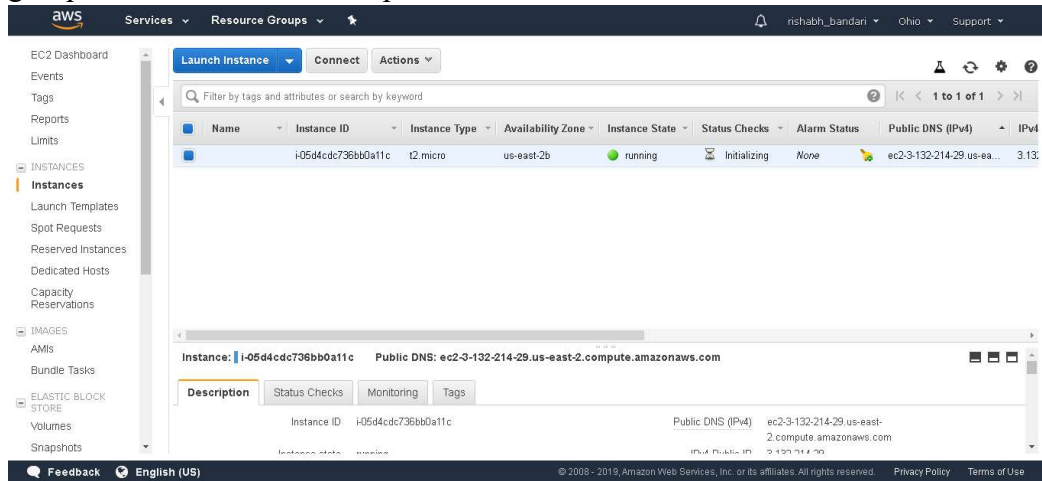


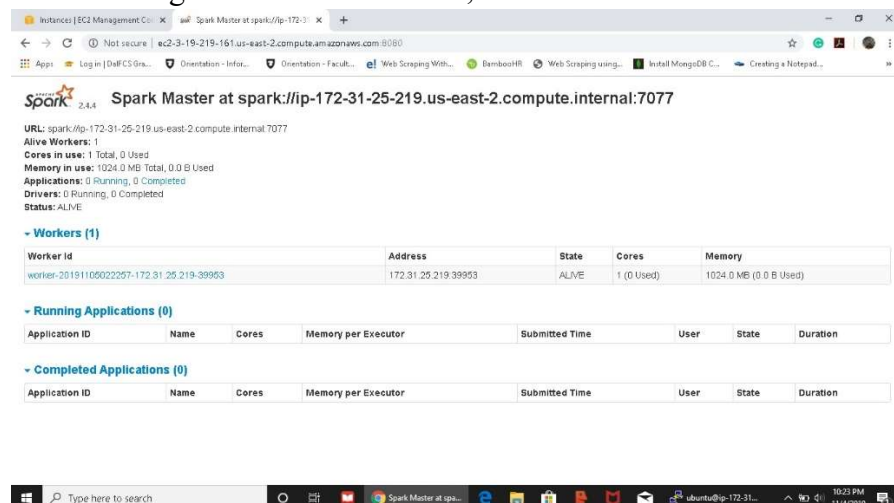
DATA MANAGEMENT WAREHOUSING AND ANALYTICS – ASSIGNMENT 2

A. Cluster Setup

- The cloud setup was done with the help of Amazon Elastic Compute Cloud (EC2). The steps in Lab: Tutorial 3 was followed to create the AWS instance. Finally, the security groups were edited to enable ports 8080 and 8081.



- After creating the instance, the installation and setup process of Apache Spark was done with the help of Tutorial 5 from the Labs. (Version 2.4.4 was installed)
 - Using the ubuntu virtual machine, the installation file was downloaded in the new server directory and unpacked using the following commands.
 - `wget http://apache.forsale.plus/spark/spark-2.4.4/spark-2.4.4-bin-hadoop2.7.tgz`
 - `sudo tar zxvf spark-2.4.4-bin-hadoop2.7.tgz`
 - 'export PYSPARK_PYTHON=python3' was used to install python.
 - After starting the master and slave, the number of workers was checked.



