# PRACTICAL-11

## AIM:

To perform Sentiment Analysis using Twitter data, Scala and Spark.

## CODE:

### Reading data:

```python
from pyspark import SparkConf, SparkContext
import sys

conf = SparkConf().setMaster("local").setAppName("RatingsHistogram")
sc = SparkContext(conf = conf)

# Create a hardcoded RDD
numbers = sc.parallelize([1, 2, 3, 4])

# Load plain text files (e.g. CSVs) from different sources
hdfs_lines = sc.textFile("hdfs:///user/cloudera/ml-100k/u.data", minPartitions=1)
local_lines = sc.textFile("file:///home/alex/ml-100k/u.data")
s3_lines = sc.textFile("s3n://bucket/ml-100k/u.data")
CLI_arg_lines = sc.textFile(sys.argv[1])

# Create RDDs from an existing Hive repository
hive_ctx = HiveContext(sc)
hive_lines = hive_ctx.sql("SELECT name, age FROM users WHERE age > 18")
```

Notice that this can not be run with the standard Python interpreter. Instead, you use the spark-submit to submit it as a batch task, or name the Shell Pyspark

Writing of data:

● The RDD class has the method saveAsTextFile. However, this saves the representation of each variable in a string. In Python, the resulting text file will contain lines such as (1949, 111).

 ● If you want to save the data in CSV or TSV format, you can either use Python's StringIO and csv modules or simply map each variable (vector) into a single string, e.g. as follows:

```python
res.saveAsTextFile("hdfs:///user/cloudera/res_raw.txt")  # bad format

res.map(lambda row: str(row[0]) + "\t" + str(row[1])) \
   .saveAsTextFile("hdfs:///user/cloudera/res_tsv.txt")  # good format
```

