

INTRODUCTION

There are many datasets available in the twitter. Understanding of each data set is difficult. To make it easy to understand we perform some operations on it using different algorithms.

The main purpose of this project is to develop a system to store, analyze and visualize Twitter's tweets. The tasks to be performed in this phase are as follows :

- To work on the tweets related to the recently released mobile phones & their accessories and to figure out how to store them in Spark SQL.
- To write interesting analytical queries to explore and understand the data collected.
- To develop interesting visualizations of the above written queries.

Motivation:

This project is motivated to develop a system to store, analyze and visualize Twitters tweets.

Significance:

It provides a wealth of information that helps to create meaningful tweets that resonates with target audience. Compare followers with different personas, demographics, interests and consumer behaviours to see brand measures up etc. Watch individual Tweet performance, cumulative overview to compare monthly activity etc

SOFTWARE'S USED:

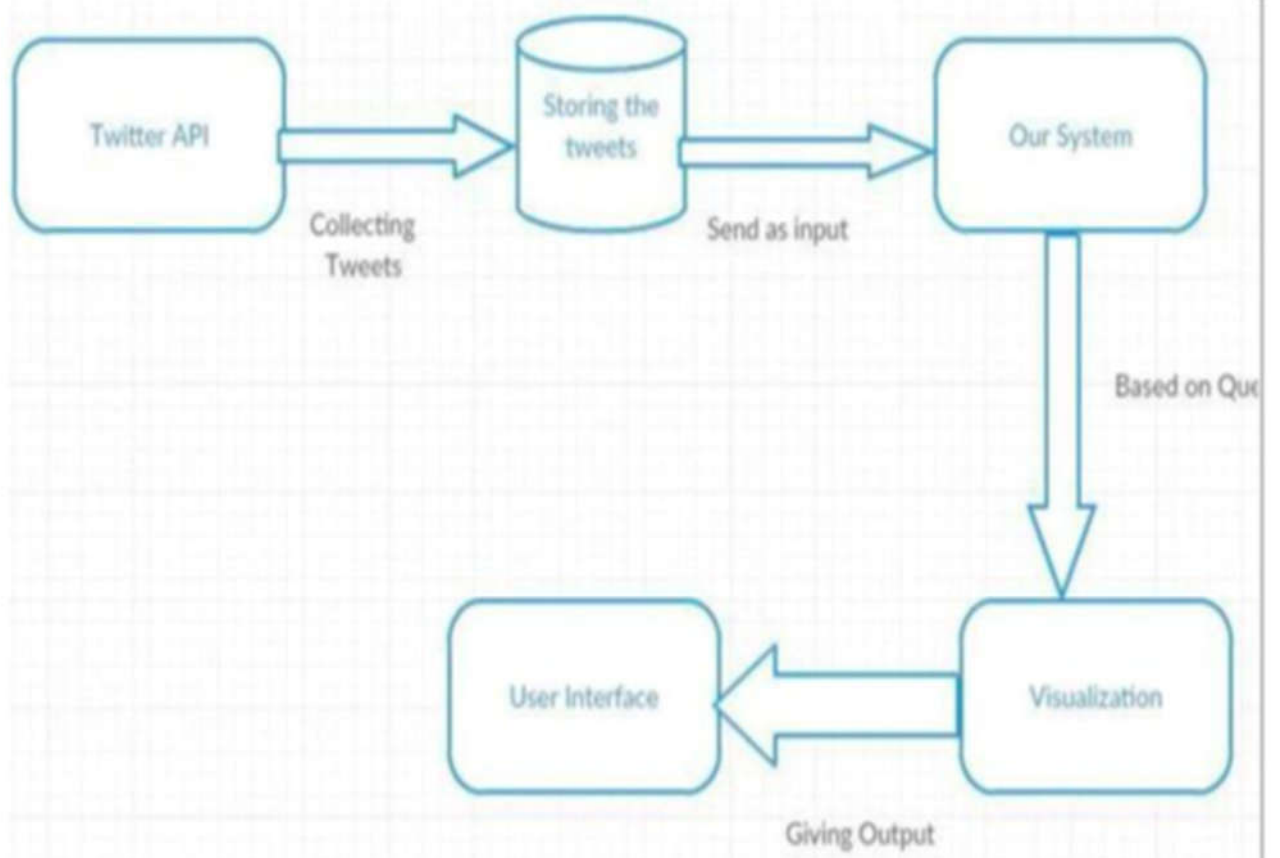
- Apache Spark and Scala
- IntelliJ IDEA
- Pycharm
- Twitter Developer Account
- Twitter4J Library for Tweets Extraction
- Tableau

DATASET : Twitter dataset (Phones / E accessories)

Detailed Description of Data set :

This data set contains data mainly about e accessories i.e mobiles . We have chosen this dataset because mobiles are used by everyone and there are different kinds of mobiles available in the market and analyzing this data makes it more interesting.

ARCHITECTURE :



IMPLEMENTATION:

Wordcount is done first in Apache Spark and Hadoop

Consumer API keys

jdY2pqlZ5uVWds4CQGIMfwxIN (**API key**)

hfjqljwzHCVG0fCLBDV9vFTJUSkXpX9uhcoy7uvFR7qJq96zGZ (**API secret key**)

Access token & access token secret

1042137305329352705-9uqx96PFyct3pTezrr29RVTSn6VXgB (Access token)

DLswmmzGGy4KRQReYF14hwgOetvbxqAP0N7mhuMwEwEg(Access token secret)
Read and write (Access level)

- To collect the tweets in JSON format, a Python program is written, the output of the program contains the tweets with all the details like the IDs, URLs, Hashtags, Created at, Text etc.

RELATED SCREENSHOTS :

Tokens and Keys Generated :

Consumer API keys:-

1) jdY2pqlZ5uVWds4CQGIMfwxIN (API key)

2) hfjq1jwzHCVG0fCLBDV9vFTJUSkXpX9uhcoy7uvFR7qJq96zGZ (API secret key)

Access token & access token secret:-

1) 1042137305329352705-9uqx96PFyct3pTezrr29RVTSn6VXgB (Access token)

2) DLswmmzGGy4KRQReYF14hwgOetvbxqAP0N7mhuMwEwEgm (Access token secret)
Read and write (Access level)

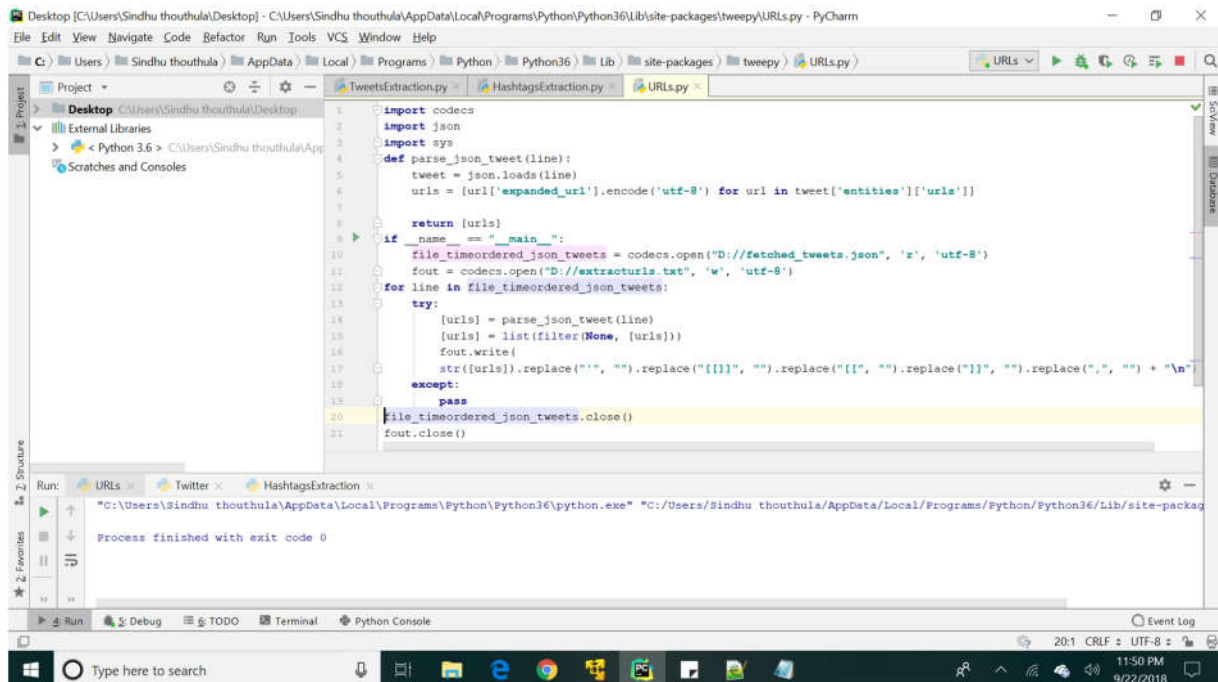
Python Code for Tweets Collection :


```
431488 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721135123181568", "id_str": "1043721135123181568", "text": "Developer Demos Web Code That Can Cause iOS 7+ Dev\n431489\n431490\n431491 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721135324454913", "id_str": "1043721135324454913", "text": "Sooriyan Poatri Tamil 108.. Use &amp; Share.. http\n431492\n431493 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721135664390144", "id_str": "1043721135664390144", "text": "RT @iatemuggles: superheroes being 198% done with\n431494\n431495 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721136150786048", "id_str": "1043721136150786048", "text": "Ki\\ulec3m so\\u00elt \\u0111\\u01b0\\u1edding \\u0111\\u1\n431496\n431497 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721135584632832", "id_str": "1043721135584632832", "text": "Kodak x Prison \\nKodak x Jail\\nKodak x Penitentiari\n431498\n431499 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721135823687680", "id_str": "1043721135823687680", "text": "RT @__BrendasS: Learned Ariana is watching Mac Mil\n431500\n431501 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721135752384512", "id_str": "1043721135752384512", "text": "I made my boyfriend watch Scott pilgrim vs the wor\n431502\n431503 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721136071237632", "id_str": "1043721136071237632", "text": "I wish I could complain like a white woman rn. Who\n431504\n431505 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721135228051456", "id_str": "1043721135228051456", "text": "eBay: IPHONE 5s (VERIZON) GRAY https://t.co/B1v\n431506\n431507 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721136347930625", "id_str": "1043721136347930625", "text": "RT @michellemalkin: https://t.co/NeFNBWnvXQ htt\n431508\n431509 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721136272338944", "id_str": "1043721136272338944", "text": "RT @TheMarkAhrens: Great example of the difference\n431510\n431511 [{"created_at": "Sun Sep 23 04:38:04 +0000 2018", "id": "1043721136285134850", "id_str": "1043721136285134850", "text": "RT @letters4trump45: \\ud83d\\udea8 Reminder\\n\\nWato\n431512\n431513 [{"created_at": "Sun Sep 23 04:38:03 +0000 2018", "id": "1043721132178857984", "id_str": "1043721132178857984", "text": "RT @dxrpingcashton: Is there any link to watch Sso\n431514\n431515 [{"created_at": "Sun Sep 23 04:38:01 +0000 2018", "id": "1043721121500200960", "id_str": "1043721121500200960", "text": "RT @Lantzlusca: Caralho que louco deixa um iPhone\n431516\n431517 [{"created_at": "Sun Sep 23 04:38:05 +0000 2018", "id": "1043721136977076224", "id_str": "1043721136977076224", "text": "RT @TaiiNeko: Apple \\u0e19\\u0e35\\u0e48\\u0e21\\u0e35\n431518\n431519 [{"created_at": "Sun Sep 23 04:38:05 +0000 2018", "id": "1043721136851378176", "id_str": "1043721136851378176", "text": "##\\u062a\\u0648\\u062b\\u064a\\u0642_\\u0633\\u0646\\u0627\n431520\n431521 [{"created_at": "Sun Sep 23 04:38:05 +0000 2018", "id": "1043721136775856129", "id_str": "1043721136775856129", "text": "Please watch The Hunting Ground", "source": "\\u003ca\n431522\n431523
```

