

# IST 664 – Natural Language Processing

## Homework 1

### Abstract

Corpus is a collection of written text. For the purpose of this report we are using the corpus text obtained from news/message boards/blogs about Corona Virus. Analysis of this corpus can lead to interesting findings about the impact of COVID-19. There are various ways to analyze a text corpus. The problem here is to understand people's views and their concerns, and things they care about.

### Data

The given data is the free dataset from news/message boards/blogs about CoronaVirus. This dataset has four months data – 5.2 M posts. The time frame of the data is December 2019 – March 2020. The posts are in English mentioning at least one of the following: "Covid" or CoronaVirus or "Corona Virus".

### Data Preprocessing

We first read the data from both the JSON files and extract the required fields 'facebook', 'title', 'published', 'replies\_count', 'author', 'url', 'country', 'text'. The required data is stored into a CSV file for further analysis.

For the purpose of corpus analysis, the text data present in text field would be the only one useful. For this reason, we extract the data and change it to a string. The following lines of code are used for the same:

```
#reading the data from csv
import pandas as pd
covid = pd.read_csv('test.csv')

#changing text data to string
str1 = ''.join(list(covid['text']))
```

### Tokenization

The assembly language is a language of numbers unlike human languages which contains words. Thus, the only thing computer can do with the human language is counting. To initiate this process of counting the text is divided or is grouped into words. This process of dividing/grouping the text is known as tokenization.

There are plenty tokenizers provided by NLTK which are available for us to use. For tokenizing the text, we can use the RegexpTokenizer available in the NLTK's tokenize library, which is a Regular Expression tokenizer. The text data contains symbols like \*, :, etc. which are of no use

thus, we use RegularExpression tokenizer to consider only the words and numbers in the text. The following lines of code were used for the same:

```
#tokenizing using RegexpTokenizer by considering only alphabets and numbers  
  
import nltk  
  
tokenizer = nltk.RegexpTokenizer('\w+')  
r_tokens = tokenizer.tokenize(str1)
```

On evaluating the tokens, we see that the tokenizer is considering same words written in different cases as different tokens. For example, “The” and “the” are being considered as different tokens. Thus, to resolve this problem we can change all of the tokens to lowercases. These will affect the frequencies of the tokens and also further analysis of those tokens. This is implemented using the following line of code:

```
#changing all the tokens to lower case  
  
tokens = [word.lower( ) for word in r_tokens]
```

### ***Removing stop words***

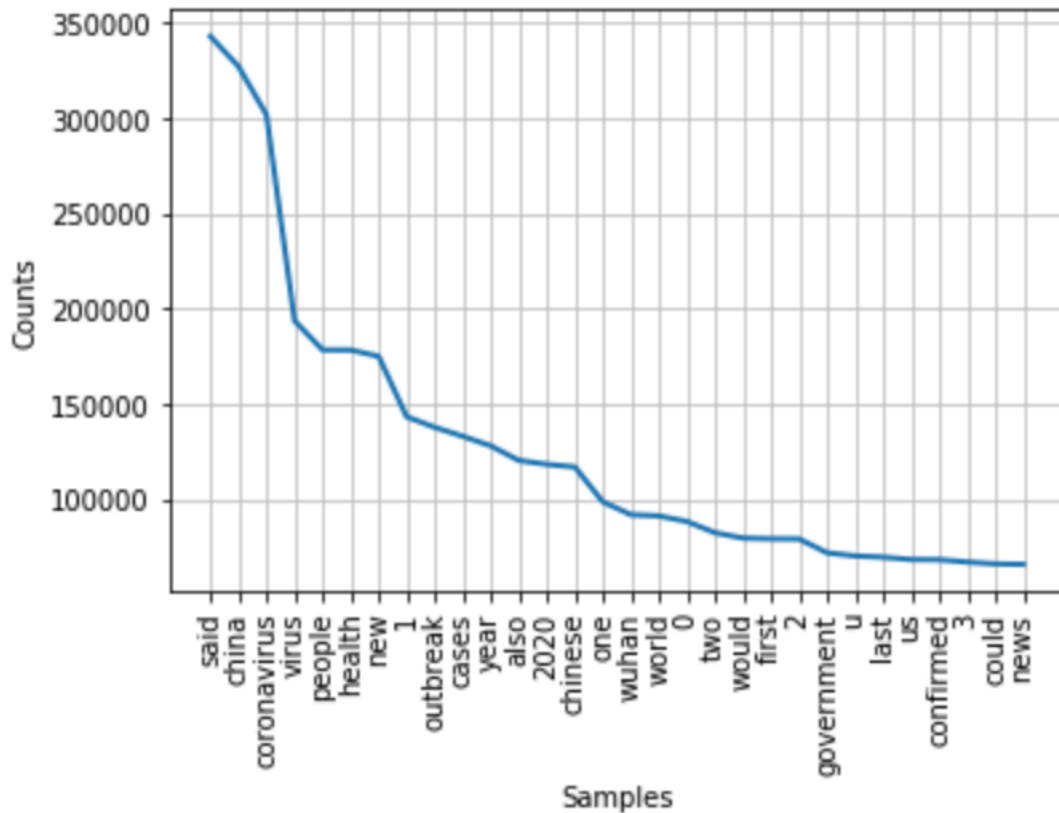
In computing, stop words are words which are filtered out before processing of natural language data. Stop words are generally the most common words in a language; there is no single universal list of stop words used by all natural language processing tools, and indeed not all tools use such a list. Some tools avoid removing stop words to support phrase search. (Reference Wikipedia)

