 N=0
151 for i in range(len(textfile_List)):
152     for j in range(len(textfile_List[i])):
153         N=N+1

```

Figure 8. Code for N which is the number of documents in the dataset


```

178 #computing the weight of word k in document i : aik
179 array_weight=[]
180 for i in range(len(textfile_List)):
181     for j in range(len(textfile_List[i])):
182         array_weight.append(dict.fromkeys(uniwd,0))
183
184     index=-1
185     for i in range(len(textfile_List)):
186         for j in range(len(textfile_List[i])):
187             index=index+1
188             for word in uniwd:
189                 array_weight[index][word]=dict_f[textfile_List[i][j]][word]*math.log(N/dict_n[word])
190

```

Figure 9. Code for the weight of word

```

168 #computing the document-by-word matrix
169 document_by_word=[]
170 for i in range(len(textfile_List)):
171     for j in range(len(textfile_List[i])):
172         document_by_word.append(dict.fromkeys(uniwd,0))
173
174     index2=-1
175     denominator=0
176     for i in range(len(textfile_List)):
177         for j in range(len(textfile_List[i])):
178             index2=index2+1
179             for word in uniwd:
180                 denominator=math.sqrt(float(array_weight[index2][word]))+denominator
181             for word in uniwd:
182                 document_by_word[index2][word]=array_weight[index2][word]/denominator
183             denominator=0

```

Figure 10. Code for the document-by-word list within dictionaries

```

185 #Creating final matrix and save into npz file
186 df_matrix=pd.DataFrame(document_by_word)
187 final_matrix=df_matrix.to_numpy()
188 np.savez('train-20ng.npz',X=final_matrix)

```

Figure 11. Code for the final matrix

Test

Following figures are examples of testing to see if the code is working properly.

- 1) Test for reading the text files and split the document text into words

```
Print subdirect_List
['alt.atheism', 'comp.graphics', 'rec.motorcycles', 'soc.religion.christian', 'talk.politics.misc']
```

Figure 12. Result of subdirect_List

