Akash Bharti
B00837970

# CSCI 5408 - Data Management, Warehousing and Analytics
# Assignment 3

### A. Cluster Setup:

I create an account on AWS [1]  and then created an instance of EC2. Figure 1 shows the console of EC2 and running instances. I then installed Apache Spark [2]  on the EC2 instance. Figure 2 shows the spark-shell running on the EC2. I then installed MongoDB [3] on the instance. Figure 3 shows the mongo shell running on the instance.
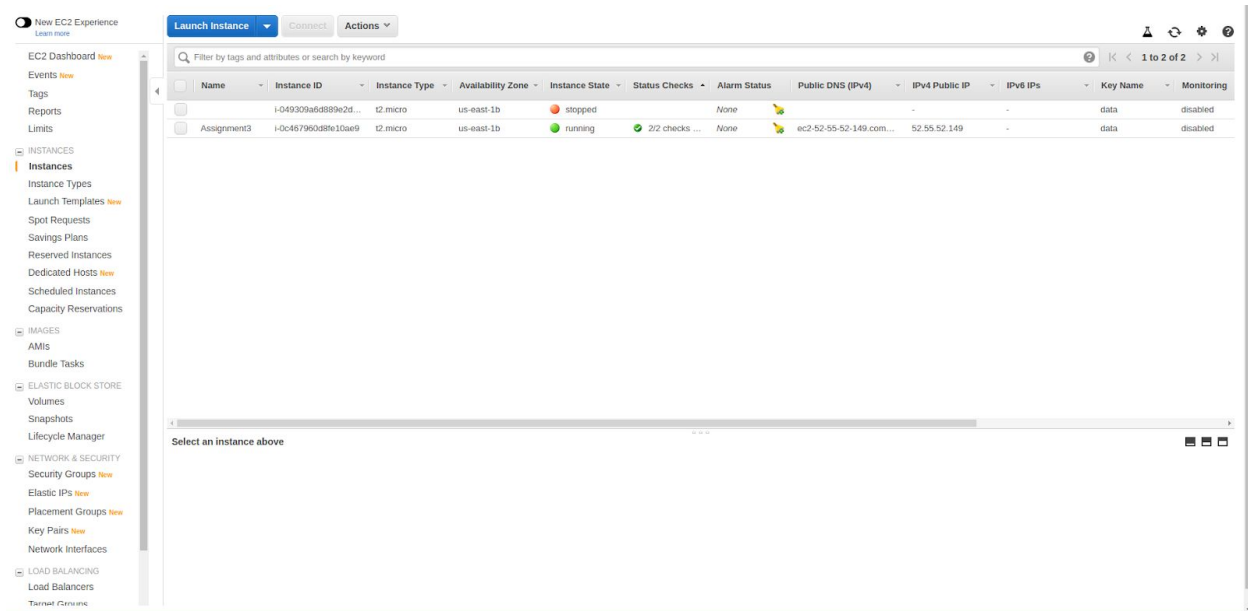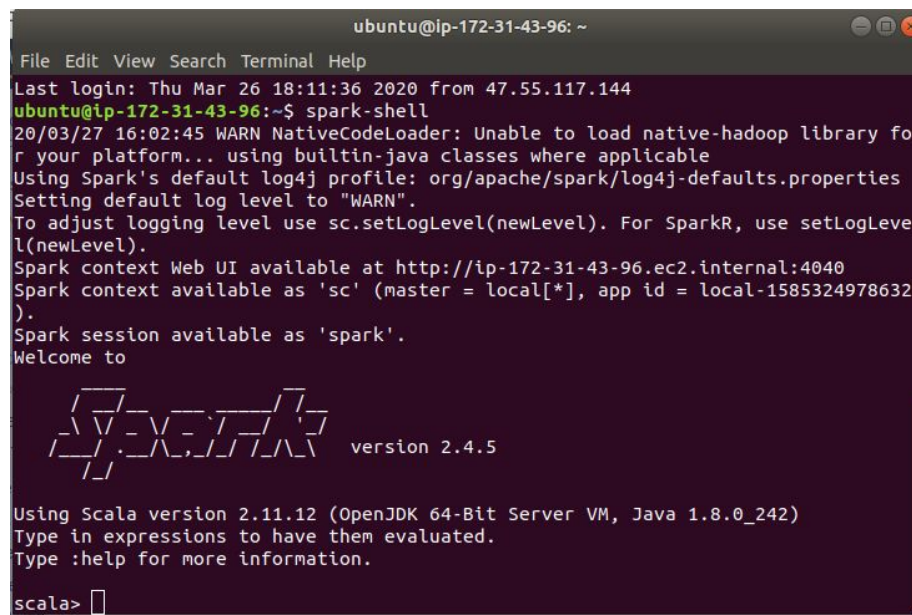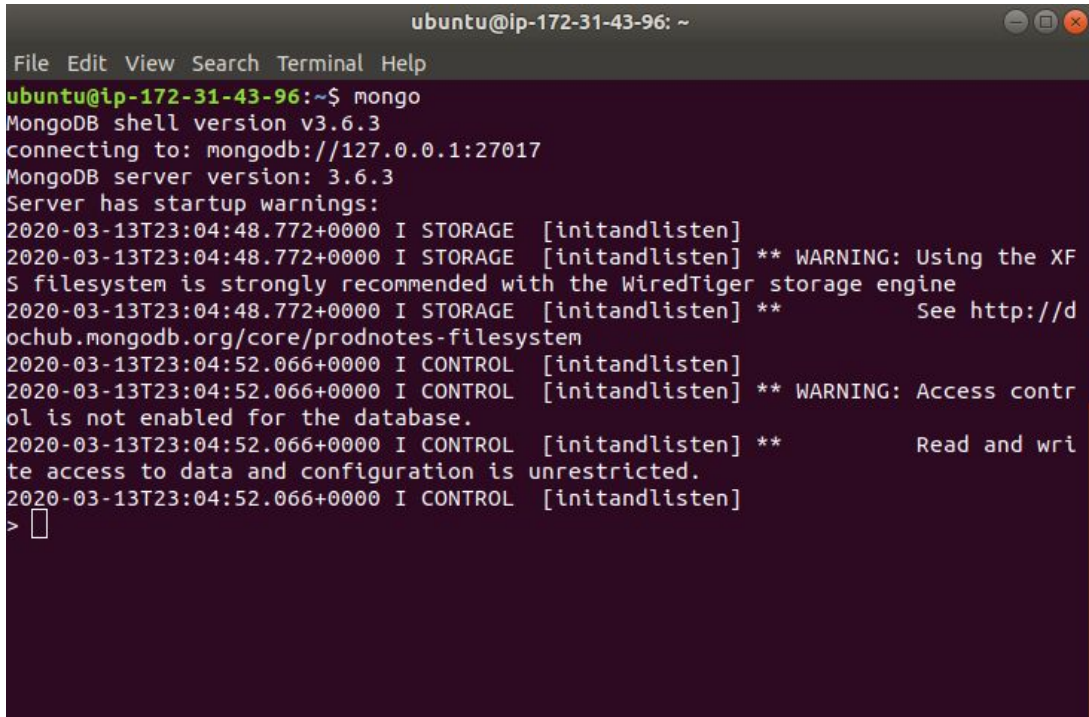


