

COVID-19 TWITTER ANALYSIS
COURSE: BIG DATA PROGRAMMING
TEAM 4

FINAL REPORT

Vidyullatha Lakshmi Kaza- 8
Aparna Manda- 11
Lohitha Yenugu- 19

Introduction:

Data analysis on tweets pertaining to COVID19. The entire world is shutdown due to the virus and we wanted to know the public opinion on this situation. So, we collected their opinion through tweets. Collected real-time tweets talking about the corona virus with keywords- COVID19, Corona and performed analysis using big data technologies- Map Reduce, Hive, Cassandra and Map Reduce sentimental analysis.

Background:

Analyzed twitter data- json structure to extract useful attributes

<https://developer.twitter.com/en/docs/tweets/data-dictionary/overview/intro-to-tweet-json>

Twitter data preprocessing to remove characters like spaces, new lines and commas that might cause issues during csv encoding.

Features developed:

1. Twitter data collection on COVID-19
2. Map Reduce to count tweets by each user
3. Data Analysis of tweets using Hive
4. Sentiment Analysis of tweets using Map Reduce
5. Data Analysis of tweets using Cassandra
6. Twitter Data Analysis using Spark SQL

Data Collection:

Python Code to extract live stream twitter data as CSV:



The screenshot shows a code editor window with the following details:

- Title Bar:** TwitterDataExtraction.py
- File Path:** ~/Desktop/BigDataProgramming
- Code Content:** Python script for extracting live stream Twitter data. It includes functions for handling status updates and errors, and a main block for initializing the API and streaming data.

```
def on_status(self, status):
    print(status.id_str)
    # if "retweeted_status" attribute exists, flag this tweet as a retweet.
    is_retweet = hasattr(status, "retweeted_status")

    # check if text has been truncated
    if hasattr(status, "extended_tweet"):
        text = status.extended_tweet["full_text"]
    else:
        text = status.text

    # check if this is a quote tweet.
    is_quote = hasattr(status, "quoted_status")
    quoted_text = ""
    if is_quote:
        # check if quoted tweet's text has been truncated before recording it
        if hasattr(status.quoted_status, "extended_tweet"):
            quoted_text = status.quoted_status.extended_tweet["full_text"]
        else:
            quoted_text = status.quoted_status.text

    # remove characters that might cause problems with csv encoding
    remove_characters = [",","\n"]
    for c in remove_characters:
        text.replace(c, " ")
        quoted_text.replace(c, " ")

    with open("/home/lohittha/Desktop/BigDataProgramming/COVID19_Data.json", "a", encoding='utf-8') as f:
        f.write("%s,%s,%s,%s,%s\n" % (status.created_at, status.user.screen_name, is_retweet, is_quote, text, quoted_text))

def on_error(self, status_code):
    print("Encountered streaming error (", status_code, ")")
    sys.exit()

if __name__ == "__main__":
    # complete authorization and initialize API endpoint
    auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
    auth.set_access_token(access_token, access_token_secret)
    api = tweepy.API(auth)

    # initialize stream
    streamListener = StreamListener()
    stream = tweepy.Stream(auth=api.auth, listener=streamListener, tweet_mode='extended')
    with open("/home/lohittha/Desktop/BigDataProgramming/COVID19_Data.json", "w", encoding='utf-8') as f:
        f.write("date,user,is_retweet,is_quote,text,quoted_text\n")
    tags = ["covid19", "covid-19", "COVID-19", "COVID19", "corona", "CORONA"]
    stream.filter(track=tags)
```

Python Code to extract live stream twitter data as JSON:

```
TwitterDataExtractioncsv.py                                     COVID19_Data1.json
is_retweet = hasattr(status, "retweeted_status")
if hasattr(status,"extended_tweet"):
    text = status.extended_tweet["full_text"]
else:
    text = status.text

# check if this is a quote tweet.
is_quote = hasattr(status, "quoted_status")
quoted_text = ""
if is_quote:
    # check if quoted tweet's text has been truncated before recording it
    if hasattr(status.quoted_status,"extended_tweet"):
        quoted_text = status.quoted_status.extended_tweet["full_text"]
    else:
        quoted_text = status.quoted_status.text

# remove characters that might cause problems with csv encoding
remove_characters = [",", "\n"]
for c in remove_characters:
    text.replace(c, " ")
    quoted_text.replace(c, " ")

with open("/home/lohitash/Desktop/BigDataProgramming/COVID19_Data1.json", "a", encoding='utf-8') as f:
    f.write('{"id": "%s", "date": "%s", "tweetuser": "%s", "isretweet": "%s", "isquote": "%s", "tweetttext": "%s", "quotedtext": "%s"}\n' % (status.id_str,status.created_at,status.user.screen_name,is_retweet,is_quote,text,quoted_text))

def on_error(self, status_code):
    print("Encountered streaming error (", status_code, ")")
    sys.exit()

if __name__ == "__main__":
    # complete authorization and initialize API endpoint
    auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
    auth.set_access_token(access_token, access_token_secret)
    api = tweepy.API(auth)

    # initialize stream
    streamListener = StreamListener()
    stream = tweepy.Stream(auth=api.auth, listener=streamListener,tweet_mode='extended')
    with open("/home/lohitash/Desktop/BigDataProgramming/COVID19_Data1.json", "w", encoding='utf-8') as f:
        f.write("id,d,usr,\n")
    tags = ["covid19", "covid19", "COVID-19", "COVID19", "corona", "CORONA"]
    stream.filter(track=tags)
```

Collected Tweets:

```
date,user,is_retweet,is_quote,text,quoted_text
2020-04-15 02:06:43,kimpetty64,False,True,@@@ and just like Corona beer, Meanwhile, Trump supporters burn their 8-tracks cursing out Keith Moon, Pete Townshend and Roger Daltrey.

Blaming WHO
2020-04-15 02:06:43,2021Diary,True,False,RT @PinkeshOfficial: While PM Modi having serious discussion with states on Covid19, Aditya Thackeray is busy with his mobile.

Anjana was...
2020-04-15 02:06:43,rdc_south,True,False,RT @chennaiacorp: Here's the Graphical Representation of total COVID-19 positive cases in Chennai as on 14-04-2020.

#covid19chennai
#GCC_#,
2020-04-15 02:06:42,openletterbot,False,False, support Patrick by signing "Support the USPS!" and I'll deliver a copy to your officials too: https://t.co/Gj02cQ65qq

↳ Last delivered to @RepStephenLynch, @SenMarkey and @SenWarren #MA08 #APoli #MAPolis #COVID19 https://t.co/SSgHHzjoi,
2020-04-15 02:06:43,ahernandez85b,True,True,RT @KlemenCari: It's been 35 years and there is still no vaccine for AIDS. https://t.co/X8yKXTStG,As we reopen #Ohio, people will have to be ve
2020-04-15 02:06:43,luca_s_bhmg,True,True,RT @AbdelDeuxFois: Sans promo ni rien. Le monde chico ☺,ALERTE INFO - 36,7 millions de téléspectateurs ont suivi l'allocution d'Emmanuel #Macron
2020-04-15 02:06:43,celinev,True,False,RT @MonsieurQB: Aquí los análisis de casos de COVID-19 en México y países de América, Europa ue es KR BR DE JP FR PE AR IC co Y...
2020-04-15 02:06:43,MrsSol54252546,True,False,RT @islamramahdotco: Kiai Anwar Zahid: Patuhui Protokoler covid-19

"Kita ajak seluruh umat agar taat dan patuh instruksi pemerintah dan pat...
2020-04-15 02:06:43,drjawswantpatil,False,False,RT @Bachchan # Homeopathy can beat Corona. Want to see the results allow homoeopaths to treat and see the change.
2020-04-15 02:06:43,xanxxz,True,False,RT @gunawan_anas: Jujurnya pemerintah masih setengah.. selain data saat ini, pmrintah harus terbuka mengenai forecasting versi mereka akan...
2020-04-15 02:06:43,WUNNA_1,True,False,RT @21savage: Bang outside I hang outside
don't come out da house cuz corona outside,
2020-04-15 02:06:43,abnesdad,True,False,RT @TrumpHoodles: Michigan kid videotaping his dad who hates the Michigan governor

Funny as heck @@@#
Qute sounds like my mom when it co...
2020-04-15 02:06:43,respirovondy,True,False,RT @crushdobb20: Flávia Pavanelli fez foto, fez preenchimento, arrumou cabelo, unhas, encontrou amigos em casa e agora vem falar q acha q...
2020-04-15 02:06:43,carterpillar62,True,False,RT @franjuero: Un reportaje de la Televisión Italiana del 2015 donde habla de que China experimenta con un virus SARS insertandole proteína...
2020-04-15 02:06:43,kak7742001,True,False,RT @BastienParisot: ● THREAD : Une proche m'envoie ce soir ces photos. Elle travaille dans un hôpital public de la région parisienne. Des b...
2020-04-15 02:06:43,FloreSgo,True,True,RT @AbdelDeuxFois: Sans promo ni rien. Le monde chico ☺,ALERTE INFO - 36,7 millions de téléspectateurs ont suivi l'allocution d'Emmanuel #Macron h
2020-04-15 02:06:43,QAlwaysWins,True,False,RT @21savage: Bang outside I hang outside
don't come out da house cuz corona outside,
2020-04-15 02:06:43,poppyp_yifan,True,True,RT @MKee3z: สงสัย ประเทศไทยต้องใจเย็นๆและ บริษัทก็ใจเย็นๆและ แบบอย่างเดียวกันนี่ใช่ไหม แต่ก็ต้องมาระบุให้ความตื่นเต้นมากขึ้นเรื่อยๆนะ... Thank you our friends from #Thailand. Thailand has long been a popu
#nnnevyy
2020-04-15 02:06:43,ikanbadutz,False,False,RT @21savage: Bang outside I hang outside

Dan semua kembali normal, yes @@#
2020-04-15 02:06:43,babyiyis_,True,False,RT @21savage: Bang outside I hang outside
```

