A. Sentiment Analysis:

1. I have considered the twitter data that I have obtained in previous assignment. I have included the python file to generate the tweets in the Sentiment Analysis folder with the file name as 'twitter.py'. I have included the csv file as 'tweets.csv' which I have obtained in the last assignment. Now to consider only the tweets message, I wrote a python script 'sentimentanalysis.py' which considers only tweet messages from 'tweets.csv'.

2. I wrote a function 'clean' to clean the twitter data. It was completed in the previous assignment. I have done cleaning tasks such as removing hashtags, special characters, URLs, and emoticons. This cleaning function is included in the 'twitter.py' file.

3. I wrote a script 'sentimentanalysis.py' which includes creation of bag-of-words for each tweet. I have removed all the duplicate tweets and generated a new csv file as 'noduplicatetweetscsv.csv'. I have used this file which contains the tweets and successfully generated the bag-of-words and moved to the next step for the sentiment analysis. The below screenshot shows the bag-of-words for each tweet in a new line. The following bag-of-words can be generated by running 'sentimentanalysis.py'



4. I downloaded the positive words and negative words from an online source which they have provided for general research usage [1]. The files are included in the Sentiment Analysis folder with file names as 'positive-words.txt' [3] and 'negative-words.txt' [4]. After obtaining
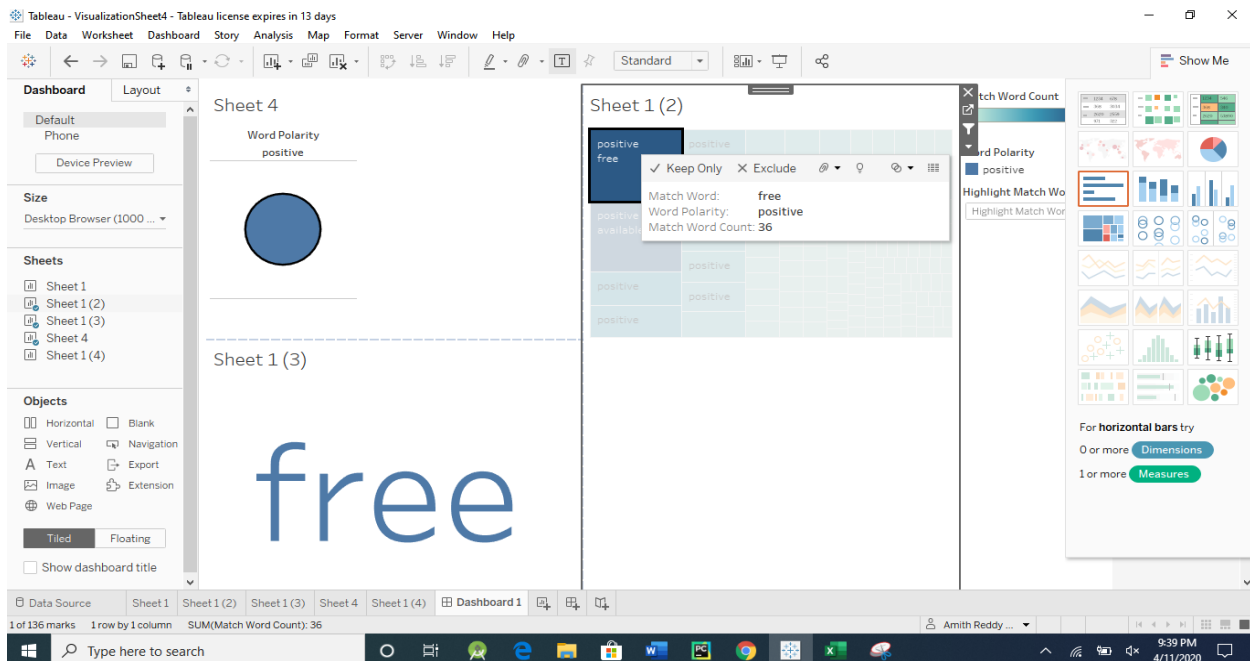
the positive and negative words, I have compared the bag-of-words with these positive and negative words and found the matches. The code for this is included in Sentiment Analysis folder as 'sentimentanalysis.py'.

