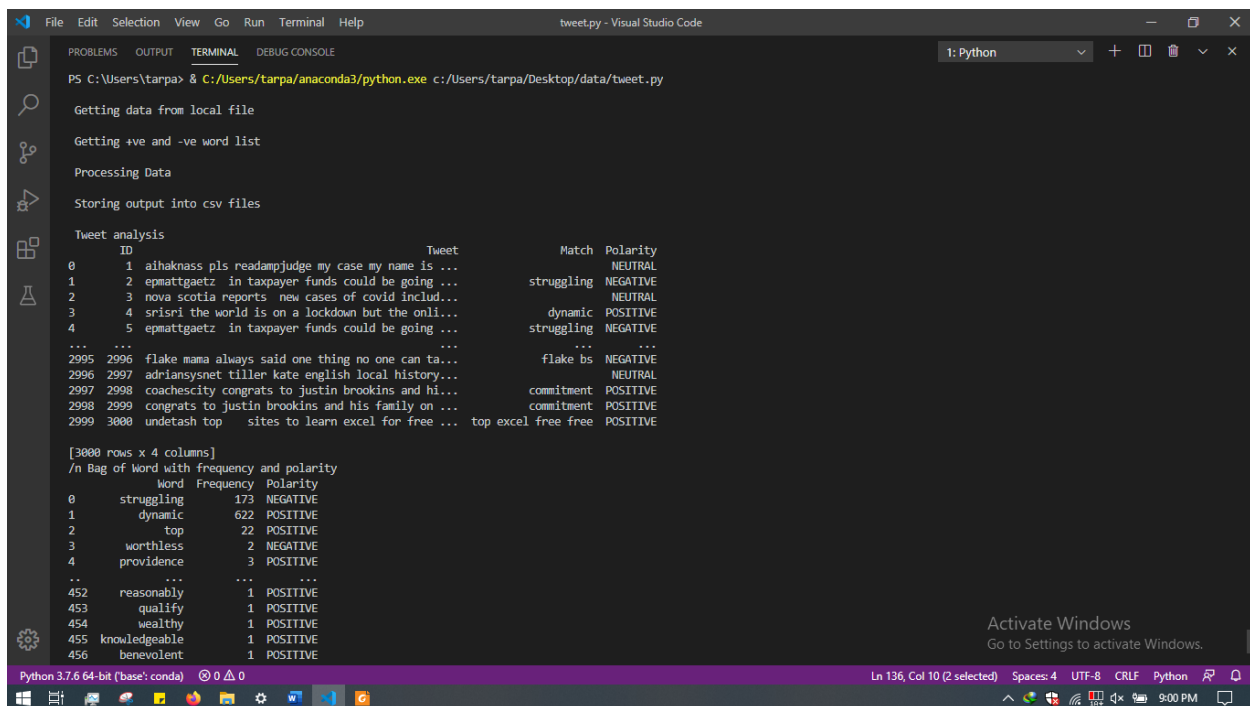# Assignment 4

## A. Sentiment Analysis

- "Tweet.py" python script fetched data and process them as given instructions.
- "getdata()" function fetched all the tweets then cleans it and store it in a text file – "data.txt", which is commented in the submitted script.
- Script then reads the cleaned data from "data.txt" and process it as per given instructions.
- Text file of positive and negative words list is obtained from links [1][2].
- Each tweet text is converted into list of words and then compared to positive and negative word list, if the count of positive words is greater than negative word count then tweet's polarity is "POSITIVE".
- Similarly, if negative word count is greater than positive word count then it is marked as "NEGATIVE" tweet. If both count is equal, then tweet is marked as "NEUTRAL".
- At the same time, frequency of each word is recorded into a dictionary – bag of word, which is then converted into dataframe to print it properly.