Dataset:

Collected real-time tweets using twitter streaming api- tweepy. Extracted features - tweet_date, tweet_user, tweet_text, is_retweet, is_quote, quoted_text.

tweet_date – timestamp

tweet_user – text

tweet_text – text

is_retweet – boolean

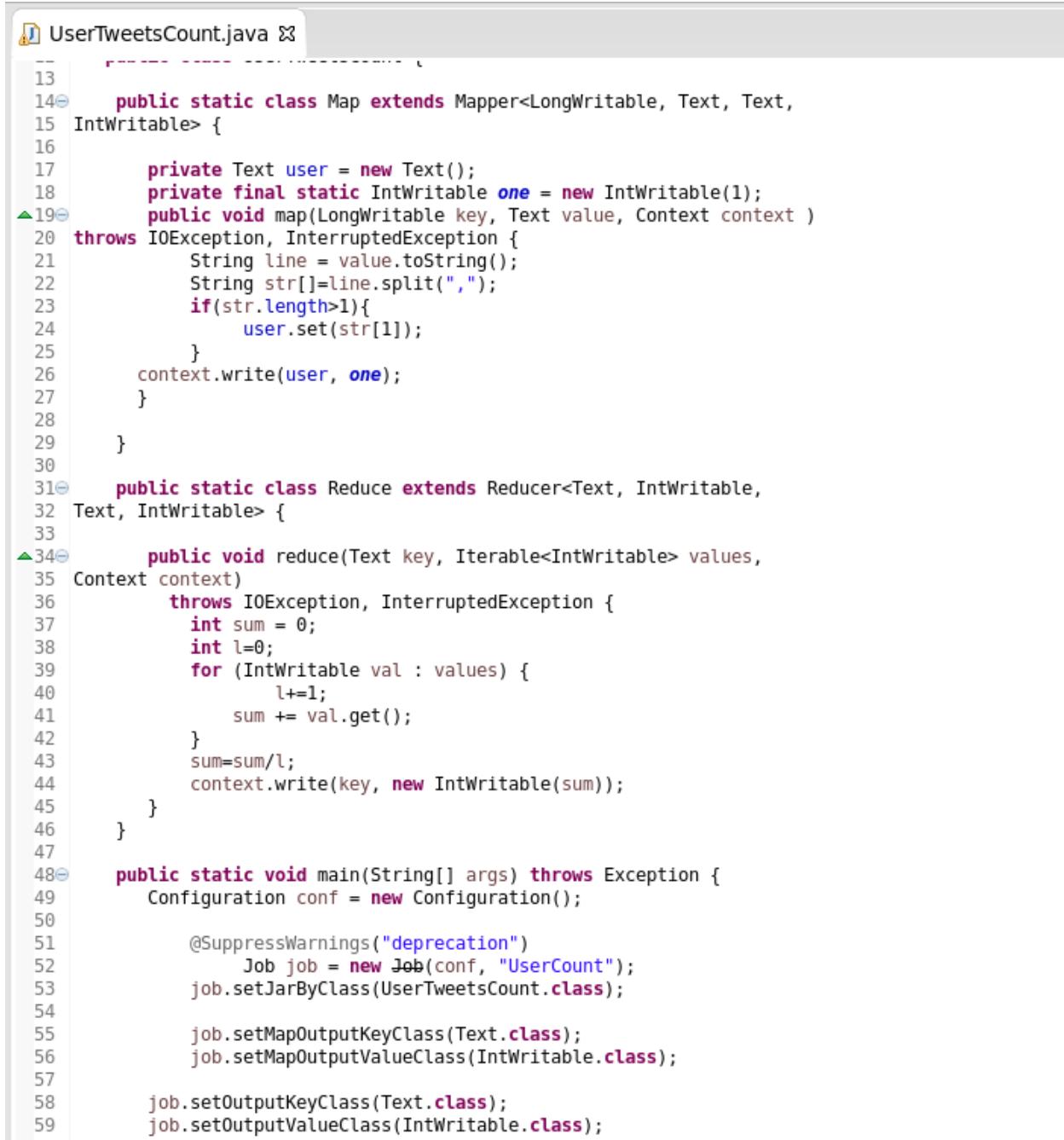
is_quote – boolean

quoted_text – text

```
cqlsh> desc bigdata.twitter_data;
CREATE TABLE bigdata.twitter_data (
    date timestamp,
    user text,
    is_quote boolean,
    is_retweet boolean,
    quoted_text text,
    text text,
    PRIMARY KEY (date, user)
) WITH CLUSTERING ORDER BY (user ASC)
    AND bloom_filter_fp_chance = 0.01
    AND caching = {'keys': 'ALL', 'rows_per_partition': 'NONE'}
    AND comment = ''
    AND compaction = {'class': 'org.apache.cassandra.db.compaction.SizeTieredCompactionStrategy', 'max_threshold': '32',
'min_threshold': '4'}
    AND compression = {'chunk_length_in_kb': '64', 'class': 'org.apache.cassandra.io.compress.LZ4Compressor'}
    AND crc_check_chance = 1.0
    AND dclocal_read_repair_chance = 0.1
    AND default_time_to_live = 0
    AND gc_grace_seconds = 864000
    AND max_index_interval = 2048
    AND memtable_flush_period_in_ms = 0
    AND min_index_interval = 128
    AND read_repair_chance = 0.0
    AND speculative_retry = '99PERCENTILE';
CREATE CUSTOM INDEX user_idx ON bigdata.twitter_data (user) USING 'org.apache.cassandra.index.sasi.SASIIndex' WITH OPTIONS = {'analyzer_class': 'org.apache.cassandra.index.sasi.analyzer.StandardAnalyzer', 'case_sensitive': 'false'};
```


Analysis of data/Implementation/Results:

Use Case 1: Map Reduce to count tweets by each user



The screenshot shows a Java code editor with the file `UserTweetsCount.java` open. The code implements a MapReduce job to count tweets by user. It includes a Mapper class that processes each tweet line and sets the user as the key and a value of 1. The Reducer class then sums up the values for each user and outputs the total count.

```
13
14  public static class Map extends Mapper<LongWritable, Text, Text,
15  IntWritable> {
16
17      private Text user = new Text();
18      private final static IntWritable one = new IntWritable(1);
19  public void map(LongWritable key, Text value, Context context )
20 throws IOException, InterruptedException {
21     String line = value.toString();
22     String str[]=line.split(",");
23     if(str.length>1){
24         user.set(str[1]);
25     }
26     context.write(user, one);
27 }
28
29 }
30
31  public static class Reduce extends Reducer<Text, IntWritable,
32  Text, IntWritable> {
33
34  public void reduce(Text key, Iterable<IntWritable> values,
35 Context context)
36 throws IOException, InterruptedException {
37     int sum = 0;
38     int l=0;
39     for (IntWritable val : values) {
40         l+=1;
41         sum += val.get();
42     }
43     sum=sum/l;
44     context.write(key, new IntWritable(sum));
45   }
46 }
47
48  public static void main(String[] args) throws Exception {
49     Configuration conf = new Configuration();
50
51     @SuppressWarnings("deprecation")
52     Job job = new Job(conf, "UserCount");
53     job.setJarByClass(UserTweetsCount.class);
54
55     job.setMapOutputKeyClass(Text.class);
56     job.setMapOutputValueClass(IntWritable.class);
57
58     job.setOutputKeyClass(Text.class);
59     job.setOutputValueClass(IntWritable.class);
```

Starting HDFS namenode and datanodes:

```
lohit@ubuntu:~/hadoop/hadoop-2.8.1/sbin$ start-dfs.sh
20/04/14 20:23:28 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Starting namenodes on [localhost]
localhost: starting namenode, logging to /home/lohit/hadoop/hadoop-2.8.1/logs/hadoop-lohit-namenode-ubuntu.out
localhost: starting datanode, logging to /home/lohit/hadoop/hadoop-2.8.1/logs/hadoop-lohit-datanode-ubuntu.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /home/lohit/hadoop/hadoop-2.8.1/logs/hadoop-lohit-secondarynamenode-ubuntu.out
20/04/14 20:23:55 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
lohit@ubuntu:~/hadoop/hadoop-2.8.1/sbin$ start-yarn.sh
starting yarn daemons
starting resourcemanager, logging to /home/lohit/hadoop/hadoop-2.8.1/logs/yarn-lohit-resourcemanager-ubuntu.out
localhost: starting nodemanager, logging to /home/lohit/hadoop/hadoop-2.8.1/logs/yarn-lohit-nodemanager-ubuntu.out
```

