



**Name:** Bhaumikkumar Patel

**Student ID:** w1813148

7BUIS025W

Web and Social Media Analytics

**Coursework:** Social Media Assignment (2020/21)

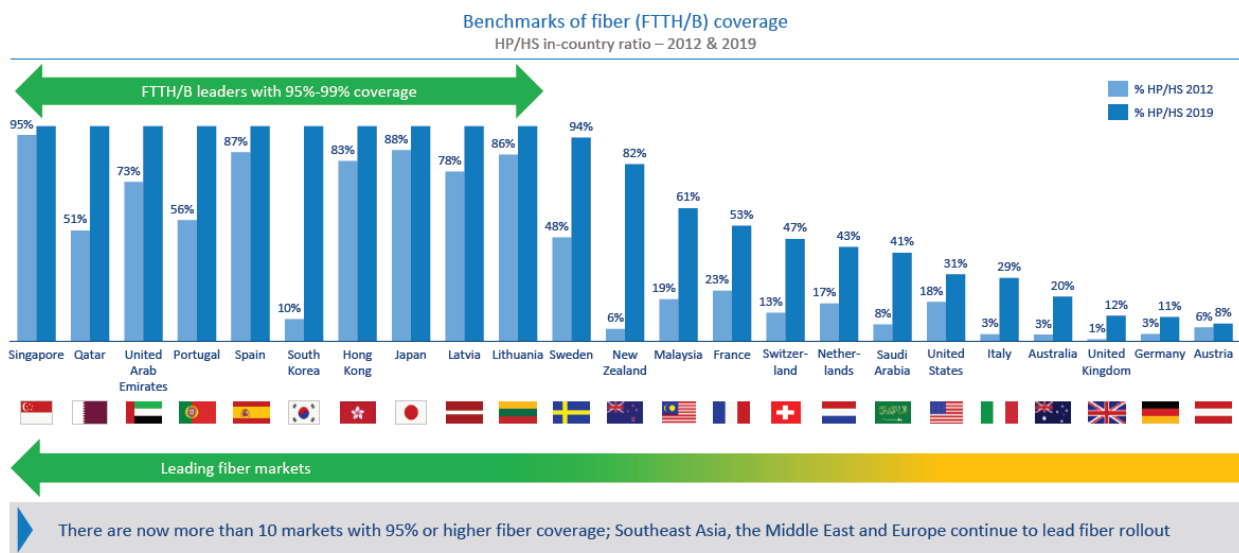
**WORD COUNT:** 1244

**Identify a brand or famous person of interest and create a list of 5 possible keywords you could use to identify relevant tweets from Twitter.**

I used the main keyword related to the Optical Fiber (Fibre in British English), it is a transparent and flexible made by glass of silica or plastic. Its diameter like slightly thicker than human hair. Optical Fibre are used for Telecom communication were transmitting over long distance as well as get higher bandwidth data transfer then electric cables. Fibers are sent data signals from one end to other end with less loss compared to the metal wires.

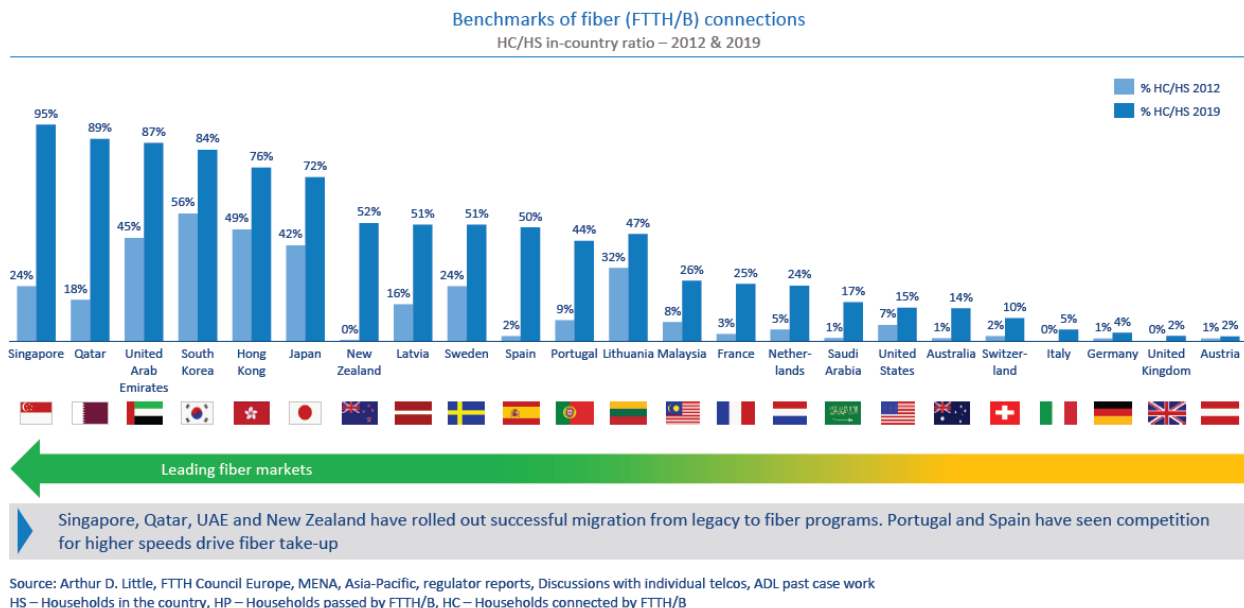
The first Optical fibre data transmission system demonstrated by German physicist Manfred Borner at Telefunken Research Labs in Ulm in 1965. The latest data transmission over fibre optic cable is 500 Gbit/s. [1] From around the world there are with 95% or higher fibre coverage; southeast Asia, the middle east and Europe continue to lead fibre rollout. On the other side Singapore, Qatar, UAE and New Zealand have rolled out successful migration. Portugal and Spain have seen competition for higher speeds drive fibre take up. Below figure explain the coverage and connection from various countries.

Figure 1. FTTH/B households Passed [2]



Source: Arthur D. Little, FTTH Council Europe, MENA, Asia-Pacific, regulator reports, Discussions with individual telcos, ADL past case work  
HS – Households in the country, HP – Households passed by FTTH/B, HC – Households connected by FTTH/B

Figure 2. Fibre (FTTH/B) households connected [2]



I used five keywords as follows “fibre”, “fiberinternet”, “Internet”, “Broadband”, “Gigabyte”. These keywords are used by the service providers to explain new technology to the customer as well as the technicians tweets their working ethics.