- Finally, the 'sudo ./spark-2.4.0-bin-hadoop2.7/bin/pyspark' command was run to start the python shell

B. Data Extraction Process

The extraction of tweets was done with the help of Python. The Tweepy API was used in retrieving and cleaning the data. In the case of the NewsAPI, the open source python code was used to extract the articles. The jupyter notebook was used for all the processes. The tweepy API gives access to the RESTful API of Twitter.

- Twitter Data Extraction & Transformation:
 1. Firstly, with the help of the Twitter developer account, a Twitter application was created. This application gives us the required Consumer Keys and Access Tokens that will be required to access the API.
 2. Tweepy was installed with the pip install command.
 3. Search API: This allows us to retrieve historical data of tweets. In simple words, it works and retrieves the same way the twitter search bar searches for tweets and related keywords. The major function used in the program uses just the Consumer Key, Consumer Secret, Access Token and Access Token Secret as parameters. As seen in the python file attached for search tweets, the program accesses the API, retrieves the tweets and stores them in an output csv file. I have only included the timestamps, tweet text and username in this retrieval process. Attention was given to excluding out special characters, links and other unwanted text.
 4. Stream API: Streaming of tweets involves retrieval of live Twitter data. The streaming API is used to push data into a particular session on a real time basis. The tweepy.Stream creates a streaming session and sends data to the stream listener. The method, on_data receives the messages depending on the features required. This the streaming of data was done with the help of the StreamListener, Stream object and the Stream Connection.

(Please find the respective ipynb files attached as part of the complete zip file.)

- News API Data Extraction
 1. A developer account was first created in newsapi.org and the API key was generated.
 2. The newsapi python library was installed using - pip install newsapi-python
 3. Using the open source modules, the data was extracted based on the keyword. I had restricted the source to CBC News. (cbc.ca)
 4. The request parameter q is used to search for specific keywords or phrases.
 5. Since a free version of the news API was used, a request limitation was encountered. Hence, only the first 100 records of data were retrieved.

C. Cleaning Process

- The data retrieved from the News API did not require much cleaning as it was mostly relevant and only limited records retrieved.
- Much of the cleaning was done while streaming and searching for the tweets with the help of the tweepy library.

- While searching for the tweets, the 'encoding="utf-8"' was while creating the output file. This removes all the unwanted characters and links.
- In the streaming process,

links = re.sub(r"https\S+", " ", tweet) was used to remove the unwanted links.

chars = re.sub(r'[^\x00-\x7f]+' , ' ', links) was used to remove the special characters.

- Finally, all the data was collected and structured manually with the help of Notepad and Excel.

(Please find attached the respective files for retrieved tweets.)

D. Data Processing

- The cleaned tweets were uploaded to Spark with the help of the scp command, which allowed for copying files between two different directory locations (Local to remote system).

```

Daniel@LAPTOP-UUDUGONS MINGW64 ~/Desktop/MACS/M1/AWS (master)
$ scp -i keypair_1.pem streamRT.txt ubuntu@ec2-3-19-219-161.us-east-2.compute.amazonaws.com:server
The authenticity of host 'ec2-3-19-219-161.us-east-2.compute.amazonaws.com (3.19.219.161)' can't be established.
ECDSA key fingerprint is SHA256:gmkX9axOnrdR+AxDMxkCOHufxa/+jXzTwKttFXt69Y.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-3-19-219-161.us-east-2.compute.amazonaws.com,3.19.219.161' (ECDSA) to the list of known hosts.
streamRT.txt                                100% 326KB 975.3KB/s   00:00

Daniel@LAPTOP-UUDUGONS MINGW64 ~/Desktop/MACS/M1/AWS (master)
$ scp -i keypair_1.pem searchlistRT.txt ubuntu@ec2-3-19-219-161.us-east-2.compute.amazonaws.com:server
searchlistRT.txt                            100% 347KB 690.6KB/s   00:00

Daniel@LAPTOP-UUDUGONS MINGW64 ~/Desktop/MACS/M1/AWS (master)
$ scp -i keypair_1.pem NewsData.txt ubuntu@ec2-3-19-219-161.us-east-2.compute.amazonaws.com:server
NewsData.txt                                100% 21KB 68.9KB/s   00:00

Daniel@LAPTOP-UUDUGONS MINGW64 ~/Desktop/MACS/M1/AWS (master)
$ |