```
Print txtfile_list
[['49968', '51068', '51119', '51120', '51121', '51122', '51123', '51124', '51125', '51126', '51127', '51128', '51130', '51131', '51132', '51133', '51134', '51135', '51136',
'51139', '51140', '51141', '51142', '51143', '51144', '51145', '51146', '51147', '51148', '51149', '51150', '51151', '51152', '51153', '51154', '51155', '51156', '51157',
'51158', '51159', '51160', '51161', '51162', '51163', '51164', '51165', '51169', '51170', '51171', '51172', '51173', '51174', '51175', '51176', '51177', '51178', '51179',
'51180', '51181', '51182', '51183', '51184', '51185', '51186', '51187', '51188', '51189', '51190', '51191', '51192', '51193', '51194', '51195', '51196', '51197', '51198',
'51199', '51200', '51201', '51202', '51203', '51204', '51205', '51206', '51208', '51209', '51210', '51211', '51212', '51213', '51214', '51215', '51216', '51217', '51218',
'51219', '51220', '51221', '51222', '51223', '51224', '51225', '51226', '51227', '51228', '51229', '51230', '51231', '51232', '51233', '51234', '51235', '51236', '51237',
'51238', '51239', '51240', '51241', '51242', '51243', '51244', '51245', '51246', '51247', '51249', '51250', '51251', '51252', '51253', '51254', '51255', '51256', '51258',
'51259', '51260', '51261', '51262', '51265', '51266', '51267', '51268', '51269', '51270', '51271', '51272', '51273', '51274', '51275', '51276', '51277', '51278', '51279',
'51280', '51281', '51282', '51283', '51284', '51285', '51286', '51287', '51288', '51290', '51291', '51292', '51293', '51294', '51295', '51296', '51297', '51298', '51299',
'51300', '51301', '51302', '51303', '51304', '51305', '51306', '51307', '51308', '51309', '51310', '51311', '51312', '51313', '51314', '51315', '51316', '51317', '51318',
'51319', '51320', '52499', '52909', '52910', '53055', '53056', '53057', '53058', '53059', '53062', '53064', '53065', '53066', '53067', '53069', '53070', '53071', '53072',
'53073', '53075', '53078', '53081', '53082', '53083', '53085', '53086', '53087', '53090', '53093', '53094', '53095', '53096', '53097', '53098', '53099', '53106', '53108',
'53110', '53111', '53112', '53113', '53114', '53117', '53118', '53120', '53121', '53122', '53123', '53124', '53125', '53126', '53127', '53130', '53131', '53132', '53133',
'53134', '53135', '53136', '53137', '53139', '53140', '53141', '53142', '53143', '53144', '53145', '53149', '53150', '53151', '53153', '53154', '53157', '53158', '53159',
'53160', '53161', '53162', '53163', '53164', '53165', '53166', '53167', '53168', '53170', '53171', '53172', '53173', '53174', '53175', '53176', '53177', '53178', '53179',
'53180', '53181', '53182', '53183', '53184', '53185', '53186', '53187', '53188', '53190', '53191', '53192', '53193', '53194', '53195', '53196', '53197', '53198', '53199',
'53201', '53203', '53208', '53209', '53210', '53211', '53212', '53213', '53214', '53215', '53216', '53217', '53218', '53219', '53220', '53221', '53222', '53223', '53224',
'53225', '53226', '53228', '53229', '53230', '53231', '53232', '53235', '53237', '53238', '53239', '53240', '53243', '53248', '53249', '53250', '53251', '53252', '53256',
'53258', '53266', '53267', '53269', '53271', '53274', '53275', '53281', '53282', '53283', '53284', '53285', '53286', '53287', '53288', '53289', '53290', '53292', '53298',