***Explain what is meant by an API and compare and contrast the two data collection APIs available on the Twitter platform.***

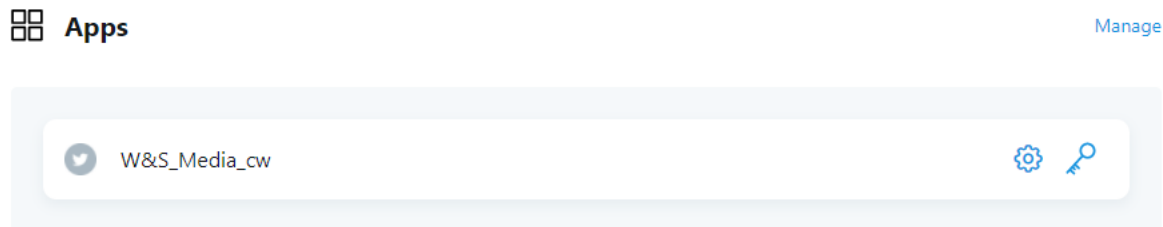
The API stands for Application Programming Interface, lets us write and read twitter data. We can use it for read profile and access our followers data, a high volume of tweets on particular subjects in specific locations.

Imagine, we are at the restaurant on the table, ordering food from the kitchen. Waiter is between kitchen and us who confirmed that kitchen staff gets all orders and meals provided to right customers. We are also free to know about items available before the order. Think about the waiter role, if we don’t have waiter in the restaurants. We need to push our self very hard to the kitchen for the food. That’s where the Application Programming Interface plays their role as a middle man between two diverse software systems or programs.

There are two types of API’s twitter offer as Historical tweets and the other one is Stream Live tweets. The diverse is innate just by their names. The Stream Live, using this tweets API we can takes continuous flow of tweets at the exact moment and record them. The Historical tweets allow us to pick up some portion of the past tweets either a specific user tweets from a certain timeframe.

**Using your suggested keywords from part (a) and your knowledge of Twitter, collect a series of Tweets surrounding your chosen brand/famous person and save them to a file. Your collected Tweets should span a minimum ONE-week period.**

The first step is to apply for the Twitter developer account and to link to my Twitter account. After getting approval, I need to register my application from the Twitter Developer Website and create an app. I use my application name is W&S\_Media\_cw.



After that we need to imports the libraries in our coding. I use replit for this project. Following are the four libraries that I am used in this program. “codecs” library allows python program to encode and decode text for diverse representations. “tweepy” library allows python program to utilize the twitter API. “os” and “sys” libraries are helps to determine the paths, directories and filenames.

```
1  import codecs
2  import tweepy
3  import os
4  import sys
```

This program used for the authentication each request via Twitter API by using specific keys. For the security reason, I hide the keys.

```
5
6  #I hidden my API keys for security reason
7  CONSUMER_KEY = "*****HIDE CONSUMER_KEY*****"
8  CONSUMER_SECRET = "*****HIDE CONSUMER_SECRET*****"
9  ACCESS_TOKEN = "*****HIDE ACCESS_TOKEN*****"
10 ACCESS_TOKEN_SECRET = "*****HIDE ACCESS_TOKEN_SECRET*****"
11
12 auth = tweepy.OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
13 auth.set_access_token(ACCESS_TOKEN, ACCESS_TOKEN_SECRET)
```

The following code will verify CONSUMER\_KEY and ACCESS\_TOKEN with the Twitter API. When this key valid and works then it will print the command “Authentication OK” but if is there any error occur in validation of this keys then it will print “Error during Authentication”

```
14
15 #Make sure that the credentials work
16 api = tweepy.API(auth)
17 try:
18     api.verify_credentials()
19     print("Authentication OK")
20 except:
21     print("Error during authentication")
22
```

Now, I need to save all the collected tweets using the file name “fibre\_tweets.txt”. I save my file with the “utf-8” encoding. I use one of the attributes “a” which is append instead of write “w” so whenever code runs then this file not overwrite but keep on writing from the last line. In some circumstances, there in no any file create by me then it will create one new file making append. Also, I used five keywords to extract the tweets from the tweeter are follows.

```
23 #to create and append the text in the file
24 f = codecs.open("fibre_tweets.txt", "a", encoding="utf-8")
25
26
27 # Define the search term and the date_since date as variables
28 keywords = ['fibre', 'fiberinternet', 'Internet', 'Broadband', 'Gigabyte']
29 |
30
```

Using the following code, I create a class that defines the tweepy StreamListener and save them in “f” that’s associated the create file name “fibre\_tewwts.txt” and the value of the tweets to collected from the Live stream. I collected 1000 tweets per day during six days period using five keywords mentioned in this file.

```

32 #tweet listener class
33 class BasicTwitterListener(tweepy.StreamListener):
34     def setup(self):
35         self.n_tweet = 0
36
37     def on_status(self, status):
38         try:
39             cleaned_tweet = status.text.replace("\n", " ")
40             if len(cleaned_tweet) < 3:
41                 return True
42             f.write(cleaned_tweet + "\n")
43             self.n_tweet += 1
44             print(self.n_tweet)
45
46             if self.n_tweet == 1000:
47                 f.close()
48                 return False
49         except Exception as e:
50             print("Exception when reading from stream:")
51             pass
52
53     def on_error(self, status_code):
54         print("Encountered error with status code: " + str(status_code))
55         return True
56
57     def on_timeout(self):
58         print("Timed out:")
59         return True
60

```

Using the following code, we must filter the collected tweets using the language only in English and utilized the keywords with the track commend

```

62 #assign the class to a variable
63 l = BasicTwitterListener()
64 l.setup()
65
66 #running the class with filter
67 live_stream = tweepy.streaming.Stream(auth, l)
68 live_stream.filter(languages=["en"], follow=None, track=keywords)

```

Running this code one-time each day during six days and 7<sup>th</sup> day I change code to collect 510 tweets, I will collect a total of 6510 tweets related to the Fibre Optics duration of 7 days.