5. The next step is to tag each "positive", "negative", or "neutral". First, I have initialized a sum variable to 0. Then, I have incremented the sum by 1 for each positive word and decremented the sum by 1 for each negative word. Now the sum variable contains the final result which helps in deciding about that particular tweet. When the sum is greater than 0, the tweet is tagged as "positive". The tweet is "neutral" if the sum is equal to 0. Further when the sum is less than 0, the tweet is tagged as "negative". The code for this is included in Sentiment Analysis folder as 'sentimentanalysis.py'.

6. For the visualization purpose, I wrote a script 'visualizationwords.py' to generate 'visualization2.csv' which contains columns such as 'Tweet_id', 'match_word', 'match_word_count', 'word_polarity'. I have downloaded and installed Tableau from the Internet. Then, the 'visualization2.csv' file is loaded into the Tableau Workbook. The 'match_word_count' column is the measure which is loaded as size and column value. The 'word_polarity' and 'match_word' columns are dimensions which are dragged onto the rows for the visualization.

I have done 11 visualizations on the Tableau Workbook and included all the screenshots and visualization sheets in Assignment folder. I found this Visualization tool very helpful to analyse the most frequently occurring words in the positive and negative tweets which I have obtained.

The first visualization is a dashboard that has three charts which clearly shows the most frequently occurring positive word and negative word. The below screenshot represents the first visualization.

The second visualization shows blocks of positive words and negative words. The size of the block tells us about the words which are most frequent and less frequent in the tweets. The highest occurring positive word is "free" with count 35 and highest occurring negative word is "break" with count 9.