'53303', '53304', '53305', '53306', '53307', '53308', '53309', '53311', '53312', '53314', '53323', '53334', '53347', '53351', '53366', '53370', '53371', '53373', '53374',
'53375', '53376', '53377', '53380', '53381', '53382', '53383', '53387', '53389', '53390', '53391', '53434', '53435', '53436', '53437', '53438', '53439', '53440', '53441',
'53442', '53443', '53445', '53449', '53459', '53460', '53465', '53466', '53467', '53468', '53471', '53477', '53478', '53483', '53509', '53510', '53512', '53515', '53518',
'53519', '53521', '53522', '53523', '53524', '53525', '53526', '53527', '53528', '53529', '53531', '53532', '53533', '53534', '53535', '53536', '53537', '53572', '53573', '53574',
'53575', '53576', '53577', '53578', '53579', '53580', '53581', '53582', '53583', '53584', '53585', '53586', '53587', '53588', '53589', '53590', '53591', '53592', '53593', '53594', '53595', '53596', '53597', '53598', '53599', '53600', '53601', '53602', '53603', '53604', '53605', '53606', '53607', '53608', '53609', '53610', '53611', '53612', '53613', '53614', '53615', '53616', '53617', '53618', '53619', '53620', '53621', '53622', '53623', '53624', '53625', '53626', '53627', '53628', '53629', '53630', '53631', '53632', '53633', '53634', '53635', '53636', '53637', '53638', '53639', '53640', '53641', '53642', '53643', '53644', '53645', '53646', '53647', '53648', '53649', '53650', '53651', '53652', '53653', '53654', '53655', '53656', '53657', '53658', '53659', '53660', '53661', '53662', '53663', '53664', '53665', '53666', '53667', '53668', '53669', '53670', '53671', '53672', '53673', '53674', '53675', '53676', '53677', '53678', '53679', '53680', '53681', '53682', '53683', '53684', '53685', '53686', '53687', '53688', '53689', '53690', '53691', '53692', '53693', '53694', '53695', '53696', '53697', '53698', '53699', '53700', '53701', '53702', '53703', '53704', '53705', '53706', '53707', '53708', '53709', '53710', '53711', '53712', '53713', '53714', '53715', '53716', '53717
```

Figure 13. Omitted result of textfile List

['from', 'mathew', 'mathew', 'mantis', 'co', 'uk', 'subject', 'alt', 'atheism', 'faq', 'atheist', 'resources', 'summary', 'books', 'addresses', 'music', 'anything', 'related', 'to', 'atheism', 'keywords', 'faq', 'atheism', 'books', 'music', 'fiction', 'addresses', 'contacts', 'expires', 'thu', 'apr', 'gmt', 'distribution', 'world', 'organization', 'mantis', 'consultants', 'cambridge', 'uk', 'supercedes', 'mantis', 'co', 'uk', 'lines', 'archive', 'name', 'atheism', 'resources', 'alt', 'atheism', 'archive', 'name', 'resources', 'last', 'modified', 'december', 'version', 'atheist', 'resources', 'addresses', 'of', 'atheist', 'organizations', 'usa', 'freedom', 'from', 'religion', 'foundation', 'darwin', 'fish', 'bumper', 'stickers', 'and', 'assorted', 'other', 'atheist', 'paraphernalia', 'are', 'available', 'from', 'the', 'freedom', 'from', 'religion', 'foundation', 'in', 'the', 'us', 'write', 'to', 'ffff', 'p', 'o', 'box', 'madison', 'wi', 'atheism', 'religion', 'evolution', 'designs', 'evolution', 'designs', 'sell', 'the', 'darwin', 'fish', 'it', 's', 'a', 'fish', 'symbol', 'like', 'the', 'ones', 'christians', 'stick', 'on', 'their', 'cars', 'but', 'with', 'feet', 'and', 'the', 'word', 'darwin', 'written', 'inside', 'the', 'deluxe', 'moulded', 