Authentication OK

1  
2  
3  
4  
5  
6

Q x

## Files



main.py

fibre\_tweets.txt

### Packager files

poetry.lock

pyproject.toml

fibre\_tweets.txt

```
6496 RT @JayElHarris: It's also extremely weird how the internet treats people existing as
"representation" and feels the need to codify existen...
6497 RT @ndekekwe: Elon Musk's SpaceX arrives Nigeria with its Starlink. The team is working
with NCC, the industry regulator for permit. Expect...
6498 RT @itsMalikEl: 2007: "don't talk to strangers on the internet!" 2021: "do any mutuals
wanna fuck?"
6499 RT @amyklobuchar: In 2021, every family in America should have access to high-speed
internet – regardless of their ZIP code.
6500 RT @desithwan: defending faramir on the internet isn't enough, i need boromir
6501 RT @MilaEric1: Bad News about Nigeria & her Government allover the Internet/social
media. #BiafraExit .
6502 RT @confusedvichar: New character: Karl Gujjar from Noida. This video series is inspired
from a very fascinating YouTube and TikTok sub-g...
6503 RT @Quartzjixler: One of the most detrimental outcomes of the internet and social media is
that now people who aren't from Eastern Pennsylv...
6504 RT @StarksGrayson: Ultron is easily the best villain in the mcu, my man spent five minutes
on internet and decided the human race needed to...
6505 The coldest video on the internet https://t.co/Y6EfMoymg6
6506 RT @JayElHarris: It's also extremely weird how the internet treats people existing as
"representation" and feels the need to codify existen...
6507 RT @cmf_wright: How does your real-life self differ from your internet persona?
#WritingCommunity #Writers #WritersCafe #Writing
6508 @existingee @heyGuhRL Wow 2021 lol everyone is so soft now wth it's a joke if being mocked
on the internet is too m... https://t.co/q4VvcJTdx7
6509 RT @totalwoke: #MothersDay is just celebration of love. Bhaidooj is celebration of
brahminical patriarchy, misogyny, racism, slavery, polar...
6510 RT @itsMalikEl: 2007: "don't talk to strangers on the internet!" 2021: "do any mutuals
wanna fuck?"
6511
```

**Using a suitable example, discuss the role of text pre-processing in the context of social media analysis. Identify TWO pre-processing steps relevant to the dataset you created in part (c) and apply them to your dataset.**

We need to follow several steps for text pre-processing in the context of social media analysis. [1] filter the language [2] Blank spaces cleaning [3] Lowercase everything [4] Removing the URL [5] Removing the Stop words [6] Removing the Emojis [7] Removing Rt (retweet) [8] Removing the selected punctuations

Now we need to create another program for the pre-processing. First of all, import the libraries. "nltk" is the new libraries we use in this program to specified for updated stop words & language detect

```
1  import nltk
2  nltk.download("stopwords")
3  import codecs
4  from langdetect import detect
5  import re
6  import sys
7  from nltk.corpus import stopwords
8
```

After that need to filter languages of the collected tweets using below comments.

```
11  #checking if the language filter worked properly
12  def find_eng_tweet(filename):
13  |   with codecs.open(filename, "r", encoding = "utf-8") as f:
14  |       lines = list(f)
15  |       for line in lines[0:]:
16  |           if len(line) > 2:
17  |               try:
18  |                   lang = detect(line)
19  |                   if lang == "en":
20  |                       yield line
21  |               except:
22  |                   pass
23
```

Now the various pre-processing methods are mention in the program with "#"

```
24  #cleaning the excessive spacing
25  def clean_t(tweet):
26  |   return tweet.strip()
27
28  #lowercase everything
29  def lowercase_t(tweet):
30  |   return tweet.lower()
31
```



```

32 #remove the url
33 def remove_url(tweet):
34     w_keep = []
35     for word in tweet.split():
36         word = word.strip()
37         if word.startswith("http"):
38             w_keep.append("")
39         else:
40             w_keep.append(word)
41     return " ".join(w_keep)
42
43 #stopwords removal
44 def remove_stopwords(tweet):
45     stop_words = set(stopwords.words('english')) #["is", "a", "an", "in", "on", "am", "are",
46     "and", "i", "the"]
47     w_keep = []
48     for word in tweet.split():
49         word = word.strip()
50         if not word in stop_words:
51             w_keep.append(word)
52     return " ".join(w_keep)
53
54 #removing emojis
55 def remove_emoji(tweet):
56     emoji_pattern = re.compile("[
57         u"\U0001F600-\U0001F64F" # emoticons
58         u"\U0001F300-\U0001F5FF" # symbols & pictographs
59         u"\U0001F680-\U0001F6FF" # transport & map symbols
60         u"\U0001F1E0-\U0001F1FF" # flags (iOS)
61         u"\U00002500-\U00002BEF" # chinese char
62         u"\U00002702-\U000027B0"
63         u"\U00002702-\U000027B0"
64         u"\U000024C2-\U0001F251"
65         u"\U0001F926-\U0001F937"
66         u"\U00010000-\U0010ffff"
67         u"\u2640-\u2642"
68         u"\u2600-\u2B55"
69         u"\u200d"
70         u"\u23cf"
71         u"\u23e9"
72         u"\u231a"
73         u"\ufe0f" # dingbats
74         u"\u3030"
75         "]+", flags=re.UNICODE)
76     return emoji_pattern.sub(r"", tweet)

```

```

77 #removing rt
78 def remove_rt(tweet):
79     w_keep = []
80     for word in tweet.split():
81         word = word.strip()
82         if word != "rt":
83             w_keep.append(word)
84     return " ".join(w_keep)
85
86 #removing the punctuations
87 def remove_punctuation(tweet):
88     punctuations = '!";:","\,<>()[]{}€$./?%^*_~'...'
89     no_punct = ""
90     for char in tweet:
91         if char not in punctuations:
92             no_punct = no_punct + char
93     return "".join(no_punct)
94
95 #calling all the other definition
96 def processed_tweets(filename):
97     for tweet in find_eng_tweet(filename):
98         tweet = clean_t(tweet)
99         tweet = lowercase_t(tweet)
100         tweet = remove_url(tweet)
101         tweet = clean_t(tweet)
102         tweet = remove_stopwords(tweet)
103         tweet = remove_emoji(tweet)
104         tweet = remove_rt(tweet)
105         tweet = remove_punctuation(tweet)
106         yield tweet
107
108 file_in = "fibre_tweets.txt"
109 file_out = "fibretweets_processed.txt"
110
111 #exporting everything on another file
112 with codecs.open(file_out, "w", encoding = "utf-8") as f:
113     for tweet in processed_tweets(file_in):
114         f.write(tweet + "\n")
115
116 print("Process Complete")
117
118 stop_words = set(stopwords.words('english'))
119 print(stop_words)