Loading input data from local to HDFS:

```
lohit@ubuntu:~/Desktop/BigDataProgramming$ hdfs dfs -put COVID19_Data.json /bigdata/
20/04/14 20:40:03 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
```

Running MapReduce:

```
lohit@ubuntu:~/Downloads$ hadoop jar UserCount.jar UserTweetsCount /bigdata/COVID19_Data.json /bigdata/tweets/
20/04/14 20:47:56 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
20/04/14 20:47:58 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8082
20/04/14 20:47:59 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
20/04/14 20:48:00 INFO input.FileInputFormat: Total input files to process : 1
20/04/14 20:48:00 WARN hdfs.DataStreamer: Caught exception
java.lang.InterruptedException
    at java.lang.Object.wait(Native Method)
    at java.lang.Thread.join(Thread.java:1252)
    at java.lang.Thread.join(Thread.java:1326)
    at org.apache.hadoop.hdfs.DataStreamer.closeResponder(DataStreamer.java:927)
    at org.apache.hadoop.hdfs.DataStreamer.endBlock(DataStreamer.java:578)
    at org.apache.hadoop.hdfs.DataStreamer.onDataStreamer(DataStreamer.java:755)
20/04/14 20:49:00 WARN hdfs.DataStreamer: Caught exception
java.lang.InterruptedException
    at java.lang.Object.wait(Native Method)
    at java.lang.Thread.join(Thread.java:1252)
    at java.lang.Thread.join(Thread.java:1326)
    at org.apache.hadoop.hdfs.DataStreamer.closeResponder(DataStreamer.java:927)
    at org.apache.hadoop.hdfs.DataStreamer.endBlock(DataStreamer.java:578)
    at org.apache.hadoop.hdfs.DataStreamer.run(DataStreamer.java:755)
20/04/14 20:49:00 INFO mapreduce.JobSubmitter: number of splits:1
20/04/14 20:49:00 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1586921057152_0003
20/04/14 20:49:01 INFO mapreduce.Job: Job tracking url: http://ubuntu:8088/proxy/application_1586921057152_0003
20/04/14 20:49:02 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1586921057152_0003
20/04/14 20:49:02 INFO mapreduce.Job: Running job: job_1586921057152_0003
20/04/14 20:49:17 INFO mapreduce.Job: Job job_1586921057152_0003 running in uber mode : false
20/04/14 20:49:17 INFO mapreduce.Job: map 0% reduce 0%
20/04/14 20:49:32 INFO mapreduce.Job: map 100% reduce 0%
20/04/14 20:49:47 INFO mapreduce.Job: map 100% reduce 100%
20/04/14 20:49:50 INFO mapreduce.Job: Job job_1586921057152_0003 completed successfully
20/04/14 20:49:50 INFO mapreduce.Job: Counters:
File System Counters
    FILE: Number of bytes read=1172059
    FILE: Number of bytes written=2617277
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=5741057
    HDFS: Number of bytes written=410786
    HDFS: Number of read operations=6
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
```

```

HDFS: Number of read operations=6
HDFS: Number of large read operations=0
HDFS: Number of write operations=2
Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=11966
    Total time spent by all reduces in occupied slots (ms)=12023
    Total time spent by all map tasks (ms)=11966
    Total time spent by all reduce tasks (ms)=12023
    Total vcore-milliseconds taken by all map tasks=11966
    Total vcore-milliseconds taken by all reduce tasks=12023
    Total megabyte-milliseconds taken by all map tasks=12253184
    Total megabyte-milliseconds taken by all reduce tasks=12311552
Map-Reduce Framework
    Map input records=52473
    Map output records=52473
    Map output bytes=1066546
    Map output materialized bytes=1172059
    Input split bytes=112
    Combine input records=0
    Combine output records=0
    Reduce input groups=24116
    Reduce shuffle bytes=1172059
    Reduce input records=52473
    Reduce output records=24116
    Spilled Records=104946
    Shuffled Maps =1
    Failed Shuffles=0
    Merged Map outputs=1
    GC time elapsed (ms)=216
    CPU time spent (ms)=3330
    Physical Memory (bytes) snapshot=302473216
    Virtual memory (bytes) snapshot=3893735424
    Total committed heap usage (bytes)=170004480
Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAP=0
    WRONG_REDUCE=0
File Input Format Counters
    Bytes Read=5740945
File Output Format Counters
    Bytes Written=410786

```

HDFS File Structure:

The screenshot shows the HDFS Web UI at the URL `localhost:50070/explorer.html#/bigdata`. The page title is "Browse Directory". The main content area displays a table of files in the `/bigdata/` directory. The table has the following columns: Permission, Owner, Group, Size, Last Modified, Replication, Block Size, and Name. There are two entries in the table:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<code>drwxr-xr-x</code>	Iohitha	supergroup	0 B	Apr 14 20:48	0	0 B	<code>tweets</code>
<code>-rw-r--r--</code>	Iohitha	supergroup	5.47 MB	Apr 14 20:40	1	128 MB	<code>COVID19_Data.json</code>

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	lohittha	supergroup	0 B	Apr 14 20:48	1	128 MB	_SUCCESS
-rw-r--r--	lohittha	supergroup	401.16 KB	Apr 14 20:48	1	128 MB	part-r-00000

Showing 1 to 2 of 2 entries

Hadoop, 2017.

Output File:

```
lohittha@ubuntu:~/Downloads$ hdfs dfs -cat /bigdata/tweets/part-r-00000 | head
20/04/14 22:29:53 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
      1
      1
111 se han recuper... 1
2 NAVOLATO 1
Gue kirimi buat Makan selama pedemi Virus Corona      1
aseguran que hay varias muertes clasificadas como "neumonia atipica" con todas las caracteristicas de ser #coronavirus. https://t.co/vpe0pZYfLM      1
very inspirational. This drug is a game changer.      1
where we are measured by our humanity not valued in silver or other coin.      1
देश के मज़रूर लोगों का धूम्रपान 1
"not responsible" and put all the responsibility on the States. The States then accepted the challenge and sourced their own supplies. NOW THIS???      1
cat: Unable to write to output stream.
lohittha@ubuntu:~/Downloads$
```

Use Case 2: Data analysis of tweets using Hive

Create Tweets Table:

```
hive> create table tweets (tweet_date STRING,user STRING,is_retweet STRING,is_quote STRING,text STRING,qouted_text STRING)
> [cloudera@quickstart ~]$ hive

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties
WARNING: Hive CLI is deprecated and migration to Beeline is recommended.
hive> create table tweets (tweet_date STRING,user STRING,is_retweet STRING,is_quote STRING,text STRING,qouted_text STRING)
> row format delimited fields terminated by "," stored as textfile;
OK
Time taken: 1.493 seconds
```

Load Twitter Data into Tweets table:

```
hive> load data local inpath "/home/cloudera/Downloads/COVID19_Data.csv" into table tweets;
Loading data to table default.tweets
Table default.tweets stats: [numFiles=1, totalSize=5740945]
OK
Time taken: 1.58 seconds
```

Query 1: Fetch top users with more number of tweets

```
hive> select user,count(1) as tweet_count from tweets where text is not null group by user order by tweet_count desc limit 10;
Query ID = cloudera_20200414212424_9bed624d-dab4-488a-891d-0dcc5d58ce86
Total jobs = 2
Launching Job 1 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1586916751037_0004, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1586916751037_0004/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1586916751037_0004
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2020-04-14 21:24:29,303 Stage-1 map = 0%, reduce = 0%
2020-04-14 21:24:49,156 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 5.51 sec
2020-04-14 21:25:10,351 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 11.04 sec
MapReduce Total cumulative CPU time: 11 seconds 40 msec
Ended Job = job_1586916751037_0004
Launching Job 2 out of 2
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1586916751037_0005, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1586916751037_0005/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1586916751037_0005
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2020-04-14 21:25:30,720 Stage-2 map = 0%, reduce = 0%
2020-04-14 21:25:49,266 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 5.59 sec
2020-04-14 21:26:06,936 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 9.05 sec
MapReduce Total cumulative CPU time: 9 seconds 50 msec
Ended Job = job_1586916751037_0005
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 11.04 sec HDFS Read: 5748913 HDFS Write: 690474 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 9.05 sec HDFS Read: 695482 HDFS Write: 201 SUCCESS
Total MapReduce CPU Time Spent: 20 seconds 90 msec
OK
openletterbot 55
Dy CM 26
HollyDeinert 21
syedsalu11 20
paragpal0920537 17
Em coletiva de imprensa realizada no @Planalto 16
JPNicholasDabot 15
Time taken: 118.462 seconds, Fetched: 10 row(s)
hive> █
```