*Figure 1: Output of tweet.py*

- These dataframes are saved as CSV files named **"bog.csv"** and **"output.csv"**.
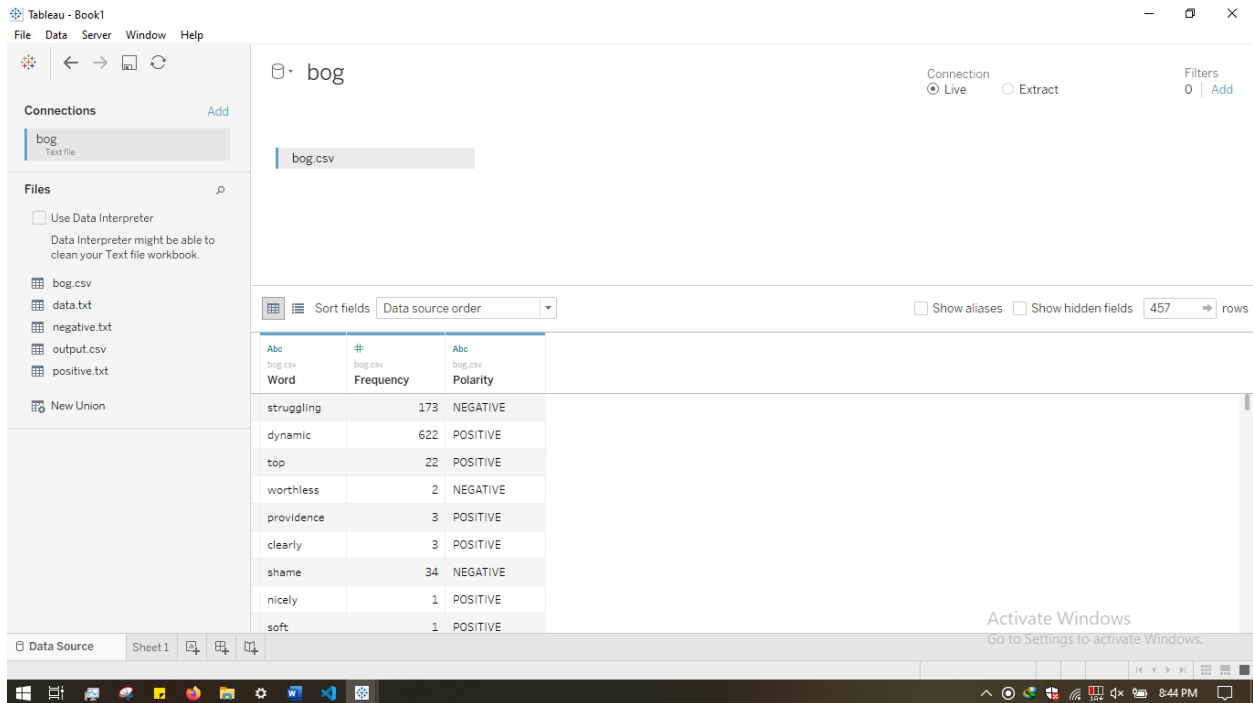- Bog.csv file is loaded as data source in Tableau to create word cloud.