```

---

For running this program, we need to install some components from the nltk using the shell commend prompt from the replit so that we can success and upgrade the pip

```
~/w1813148CWP02$ pip install numpy
Requirement already satisfied: numpy in /opt/virtualenvs/python3/lib/python3.8/site-packages (1.20.2)
WARNING: You are using pip version 19.3.1; however, version 21.1.1 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
~/w1813148CWP02$ pip install --upgrade pip
Collecting pip
  Downloading https://files.pythonhosted.org/packages/cd/6f/43037c7bcc8bd8ba7c9074256b1a11596daa15555808ec748048c1507f08/pip-21.1.1-py3-none-any.whl (1.5MB)
    | ████████████████████ | 1.6MB 3.4MB/s
ERROR: dephell 0.8.3 has requirement pip<=19.3.1,>=18.0, but you'll have pip 21.1.1 which is incompatible.
Installing collected packages: pip
  Found existing installation: pip 19.3.1
  Uninstalling pip-19.3.1:
    Successfully uninstalled pip-19.3.1
Successfully installed pip-21.1.1
~/w1813148CWP02$ pip install numpy
Requirement already satisfied: numpy in /opt/virtualenvs/python3/lib/python3.8/site-packages (1.20.2)
~/w1813148CWP02$ pip install nltk
Requirement already satisfied: nltk in /opt/virtualenvs/python3/lib/python3.8/site-packages (3.6.2)
Requirement already satisfied: click in /opt/virtualenvs/python3/lib/python3.8/site-packages (from nltk) (7.1.2)
Requirement already satisfied: regex in /opt/virtualenvs/python3/lib/python3.8/site-packages (from nltk) (2021.4.4)
Requirement already satisfied: tqdm in /opt/virtualenvs/python3/lib/python3.8/site-packages (from nltk) (4.60.0)
Requirement already satisfied: joblib in /opt/virtualenvs/python3/lib/python3.8/site-packages (from nltk) (1.0.1)
~/w1813148CWP02$ python
Python 3.8.9 (default, May 3 2021, 02:40:41)
[GCC 7.5.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> import nltk
>>> nltk.download()
~/w1813148CWP02$
```

When we run this above commands then below prompt come up with the NLTK Downloader so we press the Download button after that downloading starts. This is very important for run this code.



**Create a Python program to count the most commonly used words in your dataset and use it to generate a “word cloud”.**

```
1  import codecs
2  import re
3  from collections import Counter
4  import os
5  import sys
6
7  #importing the processed tweets
8  f = codecs.open("fibretweets_processed.txt", "r", encoding="utf-8")
9  data = f.read()
10 f.close()
11
12
13 #word counting definition
14 def word_count(str):
15     counts = Counter()
16     words = str.split()
17     for word in words:
18         if (word in counts
19             ) and not (word.startswith("@")) and not (word.startswith("#")):
20             counts[word] += 1
21         else:
22             counts[word] = 1
23     return counts
24
25
26 #tag counting definition
27 def tag_count(str):
28     counts = Counter()
29     words = str.split()
30     for word in words:
31         if (word in counts) and (word.startswith("#")):
32             counts[word] += 1
33         else:
34             counts[word] = 1
35     return counts
36
37
38 #most mentioned user counting
39 def user_count(str):
40     counts = Counter()
41     words = str.split()
42     for word in words:
43         if (word in counts) and (word.startswith("@")):
44             counts[word] += 1
45         else:
46             counts[word] = 1
47     return counts
48
```

```

50 counter = word_count(data)
51 tags = tag_count(data)
52 users = user_count(data)
53 print("\n", counter.most_common(20), "\n\n", tags.most_common(20), "\n\n",
54       | | | users.most_common(20), "\n\n")
55
56 #####
57
58 #preparation for the plot
59 word, w_counts = zip(*counter.most_common(10))
60 tag, t_counts = zip(*tags.most_common(10))
61 user, u_counts = zip(*users.most_common(10))
62
63 #importing the libraries to plot the table
64 import pandas as pd
65 import matplotlib.pyplot as plt
66
67 #creating dataframe for each counting
68 wordDF = pd.DataFrame({"word": word, "w_count": w_counts})
69 tagDF = pd.DataFrame({"tags": tag, "t_count": t_counts})
70 userDF = pd.DataFrame({"users": user, "u_count": u_counts})
71
72 #plotting tag counting
73 ax = tagDF.plot(x="tags",
74                | | | | | y=["t_count"],
75                | | | | | kind="bar",
76                | | | | | figsize=(20, 25),
77                | | | | | fontsize=12)
78 ax.set_xlabel("Tag")
79 ax.set_ylabel("Frequency")
80 plt.savefig("tag_count.png")
81
82 #plotting word counting
83 ax = wordDF.plot(x="word",
84                 | | | | | y=["w_count"],
85                 | | | | | kind="bar",
86                 | | | | | figsize=(20, 15),
87                 | | | | | fontsize=12)
88 ax.set_xlabel("Word")
89 ax.set_ylabel("Frequency")
90 plt.savefig("word_count.png")
91

```

```

92 #plotting user counting
93 ax = userDF.plot(x="users",
94                 y=["u_count"],
95                 kind="bar",
96                 figsize=(20, 15),
97                 fontsize=12)
98 ax.set_xlabel("User")
99 ax.set_ylabel("Frequency")
100 plt.savefig("user_count.png")
101
102
103 #Listing all the tags
104 def all_tags(str):
105     hashtags = []
106     words = str.split()
107     for word in words:
108         if word.startswith("#"):
109             hashtags.append(word)
110         else:
111             pass
112     return " ".join(hashtags)
113
114
115 hash_tags = all_tags(data)
116 print(hash_tags.upper())