Query 2: Retweet count

```
hive> select count(1) retweet_count from tweets where trim(lower(is_retweet)) = "true";
Query ID = cloudera_20200414213232_bfebce0c-70ca-48e4-89ed-2f15680b54ef
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1586916751037_0008, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1586916751037_0008/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1586916751037_0008
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2020-04-14 21:32:37,284 Stage-1 map = 0%,  reduce = 0%
2020-04-14 21:32:57,614 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 4.85 sec
2020-04-14 21:33:13,541 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 7.44 sec
MapReduce Total cumulative CPU time: 7 seconds 440 msec
Ended Job = job_1586916751037_0008
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1  Cumulative CPU: 7.44 sec  HDFS Read: 5749863 HDFS Write: 6 SUCCESS
Total MapReduce CPU Time Spent: 7 seconds 440 msec
OK
19731
Time taken: 54.324 seconds, Fetched: 1 row(s)
hive> █
```

Query 3: Tweets per minute

```
hive> ( select count(1) tweet_count from tweets where trim(lower(tweet_date)) like '%2020-04-15 02:07%';
NoViableAltException(290@[])
  at org.apache.hadoop.hive.ql.parse.HiveParser.statement(HiveParser.java:1028)
  at org.apache.hadoop.hive.ql.parse.ParseDriver.parse(ParseDriver.java:201)
  at org.apache.hadoop.hive.ql.parse.ParseDriver.parse(ParseDriver.java:166)
  at org.apache.hadoop.hive.ql.Driver.compile(Driver.java:522)
  at org.apache.hadoop.hive.ql.Driver.compileInternal(Driver.java:1356)
  at org.apache.hadoop.hive.ql.Driver.runInternal(Driver.java:1473)
  at org.apache.hadoop.hive.ql.Driver.run(Driver.java:1285)
  at org.apache.hadoop.hive.ql.Driver.run(Driver.java:1275)
  at org.apache.hadoop.hive.cli.CliDriver.processLocalCmd(CliDriver.java:226)
  at org.apache.hadoop.hive.cli.CliDriver.processCmd(CliDriver.java:175)
  at org.apache.hadoop.hive.cli.CliDriver.processLine(CliDriver.java:389)
  at org.apache.hadoop.hive.cli.CliDriver.executeDriver(CliDriver.java:781)
  at org.apache.hadoop.hive.cli.CliDriver.run(CliDriver.java:699)
  at org.apache.hadoop.hive.cli.CliDriver.main(CliDriver.java:634)
  at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
  at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:57)
  at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
  at java.lang.reflect.Method.invoke(Method.java:606)
  at org.apache.hadoop.util.RunJar.run(RunJar.java:221)
  at org.apache.hadoop.util.RunJar.main(RunJar.java:136)
FAILED: ParseException line 1:0 cannot recognize input near '(' 'select' 'count'
hive> select count(1) tweet_count from tweets where trim(lower(tweet_date)) like '%2020-04-15 02:07%';
Query ID = cloudera_20200414213939_1c2210e7-c840-431f-a32e-b2eb35171d20
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1586916751037_0009, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1586916751037_0009/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1586916751037_0009
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2020-04-14 21:39:45,329 Stage-1 map = 0%,  reduce = 0%
2020-04-14 21:40:04,322 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 4.46 sec
2020-04-14 21:40:21,539 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 7.14 sec
MapReduce Total cumulative CPU time: 7 seconds 140 msec
Ended Job = job_1586916751037_0009
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1  Cumulative CPU: 7.14 sec  HDFS Read: 5749880 HDFS Write: 5 SUCCESS
Total MapReduce CPU Time Spent: 7 seconds 140 msec
OK
3017
Time taken: 57.374 seconds, Fetched: 1 row(s)
hive> █
```

Query 4: Tweets on Layoffs

```
hive> select count(1) tweet_count from tweets where trim(lower(text)) like '%layoff%';
Query ID = cloudera_20200414214141_9c3cbfd4-11d9-416d-b135-e51dbd2875f3
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1586916751037_0010, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1586916751037_0010/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1586916751037_0010
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2020-04-14 21:41:58,116 Stage-1 map = 0%, reduce = 0%
2020-04-14 21:42:15,545 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 4.7 sec
2020-04-14 21:42:36,210 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 7.93 sec
MapReduce Total cumulative CPU time: 7 seconds 930 msec
Ended Job = job_1586916751037_0010
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 7.93 sec HDFS Read: 5749863 HDFS Write: 2 SUCCESS
Total MapReduce CPU Time Spent: 7 seconds 930 msec
OK
2
Time taken: 59.565 seconds, Fetched: 1 row(s)
hive> █
```

Query 5: Tweets on USA

```
hive> select count(1) tweet_count from tweets where trim(lower(text)) like '%usa%' or trim(lower(text)) like '%united states%';
Query ID = cloudera_20200414214545_ba40e4d2-1995-4b17-8f46-d647f76cdf95
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1586916751037_0011, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1586916751037_0011/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1586916751037_0011
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2020-04-14 21:45:46,527 Stage-1 map = 0%, reduce = 0%
2020-04-14 21:46:02,673 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 4.32 sec
2020-04-14 21:46:19,715 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 7.03 sec
MapReduce Total cumulative CPU time: 7 seconds 30 msec
Ended Job = job_1586916751037_0011
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 7.03 sec HDFS Read: 5750249 HDFS Write: 4 SUCCESS
Total MapReduce CPU Time Spent: 7 seconds 30 msec
OK
417
Time taken: 52.837 seconds, Fetched: 1 row(s)
hive> █
```

Use Case 3: Twitter data sentimental analysis using Map Reduce

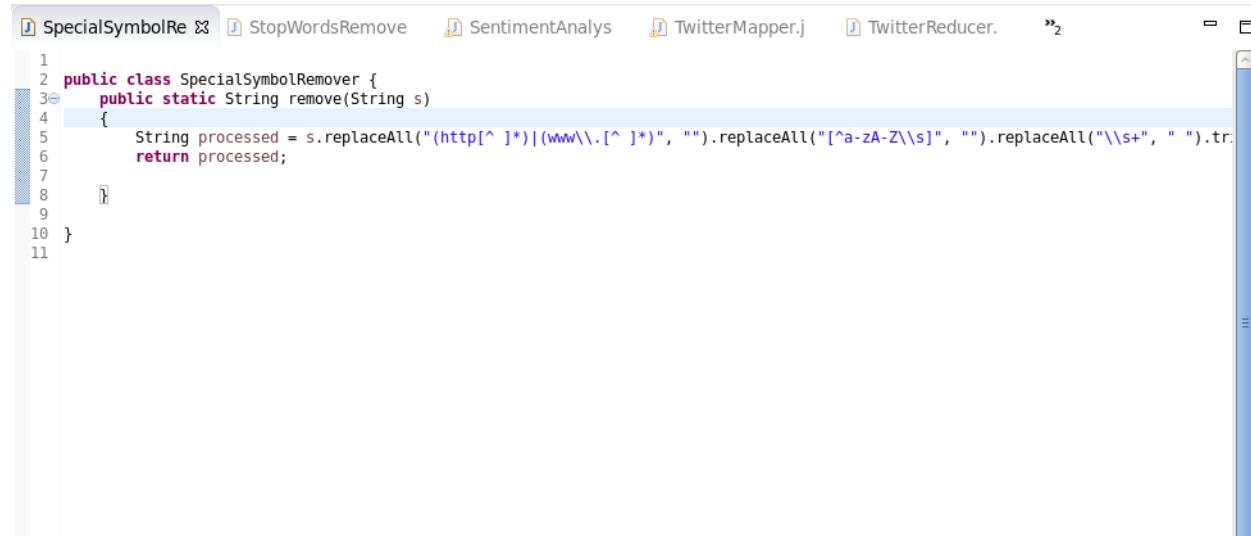
Analyzed the sensitivity of tweets. Divided tweets into 3 categories- positive, negative and neutral and displayed count for each category.

