# Assignment – 4

## A. Sentiment Analysis:

Following steps were followed to perform sentiment analysis on tweets

Step 1: Only the tweets/messages were stored in separate file (tweets.json), which was collected using python script from assignment 3. All the tweets collected using the python script was cleaned. For cleaning purposes all the punctuation, URLs, emojis, and Unicode were removed in the previous assignment. Python script file (twitter.py) can be found in the  Script folder.
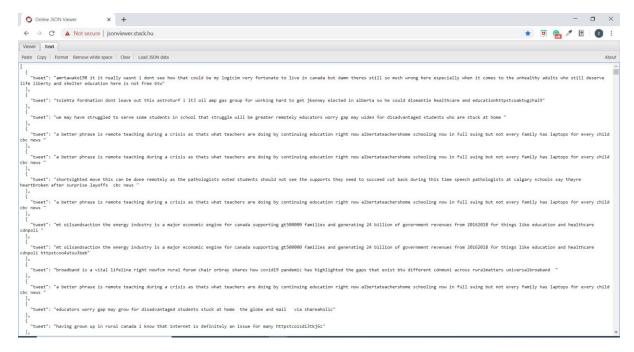


*Figure 1 tweets in tweets.json opened in http://jsonviewer.stack.hu/*

Step 2: Each tweet was read from (tweets.json) file. After that bag-of-word was created for each tweet, which shows the word and its occurrence in tweet.

Step3: Text file containing the list of words along with its polarity and intensity was downloaded from (https://sentic.net/downloads/). Senticnet5 was selected as it consists of approximately 1,00,000 words.[1]

This list was modified and stored in separate CSV file (words.csv) showing word and its polarity as we did not require intensity of the word in our analysis.

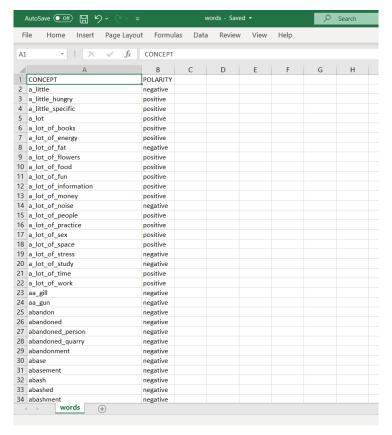*Figure 2 senticnet5.txt showing word along with polarity and intensity*



*Figure 3 words.csv containing word and its polarity*

Step 3: Now, each word from bag-of-words is matched with these words. If the word is present, then it is collected along with its polarity to compute the overall polarity of the tweet.
For each tweet, we calculated total positive frequency and negative frequency from the matching words. Now, to compute the overall frequency of the tweet:

positive frequency > negative frequency than overall polarity of tweet is "POSITIVE"
positive frequency < negative frequency than overall polarity of tweet is "NEGATIVE"
positive frequency = negative frequency than overall polarity of tweet is "NEUTRAL"

Each tweet along with the matched words, Positive Word frequency, Negative Word frequency, and Tweet Polarity is stored in the (tweet_polarity_table.csv) file it can be found in Output folder.
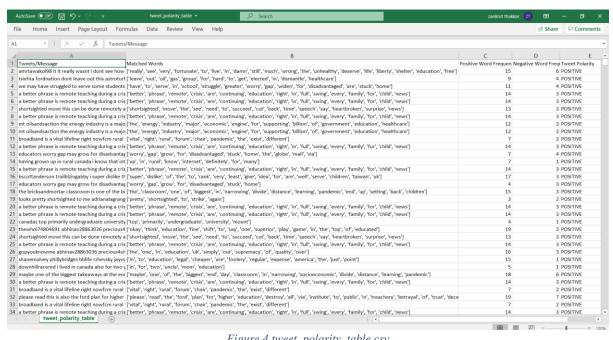