'd', 'plastic', 'fish', 'is', 'postpaid', 'in', 'the', 'us', 'write', 'to', 'evolution', 'designs', 'laurel', 'canyon', 'north', 'hollywood', 'ca', 'people', 'in', 'the', 'san', 'francisco', 'bay', 'area', 'can', 'get', 'darwin', 'fish', 'from', 'lynn', 'gold', 'try', 'mailing', 'figmo', 'netcom', 'com', 'for', 'net', 'people', 'who', 'go', 'to', 'lynn', 'directly', 'the', 'price', 'is', 'per', 'fish', 'american', 'atheist', 'press', 'aap', 'publish', 'various', 'atheist', 'books', 'critiques', 'of', 'the', 'bible', 'lists', 'of', 'biblical', 'contradictions', 'and', 'so', 'on', 'one', 'such', 'book', 'is', 'the', 'bible', 'handbook', 'by', 'w', 'p', 'ball', 'and', 'g', 'w', 'foote', 'american', 'atheist', 'press', 'pp', 'isbn', 'nd', 'edition', 'bible', 'contradictions', 'absurdities', 'atrocities', 'immoralities', 'contains', 'ball', 'foote', 'the', 'bible', 'contradicts', 'itself', 'aap', 'based', 'on', 'the', 'king', 'james', 'version', 'of', 'the', 'bible', 'write', 'to', 'american', 'atheist', 'press', 'p', 'o', 'box', 'austin', 'tx', 'or', 'cameron', 'road', 'austin', 'tx', 'telephone', 'fax', 'prometheus', 'books', 'sell', 'books', 'including', 'haught', 's', 'holy', 'horrors', 'see', 'below', 'write', 'to', 'east', 'amherst', 'street', 'buffalo', 'new', 'york', 'telephone', 'an', 'alternate', 'address', 'which', 'may', 'be', 'newer', 'on', 'older', 'is', 'prometheus', 'books', 'glenn', 'drive', 'buffalo', 'ny', 'african', 'americans', 'for', 'humanism', 'an', 'organization', 'promoting', 'black', 'secular', 'humanism', 'and', 'uncovering', 'the', 'history', 'of', 'black', 'freethought', 'they', 'publish', 'a', 'quarterly', 'newsletter', 'aah', 'examiner', 'write', 'to', 'norm', 'r', 'allen', 'jr', 'african', 'americans', 'for', 'humanism', 'p', 'o', 'box', 'buffalo', 'ny', 'united', 'kingdom', 'rationalist', 'press', 'association', 'national', 'secular', 'society', 'islington', 'high', 'street', 'holloway', 'road', 'london', 'n', 'ew', 'london', 'n', 'nl', 'british', 'humanist', 'association', 'south', 'place', 'ethical', 'society', 'lamb', 's', 'conduit', 'passage', 'conway', 'hall', 'london', 'w', 'n', 'rh', 'red', 'lion', 'square', 'london', 'w', 'n', 'rl', 'fax', 'the', 'national', 'secular', 'society', 'publish', 'the', 'freethinker', 'a', 'monthly', 'magazine', 'founded', 'in', 'germany', 'ibka', 'e', 'v', 'internationaler', 'bund', 'der', 'konfessionslosen', 'und', 'atheisten', 'postfach', 'd', 'berlin', 'germany', 'ibka', 'publish', 'a', 'journal', 'miz', 'materialien', 'und', 'informationen', 'zur', 'zeit', 'politisches', 'journal', 'der', 'konfessionslosen', 'und', 'atheisten', 'hrsg', 'ibka', 'e', 'v', 'miz', 'vertrieb', 'postfach', 'd', 'berlin', 'germany', 'for', 'atheist', 'books', 'write', 'to', 'ibdk', 'internationalen', 'b', 'ucherdienst', 'der', 'konfessionslosen', 'postfach', 'd', 'hannover', 'germany', 'telephone', 'books', 'fiction', 'thomas', 'm', 'disch', 'the', 'santa', 'claus', 'compromise', 'short', 'story', 'the', 'ultimate', 'proof', 'that', 'santa', 'exists', 'all', 'characters', 'and', 'events', 'are', 'fictitious', 'any', 'similarity', 'to', 'living', 'or', 'dead',

Figure 14. One of the omitted result of reading file

2) Test for removing stopwords from the text collection