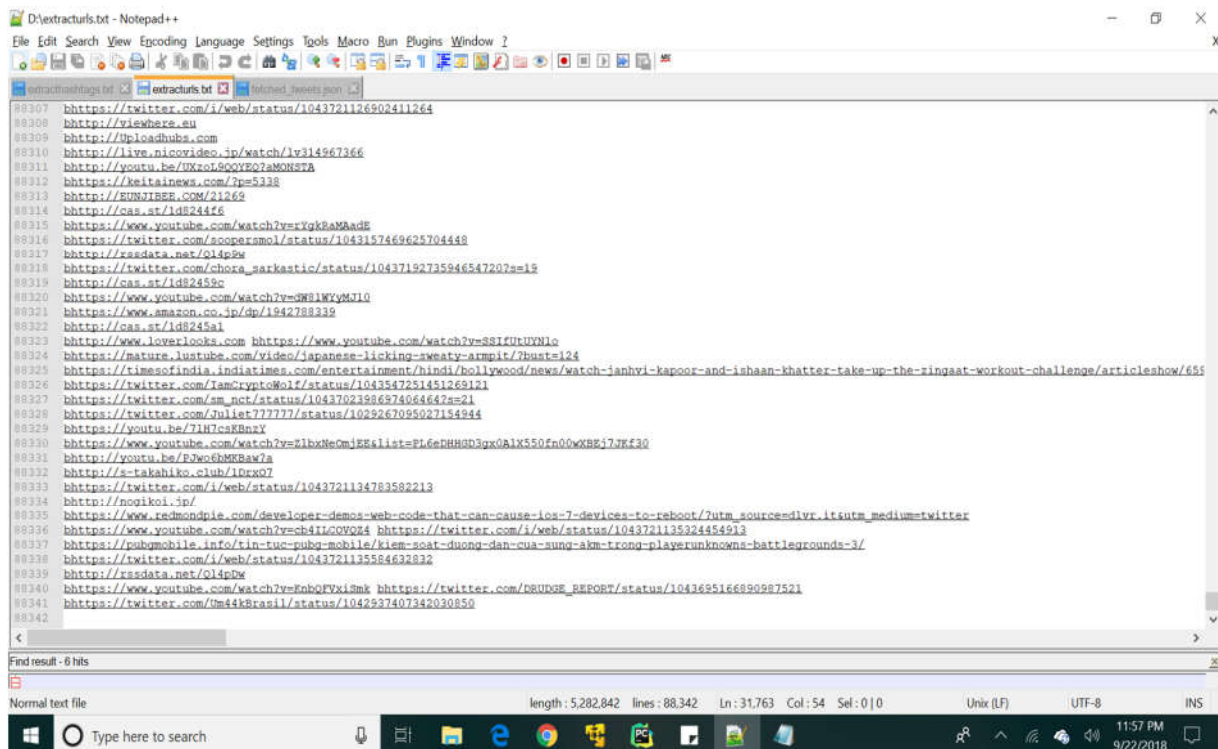
2) HASHTAGS AND URLS EXTRACTION :

To extract the hashtags and URLs from the collected Tweets, we have again run two Python programs through which we generated the files containing the hashtags and URLs alone.

Hashtags Code :



Output :



3) RUNNING THE WORDCOUNT IN APACHE HADOOP AND APACHE SPARK :

Word Count in Apache Spark :

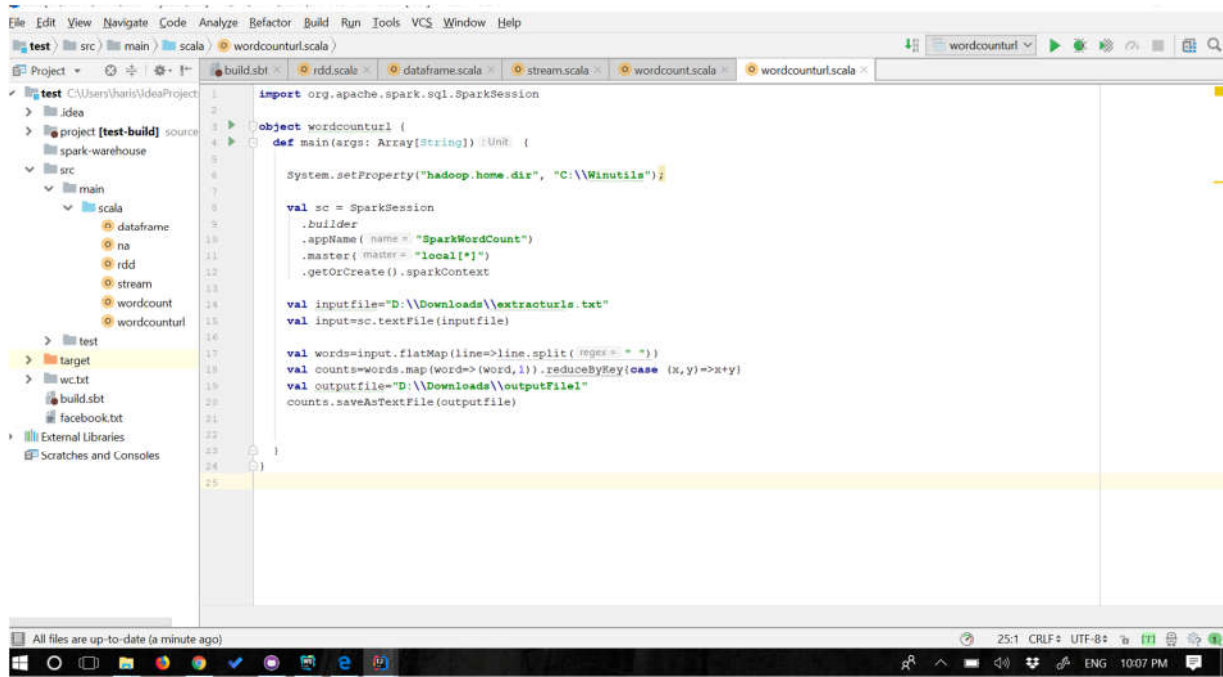
- Here we made use of IntelliJ IDEA to run the Word Count program on the extracted Hashtags and URLs in Apache Spark and Scala.
- Here in this IDEA to execute the Word Count on the extracted hashtags and URLs, Spark is integrated on Scala.
- We have given the values and input files using the following :

```
val inputfile="D:\\Downloads\\extracthashtags.txt"  
val input=sc.textFile (inputfile)
```

```
val words=input.flatMap(line=>line.split(regex= " "))  
val counts=words.map (word=>(word, 1)).reduceByKey(case  
(x,y)=>x+y)  
val outputfile="D:\\Downloads\\outputFile"  
counts.saveAsTextFile (outputfile)
```

- Here the sc.text file in Spark creates Resilient Distributed Datasets with each line considered as an element.
- RDDs are immutable and partitioned collection of records which can only be created by coarse grained operations such as the map, filter, group etc.
- Once the file is read and stored in the variable textfile, wordcount is executed with the help of the MapReduce function.
- MapReduce function is broken into two phases, the Map phase and Reduce phase where each has got a key value pair as an input and output.
- Once the reduce function groups the words with the count, the collected data is stored in the output file.

Spark :



The screenshot shows an IDE window with a project structure on the left and a Scala code editor on the right. The project structure includes a 'test' directory with sub-directories 'src' and 'main'. The 'main' directory contains a 'scala' sub-directory with files 'dataframe', 'na', 'rdd', 'stream', 'wordcount', and 'wordcounturl'. The 'wordcounturl.scala' file is open in the editor, showing the following code:

```
1 import org.apache.spark.sql.SparkSession
2
3 object wordcounturl {
4   def main(args: Array[String]): Unit = {
5
6     System.setProperty("hadoop.home.dir", "C:\\Winutils")
7
8     val sc = SparkSession
9       .builder
10      .appName("SparkWordCount")
11      .master("local[*]")
12      .getOrCreate().sparkContext
13
14     val inputFile="D:\\Downloads\\extracturis.txt"
15     val input=sc.textFile(inputFile)
16
17     val words=input.flatMap(line=>line.split(" "))
18     val counts=words.map(word=>(word,1)).reduceByKey(case (x,y)=>x+y)
19     val outputFile="D:\\Downloads\\outputFile1"
20     counts.saveAsTextFile(outputFile)
21
22   }
23 }
24
25
```

The status bar at the bottom indicates 'All files are up-to-date (a minute ago)', '25:1', 'CRLF', 'UTF-8', and 'ENG 10:07 PM'.

WordCount in Hadoop :

- To run the wordcount, we have written the code in Java language.
- In this program, the TokenizerMapper class is created with the extension of Mapper class and by overriding the map method in the Mapper class, we implemented the Map Function.
- Here the Mapper function takes a key value pair as an input and outputs the key value pair.
- The IntSumReducer class is created by extending the org.apache.hadoop.mapreduce.Reducer class and the reduce method is implemented by overriding the reduce method from the Reducer class.
- The main method in the code is to setup all the configurations and runs the map reduce job.
- In Hadoop to execute the hashtags and URLs, the .jar file is exported to cloudera.
- Using the following commands, the wordcount on hashtags and URLs is performed in Hadoop.

```
$ cat > /home/cloudera/extracturls.txt
$ hdfs dfs -mkdir /input2
$ hdfs dfs -put '/home/cloudera/input/extracturls.txt' /input2
$ hdfs dfs -cat /input2/extracturls.txt
$ hadoop jar wordcount.jar WordCount /input2 /wordcountoutput2
$ hdfs dfs -get /wordcountoutput2 /home/cloudera/output2
$ hdfs dfs -cat/output2/part-r-00000
```

WordCount Java Program – Hadoop :

The screenshot shows a VMware Workstation 14 Player window titled "cloudera-quickstart-vm-5.13.0-0-vmware - VMware Workstation 14 Player (Non-commercial use only)". Inside the VM, a Gvim editor is open with the file "WordCount.java". The code is a Java MapReduce application for word counting. It includes imports for Hadoop and Java utilities, and defines a "TokenizerMapper" and an "IntSumReducer". The code is as follows:

```

import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static class TokenizerMapper
        extends Mapper<Object, Text, Text, IntWritable>{

        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(Object key, Text value, Context context
            ) throws IOException, InterruptedException {
            StringTokenizer itr = new StringTokenizer(value.toString());
            while (itr.hasMoreTokens()) {
                word.set(itr.nextToken());
                context.write(word, one);
            }
        }
    }

    public static class IntSumReducer
        extends Reducer<Text, IntWritable, Text, IntWritable> {
        private IntWritable result = new IntWritable();

        public void reduce(Text key, Iterable<IntWritable> values,
            Context context
            ) throws IOException, InterruptedException {

            int sum = 0;
            for (IntWritable val : values) {
                sum += val.get();
            }
            result.set(sum);
            context.write(key, result);
        }
    }
}

```

The IDE interface shows a menu bar (File, Edit, Tools, Syntax, Buffers, Window, Help) and a toolbar. The status bar at the bottom of the Gvim window indicates "45,1 Top". The VMware Workstation window has a title bar with standard OS controls and a menu bar (Player, Applications, Places, System). The bottom of the screen shows the Windows taskbar with various icons and the system clock showing 8:10 PM on 9/25/2018.