```

- The text file is read, and the count of every unique word is calculated with the help of the Map-Reduce.

```

ubuntu@ip-172-31-25-219:~/server$ sudo ./spark-2.4.4-bin-hadoop2.7/bin/pyspark
Python 2.7.15+ (default, Oct 7 2019, 17:39:04)
[GCC 7.4.0] on linux2
Type "help", "copyright", "credits" or "license()" for more information.
19/11/06 13:51:05 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
19/11/06 13:51:07 WARN Utils: Service 'SparkUI' could not bind on port 4040. Attempting port 4041.
Welcome to

Spark version 2.4.4

Using Python version 2.7.15+ (default, Oct 7 2019 17:39:04)
SparkSession available as 'spark'.
>>> import sys
>>> from pyspark import SparkContext, SparkConf
>>> search = sc.textFile("searchlistRT.txt")
>>> words = search.flatMap(lambda line: line.split(" "))
>>> wordCounts = words.map(lambda word: (word, 1)).reduceByKey(lambda a,b:a +b)
>>> wordCounts.collect()
[(u',', 77380), (u'https://t.co/2nLCVw68xP', 1), (u'NuclearQuaffle', 1), (u'ESTER', 2), (u'yellow', 1), (u'Heights', 1), (u'four', 2), (u'Does', 1), (u'Thriving', 1), (u'Ohhh.Myy.Josh', 1), (u'Foundation', 3), (u'Resources', 1), (u'Schaudhary', 1), (u'Dohyon', 1), (u'Les', 2), (u'Community', 1), (u'chailbisoot', 1), (u'ItsdXPE', 1), (u'navigating', 1), (u'Height', 1), (u'Tennessee', 1), (u'Valle', 1), (u'Regional', 1), (u'Seno001', 1), (u'Bandybrexiteer', 1), (u'Sugar', 2), (u'bringing', 1), (u'DiscerningCritt', 1), (u'NACEports', 1), (u'gregrachac', 1), (u'Bigchiefkenkey', 1), (u'CGPA', 1), (u'Jeremywoorbyn', 1), (u'protest', 1), (u'29-year-old', 2), (u'RodneyReed', 1), (u'Paul', 3), (u'shows', 1), (u'JerrellLittle32', 1), (u'ility', 1), (u'Matthew', 1), (u'272', 1), (u'admitted', 1), (u'https://t.co/OkyyvpS8e0', 1), (u'formula_provent', 1), (u'FINSTA', 1), (u'Boise', 1), (u'sales', 1), (u'Smmrbates', 7), (u'SFFMemphis', 10), (u'Besakjesphoto', 1), (u'contributes', 1), (u'Presidential', 1), (u'Whose', 2), (u'Church', 2), (u'IrishTimes', 1), (u'here', 3), (u'dorm', 12), (u'brainwashes', 1), (u'departments', 1), (u'NathanielWanja3', 1)

```

- The last step involved performing a frequency count of the following substrings or words.

Word	Count
Education	378
Canada	495
University	1791
Dalhousie	2
expensive	3
good school/schools	20
bad school/schools	0
faculty	13
Computer Science	2
Graduate	586

It can be inferred that the word University, had the highest frequency whereas there were no tweets with the phrase bad schools.

E. References

- [1] ritvikmath, "GitHub," GitHub, [Online]. Available: <https://github.com/ritvikmath/ScrapingData/blob/master/Scraping%20Twitter%20Data.ipynb>. [Accessed 25th October 2019].
- [2] J. Roesslein, "Tweepy.org," [Online]. Available: http://docs.tweepy.org/en/v3.4.0/streaming_how_to.html. [Accessed 28 October 2019].
- [3] Harrison, "pythonprogramming," [Online]. Available: <https://pythonprogramming.net/twitter-api-streaming-tweets-python-tutorial/>. [Accessed 26 October 2019].
- [4] "StackOverflow," [Online]. Available: <https://pythonprogramming.net/twitter-api-streaming-tweets-python-tutorial/>. [Accessed 30 October 2019].
- [5] "Python Examples," [Online]. Available: <https://pythonexamples.org/pyspark-word-count-example/>. [Accessed 2 November 2019].
- [6] "ApacheSpark.org," [Online]. Available: <https://spark.apache.org/examples.html>. [Accessed 3 November 2019].
- [7] "Tweepy," [Online]. Available: <http://docs.tweepy.org/en/v3.4.0/index.html>. [Accessed October 2019].