```
73 #testing stopwords
74 for word in stopwords:
75     if word in wd_alt[0]:
76         print("There is still a stopword in the list")
77     print("Test is finished")
78     for i in range(len(wd_talk))

Run: main x
C:\Users\raonm\anaconda3\envs\pythonProject\python.exe C:/Users/raonm/PycharmProjects/pythonProject/main.py
Test is finished
```

Figure 15. Example of testing to see if stopwords are removed from the file

3) Code for performing word stemming to remove the word suffix

```
97 #Testing stemming
98 print(wd_alt[0])
99
100 #creating set for unique word appear even once in one of files in dataset
101 # for alt.atheism

Run: main x
C:\Users\raonm\anaconda3\envs\pythonProject\python.exe C:/Users/raonm/PycharmProjects/pythonProject/main.py
['mathew', 'mathew', 'manti', 'uk', 'subject', 'alt', 'atheism', 'faq', 'atheist', 'resourc', 'summari', 'book', 'address', 'music', 'atheism', 'keyword', 'faq', 'atheism',
'book', 'music', 'fiction', 'address', 'contact', 'expir', 'thu', 'apr', 'gmt', 'distribut', 'organ', 'manti', 'consult', 'cambridg', 'uk', 'supersed', 'manti', 'uk',
'line', 'archiv', 'atheism', 'resourc', 'alt', 'atheism', 'archiv', 'resourc', 'modifi', 'decemb', 'version', 'atheist', 'resourc', 'address', 'atheist', 'organ', 'usa',
'freedom', 'religion', 'foundat', 'darwin', 'fish', 'bumper', 'sticker', 'assort', 'atheist', 'paraphernalia', 'freedom', 'religion', 'foundat', 'write', 'ffrf', 'box',
'madison', 'wi', 'telephon', 'evolut', 'design', 'evolut', 'design', 'sell', 'darwin', 'fish', 'fish', 'symbol', 'christian', 'stick', 'car', 'feet', 'word', 'darwin',
'written', 'delux', 'mould', 'plastic', 'fish', 'postpaid', 'write', 'evolut', 'design', 'laurel', 'canyon', 'north', 'hollywood', 'peopl', 'san', 'francisco', 'bay',
'area', 'darwin', 'fish', 'lynn', 'gold', 'mail', 'figmo', 'netcom', 'net', 'peopl', 'lynn', 'price', 'fish', 'american', 'atheist', 'press', 'aap', 'publish', 'atheist',
'book', 'critiqu', 'bibl', 'list', 'bible', 'contradict', 'book', 'bibl', 'handbook', 'ball', 'foot', 'american', 'atheist', 'press', 'isbn', 'edit', 'bibl', 'contradict',
'absurd', 'atroc', 'immor', 'ball', 'foot', 'bibl', 'contradict', 'aap', 'base', 'king', 'jame', 'version', 'bibl', 'write', 'american', 'atheist', 'press', 'box', 'austin',
'tx', 'cameron', 'road', 'austin', 'tx', 'telephon', 'fax', 'prometheu', 'book', 'sell', 'book', 'includ', 'haught', 'holi', 'horror', 'write', 'east', 'amherst', 'street',
'buffalo', 'york', 'telephon', 'altern', 'address', 'newer', 'older', 'prometheu', 'book', 'glenn', 'drive', 'buffalo', 'ny', 'african', 'american', 'human', 'organ',
'promot', 'black', 'secular', 'human', 'uncov', 'histori', 'black', 'freethought', 'publish', 'quarterli', 'newslett', 'aah', 'examin', 'write', 'norm', 'allen', 'jn',
'african', 'american', 'human', 'box', 'buffalo', 'ny', 'unit', 'kingdom', 'rationalist', 'press', 'associ', 'nation', 'secular', 'societi', 'islington', 'high', 'street',
'holloway', 'road', 'london', 'ew', 'london', 'nl', 'british', 'humanist', 'associ', 'south', 'place', 'ethic', 'societi', 'lamb', 'conduit', 'passag', 'conway', 'hall',
'london', 'wc', 'rh', 'red', 'lion', 'suar', 'london', 'wc', 'rl', 'fax', 'nation', 'secular', 'societi', 'publish', 'freethink', 'monthli', 'magazin', 'found', 'germani',
'ibka', 'international', 'bund', 'der', 'konfessionslosen', 'und', 'atheisten', 'postfach', 'berlin', 'germani', 'ibka', 'publish', 'journal', 'miz', 'materialien', 'und',
'informationen', 'zur', 'zeit', 'politisch', 'journal', 'den', 'konfessionslosen', 'und', 'atheisten', 'hrsg', 'ibka', 'miz', 'vertrieb', 'postfach', 'berlin', 'germani',
```

Figure 16. Example of testing to see the result of stemming

4) Test for computing the frequency of each word in each document