*Figure 4 tweet_polarity_table.csv*

## Visualisation:

Step 4: For Visualisation of data using tableau. All the matching words along with its total occurrences in all the tweets and polarity are collected and stored in separate file (world_cloud_output.csv), that can be found in Output folder. As our file (words.csv) consists of many words, there are many words matching found from the tweets. The polarity of each matched word is based on the data we downloaded.
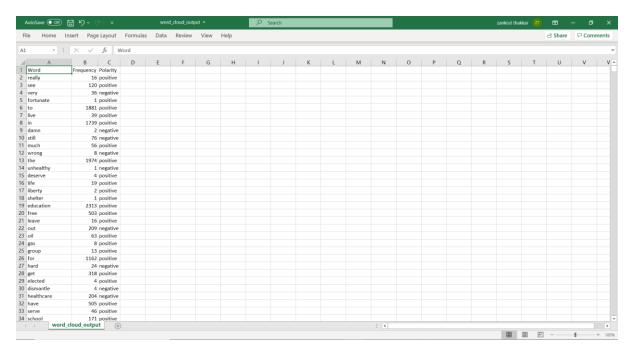
*Figure 5 word_cloud_output.csv*

**Loading csv to tableau:**

Click on File → Text File → Select the CSV file you want to load.



*Figure 6 Loading data into Tableau*

**Word Cloud:**

All the positive words are shown in green colour and negative words are shown in red colour. And the size of each word depends upon the frequency of word.[9]
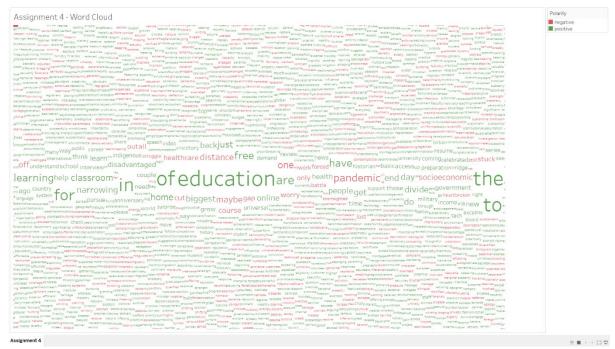


*Figure 7 Word Cloud in Tableau*

Python script file (sentiment_analysis_tweets.py) can be found in the Script folder.

## B. Semantic Analysis:

Following steps were followed to perform semantic analysis on the news data:

Step 1: The news data was collected and cleaned using the python script from assignment 3 from (https://newsapi.org/). In the previous assignment for cleaning purpose all the punctuation, URLs, emojis, Unicode were removed.[5] Python script file (news.py) can be found in the  Script folder.

Step 2: Using the python script only the "title", "description", and "content" of the article was collected and stored in different files. Here, there are 140 articles and for each article separate text file is created.
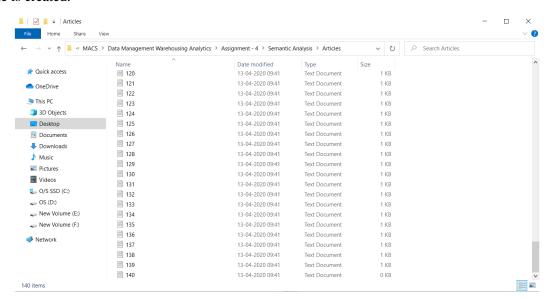


*Figure 8 Separate text file for each article*

### Part A:

Step 3: Searched each text file and took count of files in which "canada", "university", "dalhousie university", "business", and "halifax" are present. Then computed the Total number of documents/ document containing term (N/df). Also, after that $\log_{10}(N/df)$ was calculated using math.log function. Moreover, output was rounded to two decimal points. All the data was inserted in pretty table and converted to output (sematic A.csv) file.
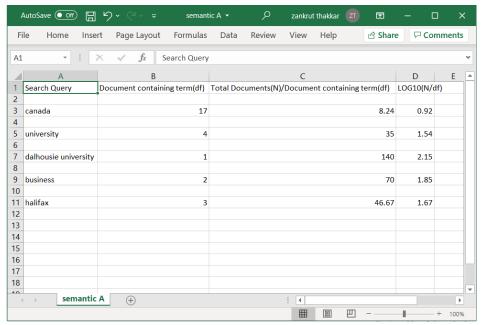
*Figure 9 semantic A.csv*