The third visualization helps us in visualizing the most frequently occurring positive and negative words in form of a word cloud. When we hover over the words, it displays the word count of the word and word polarity which is extremely useful for analysing the tweets further.
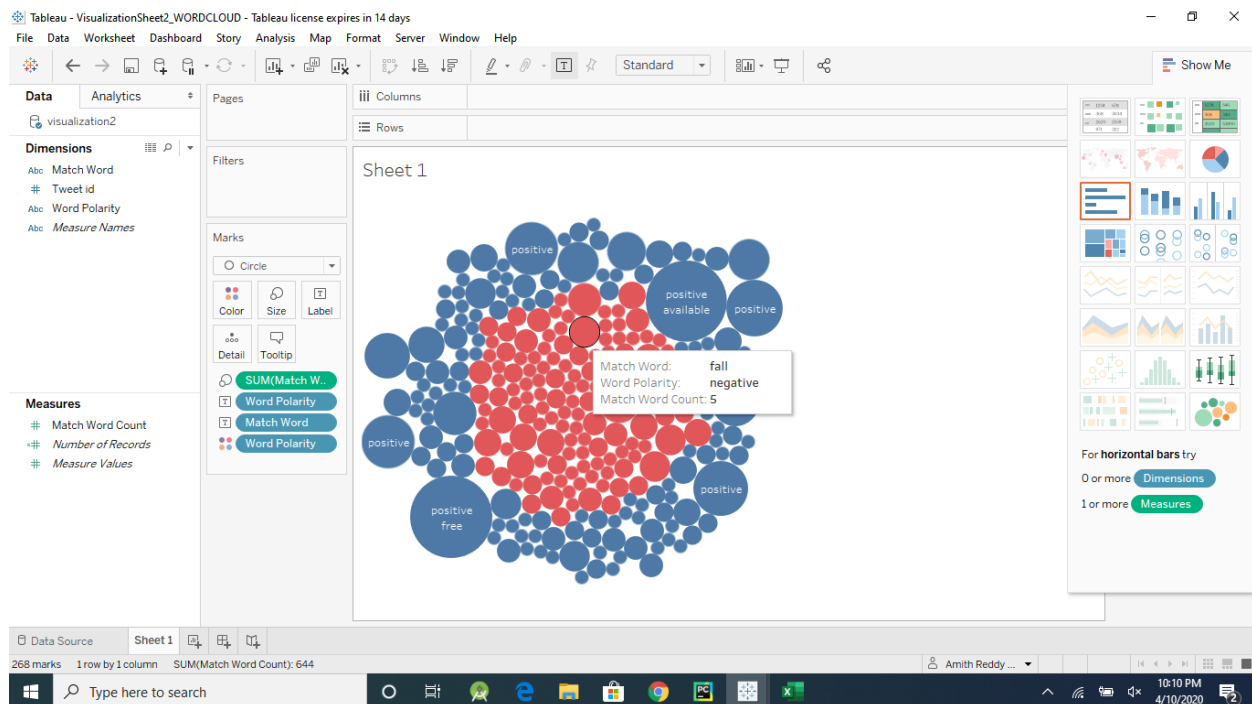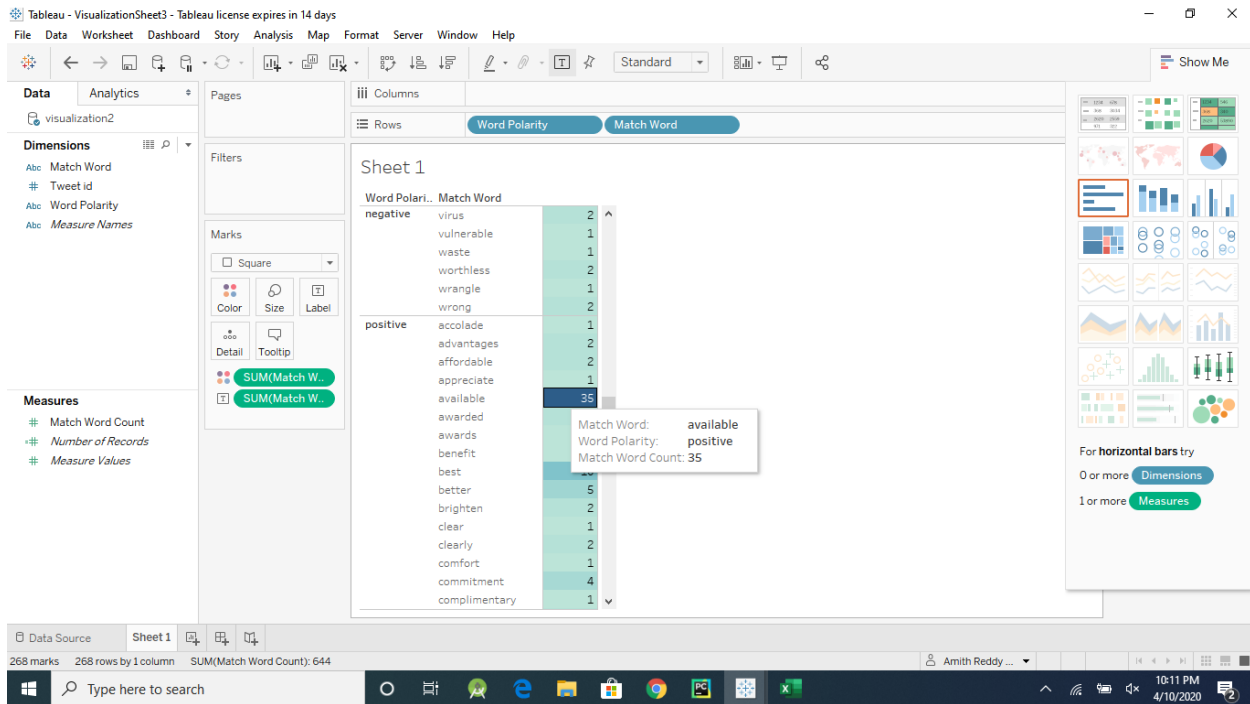
The fourth visualization is a bar graph which shows list of positive words and negative words along with their frequency count.
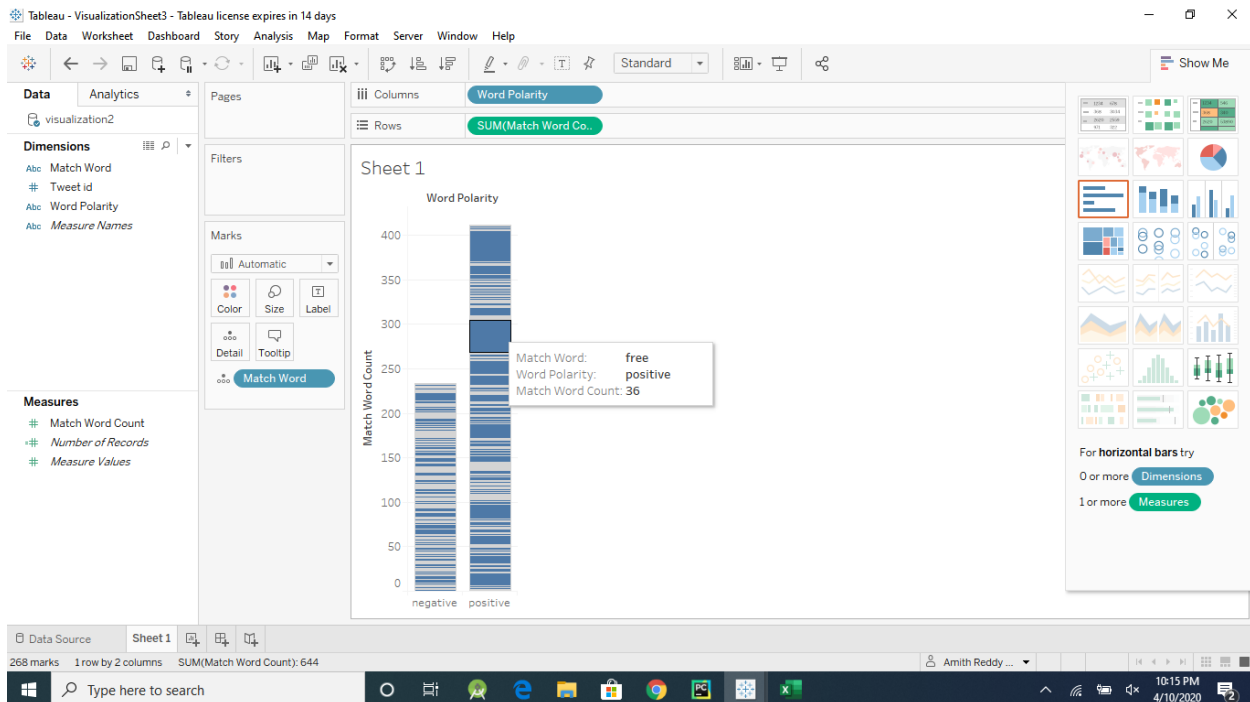


