Juhilkumar P. Zalavadiya
B00872280

## Problem #1:

## Task 1: Cluster Setup - Apache Spark Framework on GCP

1. I opened my google cloud account [1].
2. Searched for VM instance and clicked on Create button.
3. Selected E2-medium and Ubuntu 20.04 LTS as per instructions.
4. Enabled HTTP and HTTPS traffic and permitted to allow all protocols and ports.
5. Opened terminal and installed JDK, scala, and git with the command: sudo apt install default-jdk scala git –y
6. Created a new folder Apache_Spark and installed apache-spark [2], and configured the path.
7. Started the master node with the command: start-master.sh

8. After that, I started the slave node using the command: start-slave.sh <url>

9. Opened <VM adress>:8080 to confirm the working of the cluster.

10. I installed pyspark for future use.

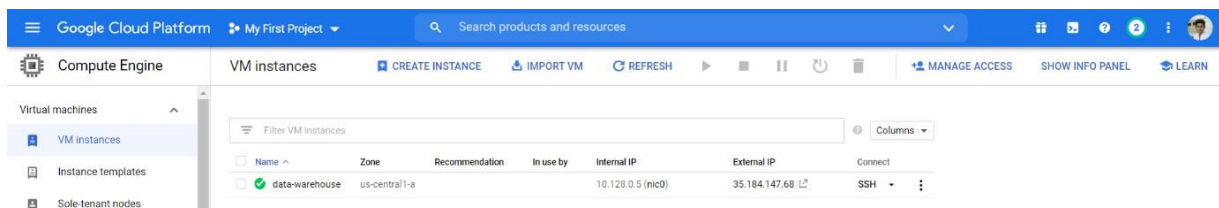   ➢ Figure 1 shows the VM instance used for the apache spark cluster.



*Figure 1 VM instance*

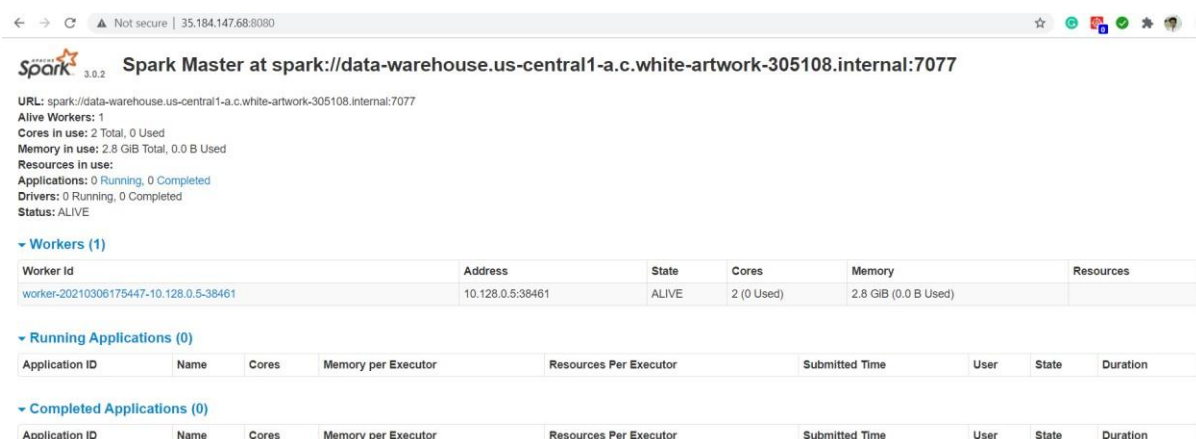   ➢ Figure 2 shows the spark dashboard with master and slave



*Figure 2 spark dashboard*

Juhilkumar P. Zalavadiya
B00872280

## Task 2:Data Extraction from Twitter:
- At first, I have created a Twitter developer account
- To extract data from Twitter, I have used tweepy library [4].
- In tweepy, I have used search and stream API for searching previously generated data and live streaming data.

## 1) Stream API:
- Stream API is  helpful in extracting live tweets.
- I have used the MongoClient library to connect with the MongoDB remote database [3].
- Also, I have authenticated the user based on Twitter-generated keys.
- I have used StreamListner class to search real-time tweets of the required keywords.
- In that class, I have given the path for the database and collection.
- I have defined a function on_status() to fetch tweet text with some metadata and inserted it into the database with the help of the insert() function provided by StreamListener class.
- To make 40 separate files for this data, I have used two counters: count and fileName. To change the file, I have used an if condition that checks for the count, and if the count is divisible by 100, it creates a new file and flushes 100 tweets in that file.
- I have used the fileName counter to give different file names in the database.
- One more if the condition will check for the counter. As mentioned in the assignment requirement, I have given an If the condition checks for the count value. If it is 4000, then the program will stop automatically.
- In the case of any error, I have handled it with the on_error() function.
- After running this program on GCP, files will automatically upload to the MongoDb remote database.