*Figure 1*



*Figure 2*

*Figure 3*

**B. (i)Twitter Data Extraction and Transformation:**

I created the twitter developer account and it was approved in 3 days. I used the tweepy [5] python package to extract the tweets from twitter APIs [4]. I used both the streaming and REST APIs to extract the tweets. I used the following scripts to extract the data:

1. *twitterStreamingApi.py* - This script is used to extract tweets for keywords "Canada" and "University". The script uses twitter streaming API to extract the tweets. The scripts stores the extracted tweets in file tweets.json.

2. *twitterRestApi.py* - This script is used to extract tweets for keywords "Dalhousie University", "Halifax" and "Canada Education". The script uses twitter REST API to extract the tweets. The scripts stores the extracted tweets in file tweets.json.

3. Clean_tweets.py - This script cleans the extracted tweets and then stores them in MongoDB cluster. The script uses regular expressions to clean the tweets. The script also saves the tweets in file *input_for_word_count.txt* which will be later used to process the word counts for different keywords.

Figure 4 shows the tweets stored in MongoDb cluster.

(ii) News Articles Data Extraction and Transformation:

The *news_api.py* script extracts the news articles from News API [6]  for all the keywords required, cleans them, stores them on the MongoDB cluster and also saves them in news_articles.json file. Figure 5 shows the contents of the news_articles.json file.