```
{
  'pioneer': 0, 'wilder': 0, 'bloc': 1, 'famou': 3, 'continuum': 0, 'nirvana': 0, 'atroc': 0, 'oster': 0, 'jato': 0, 'mrc': 0, 'coterill': 0, 'hparc': 0, 'traditionalist': 0,
  'tfarrel': 0, 'age': 0, 'nikol': 0, 'sach': 0, 'rigbi': 0, 'blunt': 0, 'whibbard': 0, 'coupl': 0, 'profici': 0, 'slaughter': 0, 'jmc': 0, 'deanza': 0, 'ofa': 0, 'floppi': 0,
  'litttl': 0, 'combust': 0, 'fulful': 0, 'ema': 0, 'leno': 0, 'havoc': 0, 'saturn': 0, 'bernoulli': 0, 'pertain': 0, 'dagger': 0, 'vgx': 0, 'peripheri': 0, 'idiot': 0,
  'hip': 0, 'ipmo': 0, 'columbia': 0, 'eal': 0, 'proceed': 0, 'appeas': 0, 'vituper': 0, 'curcuit': 0, 'milwaulke': 0, 'lambert': 0, 'rambl': 0, 'bend': 0, 'wbg': 0,
  'hermeneia': 0, 'unli': 0, 'merkun': 0, 'charism': 0, 'admin': 0, 'horsemen': 0, 'csfr': 0, 'harel': 0, 'ravi': 0, 'national': 0, 'idiosyncrat': 0, 'origin': 4, 'ey': 0,
  'reactionari': 0, 'aim': 0, 'batdud': 0, 'minimum': 0, 'myth': 1, 'nonlinear': 0, 'arai': 0, 'delauney': 0, 'white': 0, 'test': 3, 'emd': 0, 'uniac': 0, 'craze': 0, 'cvax': 0,
  'gj': 0, 'yuv': 0, 'lane': 0, 'demograph': 0, 'germani': 0, 'usc': 0, 'hypercard': 0, 'twin': 0, 'onramp': 0, 'mention': 2, 'extrateritori': 0, 'ufo': 0, 'bandwith': 0,
  'yearli': 0, 'heidi': 0, 'maxin': 1, 'dinkl': 0, 'unuseful': 0, 'zurlo': 0, 'noncar': 0, 'def': 0, 'hkuxb': 0, 'msl': 0, 'indirectli': 0, 'poli': 0, 'micron': 0, 'conferr': 0,
  'fjssl': 0, 'schechner': 0, 'skcla': 0, 'otago': 0, 'pleasur': 0, 'chasiabl': 0, 'stein': 0, 'kelli': 0, 'glide': 0, 'quantiz': 0, 'discourag': 2, 'interven': 0, 'int': 0,
  'hadith': 0, 'sieferman': 0, 'doodad': 0, 'stallion': 0, 'necessar': 0, 'maintan': 0, 'uta': 0, 'conjoin': 0, 'reluctantli': 0, 'retina': 0, 'deflat': 0, 'sandvik': 0,
  'fnq': 0, 'roast': 0, 'allegori': 0, 'wrench': 0, 'hitl': 0, 'cannib': 1, 'oddjeb': 0, 'kenworth': 0, 'molest': 0, 'militia': 0, 'calt': 0, 'stagger': 0, 'gaia': 0, 'ripbc': 0,
  'qvabj': 0, 'geenwood': 0, 'nook': 0, 'innd': 0, 'clariti': 0, 'bundi': 0, 'poetic': 0, 'jsledd': 0, 'levin': 0, 'gyro': 0, 'gamer': 0, 'judson': 0, 'bd': 0, 'recis': 0,
  'kilroy': 0, 'egl': 0, 'compel': 0, 'shaki': 0, 'jogl': 0, 'beyt': 0, 'unsteadl': 0, 'alaramon': 0, 'noturn': 0, 'jadaley': 0, 'vaytek': 0, 'corel': 0, 'automata': 0,
  'gatt': 0, 'possibl': 1, 'degener': 0, 'billion': 0, 'exempl': 1, 'locker': 0, 'plasma': 0, 'photograph': 0, 'safeguard': 0, 'shape': 0, 'uark': 0, 'mip': 0, 'freewan': 0,
  'scottsdal': 0, 'hitler': 0, 'ahpcrc': 0, 'gyroscop': 0, 'zeppelin': 0, 'peck': 0, 'virago': 0, 'pdfa': 0, 'straightway': 0, 'liaison': 0, 'requin': 1, 'keg': 0,
  'jbuddenberg': 0, 'unfortunat': 0, 'wrenc': 0, 'tniy': 0, 'christi': 0, 'pasqu': 0, 'devin': 0, 'favourit': 0, 'crown': 0, 'terraman': 0, 'condom': 1, 'personifi': 0,
  'bowman': 0, 'tendentl': 0, 'aix': 0, 'lineberri': 0, 'loc': 0, 'mutton': 0, 'magistr': 0, 'aar': 0, 'cix': 0, 'deros': 0, 'uki': 0, 'lube': 0, 'salett': 0, 'hysteria': 0,
  'newbi': 0, 'plausibl': 0, 'minitel': 0, 'pekka': 0, 'humanist': 0, 'uaccess': 0, 'microimag': 0, 'rusti': 0, 'kr': 0, 'coyot': 0, 'poorer': 0, 'univer': 0, 'elicit': 0,
  'planter': 0, 'outward': 0, 'zakat': 0, 'denmark': 0, 'bcm': 0, 'muhammadi': 0, 'illustr': 0, 'incompat': 1, 'gulp': 0, 'unbibl': 0, 'catnip': 0, 'kohn': 0, 'vice': 0,
  'balgley': 0, 'collabor': 0, 'strickl': 0, 'simi': 0, 'sbc': 0, 'underco': 0, 'blasphem': 0, 'gravi': 0, 'dual': 0, 'bkyast': 0, 'gonzaga': 0, 'describeth': 0,
  'preternatur': 0, 'reason': 5, 'rz': 0, 'nonexclus': 0, 'jump': 0, 'vtl': 0, 'erc': 0, 'bribe': 0, 'gener': 0, 'restrict': 1, 'drakula': 0, 'sculptura': 0, 'trimest': 0,
  'scratch': 0, 'ken': 0, 'haefner': 0, 'susan': 0, 'arrang': 0, 'warmong': 0, 'hki': 0, 'kevin': 0, 'trw': 0, 'sidecar': 0, 'reson': 0, 'erasmu': 0, 'miner': 0,
  'intertestament': 0, 'nnp': 0, 'centurion': 0, 'proseletyz': 0, 'icgeb': 0, 'shirley': 0, 'ukma': 0, 'xtian': 0, 'coastlin': 0, 'mmlson': 0, 'music': 0, 'newscast': 0,
  'architect': 0, 'nadja': 0, 'harpo': 0, 'nazaren': 0, 'unshroud': 0, 'flynn': 0, 'thakkar': 0, 'covington': 0, 'borean': 0, 'fare': 0, 'actionn': 0, 'enabl': 0, 'level': 1,
}
```

Figure 17. One of the omitted result of the frequency of words in file

5) Test for computing document frequency