*Figure 2: Load csv into Tableau*

- In Tableau word cloud, frequency of words is selected as size and polarity as color of text in chart, as shown in image. [3][4]
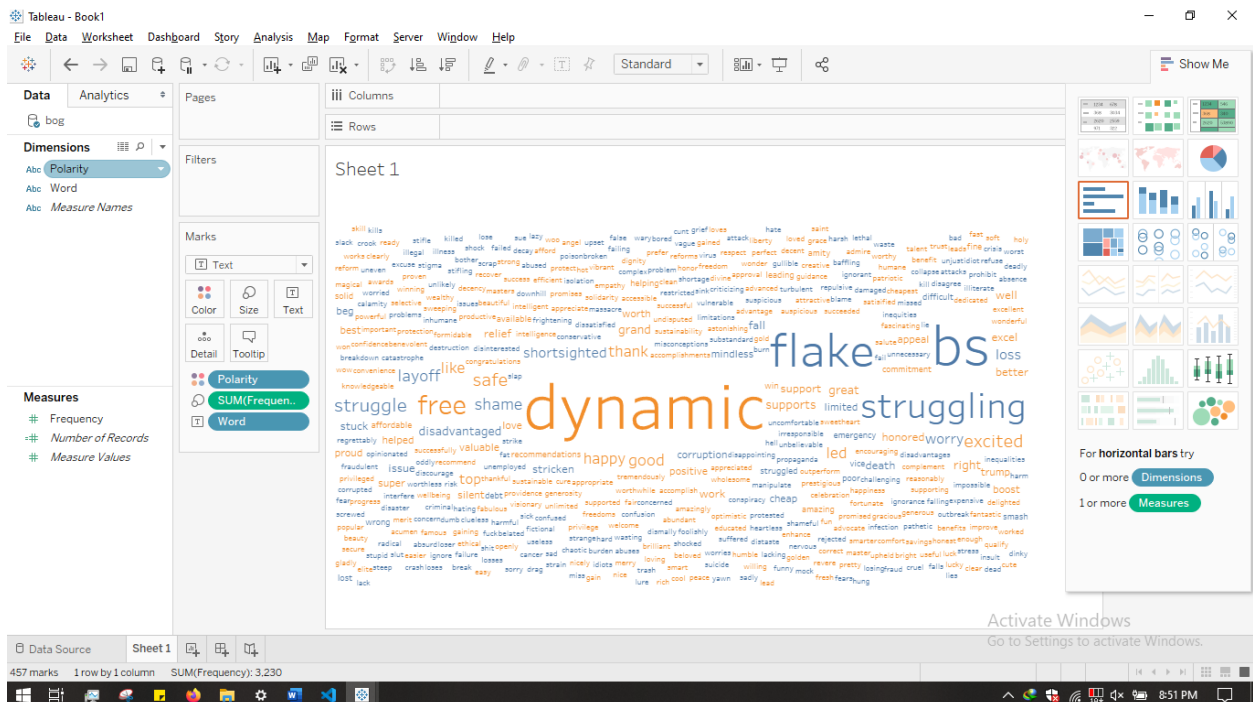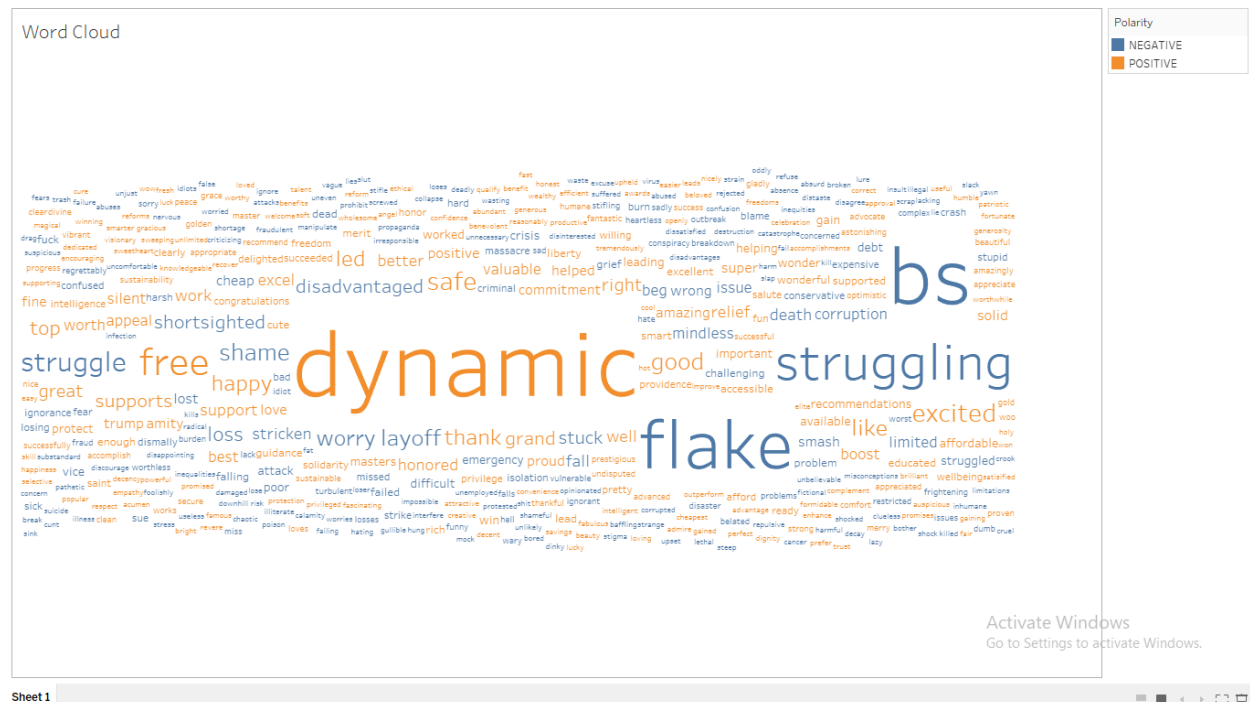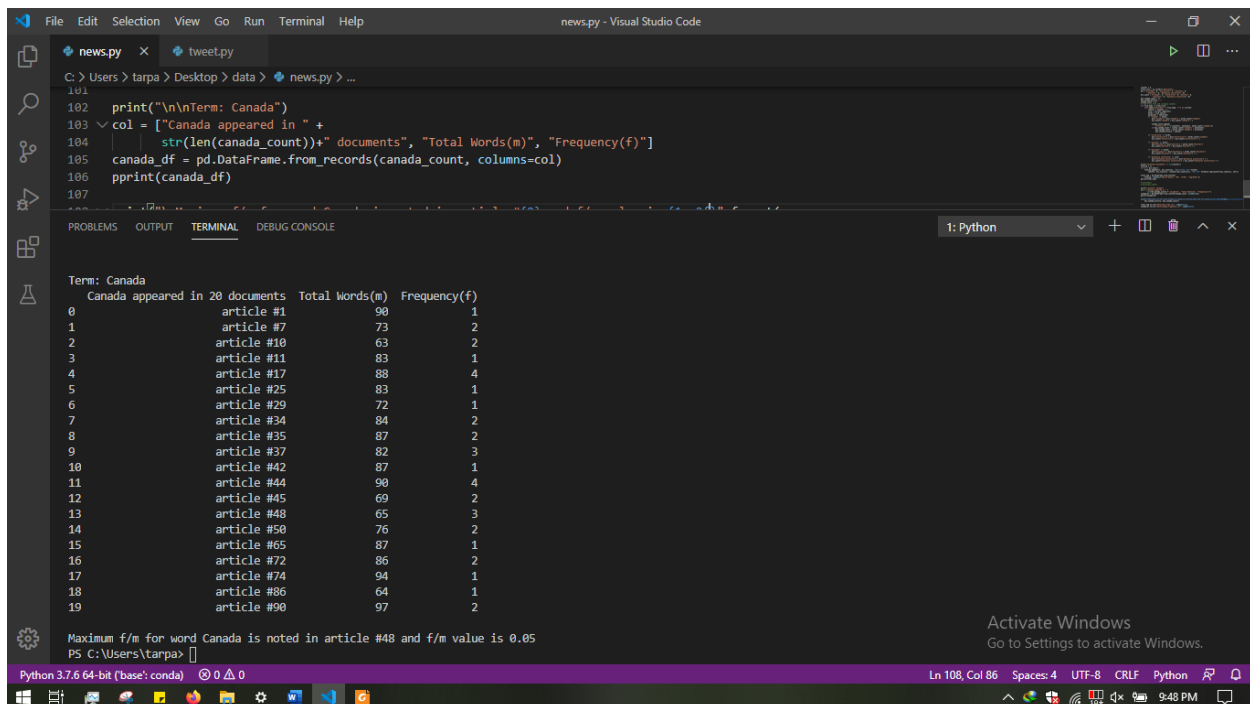