The fifth visualization is a bubble graph which clearly differentiates positive and negative words with blue and red colour. It also displays the count of each word.
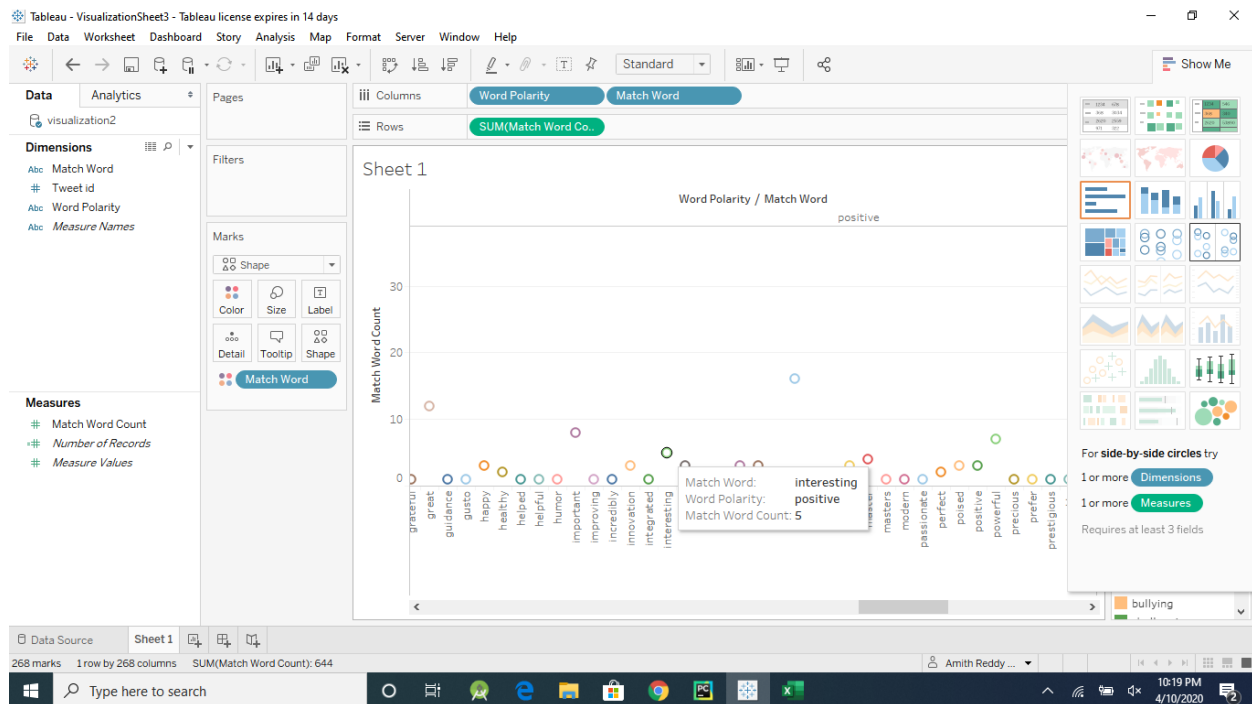
The sixth visualization is a scroll sheet where you can scroll to check all the positive words and negative words along with their count.
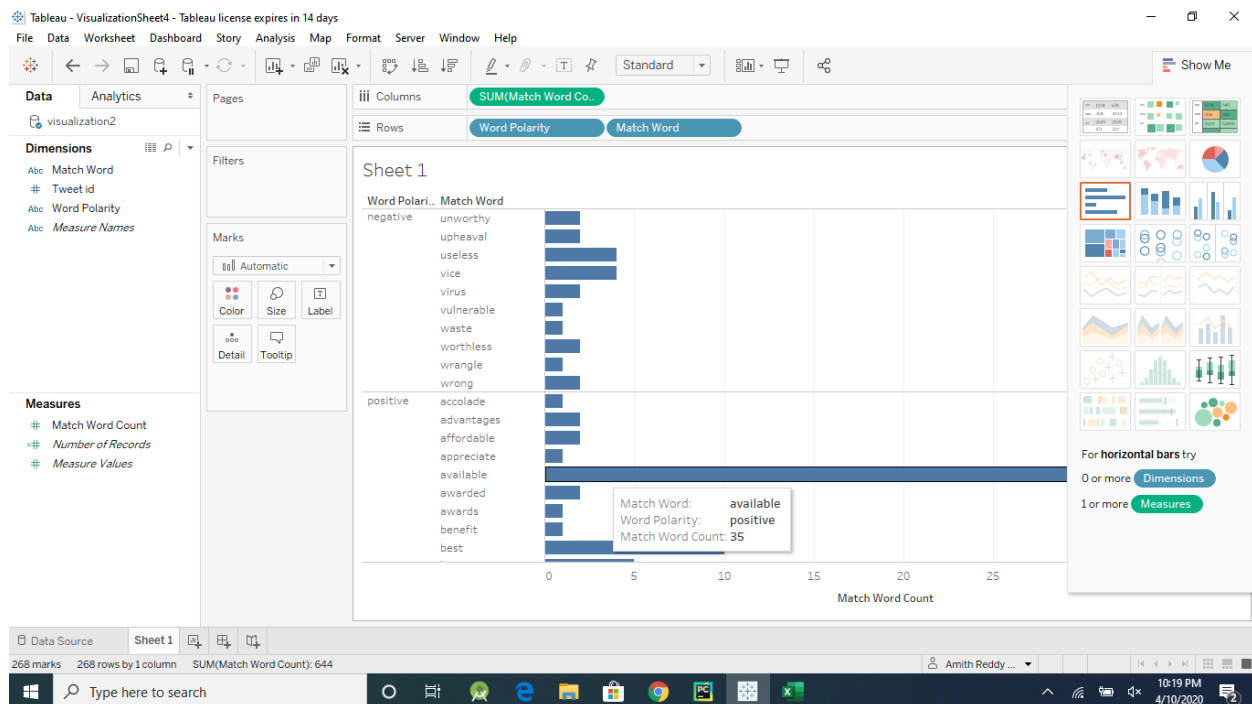


