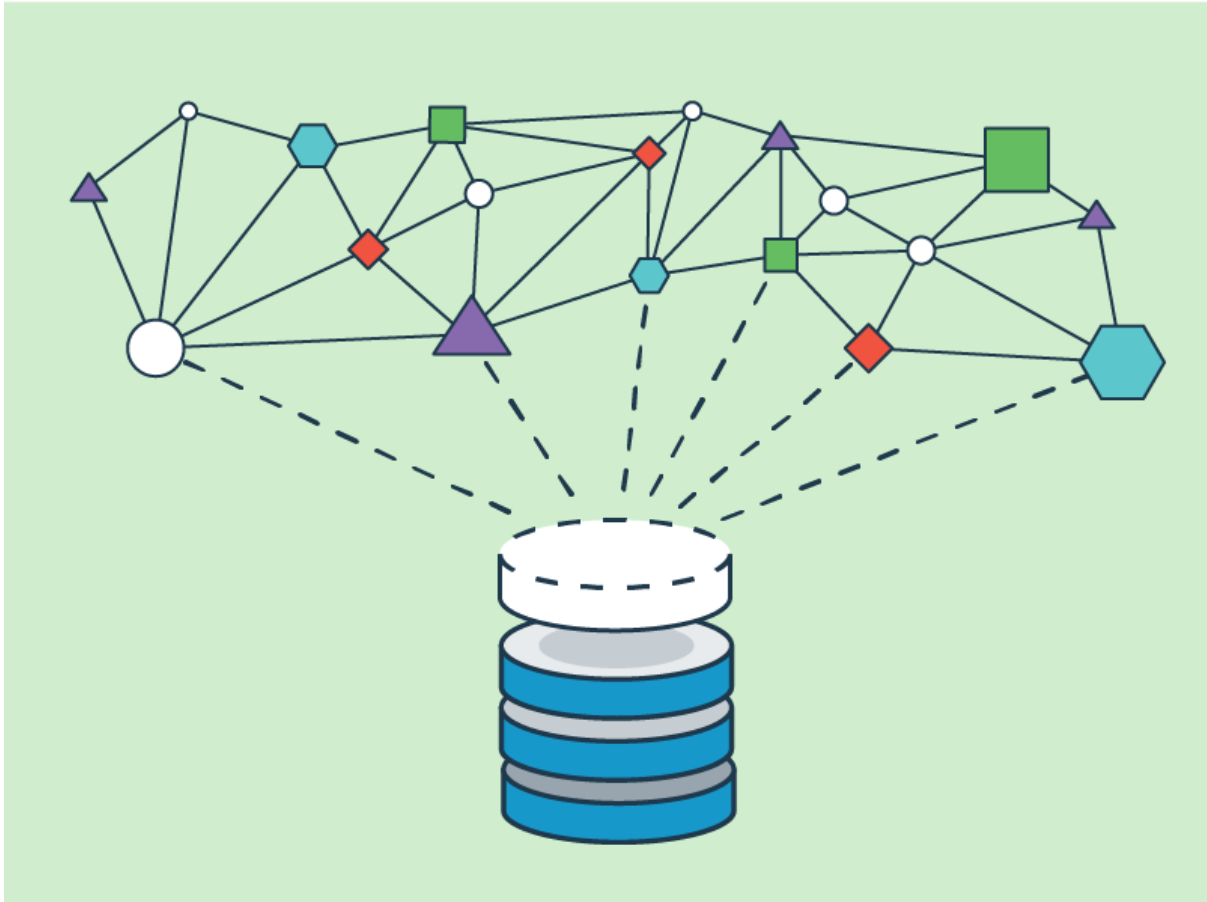# ASSIGNMENT 4 REPORT

*CSCI 5408 Data Management, Warehousing and Analytics*

# ANAND BHADANIA

B00837967

# INTRODUCTION

Assignment 4 was given with several tasks which include the following steps:

1. Sentiment Analysis
2. Semantic Analysis

# SENTIMENT ANALYSIS

To perform this task, I've used the python script which was used to extract tweet from twitter api in Assignment 3. I modified that script to collect only the tweet's message and ignored the metadata. The extracted tweets are in the clean format which I cleaned using Regular Expression (RE) in the Python script. After that, to obtain the polarity for Positive and Negative words, I downloaded files from online sources to get different positive and negative words [1][2].