Data: Loaded input file to HDFS and displayed

```
[cloudera@quickstart ~]$ hdfs dfs -cat sentimentanalysis/input/COVID19Data.txt |head
date,user,is_retweet,is_quote,text,quoted_text
2020-04-30 03:17:42,villafranca_0,0,True,False,RT @alegrant5: Pregunta tal vez tonta: ning n politico se ha contagiado con #covid19???\[1\]
2020-04-30 03:17:42,HappyHobbit4,True,False,RT @NanaJudie8484: @FredfromFlorida @GovRonDeSantis Now this is the best Governor ever!
please all Take note...
this is how our GOP
@Gov...
2020-04-30 03:17:42,antinashit,True,False,RT @anavas4: Congresistas estadounidenses instan a Bukele a no usar COVID-19 como pretexto para socavar la Constituci n https://t.co/bMKltj...,
2020-04-30 03:17:42,AKPoliticalBeat,False,False,Top story: 'A Terrible Price': The Deadly Racial Disparities of Covid-19 in America https://t.co/dFpzy0lRqc, see more https://t.co/mUz6VxxLB7,
2020-04-30 03:17:42,gromero2019,True,False,RT @rosaicela_: El uso de cubrebocas puede disminuir las probabilidades de contagio por #COVID19, si necesitas salir no olvides colocarlo d...
cat: Unable to write to output stream.
[cloudera@quickstart ~]$
```

Code:

Created class- SpecialSymbolRemover for data preprocessing to remove urls in the tweets as they do not contribute to sensitivity of tweets.



```
SpecialSymbolRe StopWordsRemove SentimentAnalysis TwitterMapper.j TwitterReducer. »2
1
2 public class SpecialSymbolRemover {
3     public static String remove(String s)
4     {
5         String processed = s.replaceAll("(http[^ ]*)|(www\\.[^ ]*)", "").replaceAll("[^a-zA-Z\\s]", "").replaceAll("\\s+", " ");
6         return processed;
7     }
8
9
10 }
```

Created class- StopWordsRemover for data preprocessing to remove stop words from data. Stop words like- a, to, pronouns are removed as they do not define tweet sensitivity.

```
SpecialSymbolRe StopWordsRemove SentimentAnalysis TwitterMapper.j TwitterReducer. »2
3 public class StopWordsRemover {
4
5 public static String[] stopWords = {"a", "as", "able", "about", "above", "according", "accordingly",
6 "across", "actually", "after", "afterwards", "again", "against", "aint", "all", "allow", "allows", "almost",
7 "alone", "along", "already", "also", "although", "always", "am", "among", "amongst", "an", "and", "another",
8 "any", "anybody", "anyhow", "anyone", "anything", "anyway", "anyways", "anywhere", "apart", "appear",
9 "appreciate", "appropriate", "are", "arent", "around", "as", "aside", "ask", "asking", "associated", "at",
10 "available", "away", "awfully", "be", "became", "because", "become", "becomes", "becoming", "been",
11 "before", "beforehand", "behind", "being", "believe", "below", "beside", "besides", "best", "better",
12 "between", "beyond", "both", "brief", "but", "by", "cmon", "cs", "came", "can", "cant", "cannot", "cant",
13 "cause", "causes", "certain", "certainly", "changes", "clearly", "co", "com", "come", "comes", "concerning",
14 "consequently", "consider", "considering", "contain", "containing", "contains", "corresponding", "could",
15 "couldnt", "course", "currently", "definitely", "described", "despite", "did", "didnt", "different", "do",
16 "does", "doesnt", "doing", "dont", "done", "down", "downwards", "during", "each", "edu", "eg", "eight",
17 "either", "else", "elsewhere", "enough", "entirely", "especially", "et", "etc", "even", "ever", "every",
18 "everybody", "everyone", "everything", "everywhere", "ex", "exactly", "example", "except", "far", "few",
19 "ff", "fifth", "first", "five", "followed", "following", "follows", "for", "former", "formerly", "forth",
20 "four", "from", "further", "furthermore", "get", "gets", "getting", "given", "gives", "go", "goes", "going",
21 "gone", "got", "gotten", "greetings", "had", "hadnt", "happens", "hardly", "has", "hasnt", "have", "havent",
22 "having", "he", "hes", "hello", "help", "hence", "her", "here", "heres", "hereafter", "hereby", "herein",
23 "hereupon", "hers", "herself", "hi", "him", "himself", "his", "hither", "hopefully", "how", "howbeit",
24 "however", "i", "id", "ill", "im", "ive", "ie", "if", "ignored", "immediate", "in", "inasmuch", "inc",
25 "indeed", "indicate", "indicated", "indicates", "inner", "insofar", "instead", "into", "inward", "is",
26 "isnt", "it", "itd", "itll", "its", "its", "itself", "just", "keep", "keeps", "kept", "know", "knows",
27 "known", "last", "lately", "later", "latter", "latterly", "least", "less", "lest", "let", "lets", "like",
28 "liked", "likely", "little", "look", "looking", "looks", "ltd", "mainly", "many", "may", "maybe", "me",
29 "mean", "meanwhile", "merely", "might", "more", "moreover", "most", "mostly", "much", "must", "my",
30 "myself", "name", "namely", "nd", "near", "nearly", "necessary", "need", "needs", "neither", "never",
```

```
31     "nevertheless", "new", "next", "nine", "no", "nobody", "non", "none", "noone", "nor", "normally", "not",
32     "nothing", "novel", "now", "nowhere", "obviously", "of", "off", "often", "oh", "ok", "okay", "old", "on",
33     "once", "one", "ones", "only", "onto", "or", "other", "others", "otherwise", "ought", "our", "ours",
34     "ourselves", "out", "outside", "over", "overall", "own", "particular", "particularly", "per", "perhaps",
35     "placed", "please", "plus", "possible", "presumably", "probably", "provides", "que", "quite", "qv",
36     "rather", "rd", "re", "really", "reasonably", "regarding", "regardless", "regards", "relatively",
37     "respectively", "right", "said", "same", "saw", "say", "saying", "says", "second", "secondly", "see",
38     "seeing", "seem", "seemed", "seeming", "seems", "seen", "self", "selves", "sensible", "sent", "serious",
39     "seriously", "seven", "several", "shall", "she", "should", "shouldnt", "since", "six", "so", "some",
40     "somebody", "somehow", "someone", "something", "sometime", "sometimes", "somewhat", "somewhere", "soon",
41     "sorry", "specified", "specify", "specifying", "still", "sub", "such", "sup", "sure", "ts", "take", "taken",
42     "tell", "tends", "th", "than", "thank", "thanks", "thanx", "that", "thats", "thats", "the", "their",
43     "theirs", "them", "themselves", "then", "thence", "there", "theres", "thereafter", "thereby", "therefore",
44     "therein", "theres", "thereupon", "these", "they", "theyd", "theyll", "theyre", "theyve", "think", "third",
45     "this", "thorough", "thoroughly", "those", "though", "three", "through", "throughout", "thru", "thus", "to",
46     "together", "too", "took", "toward", "towards", "tried", "tries", "truly", "try", "trying", "twice", "two",
47     "un", "under", "unfortunately", "unless", "unlikely", "until", "unto", "up", "upon", "us", "use", "used",
48     "useful", "uses", "using", "usually", "value", "various", "very", "via", "viz", "vs", "want", "wants",
49     "was", "wasnt", "way", "we", "wed", "well", "were", "weve", "welcome", "well", "went", "were", "werent",
50     "what", "whats", "whatever", "when", "whence", "whenever", "where", "wheres", "whereafter", "whereas",
51     "whereby", "wherein", "whereupon", "wherever", "whether", "which", "while", "whither", "who", "whos",
52     "whoever", "whole", "whom", "whose", "why", "will", "willing", "wish", "with", "within", "without", "wont",
53     "wonder", "would", "would", "wouldnt", "yes", "yet", "you", "youd", "youll", "youre", "youve", "your",
54     "yours", "yourself", "yourselves", "zero" };
```

```
SpecialSymbolRe StopWordsRemove SentimentAnalys TwitterMapper.j TwitterReducer. »2
50     was , would , way , we , wen , well , were , were , welcome , well , went , were , went ,
51     "what" , "whats" , "whatever" , "when" , "whence" , "whenever" , "where" , "wheres" , "whereafter" , "whereas" ,
52     "whereby" , "wherein" , "whereupon" , "wherever" , "whether" , "which" , "while" , "whither" , "who" , "whos" ,
53     "whoever" , "whole" , "whom" , "whose" , "why" , "will" , "willing" , "wish" , "with" , "within" , "without" , "wont" ,
54     "wonder" , "would" , "wouldn't" , "yes" , "yet" , "you" , "youd" , "youll" , "youre" , "youve" , "your" ,
55     "yours" , "yourself" , "yourselves" , "zero" };
56
57
58
59
60
61 static String remove(String s)
62 {
63     String ignoreCase = "(?i)";
64
65
66     String processed= s.replaceAll("(\\b"+ignoreCase+StringUtil.join(stopWords,"\\b|\\b")+"\\b)","");
67     .replaceAll("\\s+", " ");
68
69     return processed;
70
71 }
72
73
74
75
76
77 }
```

Main class- SentimentAnalysis

Created map, reduce jobs and added afinn-111 file to cache. Afinn-111 consists of sentiment values for various words describing emotions. Words describing emotions like anger, hatred have negative sentiment values. Words describing emotions like happiness, excitement have positive sentiment values. Words that describe neutral emotions hold sentiment value as ZERO.

```

34     job.setOutputFormatClass(TextOutputFormat.class);
35
36
37
38     try{
39
40         //adding file to distributed cache
41         job.addCacheFile(new URI("sentimentanalysis/input/AFINN-111.txt"));
42     }catch(Exception e)
43     {
44         System.out.println("file not added my dear");
45         System.exit(1);
46     }
47
48
49     FileInputFormat.addInputPath(job, new Path(args[0]));
50     FileOutputFormat.setOutputPath(job, new Path(args[1]));
51
52     System.exit(job.waitForCompletion(true)?0:1);
53
54
55 }
56
57
58
59
60 }
61