```
{
  'af': 10, 'doubtless': 3, 'weatherhead': 1, 'mistransl': 4, 'mess': 25, 'eola': 2, 'spuriou': 1, 'munsch': 1, 'davi': 28, 'geometri': 20, 'cow': 2, 'barrier': 5, 'indiana':
  37, 'bummin': 1, 'gradat': 2, 'handshak': 1, 'mysteri': 29, 'screen': 50, 'lsd': 9, 'agan': 2, 'hinder': 2, 'macalstr': 8, 'supercomput': 4, 'isaiah': 20, 'coon': 2,
  'stuba': 1, 'erwin': 3, 'excess': 16, 'tice': 1, 'cafeteria': 2, 'cvax': 3, 'joon': 2, 'benett': 3, 'drown': 3, 'believ': 180, 'loec': 1, 'chavey': 1, 'increas':
  76, 'flavor': 4, 'devoid': 6, 'baud': 4, 'yusef': 1, 'identif': 5, 'qla': 1, 'resourc': 46, 'pup': 2, 'gentl': 7, 'dcrt': 6, 'zowi': 1, 'diff': 2, 'teledetect': 1,
  'potteri': 3, 'absenc': 2, 'leonard': 11, 'racelist': 1, 'lwb': 1, 'mosaic': 6, 'sfu': 5, 'outcom': 16, 'williamson': 3, 'talud': 1, 'quieter': 3, 'mathematik': 1,
  'mouth': 40, 'colorview': 2, 'driver': 78, 'penicilin': 1, 'zephyr': 9, 'mash': 2, 'wr': 1, 'robert': 149, 'smuggl': 1, 'influenc': 35, 'esqu': 3, 'wondrou': 3, 'jtozer': 1,
  'steinberg': 1, 'fali': 1, 'gradient': 5, 'undefend': 1, 'henriksen': 1, 'distraught': 1, 'baggi': 2, 'ucr': 1, 'onan': 1, 'consortium': 4, 'waist': 1, 'eldridg': 1, 'box':
  67, 'age': 66, 'cbrook': 5, 'inflect': 3, 'crumbl': 2, 'hollasch': 2, 'undertaken': 1, 'kindergarten': 2, 'absenc': 28, 'fehr': 1, 'soni': 4, 'veino': 1, 'uncontrovert': 1,
  'quantiz': 5, 'malpractic': 4, 'kicker': 1, 'drake': 1, 'uencod': 5, 'stabilis': 1, 'ervan': 1, 'lkin': 2, 'icon': 4, 'tuition': 4, 'disappear': 23, 'prfi': 1, 'cogn': 2,
}
```

Figure 18. The omitted result of the document frequency

6) Design for computing document-by-word matrix

	hsu	unlov	folkish	grizzli	rumor	...	romeo	caputo	poj	pacifist	cactu
0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
...
2721	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
2722	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
2723	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
2724	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
2725	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
[2726 rows x 23130 columns]											

Figure 19. The result of document-by-word matrix as DataFrame for precise view