The seventh visualization shows the word count of all positive and negative words. It clearly shows that the total positive words count is nearly double the total negative words count.

The eighth visualization is almost similar to a bar graph. The only difference is that it shows the count of positive words and negative words in form of a bubble. It looks very clean and precise.



The nineth visualization shows a segregation between positive words and negative words in the form of a tilted bar graph.
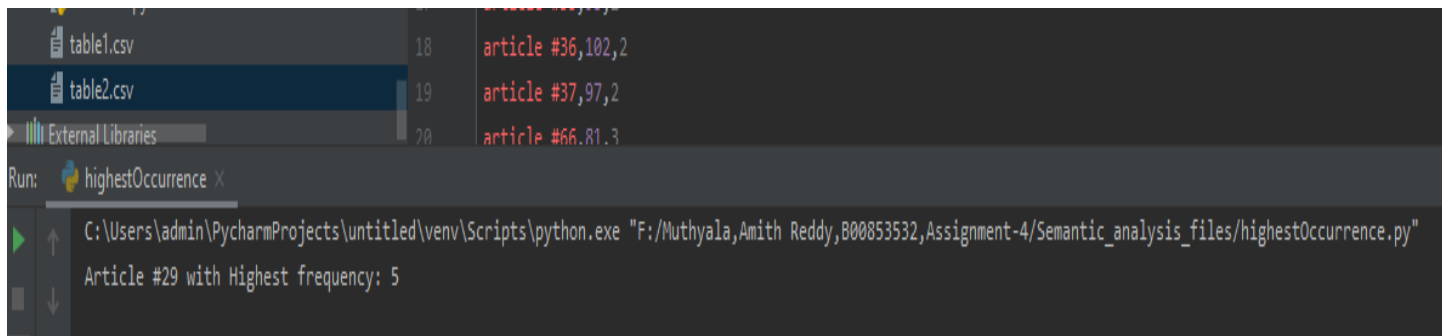
B. Semantic Analysis:

7. I wrote a well-formed script to clean and transform the news articles in python. I have stored these news articles in a file named 'news.csv'. The script to generate this 'news.csv' is included in the Semantic Analysis files folder as 'newsExtraction.py'. I have removed URLs, special characters, and emoticons.

8. I have considered each chunk of text as a document. I wrote a script 'semantic.py' to generate each news article as a separate csv document. All the news articles are named as 'file' followed by numbers that are in a sequence. These files only contain one news article per document.

9. Each news file only contains "title", "description", and the news "content". The script to generate these required columns is included as 'semantic.py'. The eighth and nineth steps are combinedly done to maintain the flow.

10. COMPUTE TF-IDF (term frequency-inverse document frequency)

    a. I have stored all the news articles in different csv files. I have included these files in the Semantic Analysis folder. In my case N=500. Now, I have used the search query "Canada", "University", "Dalhousie University", "Halifax", "Business", and searched in how many documents these words have appeared. The python script for this searching is included as 'search.py'. It takes the search query as specified and creates a table 'table1.csv' with columns "Search Query", "Document containing df", "N/df", "log10(N/df)". The below screenshot displays the number of documents which contain the search query. I have included the screenshot as 'table-1' and the csv file 'table-1.csv' is included in the Semantic Analysis folder.

| Total Documents | 500 | | |
|---|---|---|---|
| Search Query | Document containing term(df) | Total Documents(N)/ number of documents term appeared(df) | Log10(N/df) |
| Canada | 91 | 5.49 | 0.74 |
| University | 33 | 15.15 | 1.18 |
| Dalhousie University | 11 | 45.45 | 1.66 |
| Halifax | 24 | 20.83 | 1.32 |
| Business | 4 | 125 | 2.1 |

    b. The next step is to find the document which has the highest occurrence of the word "Canada". This can be done by performing frequency count of the word per document. I wrote a python script 'highestOccurrence.py' to generate 'table-2.csv' and also displays the document which has the highest occurrence of the word "Canada". The below screenshot displays the article with highest occurrence of the word "Canada".

```
            table1.csv           18   article #36,102,2
            table2.csv           19   article #37,97,2
         External Libraries       20   article #66,81,3
Run:      highestOccurrence ×

    C:\Users\admin\PycharmProjects\untitled\venv\Scripts\python.exe "F:/Muthyala,Amith Reddy,B00853532,Assignment-4/Semantic_analysis_files/highestOccurrence.py"
    Article #29 with Highest frequency: 5
```

The below table shows frequency count of the word "Canada" per document.