```

Mapper class- TwitterMapper

In the mapper class, loaded file from the cache and created map out of it. In the map function, for each tweet sentiment is calculated by summing each word sentiment. The output of map function is (sentiment_category, 1). Example: (positive, 1)

```

12 public class TwitterMapper extends Mapper<LongWritable,Text,Text,Text> {
13
14     Map<String,String> dictionary = null;
15
16
17     @Override
18     protected void setup(Context context) throws IOException,InterruptedException
19     {
20         dictionary = new HashMap<String,String>();
21
22
23         URI[] cacheFiles = context.getCacheFiles();
24
25         if (cacheFiles != null && cacheFiles.length > 0)
26         {
27             try
28             {
29                 String line ="";
30                 FileSystem fs = FileSystem.get(context.getConfiguration());
31                 Path path = new Path(cacheFiles[0].toString());
32                 BufferedReader reader = new BufferedReader(new InputStreamReader(fs.open(path)));
33
34                 while((line = reader.readLine())!=null)
35                 {
36                     String []tokens = line.split("\t");
37                     dictionary.put(tokens[0], tokens[1]);
38                 }
39             }

```

```
SpecialSymbolRe StopWordsRemove SentimentAnalys *TwitterMapper. TwitterReducer. »2
56 @Override
57 public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException
58 {
59     long sentiment_value = 0;
60     try{
61         String quoted_text=null,tweet_text=null;
62         String[] tweet = value.toString().split("-");
63         if(tweet.length==6)
64             quoted_text = tweet[5];
65         if(quoted_text == null && tweet.length==5)
66             tweet_text = tweet[4];
67         String processedTweet= null;
68
69
70
71         if(quoted_text!=null)
72         {
73
74             processedTweet = SpecialSymbolRemover.remove(quoted_text);
75
76             processedTweet = StopWordsRemover.remove(processedTweet);
77
78             String []words = processedTweet.split(" ");
79
80             for(String temp:words)
81             {
82                 if(dictionary.containsKey(temp))
83             }
84
85
86
87
88
89
90
91
92
93 }
```

The screenshot shows a Java IDE interface with multiple tabs at the top. The active tab is 'TwitterMapper.java'. The code in the editor is as follows:

```
1 SpecialSymbolRe 2 StopWordsRemove 3 SentimentAnalys 4 *TwitterMapper. 5 TwitterReducer. 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108
```

The code implements a Mapper class for Hadoop. It processes a tweet by removing special symbols and stop words, then splits the text into words. For each word, it checks if it is in a dictionary and adds its value to a sentiment score.

```

102     if(dictionary.containsKey(temp))
103     {
104         sentiment_value+=Long.parseLong(dictionary.get(temp));
105     }
106 }
107
108 }
109
110 if(quoted_text!=null || tweet_text!=null)
111 {
112     if(sentiment_value>0)
113         context.write(new Text("Positive"),new Text("1"));
114     else if(sentiment_value<0)
115         context.write(new Text("Negative"),new Text("1"));
116     else
117         context.write(new Text("Neutral"),new Text("1"));
118 }
119
120
121
122 }catch(Exception e)
123 {
124     e.printStackTrace();
125 }
126
127 }
128
129 }
130 }
```

Reducer class- TwitterReducer

In the reducer class, each tweet category is summed for the final tweet count for each category.

```

1 import java.util.ArrayList;
2 import java.util.HashMap;
3 import java.util.List;
4
5 import org.apache.hadoop.io.LongWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Reducer;
8 import org.apache.hadoop.mapreduce.Reducer.Context;
9
10
11 public class TwitterReducer extends Reducer<Text, Text, Text, Text>
12 {
13
14     private Text res = new Text();
15
16     public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException
17     {
18         int counter = 0;
19
20         for (Text val : values)
21         {
22             counter = counter + 1;
23         }
24
25         context.write(key, new Text(String.valueOf(counter))); //Generating the reduced output.
26     }
27 }
28
29
30 }
```

AFINN-111.txt

```
2447 wacky -3
2450 worse -3
2451 worsen -3
2452 worsened -3
2453 worsening -3
2454 worsens -3
2455 worshiped 3
2456 worst -3
2457 worth 2
2458 worthless -2
2459 worthy 2
2460 wow 4
2461 wowow 4
2462 wwwww 4
2463 wrathful -3
2464 wreck -2
2465 wrong -2
2466 wronged -2
2467 wtf -4
2468 yeah 1
2469 yearning 1
2470 yeeees 2
2471 yes 1
2472 youthful 2
2473 yucky -2
2474 yummy 3
2475 zealot -2
2476 zealots -2
2477 zealous 2
```

Output:

```
cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System Thu Apr 30, 11:19 AM cloudera
cloudera@quickstart:~ File Edit View Search Terminal Help
[cloudera@quickstart ~]$ hadoop jar /home/cloudera/Downloads/SentimentAnalysis.jar SentimentAnalysis sentimentanalysis/input/COVID19_Data.txt sentimentanalysis/
output1
20/04/30 11:15:04 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/04/30 11:15:05 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
20/04/30 11:15:05 INFO input.FileInputFormat: Total input paths to process : 1
20/04/30 11:15:06 INFO mapreduce.JobSubmitter: number of splits:1
20/04/30 11:15:06 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1588268869380_0002
20/04/30 11:15:07 INFO impl.YarnClientImpl: Submitted application application_1588268869380_0002
20/04/30 11:15:07 INFO mapreduce.Job: The url to track the job: http://quickstar.cloudera:8088/proxy/application_1588268869380_0002/
20/04/30 11:15:07 INFO mapreduce.Job: Running job: job_1588268869380_0002
20/04/30 11:15:21 INFO mapreduce.Job: Job job_1588268869380_0002 running in uber mode : false
20/04/30 11:15:21 INFO mapreduce.Job: map 0% reduce 0%
20/04/30 11:15:34 INFO mapreduce.Job: map 100% reduce 0%
20/04/30 11:15:47 INFO mapreduce.Job: map 100% reduce 100%
20/04/30 11:15:48 INFO mapreduce.Job: Job job_1588268869380_0002 completed successfully
20/04/30 11:15:48 INFO mapreduce.Job: Counters: 49
File System Counters
FILE: Number of bytes read=11768
FILE: Number of bytes written=313197
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=5771665
```

cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System Thu Apr 30, 11:20 AM cloudera
File Edit View Search Terminal Help

```
HDFS: Number of bytes written=36
HDFS: Number of read operations=7
HDFS: Number of large read operations=0
HDFS: Number of write operations=2

Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=11276
    Total time spent by all reduces in occupied slots (ms)=11046
    Total time spent by all map tasks (ms)=11276
    Total time spent by all reduce tasks (ms)=11046
    Total vcore-milliseconds taken by all map tasks=11276
    Total vcore-milliseconds taken by all reduce tasks=11046
    Total megabyte-milliseconds taken by all map tasks=11546624
    Total megabyte-milliseconds taken by all reduce tasks=11311104

Map-Reduce Framework
    Map input records=52473
    Map output records=965
    Map output bytes=9832
    Map output materialized bytes=11768
    Input split bytes=151
    Combine input records=0
    Combine output records=0
    Reduce input groups=3
    Reduce shuffle bytes=11768
    Reduce input records=965
    Reduce output records=3
    Spilled Records=1930
    Shuffled Maps =1
    Failed Shuffles=0
    Merged Map outputs=1

Map output materialized bytes=11768
Input split bytes=151
Combine input records=0
Combine output records=0
Reduce input groups=3
Reduce shuffle bytes=11768
Reduce input records=965
Reduce output records=3
Spilled Records=1930
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=288
CPU time spent (ms)=4550
Physical memory (bytes) snapshot=352935936
Virtual memory (bytes) snapshot=3016089600
Total committed heap usage (bytes)=226365440

Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0

File Input Format Counters
Bytes Read=5740945
File Output Format Counters
Bytes Written=36
```

```

cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System cloudera@quickstart:~ Thu Apr 30, 1:28 PM
File Edit View Search Terminal Help
Spilled Records=85570
Shuffled Maps =2
Failed Shuffles=0
Merged Map outputs=2
GC time elapsed (ms)=1097
CPU time spent (ms)=35050
Physical memory (bytes) snapshot=593186816
Virtual memory (bytes) snapshot=4519288832
Total committed heap usage (bytes)=439164928
Shuffle Errors
BAD ID=0
CONNECTION=0
IO ERROR=0
WRONG LENGTH=0
WRONG MAP=0
WRONG REDUCE=0
File Input Format Counters
Bytes Read=172927823
File Output Format Counters
Bytes Written=42
[cloudera@quickstart ~]$ hdfs dfs -ls sentimentanalysis/final_output
Found 2 items
-rw-r--r-- 1 cloudera cloudera 0 2020-04-30 13:27 sentimentanalysis/final_output/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 42 2020-04-30 13:27 sentimentanalysis/final_output/part-r-00000
[cloudera@quickstart ~]$ hdfs dfs -cat sentimentanalysis/final_output/part-r-00000
00
Negative 3889
Neutral 35802
Positive 3094
[cloudera@quickstart ~]$ 

```

The terminal window shows Hadoop job statistics and file system operations. It includes metrics like Spilled Records (85570), Shuffled Maps (2), Failed Shuffles (0), and GC time elapsed (1097 ms). It also lists file input and output format counters, such as Bytes Read (172927823) and Bytes Written (42). The user then lists files in the sentimentanalysis/final_output directory, which contains two files: _SUCCESS and part-r-00000. Finally, the user reads the contents of the part-r-00000 file.