```

[[0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 ...
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]
 [0. 0. 0. ... 0. 0. 0.]]

```

Figure 20. Corresponding matrix(numpy array) of DataFrame above

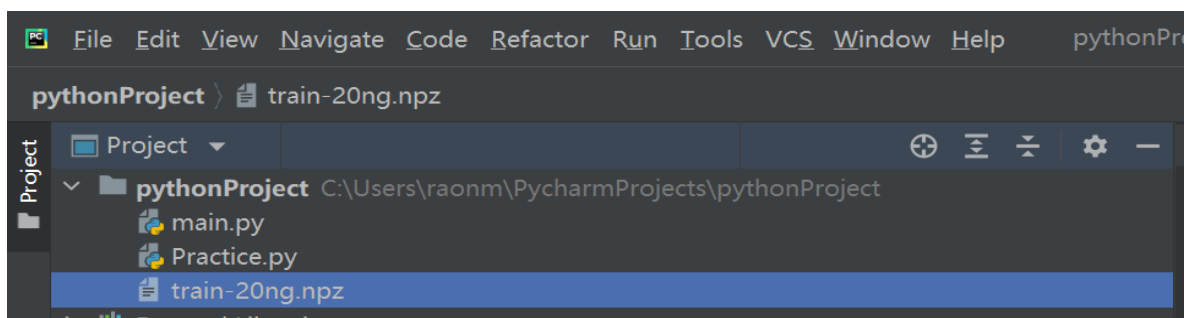


Figure 21. The result of saving final document-by-word matrix

Conclusion

The feature generation done in this lab, which is generating the features from the text collections is important in text analysis. For example, it can be used for extracting important words in particular document, or comparing similarity between documents. In the process of feature generation, various lists are used to save names of subdirectories, and files, and words that are split from each document. In this lab, when stopwords are removed from the words in list, 'List Comprehension' should be used. Since using any type of loops to work for list in python takes risk because the result not expected might be got. For example, if a list contained two 'e' only first 'e' will be removed. When unique words are obtaining, 'set' should be used because it automatically deletes duplicate words. Then, this set will be convert into list to keep index of words. For steps of computing TF-IDF, every file should be gone through by using proper lists and indexing of it. After all, for document-by-word matrix, list containing inner dictionaries for each file of A_{ik} will be converted into DataFrame and then that DataFrame will be converted into numpy array. Therefore, A_{ik} values can be stored into 2D array. Moreover, if the list of unique word is printed above the final 2D array, word can be seen more easily.