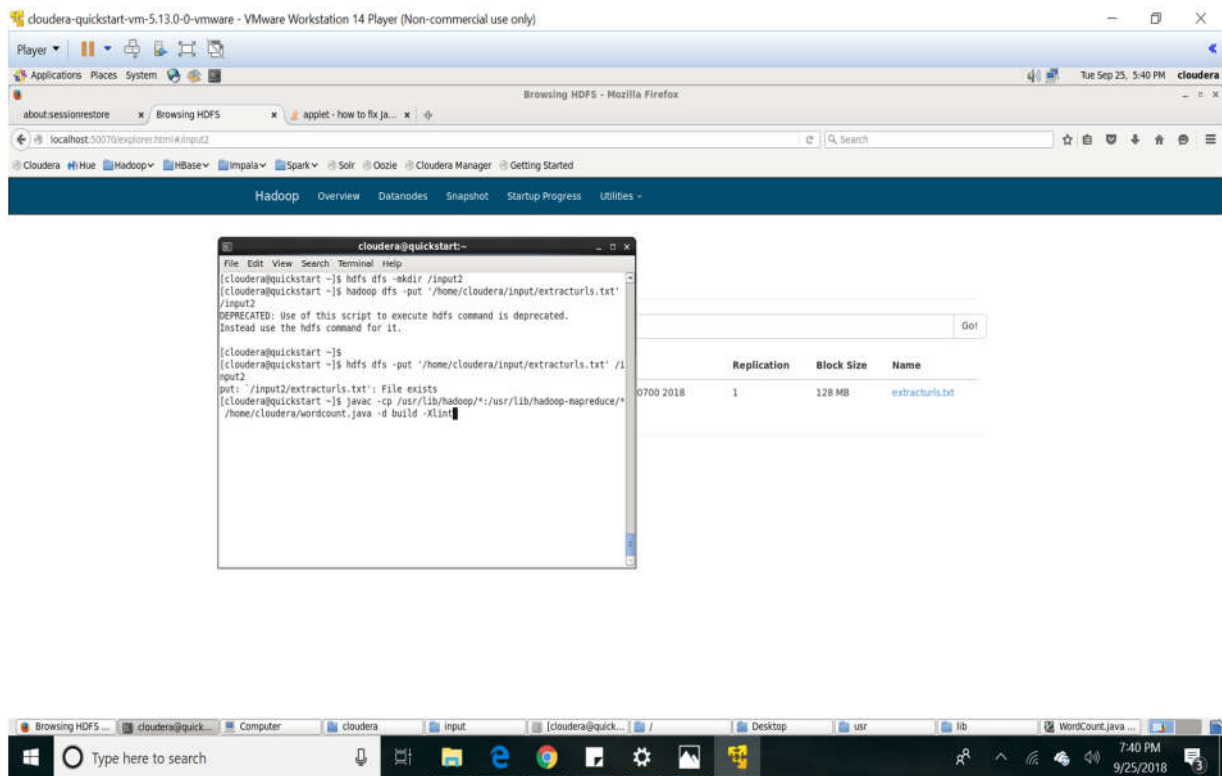
The screenshot shows a VMware Workstation 14 Player window titled "cloudera-quickstart-vm-5.13.0-0-vmware - VMware Workstation 14 Player (Non-commercial use only)". Inside the VM, a Java IDE (likely IntelliJ IDEA) is open, displaying the file "WordCount.java (-) - GVIM". The code is a Java MapReduce application for word counting. The main method sets up a Hadoop job with the following configuration:

```

public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "word count");
    job.setJarByClass(WordCount.class);
    job.setMapperClass(TokenizerMapper.class);
    job.setCombinerClass(IntSumReducer.class);
    job.setReducerClass(IntSumReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}

```

The IDE interface includes a menu bar (File, Edit, Tools, Syntax, Buffers, Window, Help), a toolbar with various icons, and a status bar at the bottom showing "45,1" and "Bot". The Windows taskbar at the very bottom shows the time as 8:11 PM on 9/25/2018.



Output : map 100% reduce 100%

```
cloudera@quickstart:~$ cat /etc/passwd | wc -l
18/09/25 17:42:24 INFO mapreduce.Job: The url to track the job: http://quickstar
18/09/25 17:42:24 INFO mapreduce.Job: Running job: job_1537917063927_0002
18/09/25 17:42:39 INFO mapreduce.Job: Job job_1537917063927_0002 running in uber
mode : false
18/09/25 17:42:39 INFO mapreduce.Job: map 0% reduce 0%
18/09/25 17:42:54 INFO mapreduce.Job: map 100% reduce 0%
18/09/25 17:43:11 INFO mapreduce.Job: map 100% reduce 100%
18/09/25 17:43:13 INFO mapreduce.Job: Job job_1537917063927_0002 completed succe
ssfully
18/09/25 17:43:13 INFO mapreduce.Job: Counters: 49
File System Counters
  FILE: Number of bytes read=3471387
  FILE: Number of bytes written=7229455
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=5282961
  HDFS: Number of bytes written=3247694
  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)=12932
Total time spent by all reduces in occupied slots (ms)=14936
Total time spent by all map tasks (ms)=12932
Total time spent by all reduce tasks (ms)=14936
Total vcore-millisecods taken by all map tasks=12932
Total vcore-millisecods taken by all reduce tasks=14936
Total megabyte-millisecods taken by all map tasks=13242368
Total megabyte-millisecods taken by all reduce tasks=15294464
Map-Reduce Framework
  Map input records=88341
  Map output records=96648
  Map output bytes=5671927
  Map output materialized bytes=3471387
  Input split bytes=119
  Combine input records=96648
  Combine output records=55273
  Reduce input groups=55273
  Reduce shuffle bytes=3471387
  Reduce input records=55273
  Reduce output records=55273
  Spilled Records=110546
  Shuffled Maps=1
```

- Initially collected the tweets in JSON format for which a Python program is written, the output of the program contains the tweets with all the details like the IDs, URLs, Hashtags, Created at, Text etc.
- The twitter data is collected on the concept based on to analyze & visualize the data regarding various phone/e-accessories.
- The extracted JSON tweets are persisted into the Apache Spark in the form of tables.
- Query written in Scala language will be sent to spark server and the outputs files are stored in the form of CSV/JSON files.
- These CSV/JSON output files are used to visualize the data using Bar Graphs, Pie Charts through Tableau.
- Key-words used in the tweets extraction are as follows :


iphone, iphonex, iPhoneXs, iPhoneXR, iPhoneXr, iPhone, #iPhone, AirPods, mobile, watch, technology, Accessories, Mac, iOS, update, music, latest etc.

ANALYTICAL QUERIES:

The following are the 10 queries on which we performed the visualizations.

```
build.sbt x queries.scala x
30 "WHEN text like '%technology%' THEN 'TECHNOLOGY'" +
31 "WHEN text like '%Accessories%' THEN 'ACCESSORIES'" +
32 "END AS phoneType from tweets where text is not null")
33 disCat.createOrReplaceTempView( viewName = "disCat2")
34 val disCat3 = sqlContext.sql( sqlText = "SELECT user.name as UserName,user.location as loc,text,created_at," +
35 "CASE WHEN text like '%iPhoneX%' OR text like '%Iphonex%' OR text like '%Iphonex%' OR text like '%Iphonex%' THEN 'Iphone X'" +
36 "WHEN text like '%iphone7%' OR text like '%iphone7plus%' OR text like '%iPHONE7%' OR text like '%iPHONE 7%' OR text like '%iphone 7%' OR text like '%iPHONE7p'" +
37 "WHEN text like '%iphone8%' OR text like '%iPHONE 8%' OR text like '%iphone 8%' OR text like '%iphone8plus%' OR text like '%iPHONE8%' OR text like '%iPHONE8p'" +
38 "WHEN text like '%AirPods%' OR text like '%airpods%' THEN 'AirPods'" +
39 "WHEN text like '%watch%' OR text like '%Watch%' OR text like '%technology%' OR text like '%Technology%' THEN 'TECHNOLOGY'" +
40 "WHEN text like '%ios%' OR text like '%IOS%' OR text like '%ios%' THEN 'IOS'" +
41 "WHEN text like '%accessories%' OR text like '%ACCESSORIES%' THEN 'Accessories'" +
42 "WHEN text like '%Mac%' OR text like '%mac%' OR text like '%MAC%' THEN 'MAC'" +
43 "WHEN text like '%mobile%' OR text like '%MOBILE%' THEN 'Mobile'" +
44 "END AS phoneType from tweets where text is not null")
45 disCat3.createOrReplaceTempView( viewName = "disCat4")
46 println("Enter any one of the following query to get data")
47 println("1.Query-1:This query fetches the phone/e-accessories and their popularity based on tweets data")
48 println("2.Query-2:Which user tweeted most about which type of phone/e-accessories")
49 println("3.Query-3:Tweets from different countries about phone/e-accessories")
50 println("4.Query-4:On which day more tweets are done")
51 println("5.Query-5:This query fetches tweets count for different types of phone/e-accessories")
52 println("6.Query-6:Language mostly used for tweeting about phone/e-accessories")
53 println("7.Query-7:Number of tweets for particular date ")
54 println("8.Query-8:Tweets from verified accounts")
55 println("9.Query-9:On Which hours More Tweets Were Done")
56 println("10.Query-10:Which state is mostly having tweets about type of phone/e-accessories")
57 println("Enter any one of the following query to get data:")
58 val count = scala.io.StdIn.readLine()
59 count match {
```

```
18/12/02 17:21:10 INFO FileScanRDD: Reading File path: file:///C:/Users/Sindhu%20thouthula/Desktop/PB%20Project/Phase%20-1/Source/fetched tweets.json, range: 1
18/12/02 17:21:10 INFO Executor: Running task 8.0 in stage 0.0 (TID 8)
18/12/02 17:21:10 INFO FileScanRDD: Reading File path: file:///C:/Users/Sindhu%20thouthula/Desktop/PB%20Project/Phase%20-1/Source/fetched tweets.json, range: 1
18/12/02 17:21:10 INFO Executor: Finished task 5.0 in stage 0.0 (TID 5). 34391 bytes result sent to driver
18/12/02 17:21:10 INFO TaskSetManager: Finished task 5.0 in stage 0.0 (TID 5) in 2904 ms on localhost (executor driver) (8/11)
18/12/02 17:21:11 INFO Executor: Finished task 10.0 in stage 0.0 (TID 10). 33931 bytes result sent to driver
18/12/02 17:21:11 INFO TaskSetManager: Finished task 10.0 in stage 0.0 (TID 10) in 698 ms on localhost (executor driver) (9/11)
18/12/02 17:21:12 INFO Executor: Finished task 9.0 in stage 0.0 (TID 9). 34519 bytes result sent to driver
18/12/02 17:21:12 INFO TaskSetManager: Finished task 9.0 in stage 0.0 (TID 9) in 2101 ms on localhost (executor driver) (10/11)
18/12/02 17:21:12 INFO Executor: Finished task 8.0 in stage 0.0 (TID 8). 34466 bytes result sent to driver
18/12/02 17:21:12 INFO TaskSetManager: Finished task 8.0 in stage 0.0 (TID 8) in 2152 ms on localhost (executor driver) (11/11)
18/12/02 17:21:12 INFO TaskSchedulerImpl: Removed TaskSet 0.0, whose tasks have all completed, from pool
18/12/02 17:21:12 INFO DAGScheduler: ResultStage 0 (json at queries.scala:16) finished in 9.774 s
18/12/02 17:21:12 INFO DAGScheduler: Job 0 finished: json at queries.scala:16, took 9.850900 s
18/12/02 17:21:13 WARN Utils: Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting 'spark.debug.maxTo
Enter any one of the following query to get data
1.Query-1:This query fetches the phone/e-accessories and their popularity based on tweets data
2.Query-2:Which user tweeted most about which type of phone/e-accessories
3.Query-3:Tweets from different countries about phone/e-accessories
4.Query-4:On which day more tweets are done
5.Query-5:This query fetches tweets count for different types of phone/e-accessories
6.Query-6:Language mostly used for tweeting about phone/e-accessories
7.Query-7:Number of tweets for particular date
8.Query-8:Tweets from verified accounts
9.Query-9:On Which hours More Tweets Were Done
10.Query-10:Which state is mostly having tweets about type of phone/e-accessories
Enter any one of the following query to get data:
```



1.Query to fetch the tweets of phones/e-accessories based on the popularity.