*Figure 3: Create Word cloud*

Word Cloud



Figure 4: Word cloud

## B. Semantic Analysis

- Created a script named "news.py", that reads "news.json" file. It fetches title, content and description of every news article and converts into one string,
- This string is then cleaned and stored into new file in "news" folder.
- After creating new file for every news article, script reads each file one by one and do the processing as per given instructions.
- Script finds the number of text files in "news" folder and records the frequency of words given in instruction.
- At the same time, it records the article name in which "Canada" appears, and records it in separate list, to create Canada frequency count per document table, as shown in following image.



*Figure 5: Output of news.py*

- Outputs are stored as CSV files named – **"news_tfdf.csv"** and **"Canada_frequency.csv"** that can be found in "data" folder.

**References:**

[1] *Ptrckprry.com*, 2020. [Online]. Available: http://ptrckprry.com/course/ssd/data/negative-words.txt. [Accessed: 14- Apr- 2020]

[2] *Ptrckprry.com*, 2020. [Online]. Available: http://ptrckprry.com/course/ssd/data/positive-words.txt. [Accessed: 14- Apr- 2020]

[3] "Word Clouds in Tableau: Quick & Easy.", *Medium*, 2020. [Online]. Available: https://towardsdatascience.com/word-clouds-in-tableau-quick-easy. [Accessed: 14- Apr- 2020]

[4] "Thanks for choosing a free trial of Tableau Desktop.", *Tableau Software*, 2020. [Online]. Available: https://www.tableau.com/en-ca/products/desktop/download. [Accessed: 14- Apr- 2020]