**Part B:**

Step 4: After the above-mentioned step each file was searched and counted number of times word "canada" was present in it by appending the data into list and counting word using count() method provided by list. The output was inserted into pretty table and converted to output (semantic B.csv) file.
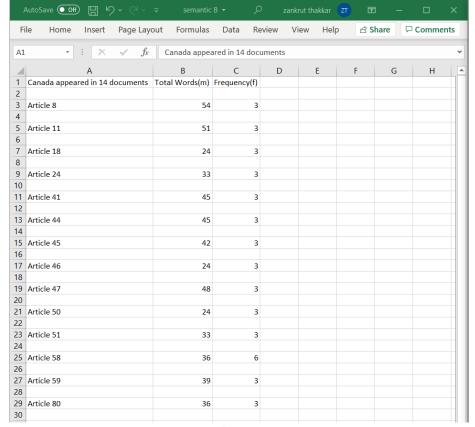


*Figure 10 semantic B.csv*

**Part C:**

Step 5: After that relative frequency = Total words (m) / Frequency (f) was computed for the articles to in which "canada" was present. And article with highest relative frequency was taken and stored in output (semantic C.csv) file.



*Figure 11 semantic C.csv*

Python script file (semantic_analysis_news.py) can be found in the  Script folder. All the three output files semantic A.csv, semantic B.csv, and semantic C.csv can be found in Output folder.

# References:

[1] "Downloads « SenticNet", *Sentic.net*, 2020. [Online]. Available: https://sentic.net/downloads/. [Accessed: 13- Apr- 2020].

E Cambria, S Poria, D Hazarika, K Kwok. SenticNet 5: Discovering conceptual primitives for sentiment analysis by means of context embeddings. In: AAAI, pp. 1795-1802 (2018)

[2] "Docs," *Twitter.com*, 2020. [Online]. Available: https://developer.twitter.com/en/docs.  [Accessed: 27-Mar-2020]

[3] "Tweepy Documentation — tweepy 3.8.0 documentation", *Docs.tweepy.org*, 2020. [Online]. Available: http://docs.tweepy.org/en/latest/.  [Accessed: 13- Apr- 2020]

[4] "apache/spark", *GitHub*, 2020. [Online]. Available: https://github.com/apache/spark/blob/master/examples/src/main/python/wordcount.py.  [Accessed: 13- Apr- 2020]

[5] "News API - A JSON API for live news and blog articles", *Newsapi.org*, 2020. [Online]. Available: https://newsapi.org/.  [Accessed: 13- Apr- 2020]

[6]"Python: Count the occurrences of each word in a given sentence - w3resource", *w3resource*, 2020. [Online]. Available: https://www.w3resource.com/python-exercises/string/python-data-type-string-exercise-12.php. [Accessed: 13- Apr- 2020].

[7] J. Bodnar, "Python PrettyTable tutorial - generating tables in Python with PrettyTable", *Zetcode.com*, 2020. [Online]. Available: http://zetcode.com/python/prettytable/. [Accessed: 13- Apr- 2020].

[8] C. line, A. Smith, J. Böcker and A. Klinke, "Convert Python Pretty table to CSV using shell or batch command line", *Stack Overflow*, 2020. [Online]. Available: https://stackoverflow.com/questions/32128226/convert-python-pretty-table-to-csv-using-shell-or-batch-command-line.  [Accessed: 13- Apr- 2020].

[9] "Creating a Word Cloud | Tableau Software", *Kb.tableau.com*, 2020. [Online]. Available: https://kb.tableau.com/articles/howto/creating-a-word-cloud . [Accessed: 13- Apr- 2020].