This query is written to analyze the tweets that are made by the users on the particular key-words which would reflect the popularity of the phones/e-accessories. It would result the count i.e., how many times the phone/e-accessories would appear in the tweets made.

Query-Code :


```

/*-----Query 1: This query fetches the phones and its popularity based on tweets data-----*/
val textFile = sc.textFile( path = "C:\\Users\\Sindhu thouthula\\Desktop\\PB Project\\Phase -1\\Source\\fetched_tweets.json")
val iphone = (textFile.filter(line => line.contains("iphone")).count())
val iPhoneX = (textFile.filter(line => line.contains("iPhoneX")).count())
val iPhoneXs = (textFile.filter(line => line.contains("iPhoneXs")).count())
val iPhoneXR = (textFile.filter(line => line.contains("iPhoneXR")).count())
val Mac = (textFile.filter(line => line.contains("Mac")).count())
val iOS = (textFile.filter(line => line.contains("iOS")).count())
val AirPods = (textFile.filter(line => line.contains("AirPods")).count())
val mobile = (textFile.filter(line => line.contains("mobile")).count())
val watch = (textFile.filter(line => line.contains("watch")).count())
val technology = (textFile.filter(line => line.contains("technology")).count())
val Accessories = (textFile.filter(line => line.contains("Accessories")).count())
println("*****")
println("Number of tweets on different types of phones")
println("*****")
println("iphone : %s".format(iphone))
println("iPhoneX : %s".format(iPhoneX))
println("iPhoneXs : %s".format(iPhoneXs))
println("iPhoneXR : %s".format(iPhoneXR))
println("Mac : %s".format(Mac))
println("iOS : %s".format(iOS))
println("AirPods : %s".format(AirPods))
println("mobile : %s".format(mobile))
println("watch : %s".format(watch))
println("technology : %s".format(technology))
println("Accessories : %s".format(Accessories))

```

Output :

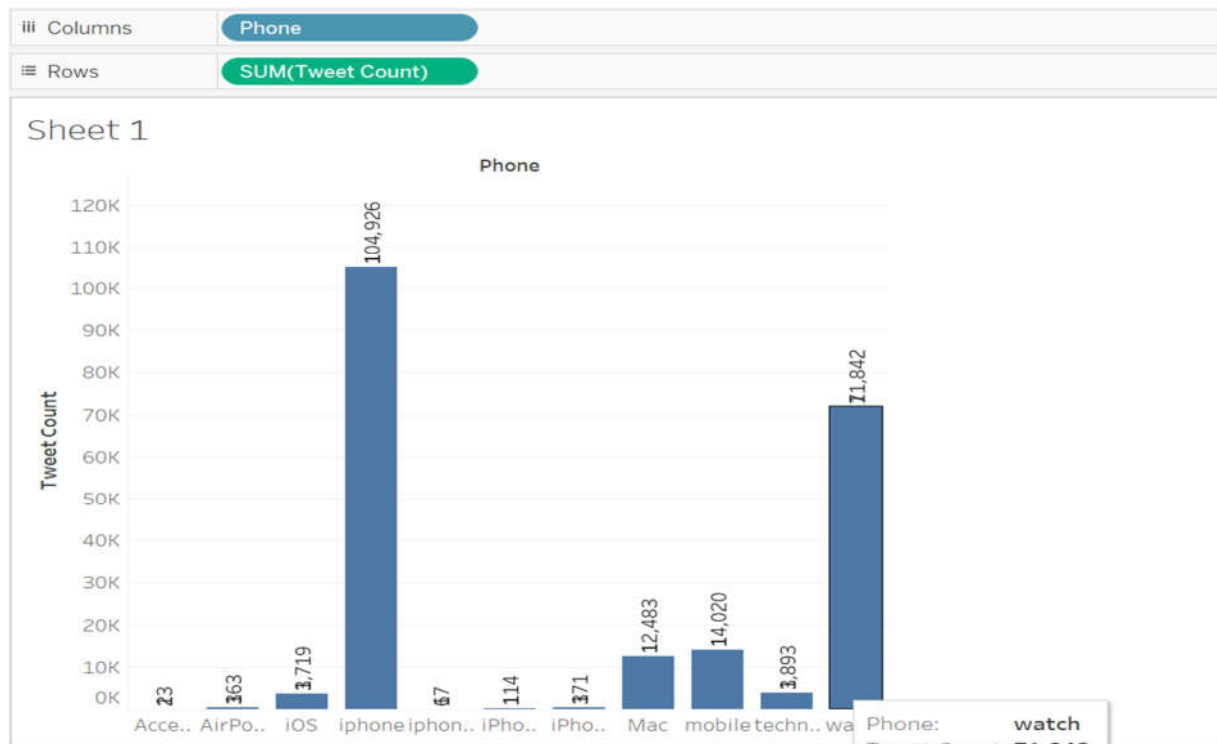
```

*****
Number of tweets on different types of phones
*****
iphone : 104926
iPhoneX : 67
iPhoneXs : 371
iPhoneXR : 114
Mac : 12493
iOS : 3719
AirPods : 363
mobile : 14020
watch : 71842
technology : 3893
Accessories : 23
18/12/02 17:22:30 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:22:30 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-RF770CO:4040
18/12/02 17:22:30 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:22:30 INFO MemoryStore: MemoryStore cleared
18/12/02 17:22:30 INFO BlockManager: BlockManager stopped
18/12/02 17:22:30 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:22:30 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:22:30 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:22:30 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:22:30 INFO ShutdownHookManager: Deleting directory C:\Users\Sindhu thouthula\AppData\Local\Temp\spark-d87d2d92-a575-4b59-97f3-0d8e487d57cd

Process finished with exit code 0

```





2. Query for finding which user tweeted more about the type of phone/e-accessories.

This query is written to find the user that most tweeted about the phone/e-accessories so that it would result the count of how many times a user tweeted at most for each kind particularly.

Query-Code :

```

val r1 = sqlContext.sql( sqlText = "SELECT UserName,'IPHONE' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='IPHONE' " +
    "group by UserName order by count desc limit 1")
val r2 = sqlContext.sql( sqlText = "SELECT UserName,'IPHONEX' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='IPHONEX' " +
    "group by UserName order by count desc limit 1")
val r3 = sqlContext.sql( sqlText = "SELECT UserName,'IPHONEXS' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='IPHONEXS' " +
    "group by UserName order by count desc limit 1")
val r4 = sqlContext.sql( sqlText = "SELECT UserName,'IPHONEXR' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='IPHONEXR' " +
    "group by UserName order by count desc limit 1")
val r5 = sqlContext.sql( sqlText = "SELECT UserName,'MAC' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='MAC' " +
    "group by UserName order by count desc limit 1")
val r6 = sqlContext.sql( sqlText = "SELECT UserName,'IOS' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='IOS' " +
    "group by UserName order by count desc limit 1")
val r7 = sqlContext.sql( sqlText = "SELECT UserName,'AIRPODS' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='AIRPODS' " +
    "group by UserName order by count desc limit 1")
val r8 = sqlContext.sql( sqlText = "SELECT UserName,'MOBILE' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='MOBILE' " +
    "group by UserName order by count desc limit 1")
val r9 = sqlContext.sql( sqlText = "SELECT UserName,'WATCH' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='WATCH' " +
    "group by UserName order by count desc limit 1")
val r10 = sqlContext.sql( sqlText = "SELECT UserName,'TECHNOLOGY' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='TECHNOLOGY' " +
    "group by UserName order by count desc limit 1")
val r11 = sqlContext.sql( sqlText = "SELECT UserName,'ACCESSORIES' as phoneType,count(*) as count FROM disCat2 WHERE phoneType='ACCESSORIES' " +
    "group by UserName order by count desc limit 1")

val rdd1 = r1.union(r2).union(r3).union(r4).union(r5).union(r6).union(r7).union(r8).union(r9).union(r10).union(r11)

println("*****")
println("Which user tweeted more on which type of phone")
println("*****")
rdd1.show()

```

Output :