```

Word count , tag count and user count Out put

```

[('internet', 1198), ('-', 577), ('net', 434), ('via', 431), ('please', 413), ('official', 355), ('send', 351), ('activity', 348), ('related', 346), ('email', 345), ('mail', 344), ('address', 343), ('inquiries', 340), ('bis', 340), ('supporting', 340), ('gtid@131labeln', 340), ('people', 269), ('amp', 269), ('2021', 258), ('wanna', 225)]

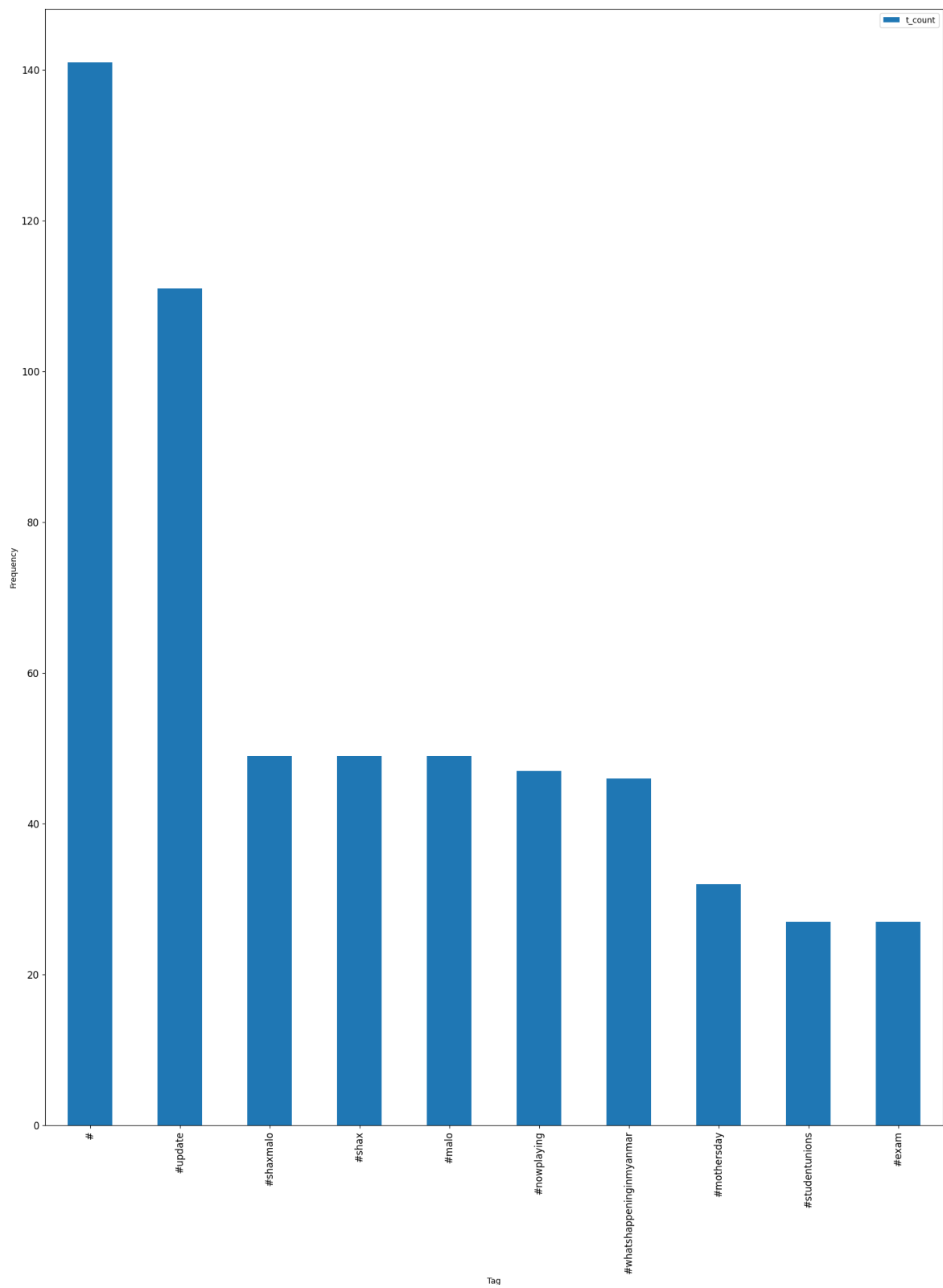
[('#', 141), ('#update', 111), ('#shaxmalo', 49), ('#shax', 49), ('#malo', 49), ('#nowplaying', 47), ('#whatshappeninginmyanmar', 46), ('#mothersday', 32), ('#studentunions', 27), ('#exam', 27), ('#astro', 26), ('#mdcat', 26), ('#dogecoin', 25), ('#helixstudios', 23), ('#mellow975xmss', 21), ('#affiliatemarketing', 17), ('#cardano', 16), ('#radio', 16), ('#cybersecurity', 16), ('#job', 15)]

[('@131label', 342), ('@itsmalikel', 148), ('@balloonwanted', 119), ('@teamviewer', 87), ('@newyorkstateag', 65), ('@prettyinbri', 61), ('@lustfulelaine', 55), ('@starksgrayson', 53), ('@jayelharris', 53), ('@shaxtw', 49), ('@teamviewerhelp', 47), ('@piichi02', 40), ('@rellromance', 39), ('@ksgupdates', 28), ('@mattdpearce', 27), ('@haqoogekhalq', 27), ('@mjibrannasir', 26), ('@astrofancafe', 23), ('@totalwoke', 23), ('@', 20)]

```

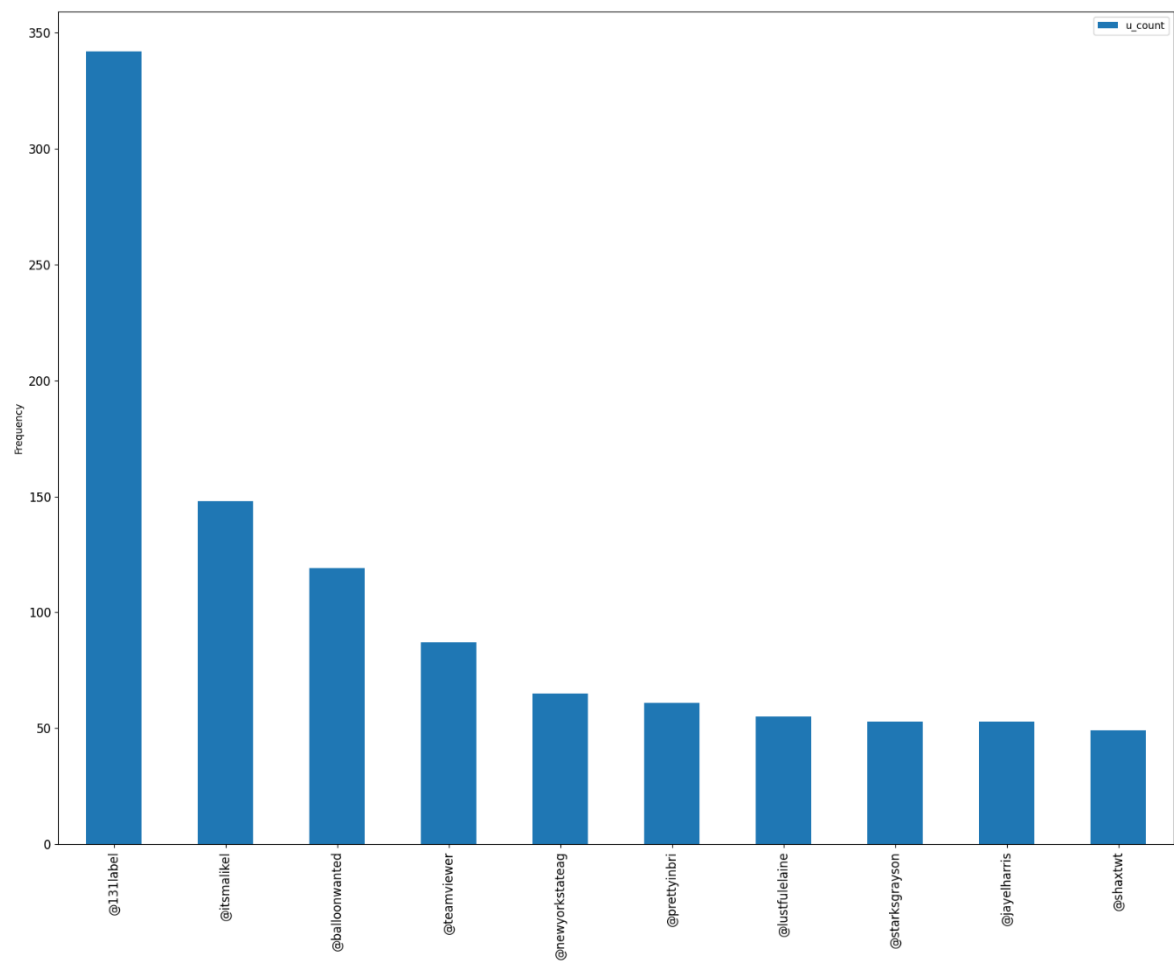
I also get graphs as a outputs of the program