Use Case 4: Twitter data sentimental analysis using Cassandra

Useful analysis on tweets are performed using complex queries and user defined functions(UDFs).

Create Table:

Twitter_data table is created with 6 columns- date, user, is_retweet, is_quote, text, quoted_text. Primary key is a composite key with columns (date,user).

```
cqlsh:bigdata> CREATE TABLE twitter_data(
...     date timestamp,
...     user varchar,
...     is_retweet boolean,
...     is_quote boolean,
...     text text,
...     quoted_text text,
...     PRIMARY KEY (date,user)
... );
```

Load Twitter data:

COVID19 csv file is loaded into twitter_data table with more than 2 lakh records.

```
C:\Windows\System32\cmd.exe -cqlsh

35, in create_timer
    File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\io\asyncorereactor.py", line 335
      In create_timer
        File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\io\asyncorereactor.py", line 335
          In create_timer
            self._connection.close()
            cls._loop.add_timer(timer)
            cls._loop.add_timer(timer)
File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\cluster.py", line 2858, in shutdown
  AsyncoreConnection.create_timer(0, partial(asyncore.dispatcher.close, self))
  self._connection.close()
  cls._loop.add_timer(timer)
  cls._loop.add_timer(timer)
A File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\io\asyncorereactor.py", line 373
  In close
    cls._loop.add_timer(timer)
File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\io\asyncorereactor.py", line 35
  In create_timer
AAAAttributeError: 'NoneType' object has no attribute 'add_timer'
  File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\io\asyncorereactor.py", line 373
  In close
    ttributeError: 'NoneType' object has no attribute 'add_timer'
ttributeError: 'NoneType' object has no attribute 'add_timer'
ttributeError: 'NoneType' object has no attribute 'add_timer'
  AsyncoreConnection.create_timer(0, partial(asyncore.dispatcher.close, self))
  cls._loop.add_timer(timer)
  AsyncoreConnection.create_timer(0, partial(asyncore.dispatcher.close, self))
  File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\io\asyncorereactor.py", line 373
  In close
    A File "C:\Program Files\apache-cassandra-3.11.6-bin\apache-cassandra-3.11.6\bin..\lib\cassandra-driver-internal-only-3.11.0-bb96859b.zip\cassandra-driver-3.11.0-bb96859b\cassandra\io\asyncorereactor.py", line 335
  In create_timer
    cls._loop.add_timer(timer)
  ttributeError: 'NoneType' object has no attribute 'add_timer'
  A  cls._loop.add_timer(timer)
ttributeError: 'NoneType' object has no attribute 'add_timer'
AttributeError: 'NoneType' object has no attribute 'add_timer'
Processed: 350000 rows; Rate:  8928 rows/s; Avg. rate:  2709 rows/s
350000 rows imported from 1 files in 2 minutes and 9.193 seconds (0 skipped).
cqsh:bigdata> select count(*) from twitter_data;

count
-----
215151

(1 rows)

Warnings :
Aggregation query used without partition key

cqsh:bigdata>
```

Analysis 1: Display the users and their tweet count

Cassandra does not support GROUPBY clause. To perform group by on the user column to count the tweets by each user, created UDFs- cumulateCounter and groupCountByUser. CumulateCounter maps each user to 1 if it is a new user else increments the existing count.

```
cqlsh:bigdata> CREATE OR REPLACE FUNCTION cumulateCounter(state map<varchar,bigint>, user varchar, count counter)
...    . . .
...    RETURNS NULL ON NULL INPUT
...    RETURNS map<varchar,bigint>
...    LANGUAGE java
...    AS '
...    if(state.containsKey(user)) {
...        state.put(user, state.get(user) + count);
...    } else {
...        state.put(user, count);
...    }
...    return state;
...    ';
cqlsh:bigdata> CREATE OR REPLACE AGGREGATE groupCountByUser(varchar, counter)
...    SFUNC cumulateCounter
...    STYPE map<varchar,bigint>
...    INITCOND {};
cqlsh:bigdata>
```

```
C:\Windows\System32\cmd.exe - cqish
```

Analysis 2: Count the number of retweets

To perform GROUPBY aggregation on is_retweet column, created UDFs-
cumulateCounter2 and groupCountByRetweet. CumulateCounter2 aggregates
retweets by their Boolean value- false and true.

```
cqlsh:bigdata> CREATE OR REPLACE FUNCTION cumulateCounter2(state map<boolean,int>, is_retweet boolean)
... RETURNS NULL ON NULL INPUT
... RETURNS map<boolean,int>
... LANGUAGE java
... AS '
... if(state.containsKey(is_retweet)) {
...     state.put(is_retweet, state.get(is_retweet) + 1);
... } else {
...     state.put(is_retweet, 1);
... }
... return state;
...
';
cqlsh:bigdata> CREATE OR REPLACE AGGREGATE groupCountByRetweet(boolean)
... SFUNC cumulateCounter2
... STYPE map<boolean,int>
... INITCOND {};
cqlsh:bigdata> SELECT groupCountByRetweet(is_retweet)
... FROM twitter_data;

bigdata.groupcountbyretweet(is_retweet)
-----
{False: 179528, True: 35623}

(1 rows)

Warnings :
Aggregation query used without partition key

cqlsh:bigdata>
```

Analysis 3: Display users and their tweets whose username starts with character 'L'

Cassandra directly does not support LIKE operator. To display user tweets whose username starts with 'L', SASI index is created on user column.

```

cqlsh:bigdata> DROP INDEX user_idx;
cqlsh:bigdata> CREATE CUSTOM INDEX user_idx ON twitter_data (user) USING
...     ... 'org.apache.cassandra.index.sasi.SASIIndex' WITH
...     ... OPTIONS = {'analyzer_class': 'org.apache.cassandra.index.sasi.analyzer.StandardAnalyzer', 'case_sensitive': 'false'};

Warnings :
SASI indexes are experimental and are not recommended for production use.

cqlsh:bigdata> SELECT user, text FROM twitter_data WHERE user LIKE 'L%' LIMIT 15;

user          | text
-----+-----
Liberbot_|           null
LiekoIchihara |        null
LauraRizo |        null
Linomopeko |        null
LauraCalvoa |        null
pymes a trav@s de EOIs para contrarrestar los efectos COVID-19.           La Secretaria General de Industria y Pyme @IndustriaGob ofrece formaci@n y asesoramiento a #
LoriSums |        null
Lartyfat1 |        null
Lloydfinger |        null
Lerato_Kea |        null
LDPhilippe | 4./ #Covid19 que % la vitesse de r@-action des #Chinois dans la gestion des #-pid-@-mies a |-t|- stup@-fiant... que la #chloroquine ... peut-@-tre le meilleur #traitemen
la meilleure pr@-vention T@l @raoult_didier https://t.co/8x0skZaXn
Lizbeth12808221 |        null
LordSriRama333 |        null
LancsSocial |        null
Liv3_Ron |        null
SE. They INCREASED to 5000. The CURVE DID NOT FLATTEN AT ALL           We were under a month long lockdown & cases of Covid-19 DID NOT DECREA
LADominator |        null

(15 rows)
cqlsh:bigdata>

```

Analysis 4: Display tweets that are tweeted after date- '2020-04-13'

```
cqlsh:bigdata> select * from twitter_data where date >= '2020-04-13' LIMIT 5 ALLOW FILTERING;
date           | user          | is_quote | is_retweet | quoted_text
-----+-----+-----+-----+-----+
2020-04-30 01:06:29.000000+0000 | 2hmonster | False | False | How Milwaukee represents the COVID-19 pandemic for African Americans Race explains why cities across the country are seeing a dispropor
tionate amount of black deaths - https://t.co/8g0lwxRCZE https://t.co/0s5L9RMwC | null
2020-04-30 01:06:29.000000+0000 | AlejoCons | True | False |
RT @publico.es: #H11 ÚLTIMA HORA | La economía española se desploma un 5 % en el primer trimestre por la emergencia de la covid-19
2020-04-30 01:06:29.000000+0000 | Caveat_E | True | False |
RT @linamichi: Good news from South Korea. Looks like recovered COVID19 patients testing positive again was due to RNA fragments released from the body | null
2020-04-30 01:06:29.000000+0000 | Chonnnl | True | False |
RT @cpnpratt: 気温は毎日毎日下がってます。寒いです。 | null
2020-04-30 01:06:29.000000+0000 | Colonia Digital | False | False |
RT @ColoniaDigital: Consulta aquí el avance del Covid-19 en tu municipio: hay diagnósticos en 823 localidades https://t.co/qFS4LBDeDz | null
(5 rows)
cqlsh:bigdata>
```