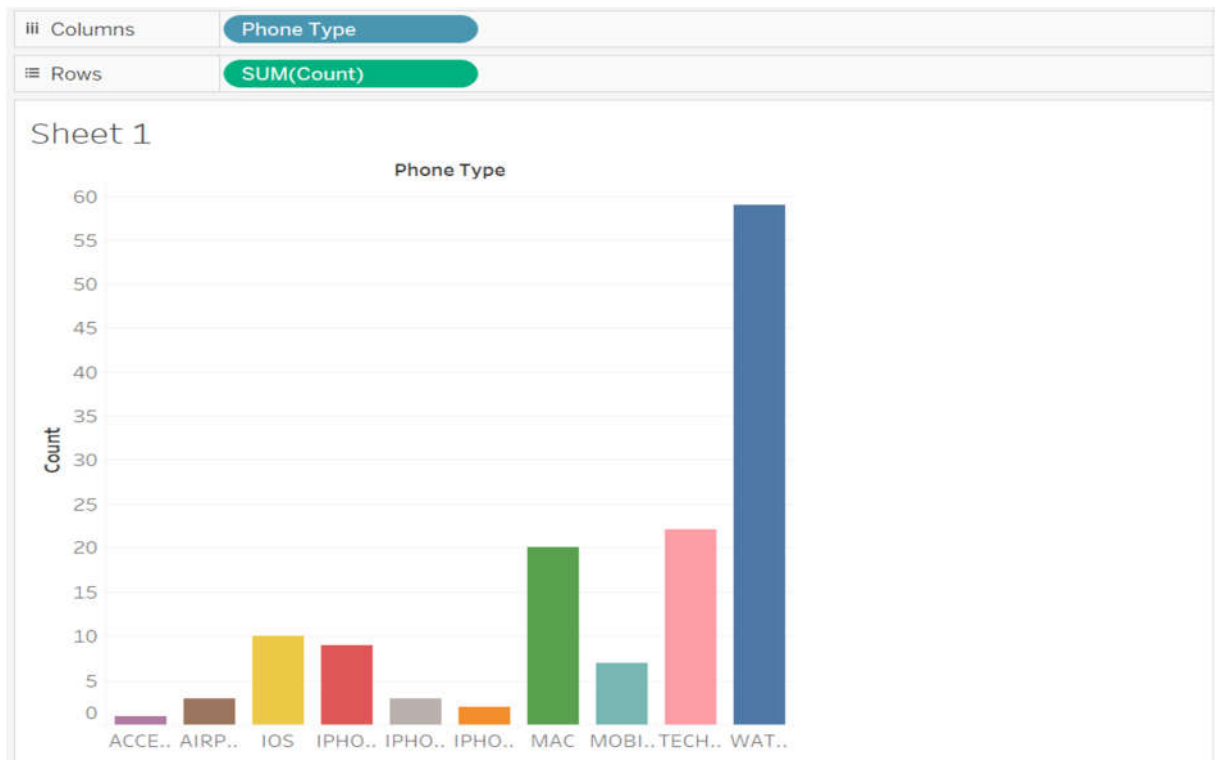
```

18/12/02 17:25:21 INFO DAGScheduler: Job 3 finished: show at queries.scala:119, took 0.240341 s
+-----+-----+-----+
|      UserName|  phoneType|count|
+-----+-----+-----+
|  Art_By_Lamarcus|    IPHONE|    9|
|      இராமன்|    IPHONEXS|    2|
|   Shane Salkey|    IPHONEXR|    3|
| jenna montemayor|        MAC|   20|
| Muhammad Farhan B...|        IOS|   10|
|      Ryan Ngala|    AIRPODS|    3|
| Bharti Airtel India|    MOBILE|    7|
|      94.9 THE BULL|    WATCH|   59|
|      ramana|    TECHNOLOGY|   22|
| Official_Tfshop|ACCESSORIES|    1|
+-----+-----+-----+

18/12/02 17:25:21 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:25:21 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-RF770C0:4040
18/12/02 17:25:21 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:25:23 INFO MemoryStore: MemoryStore cleared
18/12/02 17:25:23 INFO BlockManager: BlockManager stopped
18/12/02 17:25:23 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:25:23 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:25:23 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:25:23 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:25:23 INFO ShutdownHookManager: Deleting directory C:\Users\Sindhu\thouthula\AppData\Local\Temp\spark-f11bdba4-dc07-4da7-a022-6d8f5f5ba85f

Process finished with exit code 0

```



3.Query for fetching tweets from different countries.

This query is written to find the tweets based on the locations such that it would count how many tweets are posted about the phones/e-accessories from different countries.

Query-Code :

```
/*-----Query 3: Tweets from different countries about phones -----*/
case "3" =>
val countrytweetscount = sqlContext.sql( sqlText = "SELECT distinct place.country, count(*) as count FROM tweets where place.country is not null " + "GROUP BY
countrytweetscount.createOrReplaceTempView( viewName = "countrytweetscount")
println("*****")
println("Tweets from different countries")
println("*****")
countrytweetscount.show()
```

Output :

```
+-----+
|      country|count|
+-----+
|      United States| 1687|
|           India|   134|
| Republic of the P...|   93|
|           Canada|   73|
|         Australia|   54|
|           日本|   48|
|         Malaysia|   43|
|           Brasil|   32|
|         Indonesia|   20|
|        South Africa|   17|
|           México|   14|
|      United Kingdom|   13|
|           Nigeria|   12|
| Republic of Korea|   12|
|           Brazil|   11|
|        New Zealand|   11|
|         Sri Lanka|    9|
|          Uruguay|    9|
|         Singapore|    8|
|          Argentina|    8|
+-----+
only showing top 20 rows

18/12/02 17:26:55 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:26:55 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-RF770CO:4040
18/12/02 17:26:55 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:26:55 INFO MemoryStore: MemoryStore cleared
18/12/02 17:26:55 INFO BlockManager: BlockManager stopped
18/12/02 17:26:55 INFO BlockManagerMaster: BlockManagerMaster stopped
```

TODO

sbt shell

Terminal

Build

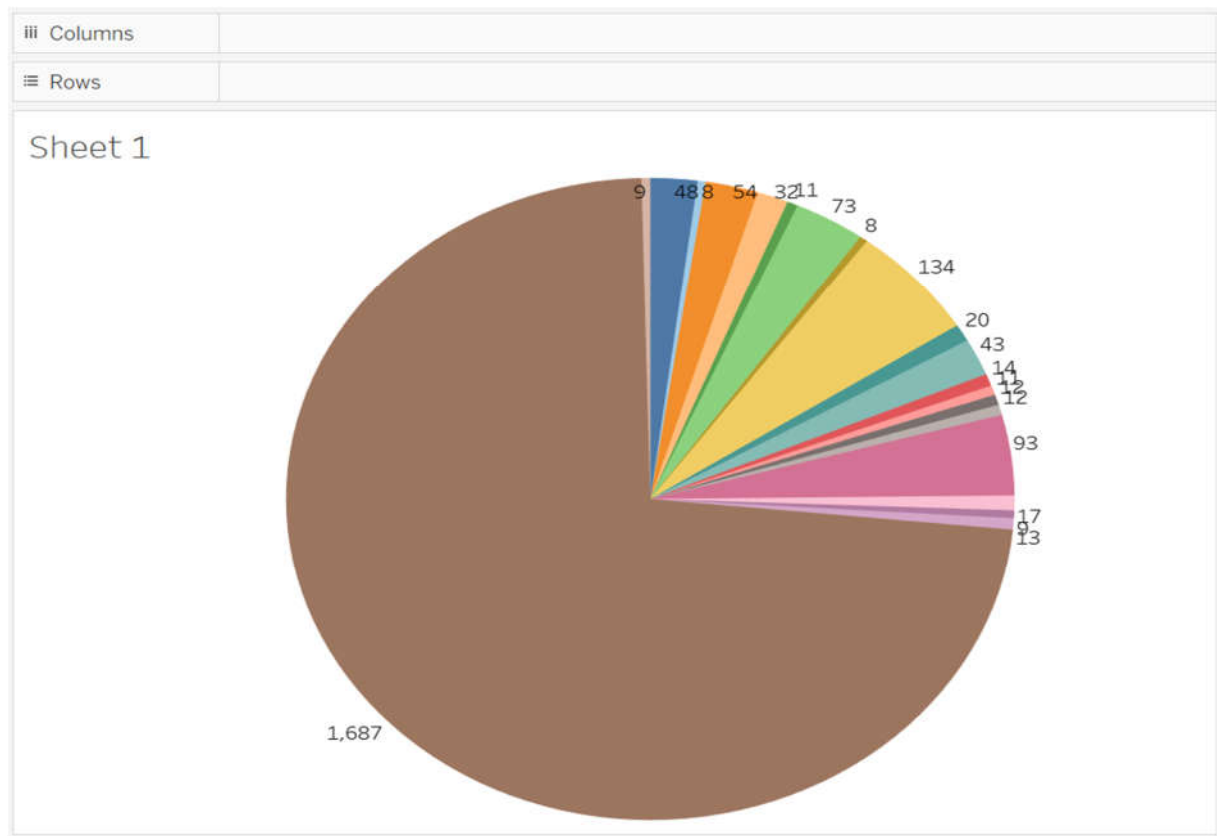
Event Log

re up-to-date (a minute ago)

2717:58 CRLF UTF-8 2 spaces

Type here to search

5:27 PM
12/2/2018



4. Query to fetch tweets to check on which day most of the tweets were made.

This query is written to find out the day on which more tweets were done so that it would count about giving us a figure.

Query-Code :


```

/*-----Query 4 : On which Day More Tweets are posted-----*/
case "4" =>
} val day_data = sqlContext.sql(sqlText = "SELECT substring(user.created_at,1,3) as day from tweets where text is not null")
day_data.createOrReplaceTempView(viewName = "day_data")

} val days_final = sqlContext.sql(
} """ SELECT Case
| when day LIKE '%Mon%' then 'WEEKDAY'
| when day LIKE '%Tue%' then 'WEEKDAY'
| when day LIKE '%Wed%' then 'WEEKDAY'
| when day LIKE '%Thu%' then 'WEEKDAY'
| when day LIKE '%Fri%' then 'WEEKDAY'
| when day LIKE '%Sat%' then 'WEEKEND'
| when day LIKE '%Sun%' then 'WEEKEND'
| else
| null
| end as day1 from day_data where day is not null""",stripMargin)

days_final.createOrReplaceTempView(viewName = "days_final")

val res = sqlContext.sql(sqlText = "SELECT day1 as Day,Count(*) as Day_Count from days_final where day1 is not null group by day1 order by count(*) desc")

println("*****")
println("On Which Day More Tweets Were Done")
println("*****")
} res.show()

```

Output :

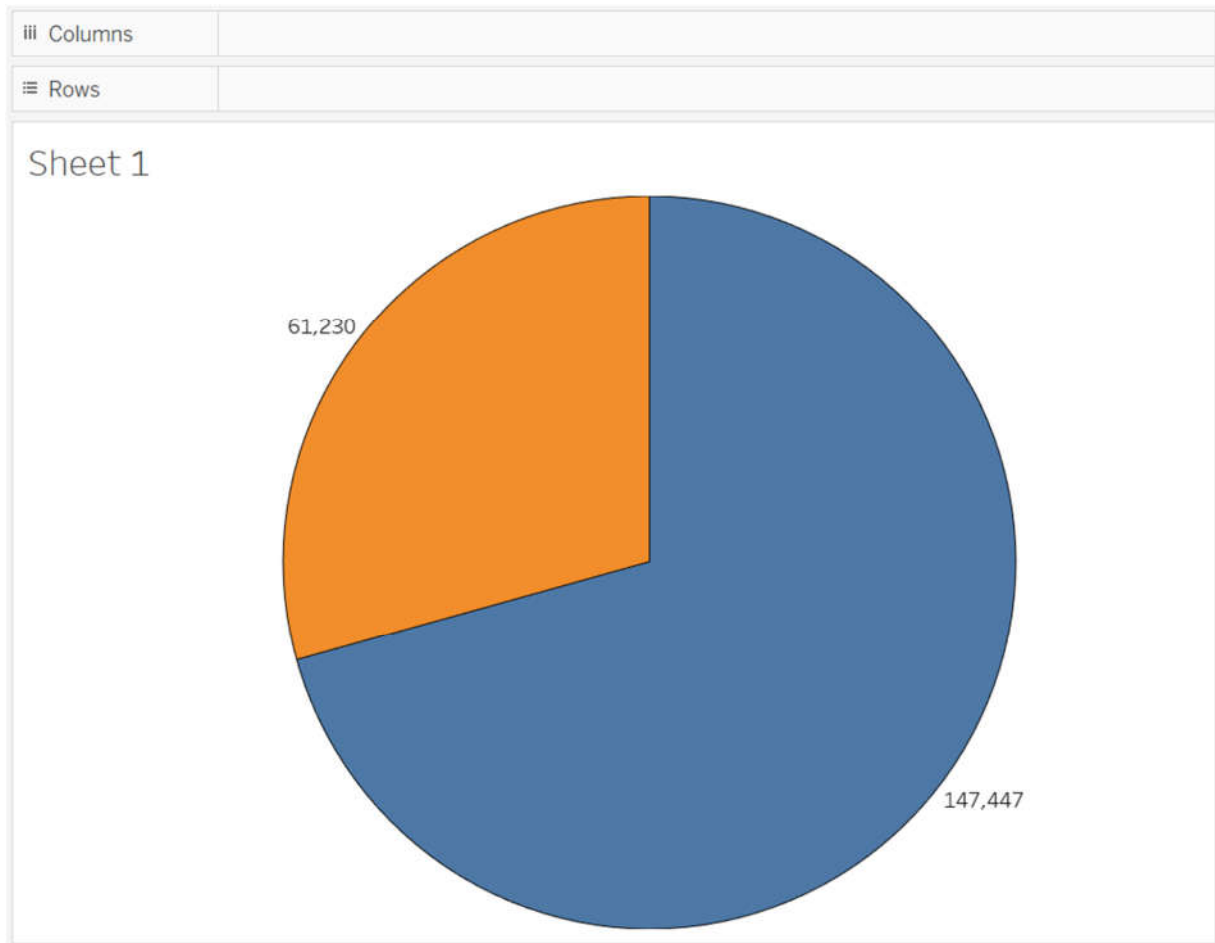
```

+-----+
| Day|Day_Count|
+-----+
|WEEKDAY| 147447|
|WEEKEND| 61230|
+-----+

18/12/02 17:28:36 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:28:36 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-RF770CO:4040
18/12/02 17:28:36 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:28:36 INFO MemoryStore: MemoryStore cleared
18/12/02 17:28:36 INFO BlockManager: BlockManager stopped
18/12/02 17:28:36 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:28:36 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:28:36 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:28:36 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:28:36 INFO ShutdownHookManager: Deleting directory C:\Users\Sindhu\thouthula\AppData\Local\Temp\spark-ecf4d02d-88b4-49e9-9492-6bc5317ea4a2

Process finished with exit code 0

```



5.Query to fetch the tweets for the various series of the phone/e-accessories.