To get a better insight of the reviews I have removed all the stop words from tokens predefined as stopwords in English language. Line of code to remove the stop words:

```
#removing stop words  
  
from nltk.corpus import stopwords  
  
stop_words = set(stopwords.words('english'))  
covid_tokens = [token for token in tokens if token not in stop_words]
```

## **Data Analysis**

We have a total of 106017 different data regarding Corona Virus from 149 countries. ‘said’, ‘china’ and ‘coronavirus’ are the most frequent words in the corpus.



The aim to clean and process the review text is to analyze it and derive meaningful insights from it. The analysis has been performed in three forms:

- i) listing out the most commonly occurring words
  - ii) listing out the most commonly occurring bigrams and
  - iii) listing out bigrams by their mutual information score
- in descending order, occurring at-least a specific number of times.

The top 50 words with highest frequencies are obtained with the help of following block of code:

```
#top 50 words by frequency
freq_dist = nltk.FreqDist(covid_tokens)
most_common_50 = freq_dist.most_common(50)

for w in most_common_50:
    print(w)
```

The FreqDist function generates a list of tuples as (word, frequency) and the most\_common(50) lists out the 50 most frequent tokens.

The 50 most frequent bigrams can be obtained using the following block of code:

```
#top 50 bigrams by frequencies

bigrams = list(nltk.bigrams(tokens))
bigram_freq_dist = nltk.FreqDist(bigrams)
bigram_most_common_50 = bigram_freq_dist.most_common(50)

for w in bigram_most_common_50:
    print(w)
```

The FreqDist function generates a list of tuples as (word1\_of\_bigram, word2\_of\_bigram, frequency) and the most\_common(50) lists out the 50 most frequent tokens.

For listing out bigrams according to their mutual information score in descending order the following code is used:

```
#top 50 bigrams by their mutual information scores (using min frequency 5)

bigram_measures = nltk.collocations.BigramAssocMeasures()
finder = nltk.BigramCollocationFinder.from_words(tokens)
finder.apply_freq_filter(5)
scored = finder.score_ngrams(bigram_measures.pmi)

for bscore in scored[:50]:
    print(bscore)
```

## Results

The following are the results obtained for top 50 tokens according to their frequencies, top 50 bigrams according to their frequencies and top 50 bigrams according to their mutual information score with minimum occurrence as 5.

('said', 343052)  
 ('china', 327291)  
 ('coronavirus', 301838)  
 ('virus', 193760)  
 ('people', 178459)  
 ('health', 178421)  
 ('new', 175110)  
 ('1', 143353)  
 ('outbreak', 137832)  
 ('cases', 133120)  
 ('year', 128090)  
 ('also', 120448)  
 ('2020', 118372)  
 ('chinese', 117007)  
 ('one', 98708)  
 ('wuhan', 91962)  
 ('world', 91307)  
 ('0', 88367)  
 ('two', 82481)  
 ('would', 79658)  
 ('first', 79299)  
 ('2', 79221)  
 ('government', 71936)  
 ('u', 70285)  
 ('last', 69631)  
 ('us', 68420)  
 ('confirmed', 68323)  
 ('3', 67094)  
 ('could', 66108)  
 ('news', 65892)  
 ('time', 65819)  
 ('spread', 64832)  
 ('million', 62400)  
 ('text', 61747)  
 ('global', 61337)  
 ('public', 61128)  
 ('000', 61126)  
 ('week', 61048)  
 ('reported', 60711)  
 ('disease', 59149)  
 ('day', 58627)  
 ('february', 58450)  
 ('feb', 56897)  
 ('ship', 56695)  
 ('2019', 55762)  
 ('5', 55127)  
 ('company', 54991)  
 ('like', 54876)  
 ('10', 53674)  
 ('including', 53504)