*Figure 5*

Figure 6 shows the news articles stored in MongoDb cluster.



*Figure 6*

**(iii) Movie Data Extraction and Transformation:**

The *movies_api.py* script extracts the movies data from OMDb API [7] for all the keywords required, cleans them, stores them on the MongoDB cluster and also saves them in movies.json file. Figure 7 shows the contents of the movies.json file.

*Figure 7*

Figure 8 shows the movies data stores in MongoDb cluster.



*Figure 8*

## C. Data Processing (Spark):

The script *word_count.py* reads cleaned tweets and news articles from file *input_for_word_count.txt* and uses pyspark python package [8] to calculate the word count of the provided keywords. The script initializes a spark context and then uses map-reduce to print the word count of the provided keyword. Figure 9 shows the output of the script.

```
20/03/27 16:55:57 WARN Utils: Your hostname, akash-Predator-PH315-51 resolves to
 a loopback address: 127.0.1.1; using 192.168.2.19 instead (on interface wlp0s20
f3)
20/03/27 16:55:57 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another
address
20/03/27 16:55:57 WARN NativeCodeLoader: Unable to load native-hadoop library fo
r 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 setLogLeve
l(newLevel).
('education', 23)
('Canada', 899)
('University', 699)
('Dalhousie', 22)
('expensive', 2)
('good school', 0)
('good schools', 0)
('bad school', 0)
('faculty', 3)
('Computer Science', 0)
('graduate', 14)
```

*Figure 9*

## Extracting movie Information:

The script *movies_processing.py* fetch the movies data from mongoDb and then displays the movie title, rating, genre and plot on the console. Figure 10 shows the output of the script.

```
akash@akash-Predator-PH315-51: ~/Documents/A3

File  Edit  View  Search  Terminal  Help
akash@akash-Predator-PH315-51:~/Documents/A3$ python movie_processing.py
Title : Big Brother Canada
Rating :
Internet Movie Database  5.5/10
Genre : Reality-TV
Plot : The format for Big Brother in Canada remains largely unchanged from the U S edition making th
em the only two version of the series thus far to follow this format HouseGuests are

Title : University
Rating :
Internet Movie Database  6.4/10
Genre : N/A
Plot : N A

Title : Halifax
Rating :
Genre : Short, Horror, Thriller
Plot : Following the death of her father Beatrice returns to her childhood home causing a series of
repressed memories to surface which sparks a nightmarish presence to haunt her and force her

Title : CBC News: Moncton Shooting
Rating :
Genre : News
Plot : N A

Title : I Hate Toronto: A Love Story
Rating :
Internet Movie Database  3.5/10
Genre : Comedy, Drama, Romance
Plot : A young man s life is struck by tragedy The only way out seems to be suicide but he decides t
o have a bit of fun first It takes 200 meaningless women to meet the one woman that could give meani
ng to his life again

Title : Vancouver
Rating :
```

*Figure 10*

References:
[1] "What Is Amazon EC2? - Amazon Elastic Compute Cloud," *Amazon.com*, 2020.
[Online]. Available:
https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/concepts.html. [Accessed:
27-Mar-2020].

[2] "Overview - Spark 2.4.5 Documentation," *Apache.org*, 2020. [Online]. Available:
https://spark.apache.org/docs/latest/. [Accessed: 27-Mar-2020].

[3] "MongoDB Documentation," *Mongodb.com*, 2020. [Online]. Available:
https://docs.mongodb.com/. [Accessed: 27-Mar-2020].

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

[5] "Tweepy Documentation — tweepy 3.8.0 documentation," *Tweepy.org*, 2020.
[Online]. Available: http://docs.tweepy.org/en/latest/. [Accessed: 27-Mar-2020].

[6] Documentation - News API, "Documentation - News API," *Newsapi.org*, 2020.
[Online]. Available: https://newsapi.org/docs. [Accessed: 27-Mar-2020].

[7] "OMDb API - The Open Movie Database," *Omdbapi.com*, 2020. [Online]. Available:
http://www.omdbapi.com/. [Accessed: 27-Mar-2020].

[8] "Welcome to Spark Python API Docs! — PySpark 2.4.5 documentation," *Apache.org*,
2020. [Online]. Available: https://spark.apache.org/docs/latest/api/python/index.html.
[Accessed: 27-Mar-2020].