This query is written to extract the tweets that made on the different series of the phone/e-accessories so that it would count for it.

Query-Code :

```

/*-----Query 5: Tweets count for different types of phone models -----*/
case "5" =>
val r1 = sqlContext.sql( sqlText = "SELECT loc,'Iphone X' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='Iphone X' " +
"group by loc order by count desc limit 10")
val r2 = sqlContext.sql( sqlText = "SELECT loc,'iphone7 Series' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='iphone7 Series' " +
"group by loc order by count desc limit 10")
val r3 = sqlContext.sql( sqlText = "SELECT loc,'iphone8 Series' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='iphone8 Series' " +
"group by loc order by count desc limit 10")
val r4 = sqlContext.sql( sqlText = "SELECT loc,'AirPods' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='AirPods' " +
"group by loc order by count desc limit 10")
val r5 = sqlContext.sql( sqlText = "SELECT loc,'TECHNOLOGY' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='TECHNOLOGY' " +
"group by loc order by count desc limit 10")
val r6 = sqlContext.sql( sqlText = "SELECT loc,'iOS' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='iOS' " +
"group by loc order by count desc limit 10")
val r7 = sqlContext.sql( sqlText = "SELECT loc,'Accessories' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='Accessories' " +
"group by loc order by count desc limit 10")
val r8 = sqlContext.sql( sqlText = "SELECT loc,'MAC' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='MAC' " +
"group by loc order by count desc limit 10")
val r9 = sqlContext.sql( sqlText = "SELECT loc,'Mobile' as phoneType,count(*) as count FROM disCat4 WHERE phoneType='Mobile' " +
"group by loc order by count desc limit 10")
val rdd1 = r1.union(r2).union(r3).union(r4).union(r5).union(r6).union(r7).union(r8).union(r9)
rdd1.createOrReplaceTempView( viewName = "rdd1")
val res = sqlContext.sql( sqlText = "SELECT phoneType, Count(*) as Count from rdd1 where phoneType is not null group by phoneType")
println("*****")
println("Model Type")
println("*****")
res.show()

```

Output :

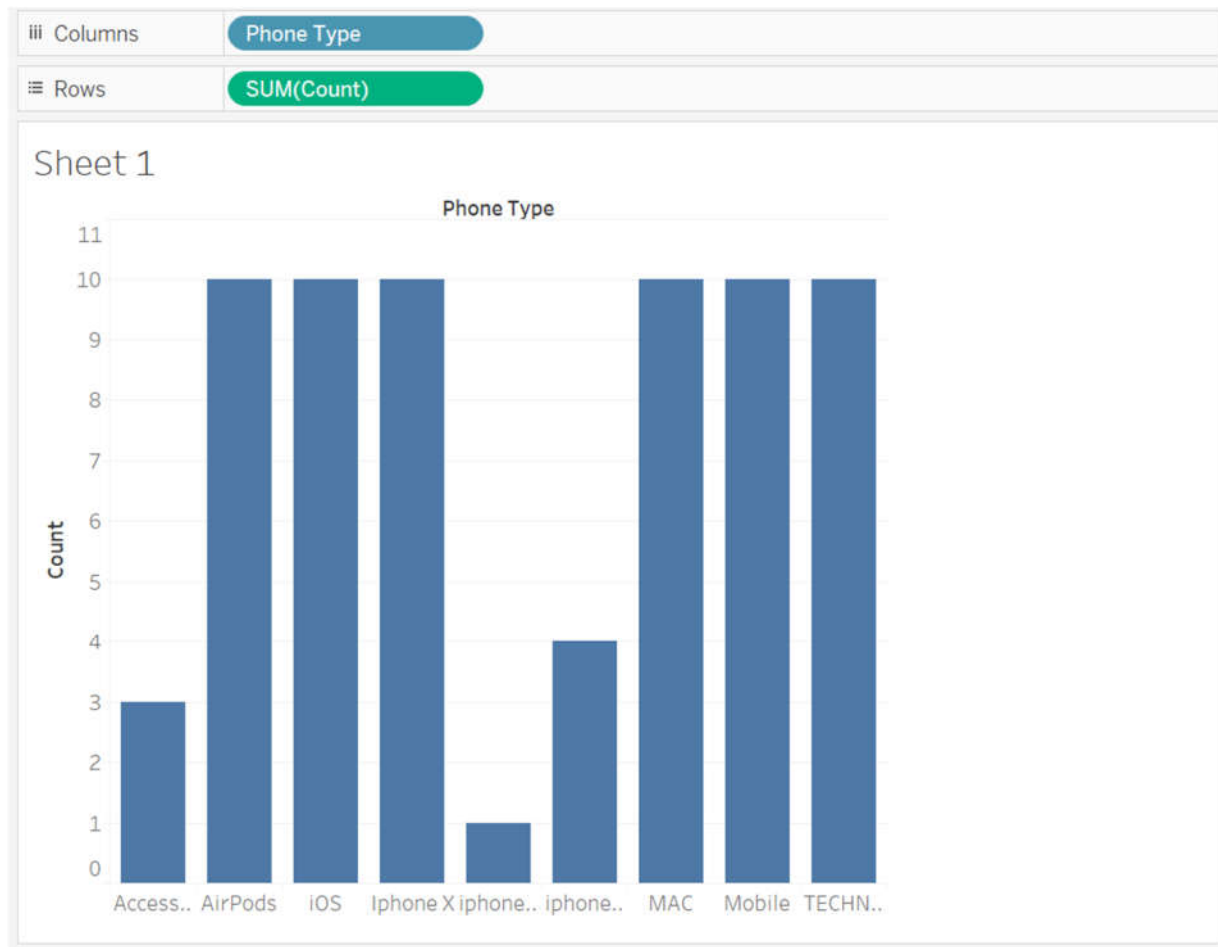
```

+-----+
| phoneType|Count|
+-----+
| iOS| 10|
| Iphone X| 10|
| AirPods| 10|
| Accessories| 3|
| MAC| 10|
| TECHNOLOGY| 10|
| Mobile| 10|
| iphone8 Series| 4|
| iphone7 Series| 1|
+-----+

18/12/02 17:31:14 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:31:14 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-RE770CO:4040
18/12/02 17:31:14 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:31:16 INFO MemoryStore: MemoryStore cleared
18/12/02 17:31:16 INFO BlockManager: BlockManager stopped
18/12/02 17:31:16 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:31:16 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:31:16 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:31:16 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:31:16 INFO ShutdownHookManager: Deleting directory C:\Users\Sindhu\thouthula\AppData\Local\Temp\spark-91623296-27d9-4328-a6f7-e08018788efe

Process finished with exit code 0

```



6. Query to fetch the languages mostly used for tweeting about the phone/e-accessories.

This query is written to analyze the language mostly used by the users so that it would count how many times the users tweeted about the phone/e-accessories in a particular language.

Query-Code :

```

/*-----Query 6 Popular languages used for tweeting tweets about phones -----*/
case "6" =>
} val langWstCount = sqlContext.sql( sqlText = "SELECT distinct id," +
  "CASE when user.lang LIKE '%en%' then 'English'" +
  "when user.lang LIKE '%ja%' then 'Japanese'" +
  "when user.lang LIKE '%es%' then 'Spanish'" +
  "when user.lang LIKE '%fr%' then 'French'" +
  "when user.lang LIKE '%it%' then 'Italian'" +
  "when user.lang LIKE '%ru%' then 'Russian'" +
  "when user.lang LIKE '%ar%' then 'Arabic'" +
  "when user.lang LIKE '%bn%' then 'Bengali'" +
  "when user.lang LIKE '%cs%' then 'Czech'" +
  "when user.lang LIKE '%da%' then 'Danish'" +

```

Executed Query Output :

language	Count
English	150413
Japanese	24997
Spanish	9803
Thai	6538
Portuguese	3993
Indonesian	2625
Korean	2141
Turkish	1901
French	1228
Arabic	982
Russian	839
Vietnamese	519
German	390
Italian	264
Chinese (Simplified)	212
Chinese (Traditio...	178
Dutch	176
Romanian	95
Polish	87
Hebrew	79

only showing top 20 rows

18/12/02 17:33:12 INFO SparkContext: Invoking stop() from shutdown hook
 18/12/02 17:33:12 INFO SparkUI: Stopped Spark web UI at <http://DESKTOP-8F770C0:4040>
 18/12/02 17:33:12 INFO MesOutputTrackerMasterEndpoint: MesOutputTrackerMasterEndpoint stopped!