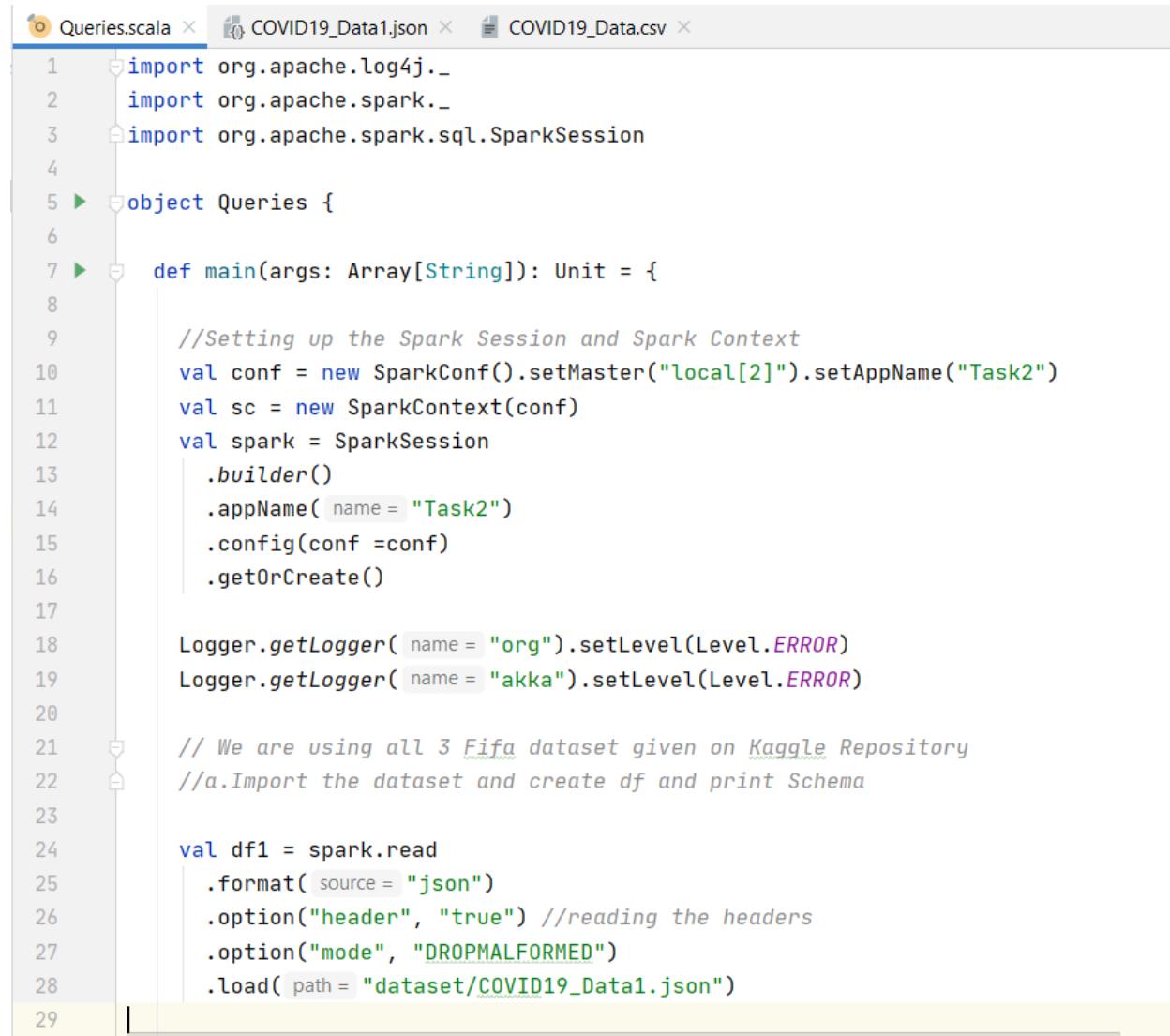
Query 5: Display table data in json format

```
cqlsh:bigdata> select json date, user, text from twitter_data LIMIT 15 ALLOW FILTERING;
[json]
-----
{"date": "2020-04-30 01:06:29.000Z", "user": "AlejoComs", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "2hmonster", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "Caveat_E", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "Charmi", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "ColoniaCapital", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "DILIPKUMAR9999", "text": "an\u00f1a an\u00f3nima an\u00f3nica an\u00f3nica an\u00f3nica an\u00f3nica an\u00f3nica..."}
 {"date": "2020-04-30 01:06:29.000Z", "user": "Fegzie_ ", "text": "#COVID19 confirmed cases have been reported in 34 states and the Federal Capital Territory"}
 {"date": "2020-04-30 01:06:29.000Z", "user": "HewallB", "text": "Help Me Retweet "}
 {"date": "2020-04-30 01:06:29.000Z", "user": "ImSandeshLain", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "KatherineKohler", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "KirstenJassies", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "MattGibbo1801", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "MelindaDaugh", "text": null}
 {"date": "2020-04-30 01:06:29.000Z", "user": "PaschalMalley", "text": null}

(15 rows)
cqlsh:bigdata>
```

Use Case 5: Twitter data sentimental analysis using Spark SQL

Load json data into SparkContext



```
Queries.scala x COVID19_Data1.json x COVID19_Data.csv x
1 import org.apache.log4j._
2 import org.apache.spark._
3 import org.apache.spark.sql.SparkSession
4
5 object Queries {
6
7   def main(args: Array[String]): Unit = {
8
9     //Setting up the Spark Session and Spark Context
10    val conf = new SparkConf().setMaster("local[2]").setAppName("Task2")
11    val sc = new SparkContext(conf)
12    val spark = SparkSession
13      .builder()
14      .appName( name = "Task2")
15      .config(conf =conf)
16      .getOrCreate()
17
18    Logger.getLogger( name = "org").setLevel(Level.ERROR)
19    Logger.getLogger( name = "akka").setLevel(Level.ERROR)
20
21    // We are using all 3 Fifa dataset given on Kaggle Repository
22    //a.Import the dataset and create df and print Schema
23
24    val df1 = spark.read
25      .format( source = "json")
26      .option("header", "true") //reading the headers
27      .option("mode", "DROPMALFORMED")
28      .load( path = "dataset/COVID19_Data1.json")
29}
```

Analysis 1: Display the count of the tweets that are retweeted

```
val Q1 = spark.sql( sqlText = "select count(1) retweet_count from Covid WHERE isretweet LIKE '%True%' ")  
Q1.show()
```

```
+-----+  
|retweet_count|  
+-----+  
|      19964|  
+-----+
```

Analysis 2: Display the top 5 users with highest number of tweets along with count

```
var Q2=spark.sql( sqlText = "select tweetuser user, count(1) tweets from Covid GROUP BY tweetuser ORDER BY tweets DESC LIMIT 5")  
Q2.show();
```

```
+-----+-----+  
|          user|tweets|  
+-----+-----+  
| "tweetuser":"auro...|    31|  
| "tweetuser":"maxp...|    29|  
| "tweetuser":"BotR...|    27|  
| "tweetuser":"PHLN...|    25|  
| "tweetuser":"Vala...|    16|  
+-----+-----+
```

Analysis 3: Display the top 5 users with highest number of quoted tweets along with count

```
val Q3 = spark.sql( sqlText = "select tweetuser user,count(1) quoted_tweets,isquote from Covid where isquote LIKE '%True%' " +  
    "|GROUP BY tweetuser,isquote ORDER BY quoted_tweets DESC LIMIT 5")  
Q3.show()
```

user	quoted_tweets	isquote
"tweetuser":"shei...	13	"isquote":"True"
"tweetuser":"EndTax"	11	"isquote":"True"
"tweetuser":"Rosa...	7	"isquote":"True"
"tweetuser":"Eliz...	6	"isquote":"True"
"tweetuser":"rwoo...	5	"isquote":"True"

Analysis 4: Display the number of tweets taking about 'United States'

```
val Q4 = spark.sql( sqlText = "select count(1) tweets_USA from Covid WHERE lower(tweettext) LIKE '%usa%' OR"  
    "| lower(tweettext) LIKE '%united states%' OR lower(tweettext) LIKE '%america%' ")  
Q4.show()
```

tweets_USA
836

Work Completed (100%):

- Collected Data
- Analyzed data using map reduce and hive
- Analyzed data using Cassandra
- Performed sentiment analysis using map reduce
- Analyzed data using Spark SQL

Responsibility:

Vidyullatha Lakshmi Kaza- 34%

Aparna Manda- 33%

Lohitha Yenugu- 33%

References:

1. <https://developer.twitter.com/en/docs/tweets/data-dictionary/overview/intro-to-tweet-json>
2. <https://cassandra.apache.org/doc/latest/cql/dml.html>
3. <https://stackoverflow.com/questions/17342176/max-distinct-and-group-by-in-cassandra>
4. https://docs.datastax.com/en/cql-oss/3.3/cql/cql_reference/cqlCreateAggregate.html