## 2) Search API:

- Search API is helpful in extracting previously generated tweets based on date. You have to specify a date, and the tweepy library will search it with the help of the Cursor() function and give you the tweets since that date.
- I have authenticated the user based on the keys given by Twitter.
- After that, I have connected MongoDB remote database with the help of the pymongo library.
- I have defined one list named "WORDS" to define keywords mentioned in the assignment requirements.
- I have defined a for loop to take each tweet with the [key, value] pair and extracted the text and metadata. After that, with the help of a predefined method insert(), I have inserted those fields into a MongoDB remote database collection.
- To store these data in separate files, I have used two counters: count and fileName. When the count is divisible by 100, this condition will store 100 tweets in one file in the MongoDb database named "myMongoTweet."
- I have used the fileName counter to give different file names in the database.
- I have defined one more if condition to terminate the program when the count value reaches 4000.

Juhilkumar P. Zalavadiya
B00872280

## Output screenshots for twitter data extraction:

### 1) Before generating the "myMongoTweet" database



*Figure 2. Before running the script*

### 2) After running the search data extraction script



*Figure 3. After running the search data extraction script*

## 3) Output of the search data extraction script on MongoDB



*Figure 4. The output of search data extraction on MongoDB*

## 4) After running the stream data extraction script



*Figure 5. After running the stream data extraction script*

Juhilkumar P. Zalavadiya
B00872280

**5) Output of the stream data extraction script on MongoDB**



*Figure 6. The output of stream data extraction on MongoDB*

**3) Cleaning process:**
- At first, I have established a connection to store data in the MongoDB remote database.
- After that, I have defined paths for a database named "myMongoTweet" where all the streaming files are located. I have used previously generated streaming data for cleaning.
- I have used the find() method to find the tweets in the file.
- With the help of for loop, I have accessed each field of tweets like id, text, location and cleaned those fields with the clean_data() function. This function will take that field as input and checks it with the regular expression. As per the assignment requirement, I have defined regular expressions to remove the emoticons, special characters, and URLs.
- After the cleaning part, I have used the insert_one() function to insert all that data into a separate file.

Juhilkumar P. Zalavadiya
B00872280

## **Output screenshots for twitter data extraction:**

### **1) Before applying the cleaning**



*Figure 7. Before applying the cleaning*

### **2) After running the cleaning script**



*Figure 8. After running the cleaning script*

Juhilkumar P. Zalavadiya
B00872280

## 3) Output of the cleaned data on MongoDB



*Figure 9. The output of the cleaned data on MongoDB*

## Problem #2

**Task 1: Data Processing using Spark - MapReduce to perform count**
- MapReduce is used for the processing of extensive data set. The task is divided into many smaller parts to make the processing efficient and fast.
- To perform this task, I have used the SparkContext library for mapping and reduce tasks. I also used the MongoClient library to connect with the MongoDB remote database.
- After the connection part, I have given the path for the database where all the cleaned data is stored in different files.
- For storing the list of searched keywords mentioned in the assignment requirement, I have defined the initial_list.
- After that, I have used one for loop to access different files in the database.
- I have used another for loop to break the string into a list of words and compare it with the keywords I want to search.
- After all these steps, I will have a list of words required for the given task. I have stored all these keywords in the initial_list.
- I used parallelized() method to convert that initial_list into Resilient Distributed Dataset.
- After that, mapping is done with the help of the map() function. Its output will be given for the reducing part. SpartkContext library has one method named reduceByKey() for the reducing part. So, I have used it and stored the output in the reduced_map list.

Juhilkumar P. Zalavadiya
B00872280

## The output of the MapReduce program:



*Figure 10. The output of the MapReduce - 1*



*Figure 11. The output of the MapReduce - 2*

Highest frequencies: emergency -> 152
Lowest frequencies: flu -> 48

Juhilkumar P. Zalavadiya
B00872280

# References:

[1]     "Google Cloud Services" , Google Cloud, 2021. [Online]. Available:
        https://cloud.google.com/ [Accessed: 5-March 2021]

[2]     "Apache Spark – Python Programming Guide" [Online]. Available:
        https://spark.apache.org/docs/0.9.0/python-programming-guide.html
        [Accessed: 6-March 2021]

[3]     "MongoDB Cloud", MongoDB,2021. [Online]. Available:
        https://www.mongodb.com/cloud [Accessed: 8-March 2021]

[4]     "Tweepy – An easy-to-use Python library for accessing the Twitter API"     [Online].
        Available: https://www.tweepy.org/
         [Accessed: 8-March 2021]