to the previous occurrence

2685:1 CRLF UTF-8 2 spaces

Type here to search

5:33 PM
12/2/2018

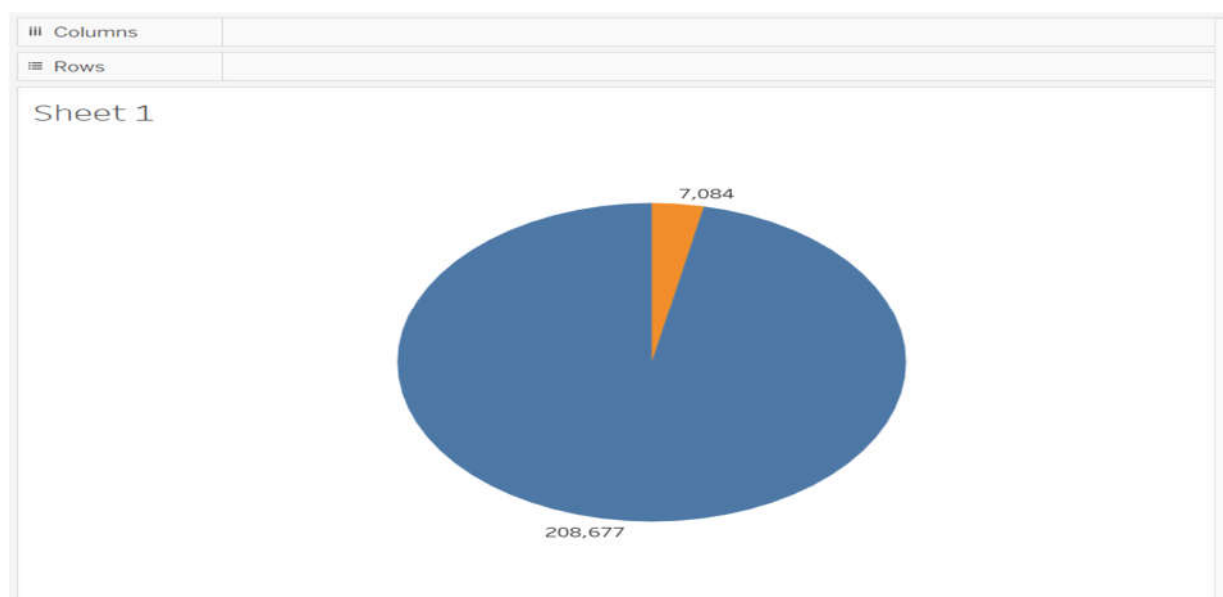

```

+-----+
|tweet_date|tweet_count|
+-----+
|Sun Sep 23|    208677|
|      null|      7084|
+-----+

18/12/02 17:34:31 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:34:31 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-RF770CO:4040
18/12/02 17:34:31 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:34:32 INFO MemoryStore: MemoryStore cleared
18/12/02 17:34:32 INFO BlockManager: BlockManager stopped
18/12/02 17:34:32 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:34:32 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:34:32 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:34:32 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:34:32 INFO ShutdownHookManager: Deleting directory C:\Users\Sindhu\thouthula\AppData\Local\Temp\spark-5455e683-72e5-4046-0798-467a136bf11c

Process finished with exit code 0

```



8. Query for fetching the tweets made from verified accounts.

This query is written to make an analysis on the number of verified users. This query counts for the tweets made from the verified accounts only and gives us the value.

Query-Code :

```

/*-----Query # Account Verification tweets -----*/
case "8" ->
val acctVerify = sqlContext.sql( sqlText = "SELECT distinct id, " +
"CASE when user.verified LIKE '%true%' THEN 'VERIFIED ACCOUNT'" +
"when user.verified LIKE '%false%' THEN 'NON-VERIFIED ACCOUNT'" +
"END AS Verified from tweets where text is not null")
acctVerify.createOrReplaceTempView( viewName = "acctVerify")
var acctVerifydata = sqlContext.sql( sqlText = "SELECT Verified, Count(Verified) as Count from acctVerify where id is NOT NULL and Verified is not null group
by Verified")
println("*****")
println("Account Verification")
println("*****")
acctVerifydata.show()

```

Output :

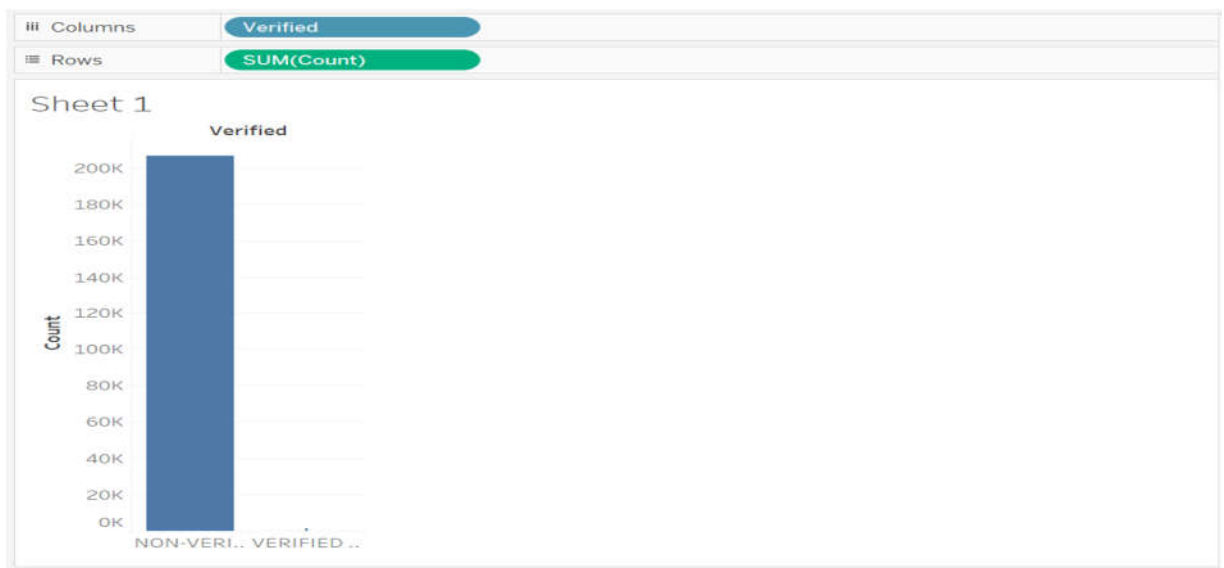
```

+-----+
| Verified | Count |
+-----+
| NON-VERIFIED ACCOUNT | 206624 |
| VERIFIED ACCOUNT | 1590 |
+-----+

18/12/02 17:36:02 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:36:02 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-BK770CO:4040
18/12/02 17:36:02 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:36:03 INFO MemoryStore: MemoryStore cleared
18/12/02 17:36:03 INFO BlockManager: BlockManager stopped
18/12/02 17:36:03 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:36:03 INFO OutputCommitCoordinatorOutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:36:03 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:36:03 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:36:03 INFO ShutdownHookManager: Deleting directory C:\Users\Gindhu\thouthula\AppData\Local\Temp\spark-Cc2af815-7cf9-4058-9dfe-did96a8b1413

Process finished with exit code 0

```



9. Query for fetching the tweets based on the hours on which most of them were made.

This query is written to analyze on which hours tweets are made like mornings, afternoon, evenings etc.

Query-Code :

```
/*-----Query 9 On Which hours More Tweets Were Done -----*/
case "9" =>
  val timehour = sqlContext.sql( sqlText = "SELECT SUBSTRING(created_at,12,2) as hour from tweets where text is not null")

  timehour.createOrReplaceTempView( viewName = "timehour")

  val timeAnalysis = sqlContext.sql(
    """ SELECT Case
        |when hour>=0 and hour <4 then 'midnight'
        |when hour>=4 and hour <7 then 'early Morning'
        |when hour>=7 and hour <12 then 'Morning'
        |when hour>=12 and hour <15 then 'afternoon'
        |when hour>=15 and hour <18 then 'evening'
        |when hour>=18 and hour <=23 then 'night'
        end as time from timehour""" .stripMargin)

  timeAnalysis.createOrReplaceTempView( viewName = "timeAnalysis")

  val res = sqlContext.sql( sqlText = "SELECT time as hour,Count(*) as tweets_count from timeAnalysis where time is not null group by time order by count(*)")


  println("*****")
  println("On Which hours More Tweets Were Done")
  println("*****")
  ...
}
```

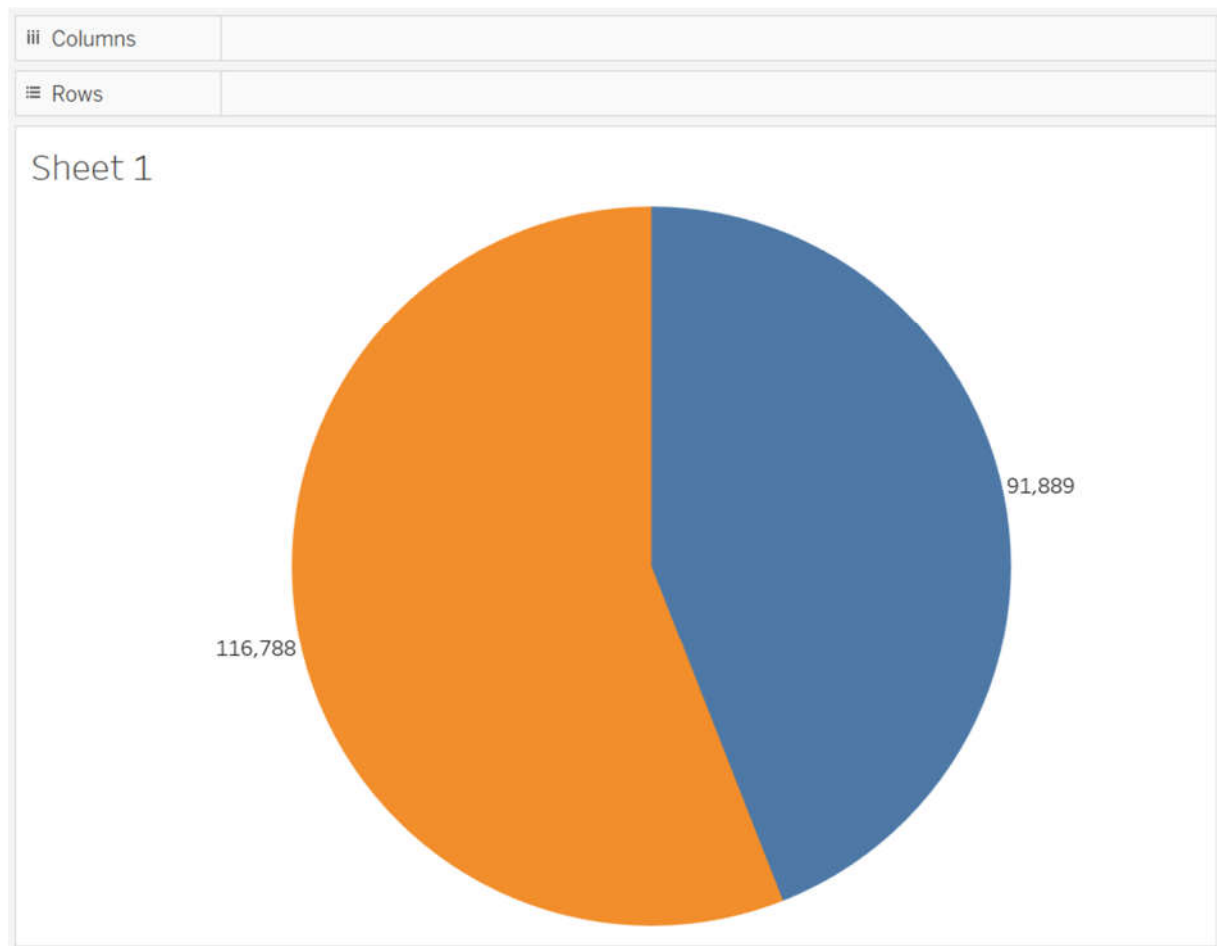
Output :

```
+-----+
|      hour|tweets_count|
+-----+
|  midnight|      116788|
|early Morning|      91889|
+-----+
```

```
18/12/02 17:37:18 INFO BlockManagerInfo: Removed broadcast_4_piece0 on DESKTOP-RF770CO:55163 in memory (size: 15.5 KB, free: 892.2 MB)
18/12/02 17:37:18 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:37:18 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-RF770CO:4040
18/12/02 17:37:18 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:37:18 INFO MemoryStore: MemoryStore cleared
18/12/02 17:37:18 INFO BlockManager: BlockManager stopped
18/12/02 17:37:18 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:37:18 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:37:18 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:37:18 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:37:18 INFO ShutdownHookManager: Deleting directory C:\Users\Sindhu\thouthula\AppData\Local\Temp\spark-212bdae2-d417-4685-81e9-67daa806a02e

Process finished with exit code 0
```





10.Query for fetching tweets based on which state has more tweets about particular type of phone/e-accessories.