Then I wrote another python script to classify each tweet as either "positive", "negative", or "neutral". The script is named "Sentiment_Analysis_Script" inside the folder named "sentiment analysis". The script generates 3 excel files which are also inside the folder mentioned above. File named "sentiment_analysis" records the polarity of each tweet along with tweet's message and the word with which the content of the tweet matched.
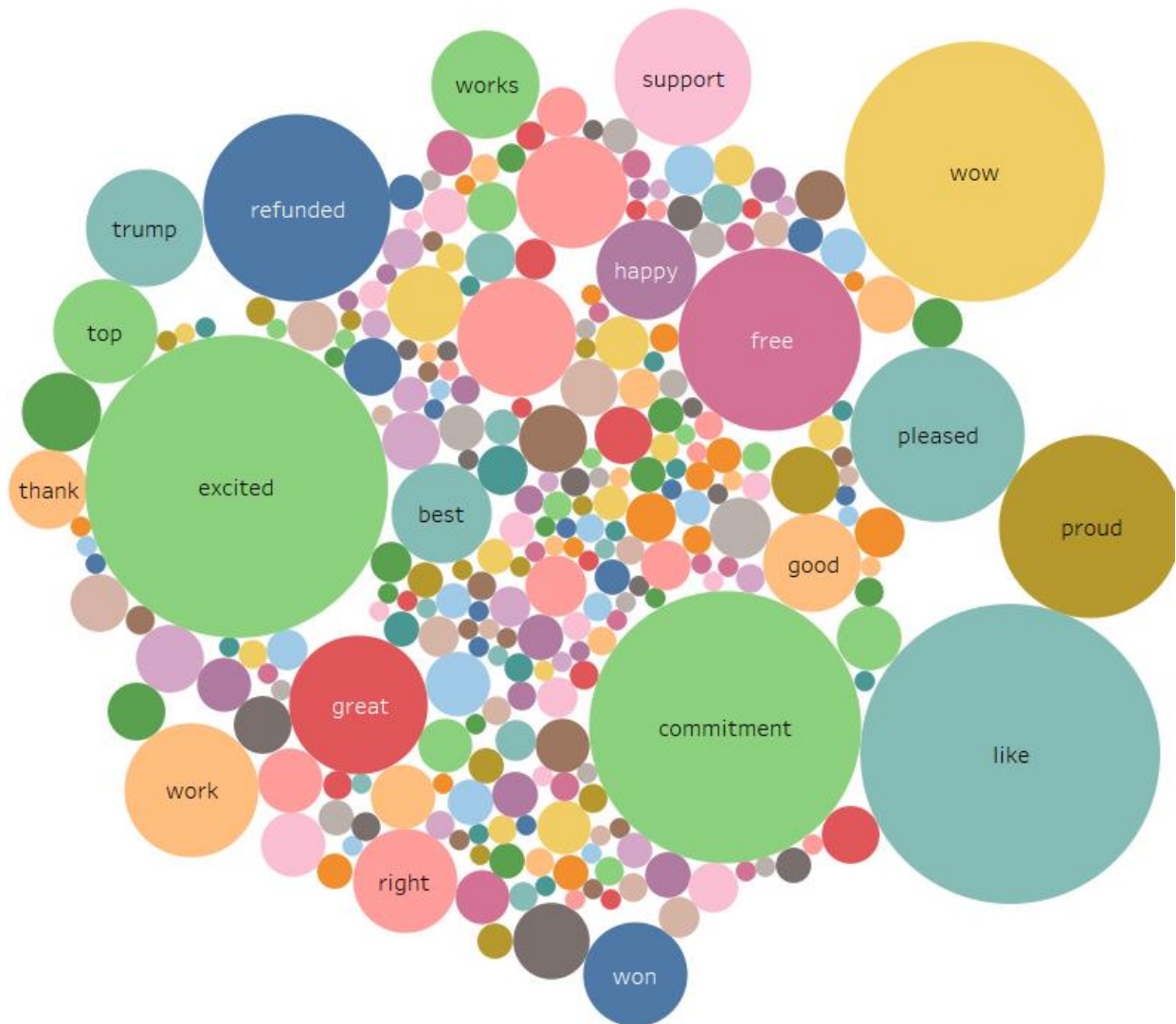
Other 2 files named "Positive_words_count" "Negative_words_count" contains the positive and negative words along with their count of occurrences in the tweets respectively.