Tag count





User Count



Total word count after pre-processing

6256



Word Cloud using the word count words via worditout.com



***Use your processed data file to produce a series of graphs or charts to summaries the following information.***

- I. The number of tweets posted per day***
- II. The number of unique users per day***
- III. The top 10 most active users over the entire period***

**Using a suitable approach, construct a LDA topic model to identify themes of discussion within your dataset.**

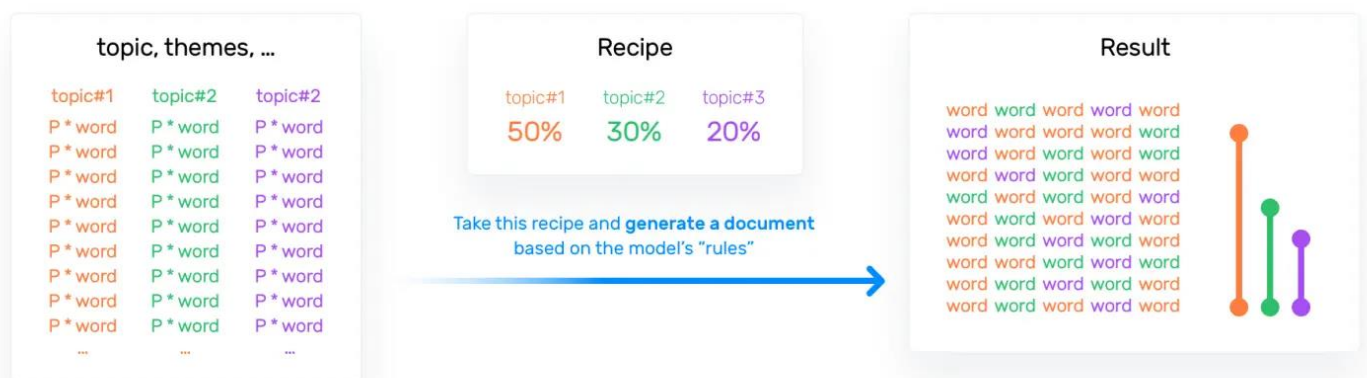
Topic modeling is a one type of machine learning technique that analyze text data automatically using cluster words for set of documents. This is unsupervised machine learning. It's a easy and quick way to analyze our data. Counting words and similar word grouping are part of the topic modeling within unstructured data.

Topic modeling refers to the process of dividing a corpus of documents in two:[3]

1. A list of the topics covered by documents in the corpus
2. Serval set of documents from the corpus by topics they cover.

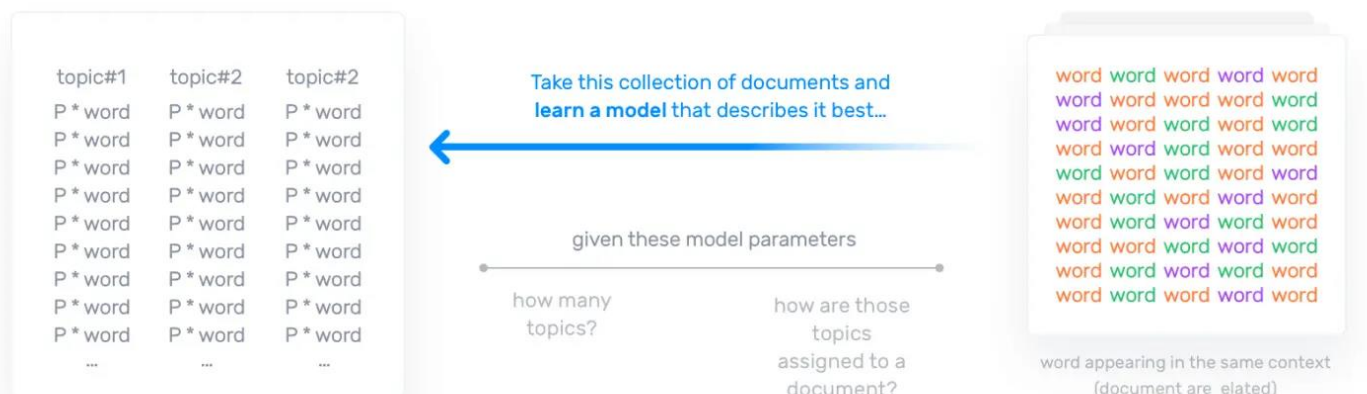
LDA assumes that topics and documents look like this:[3]

Lets assume that...



And, when LDA models a new document, it works this way:[3]

What really happens...



We, used the tweet cleaning process once but still in this step we can cleaning some of the data cleaning using R. first of all, I installed the missing packages that required to run the program

```
Install.packages(c("tm"))
```

```
Install.packages(c("topicmodels"))
```

```
text <- readLines("C:\\Users\\Bhaumik\\Desktop\\fibretweets_processed.txt",encoding="UTF-8")
```

```
doc.vec <- VectorSource(text)
```

```
doc.corpus <- Corpus(doc.vec)
```

```
doc.corpus <- tm_map(doc.corpus, function(x) iconv(enc2utf8(x), sub = "byte"))
```

```
doc.corpus <- tm_map(doc.corpus, PlainTextDocument)
```

```
doc.corpus <- tm_map(doc.corpus, content_transformer(tolower))
```