This query is written to check which state users made most number of tweets about a type of phone/e-accessories.

Query-Code :


```

-----Query 10 Which state is mostly having tweets about type of phone -----
case "10" =>
val iPhoneRDD = sqlContext.sql( sqlText = """ SELECT 'iphone' as phoneType, user.location as loc from tweets where text LIKE '%$iphone%' """)
val iPhoneXRDD = sqlContext.sql( sqlText = """ SELECT 'iPhoneX' as phoneType, user.location as loc from tweets where text LIKE '%$iPhoneX%' """)
val iPhoneXsRDD = sqlContext.sql( sqlText = """ SELECT 'iPhoneXs' as phoneType, user.location as loc from tweets where text LIKE '%$iPhoneXs%' """)
val watchRDD = sqlContext.sql( sqlText = """ SELECT 'watch' as phoneType, user.location as loc from tweets where text LIKE '%$watch%' """)
//val BreakfastRDD = sqlContext.sql( """ SELECT 'Breakfast' as MealType, SUBSTRING(created_at,12,2) as hour, user.location as loc from dfs where text = 'Breakfast' """ )
//val BrunchRDD = sqlContext.sql( """ SELECT 'Brunch' as MealType, SUBSTRING(created_at,12,2) as hour, user.location as loc from dfs where text = 'Brunch' """ )
val sql2RDD = iPhoneRDD.union(iPhoneXRDD).union(iPhoneXsRDD).union(watchRDD)
sql2RDD.createOrReplaceTempView( viewName = "sql2RDD")
val locate = sqlContext.sql(
""" SELECT phoneType, loc from sql2RDD where
|loc LIKE '%Alaska%' OR loc LIKE '%Arizona%' OR loc LIKE '%Arkansas%' OR loc LIKE '%California%' OR loc LIKE '%Colorado%' OR loc LIKE '%
|OR loc LIKE '%Florida%'
|OR loc LIKE '%Georgia%'
|OR loc LIKE '%Hawaii%'
|OR loc LIKE '%Idaho%'
|OR loc LIKE '%Illinois%'
|OR loc LIKE '%Indiana%'

```

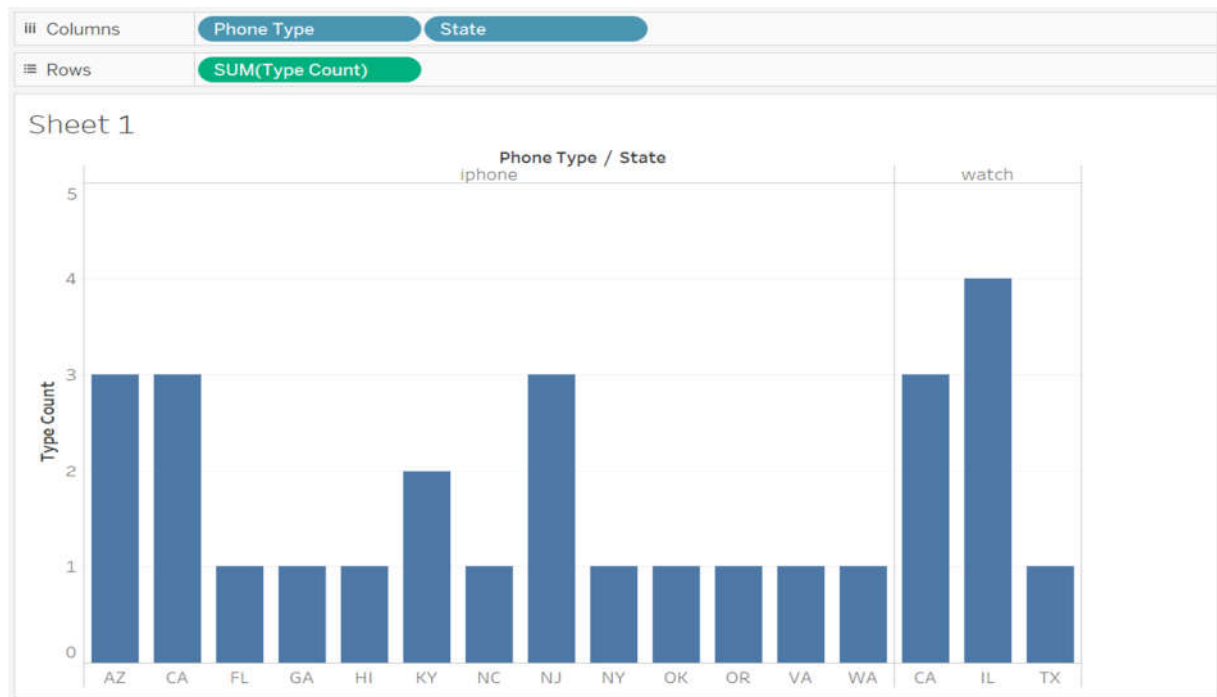
Output :

```

-----
|phoneType|state|type_count|
-----+-----+-----
|iphone|CA|3|
|watch|CA|3|
|iphone|FL|1|
|iphone|AZ|3|
|iphone|KY|2|
|iphone|VA|1|
|watch|TX|1|
|iphone|OR|1|
|watch|IL|4|
|iphone|HI|1|
|iphone|WA|1|
|iphone|OR|1|
|iphone|NC|1|
|iphone|NJ|3|
|iphone|GA|1|
|iphone|NY|1|
-----

18/12/02 17:38:51 INFO SparkContext: Invoking stop() from shutdown hook
18/12/02 17:38:51 INFO SparkUI: Stopped Spark web UI at http://DESKTOP-BF77000:4040
18/12/02 17:38:51 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
18/12/02 17:38:51 INFO MemoryStore: MemoryStore cleared
18/12/02 17:38:51 INFO BlockManager: BlockManager stopped
18/12/02 17:38:51 INFO BlockManagerMaster: BlockManagerMaster stopped
18/12/02 17:38:51 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
18/12/02 17:38:51 INFO SparkContext: Successfully stopped SparkContext
18/12/02 17:38:51 INFO ShutdownHookManager: Shutdown hook called
18/12/02 17:38:51 INFO ShutdownHookManager: Shutdown hook called

```



DOCKER :

Docker is a popular independent software container platform that allows you to build and ship your applications, along with all its environments, libraries and dependencies in containers



`docker` is configured to use the `default` machine with IP `192.168.99.100`
For help getting started, check out the docs at <https://docs.docker.com>

Start interactive shell

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$ docker pull chkkrish9/msp-pb-phase3:2
```

```
2: Pulling from chkkrish9/msp-pb-phase3
```

```
Digest: sha256:7ac07c85e32407ab52de4f42bdb405aaadb05e8cef94fe1755b989d42e5f599c
```

```
Status: Image is up to date for chkkrish9/msp-pb-phase3:2
```

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$
```



`docker` is configured to use the `default` machine with IP `192.168.99.100`
For help getting started, check out the docs at <https://docs.docker.com>

Start interactive shell

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$ docker pull chkrish9/msp-pb-phase3:2
```

```
2: Pulling from chkrish9/msp-pb-phase3
```

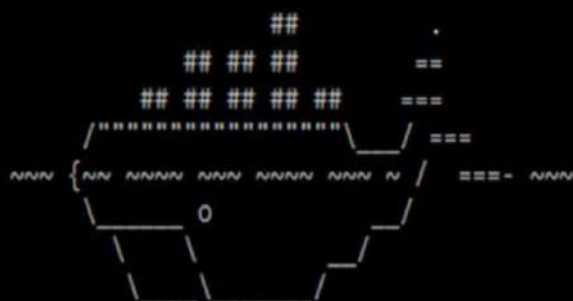
```
Digest: sha256:7ac07c85e32407ab52de4f42bdb405aaadb05e8cef94fe1755b989d42e5f599c
```

```
Status: Image is up to date for chkrish9/msp-pb-phase3:2
```

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$ docker run -it chkrish9/msp-pb-phase3:2 bash
```

```
root@863486393d93:~#
```



`docker` is configured to use the `default` machine with IP `192.168.99.100`
For help getting started, check out the docs at <https://docs.docker.com>

Start interactive shell

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$ docker pull chkkrish9/msp-pb-phase3:2
```

```
2: Pulling from chkkrish9/msp-pb-phase3
```

```
Digest: sha256:7ac07c85e32407ab52de4f42bdb405aaadb05e8cef94fe1755b989d42e5f599c
```

```
Status: Image is up to date for chkkrish9/msp-pb-phase3:2
```

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$ docker run -it chkkrish9/msp-pb-phase3:2 bash
```

```
root@c3e0aa440030:~# ls
```

```
data.json derby.log metastore_db
```

```
root@c3e0aa440030:~#
```



```
Start interactive shell
```

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$ docker pull chkkrish9/msp-pb-phase3:2
```

```
2: Pulling from chkriash9/msp-pb-phase3
```

Digest: sha256:7ac07c85e32407ab52de4f42bdb405aaadb05e8cef94fe1755b989d42e5f599c

```
Status: Image is up to date for chkrish9/msp-pb-phase3:2
```

```
Kite@DESKTOP-D128J0A MINGW64 /c/Program Files/Docker Toolbox
```

```
$ docker run -it chkrish9/msp-pb-phase3:2 bash
```

```
root@c6097b0b07ff:~# ls
```

```
data.json  derby.log  metastore.db
```

```
root@c6097b0b07ff:~# spark-shell
```

Spark context Web UI available at <http://172.17.0.2:4040>

```
Spark context available as 'sc' (master = local[*], app id = local-1558120102150).
```

```
Spark session available as 'spark'.
```

Welcome to

[illegible]

```
Using Scala version 2.11.8 (OpenJDK 64-Bit Server VM, Java 1.8.0_131)
```

Type in expressions to have them evaluated.

```
Type :help for more information.
```

```
scala>
```

```
data.json derby.log metastore_db
root@c6097b0b07ff:~# spark-shell
Spark context Web UI available at http://172.17.0.2:4040
Spark context available as 'sc' (master = local[*], app id = local-1558120102150).
Spark session available as 'spark'.
Welcome to
```



```
Using Scala version 2.11.8 (OpenJDK 64-Bit Server VM, Java 1.8.0_131)
Type in expressions to have them evaluated.
Type :help for more information.
```

```
scala> val df = spark.read.json("data.json")
df: org.apache.spark.sql.DataFrame = [_corrupt_record: string, contributors: string ... 36 more fields]
```

```
scala> import org.apache.spark.sql.SQLContext
import org.apache.spark.sql.SQLContext
```

```
scala> val sqlContext = new SQLContext(sc)
warning: there was one deprecation warning; re-run with -deprecation for details
sqlContext: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@1e08f0cd
```

```
scala> df.registerTempTable("DataTable")
warning: there was one deprecation warning; re-run with -deprecation for details
```

```
scala>
```

PROJECT MANAGEMENT :

WORK COMPLETED :

Pravalhika Kampally :

Collected Data and executed first 3 quires.

Joshmitha Tammareddy :

Executes 4 to 7 Quires and worked for Docker

Rupesh Sai Ram Doddala :

Executed 7 to 10 quires and worked for visualization.