| Term | Canada | |
| --- | --- | --- |
| Canada appeared in | Total words(m) | Frequency(f) |
| article #1 | 101 | 2 |
| article #4 | 107 | 1 |
| article #5 | 99 | 2 |
| article #6 | 98 | 1 |
| article #7 | 101 | 1 |
| article #8 | 105 | 2 |
| article #9 | 103 | 2 |
| article #12 | 91 | 2 |
| article #25 | 93 | 2 |
| article #29 | 105 | 5 |
| article #30 | 115 | 3 |
| article #31 | 106 | 3 |
| article #32 | 111 | 1 |
| article #33 | 95 | 2 |
| article #35 | 95 | 2 |
| article #36 | 102 | 2 |
| article #37 | 97 | 2 |
| article #66 | 81 | 3 |
| article #67 | 98 | 3 |

c.  To find out the article with highest relative frequency, I wrote a well-formed python script 'highestrelativefrequency.py' which is included in the Semantic Analysis folder. This python file generates a table 'table-3.csv' which has relative frequency values per document. The below screenshot displays the article with highest relative frequency (f/m).



```
▶ External Libraries          18   article #36,102,2,0.02
  Scratches and Consoles      19   article #37,97,2,0.02
                              20   article #66,81,3,0.04
Run:   highestrelativefrequency (1) ×
   C:\Users\admin\PycharmProjects\untitled\venv\Scripts\python.exe "F:/Muthyala,Amith Reddy,B00853532,Assignment-4/Semantic_analysis_files/highestrelativefrequency.py"
   Article #29 with Highest Relative frequency: 0.05
```

The below table shows the relative frequency of news articles.

| Term | Canada | | |
| --- | --- | --- | --- |
| Canada appeared in 500 documents | Total words(m) | Frequency(f) | Relative Frequency(f/m) |
| article #1 | 101 | 2 | 0.02 |
| article #4 | 107 | 1 | 0.01 |
| article #5 | 99 | 2 | 0.02 |
| article #6 | 98 | 1 | 0.01 |
| article #7 | 101 | 1 | 0.01 |
| article #8 | 105 | 2 | 0.02 |
| article #9 | 103 | 2 | 0.02 |
| article #12 | 91 | 2 | 0.02 |
| article #25 | 93 | 2 | 0.02 |
| article #29 | 105 | 5 | 0.05 |
| article #30 | 115 | 3 | 0.03 |
| article #31 | 106 | 3 | 0.03 |
| article #32 | 111 | 1 | 0.01 |
| article #33 | 95 | 2 | 0.02 |
| article #35 | 95 | 2 | 0.02 |
| article #36 | 102 | 2 | 0.02 |
| article #37 | 97 | 2 | 0.02 |
| article #66 | 81 | 3 | 0.04 |
| article #67 | 98 | 3 | 0.03 |

References:

[1] M. Hu and B. Liu, "Mining and summarizing customer reviews," Proceedings of the 2004 ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '04, 2004 [Online]. Available: https://dl.acm.org/doi/pdf/10.1145/1014052.1014073. [Accessed: 11-Apr-2020]

[2] B. Liu, M. Hu, and J. Cheng, "Opinion observer," Proceedings of the 14th international conference on World Wide Web - WWW '05, 2005 [Online]. Available: https://dl.acm.org/doi/pdf/10.1145/1060745.1060797. [Accessed: 11-Apr-2020]

[3] *Ptrckprry.com*, 2010. [Online]. Available: http://ptrckprry.com/course/ssd/data/positive-words.txt. [Accessed: 11-Apr-2020].

[4] *Ptrckprry.com*, 2010. [Online]. Available: http://ptrckprry.com/course/ssd/data/negative-words.txt. [Accessed: 11-Apr-2020].

[5] "Natural Language Processing - Semantic Analysis - Tutorialspoint," *Tutorialspoint.com*, 2020. [Online]. Available: https://www.tutorialspoint.com/natural_language_processing/natural_language_processing_semantic_analysis.htm. [Accessed: 11-Apr-2020].

[6] Abdullah Alsaeedi and Mohammad Zubair Khan, "A Study on Sentiment Analysis Techniques of Twitter Data," *ResearchGate*, 28-Feb-2019. [Online]. Available: https://www.researchgate.net/publication/331411860_A_Study_on_Sentiment_Analysis_Techniques_of_Twitter_Data. [Accessed: 11-Apr-2020].