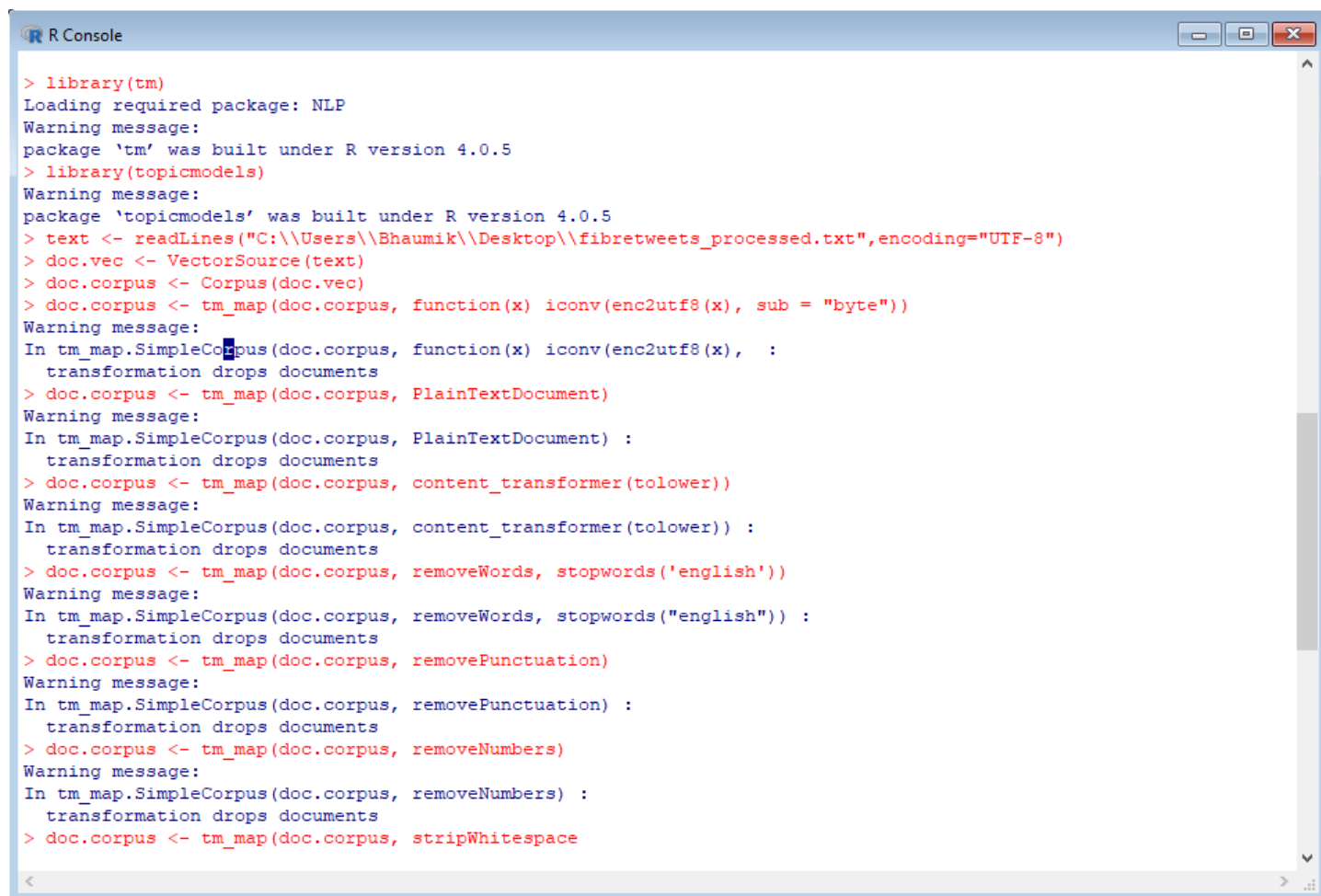
```
doc.corpus <- tm_map(doc.corpus, removeWords, stopwords('english'))
```

```
doc.corpus <- tm_map(doc.corpus, removePunctuation)
```

```
doc.corpus <- tm_map(doc.corpus, removeNumbers)
```

```
doc.corpus <- tm_map(doc.corpus, stripWhitespace)
```

This is the outputs from the R were running the above code



```
> library(tm)
Loading required package: NLP
Warning message:
package 'tm' was built under R version 4.0.5
> library(topicmodels)
Warning message:
package 'topicmodels' was built under R version 4.0.5
> text <- readLines("C:\\Users\\Bhaumik\\Desktop\\fibretweets_processed.txt",encoding="UTF-8")
> doc.vec <- VectorSource(text)
> doc.corpus <- Corpus(doc.vec)
> doc.corpus <- tm_map(doc.corpus, function(x) iconv(enc2utf8(x), sub = "byte"))
Warning message:
In tm_map.SimpleCorpus(doc.corpus, function(x) iconv(enc2utf8(x), :
transformation drops documents
> doc.corpus <- tm_map(doc.corpus, PlainTextDocument)
Warning message:
In tm_map.SimpleCorpus(doc.corpus, PlainTextDocument) :
transformation drops documents
> doc.corpus <- tm_map(doc.corpus, content_transformer(tolower))
Warning message:
In tm_map.SimpleCorpus(doc.corpus, content_transformer(tolower)) :
transformation drops documents
> doc.corpus <- tm_map(doc.corpus, removeWords, stopwords('english'))
Warning message:
In tm_map.SimpleCorpus(doc.corpus, removeWords, stopwords("english")) :
transformation drops documents
> doc.corpus <- tm_map(doc.corpus, removePunctuation)
Warning message:
In tm_map.SimpleCorpus(doc.corpus, removePunctuation) :
transformation drops documents
> doc.corpus <- tm_map(doc.corpus, removeNumbers)
Warning message:
In tm_map.SimpleCorpus(doc.corpus, removeNumbers) :
transformation drops documents
> doc.corpus <- tm_map(doc.corpus, stripWhitespace)
```

```
R Console

> doc.corpus <- tm_map(doc.corpus, removePunctuation)
Warning message:
In tm_map.SimpleCorpus(doc.corpus, removePunctuation) :
  transformation drops documents
> doc.corpus <- tm_map(doc.corpus, removeNumbers)
Warning message:
In tm_map.SimpleCorpus(doc.corpus, removeNumbers) :
  transformation drops documents
> doc.corpus <- tm_map(doc.corpus, stripWhitespace)
+ doc.corpus <- tm_map(doc.corpus, stripWhitespace)
Error: unexpected symbol in:
"doc.corpus <- tm_map(doc.corpus, stripWhitespace
doc.corpus"
> doc.corpus <- tm_map(doc.corpus, stripWhitespace)
Warning message:
In tm_map.SimpleCorpus(doc.corpus, stripWhitespace) :
  transformation drops documents
> dtm <- DocumentTermMatrix(doc.corpus)
> dtm <- removeSparseTerms(dtm, 0.98)
> x <- as.matrix(dtm)
> x <- x[which(rowSums(x) > 0),]
> rownames(x) <- 1:nrow(x)
> lda <- LDA(x, <<nooftopics>>)
Error: unexpected input in "lda <- LDA(x, <<"
> lda <- LDA(x,6)
> terms(lda, 6)
      Topic 1 Topic 2      Topic 3      Topic 4      Topic 5      Topic 6
[1,] "one"    "internet" "wanna"  "official" "net"    "via"
[2,] "like"   "amp"      ""       "send"    "people" "please"
[3,] "time"   "new"      "talk"   "label"   "update" "broadband"
[4,] "need"   "broadband" "strangers" "activity" "covid"  "new"
[5,] "get"    "best"     "internet" "related" "day"    "internet"
[6,] "may"    "dont"     "itsmalikel" "email"  "live"   "email"
> perplexity(lda)
[1] 22.69929
```