*Figure 1 sentiment_analysis CSV file | Source: Author*

For visualization, I installed tableau from the link mentioned in the lab7 under Lab-Tutorial section under course's content on brightspace [3][4]. I imported excel files named

"Positive_words_count" and "Negative_words_count" to create visualization for the most frequently occurring words in the positive and negative tweets respectively [5].



Words. Color shows details about words. Size shows sum of Count. The marks are labeled by words.

*Figure 2 Positive Words Visualisation| Source: Author*

Words. Color shows details about words. Size shows sum of Count. The marks are labeled by words.

*Figure 3 Negative Words Visualisation| Source: Author*

The visualisation for both positive and negative words and attached in pictures folders.

## SEMANTIC ANALYSIS

To perform Semantic analysis, I have used the python script, to extract news articles from the news API, used in Assignment 3. I modified that python script to generate different files for different news articles. Total files generated are 700 and all files are inside the folder named "semantic analysis". I created a python script named "Semantic_Analysis_Script" under folder "semantic analysis" which creates 2 excel files and print the article on the console which had highest relative frequency which is explained in the assignment 4 pdf file. One file is named

"_SemanticAnalysis" which stores the frequency count of occurrence of different words(Canada, University, Dalhousie University, Halifax, Business) in different news articles along with additional information. Second file named "_SemanticAnalysis2" stores the information about different articles which contained the word "Canada" and stores the number of words in the article along with the frequency of the word "Canada".



*Figure 4 : _SemanticAnalysis CSV file| Source: Author*



*Figure 5 _SemanticAnalysis2 CSV file| Source: Author*

The article with the highest relative frequency printed on the console is shown in the Figure 6.



*Figure 6 Highest relative frequency Article| Source: Author*

**\*** *For this assignment, all the scripts and images are available in the folder.*

# REFERENCES

1. 262588213843476, "negative-words.txt," *Gist*, 14-Dec-2012. [Online]. Available: https://gist.github.com/mkulakowski2/4289441. [Accessed: 08-Apr-2020].
2. 262588213843476, "positive-words.txt," *Gist*, 14-Dec-2012. [Online]. Available: https://gist.github.com/mkulakowski2/4289437. [Accessed: 08-Apr-2020].
3. "BI Tools and Tableau - CSCI5408 - Data Mgmt  Warhsng Analytics (Sec 1) - 2020 Winter," *Brightspace.com*, 2020. [Online]. Available: https://dal.brightspace.com/d2l/le/content/110354/viewContent/1624572/View. [Accessed: 08-Apr-2020].
4. "Product Trials," *Tableau Software*, 2012. [Online]. Available: https://www.tableau.com/products/trial. [Accessed: 08-Apr-2020].
5. Parul Pandey, "Word Clouds in Tableau: Quick & Easy. - Towards Data Science," *Medium*, 20-Feb-2019. [Online]. Available: https://towardsdatascience.com/word-clouds-in-tableau-quick-easy-e71519cf507a. [Accessed: 08-Apr-2020].