Top 50 tokens

(('of', 'the'), 354918)  
 (('in', 'the'), 254326)  
 (('to', 'the'), 149933)  
 (('the', 'virus'), 126370)  
 (('on', 'the'), 113739)  
 (('for', 'the'), 111931)  
 (('the', 'coronavirus'), 97023)  
 (('in', 'china'), 86148)  
 (('from', 'the'), 84900)  
 (('with', 'the'), 82180)  
 (('and', 'the'), 82076)  
 (('at', 'the'), 81219)  
 (('to', 'be'), 74476)  
 (('have', 'been'), 74318)  
 (('in', 'a'), 73525)  
 (('that', 'the'), 68879)  
 (('more', 'than'), 66486)  
 (('u', 's'), 65138)  
 (('the', 'world'), 65044)  
 (('by', 'the'), 59970)  
 (('the', 'outbreak'), 52364)  
 (('it', 's'), 51748)  
 (('will', 'be'), 51583)  
 (('has', 'been'), 51294)  
 (('china', 's'), 48313)  
 (('said', 'the'), 46911)  
 (('according', 'to'), 46475)  
 (('of', 'a'), 46348)  
 (('number', 'of'), 45618)  
 (('the', 'new'), 44666)  
 (('coronavirus', 'outbreak'), 44558)  
 (('as', 'the'), 42657)  
 (('hong', 'kong'), 42528)  
 (('the', 'first'), 40684)  
 (('he', 'said'), 40560)  
 (('it', 'is'), 39940)  
 (('to', 'a'), 39653)  
 (('novel', 'coronavirus'), 39132)  
 (('due', 'to'), 38710)  
 (('is', 'a'), 37397)  
 (('as', 'a'), 36539)  
 (('the', 'country'), 35637)  
 (('the', 'u'), 35462)  
 (('the', 'chinese'), 35160)  
 (('about', 'the'), 34608)  
 (('covid', '19'), 34453)  
 (('the', 'company'), 33698)  
 (('for', 'a'), 30624)  
 (('with', 'a'), 29616)  
 (('over', 'the'), 29584)

Top 50 bigrams

(('1084673', 'impact\_coronavirus\_on\_the\_steel\_industry'), 23.454796082453907)  
 (('1087991', 'mmc\_industrial\_supply'), 23.454796082453907)  
 (('1087993', 'mmc\_monitoring'), 23.454796082453907)  
 (('1087995', 'mmc\_missions'), 23.454796082453907)  
 (('deserthottie', 'qtanoni'), 23.454796082453907)  
 (('4346', 'bce3'), 23.454796082453907)  
 (('516bc8cb', 'b44e'), 23.454796082453907)  
 (('5384', '5421'), 23.454796082453907)  
 (('7771', '7805'), 23.454796082453907)  
 (('84883561', '联系'), 23.454796082453907)  
 (('9817', '9858'), 23.454796082453907)  
 (('aff25', 'halston'), 23.454796082453907)  
 (('aline', 'oyamada'), 23.454796082453907)  
 (('andreo', 'calonzo'), 23.454796082453907)  
 (('artha', 'ardhana'), 23.454796082453907)  
 (('asaduddin', 'owaisi'), 23.454796082453907)  
 (('b44e', '4346'), 23.454796082453907)  
 (('babulal', 'marandi'), 23.454796082453907)  
 (('bahceli', 'kktc'), 23.454796082453907)  
 (('bce3', '0659080e396b'), 23.454796082453907)  
 (('billca', 'deserthottie'), 23.454796082453907)  
 (('bissell', 'linsk'), 23.454796082453907)  
 (('boulardi', 'toti'), 23.454796082453907)  
 (('celestino', 'gallares'), 23.454796082453907)  
 (('chiwuke', 'onyeanu'), 23.454796082453907)  
 (('combizym', 'hirudoid'), 23.454796082453907)  
 (('derhal', 'istifa'), 23.454796082453907)  
 (('ekstra', 'bladet'), 23.454796082453907)  
 (('eren', 'sengezer'), 23.454796082453907)  
 (('ezzedin', 'bahader'), 23.454796082453907)  
 (('filou', 'oostende'), 23.454796082453907)  
 (('foodgrain', 'production'), 23.454796082453907)  
 (('galtung', 'dovig'), 23.454796082453907)  
 (('ghanti', 'bajao'), 23.454796082453907)  
 (('gonzalez', 'garcoa'), 23.454796082453907)  
 (('hatidze', 'muratova'), 23.454796082453907)  
 (('heikki', 'kovalainen'), 23.454796082453907)  
 (('intravascular', 'coagulopathy'), 23.454796082453907)  
 (('iodinepowerpack', 'thumbnail\_1'), 23.454796082453907)  
 (('istifa', 'etmeli'), 23.454796082453907)  
 (('jeoffrey', 'lamujon'), 23.454796082453907)  
 (('jirasat', 'wittaya'), 23.454796082453907)  
 (('joventut', 'badalona'), 23.454796082453907)  
 (('jinpíng', '习近平'), 23.454796082453907)  
 (('kalagayang', 'epidemiko'), 23.454796082453907)  
 (('kalidou', 'koulbaly'), 23.454796082453907)  
 (('kalyeena', 'makortoff'), 23.454796082453907)  
 (('kinimasa', 'mayama'), 23.454796082453907)  
 (('kkl', 'jnf'), 23.454796082453907)  
 (('kktc', 'cumhurbařkani'), 23.454796082453907)

Top 50 bigrams (pmi)

As we can see in the above images most of the reviews are about china, coronavirus, health, outbreak. Also, from bigrams we can detect that most of the people spoke about virus and its origin. Furthermore, we can also see that the corpus doesn't follow the Zipf's law as frequency of china is not twice the frequency of coronavirus. For further evaluation I suggest performing sentiment analysis on tokens and bigrams and collectively determine people's view on COVID-19 and its spread. We should tokenize each post separately and generate a sentiment for each post as positive, neutral or negative.