From this analysis, we got the most talked 6 topics and perplexity LDA value 22.69929

**Apply noun phrase recognition to your dataset and identify the top five most mentioned noun phrases. Construct a sentiment model for each of your identified noun phrases and compare and contrast the differences in both polarity and sentiment**

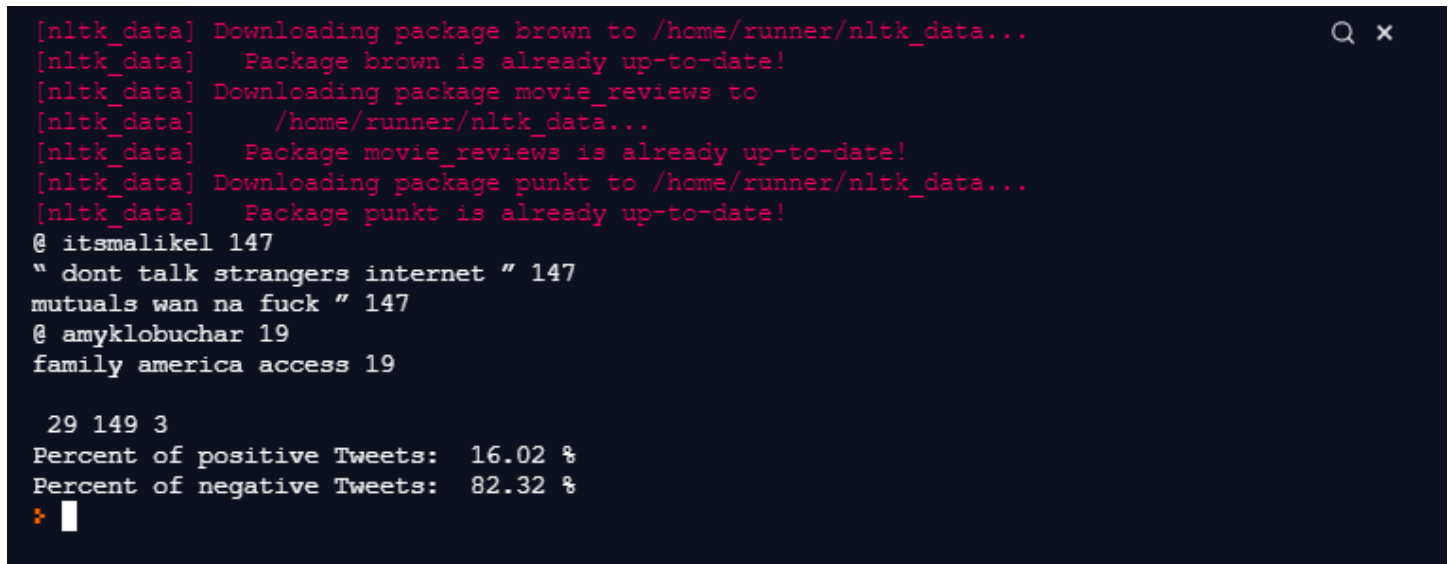
```
1  from textblob import TextBlob
2  from textblob.sentiments import NaiveBayesAnalyzer
3  import codecs
4  from textblob.np_extractors import ConllExtractor
5  import nltk
6  nltk.download('brown')
7  nltk.download('movie_reviews')
8  nltk.download('punkt')
9
10 prebuilt_classifier = NaiveBayesAnalyzer()
11 extractor = ConllExtractor()
12
13 f = codecs.open("fibretweets_processed.txt", encoding = "utf-8")
14
15 positive = 0
16 negative = 0
17 neutral = 0
18 np = {}
19
20 #defining the terms to use
21 for line in f.readlines():
22     lc_line = line.lower()
23
24     if "internet" in lc_line and "2021" in lc_line:
25         tweet = TextBlob(lc_line, analyzer = prebuilt_classifier)
26
27         for n in tweet.noun_phrases:
28             if n in np:
29                 np[n] += 1
30             else:
31                 np[n] = 1
32
33         if tweet.sentiment.p_pos >= 0.7:
34             positive += 1
35         elif tweet.sentiment.p_neg >= 0.7:
36             negative += 1
37         else:
38             neutral += 1
39
40 #sorting the noun phrases by the highest appearance
41 for w in sorted(np, key = np.get, reverse = True) [0:5]:
42     print(w, np[w])
43
```

```

44 #print out the values
45 print("\n", positive, negative, neutral)
46 total = positive + negative + neutral
47 positivePercent = positive/total*100
48 negativePercent = negative/total*100
49 print("Percent of positive Tweets: ", round(positivePercent, 2), "%")
50 print("Percent of negative Tweets: ", round(negativePercent, 2), "%")

```

Outputs of the program



```

[nltk_data] Downloading package brown to /home/runner/nltk_data...
[nltk_data]   Package brown is already up-to-date!
[nltk_data] Downloading package movie_reviews to
[nltk_data]   /home/runner/nltk_data...
[nltk_data]   Package movie_reviews is already up-to-date!
[nltk_data] Downloading package punkt to /home/runner/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
@ itsmalikel 147
" dont talk strangers internet " 147
mutuals wan na fuck " 147
@ amyklobuchar 19
family america access 19

29 149 3
Percent of positive Tweets:  16.02 %
Percent of negative Tweets:  82.32 %

```

During the coursework,

## References

[1] Technical University of Munich, Date 21 Feb 2019

[https://www.technologist.eu/new-record-data-transfer-speed-in-fiber-optic-network/#:~:text=In%20an%20intercity%20field%20experiment,s\)%20with%20a%20single%20wavelength.](https://www.technologist.eu/new-record-data-transfer-speed-in-fiber-optic-network/#:~:text=In%20an%20intercity%20field%20experiment,s)%20with%20a%20single%20wavelength.)

[2] Figure 1 and 2: Mark Jackson from Dorset (England)

<https://www.ispreview.co.uk/index.php/2020/09/uk-trails-as-10-countries-pass-95-full-fibre-broadband-cover.html>

[3] Introduction to Topic Modeling, Federico Pascual, 26<sup>th</sup> September 2019

<https://monkeylearn.com/blog/introduction-to-topic-modeling/#:~:text=Topic%20modeling%20is%20an%20unsupervised,characterize%20a%20set%20of%20documents.>