**Sentiment Analysis using Scala:**

**Spark Streaming Implementation:**

```scala
//Import the necessary packages into the Spark Program
import org.apache.spark.streaming.{Seconds, StreamingContext}
import org.apache.spark.SparkContext._
...
import java.io.File

object twitterSentiment {

def main(args: Array[String]) {
if (args.length < 4) {
System.err.println("Usage: TwitterPopularTags <consumer key> <consumer secret>
System.exit(1)
}

StreamingExamples.setStreamingLogLevels()
//Passing our Twitter keys and tokens as arguments for authorization
val Array(consumerKey, consumerSecret, accessToken, accessTokenSecret) = args.
val filters = args.takeRight(args.length - 4)

// Set the system properties so that Twitter4j library used by twitter stream
// Use them to generate OAuth credentials
System.setProperty("twitter4j.oauth.consumerKey", consumerKey)
...
System.setProperty("twitter4j.oauth.accessTokenSecret", accessTokenSecret)

val sparkConf = new SparkConf().setAppName("twitterSentiment").setMaster("loca
val ssc = new Streaming Context
val stream = TwitterUtils.createStream(ssc, None, filters)

//Input DStream transformation using flatMap
val tags = stream.flatMap { status => Get Text From The Hashtags }

//RDD transformation using sortBy and then map function
tags.countByValue()
.foreachRDD { rdd =>
val now = Get current time of each Tweet
rdd
.sortBy(_._2)
.map(x => (x, now))
//Saving our output at ~/twitter/ directory
.saveAsTextFile(s"~/twitter/$now")
}

data.print()
//Saving our output at ~/ with filenames starting like twitters
data.saveAsTextFiles("~/twitters","20000")

ssc.start()
ssc.awaitTermination()
 }
}
```
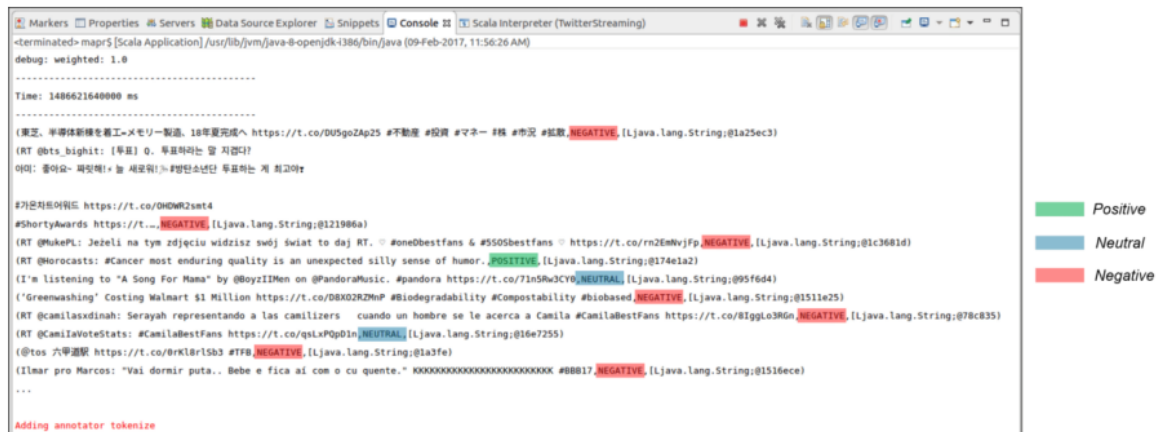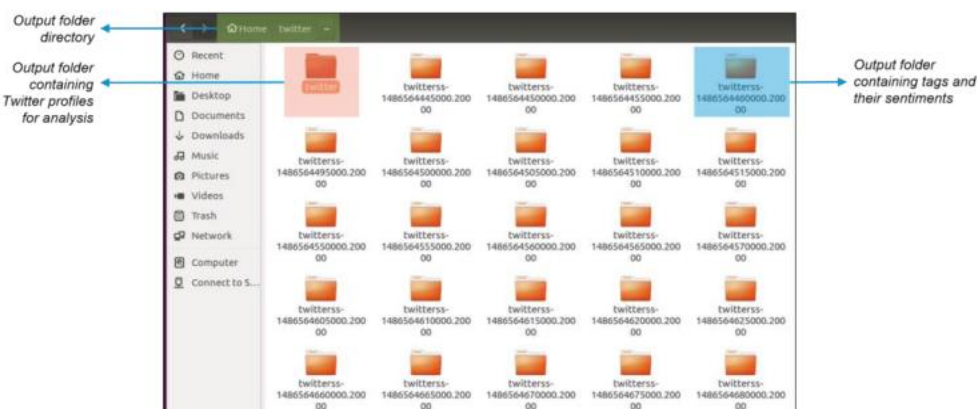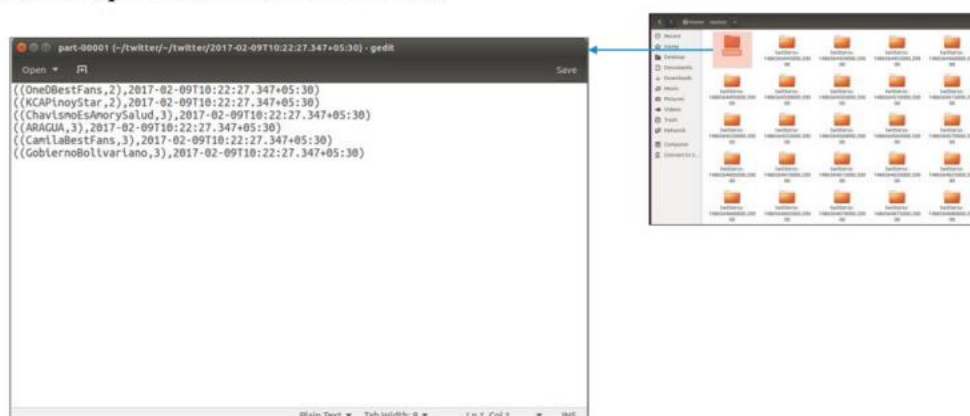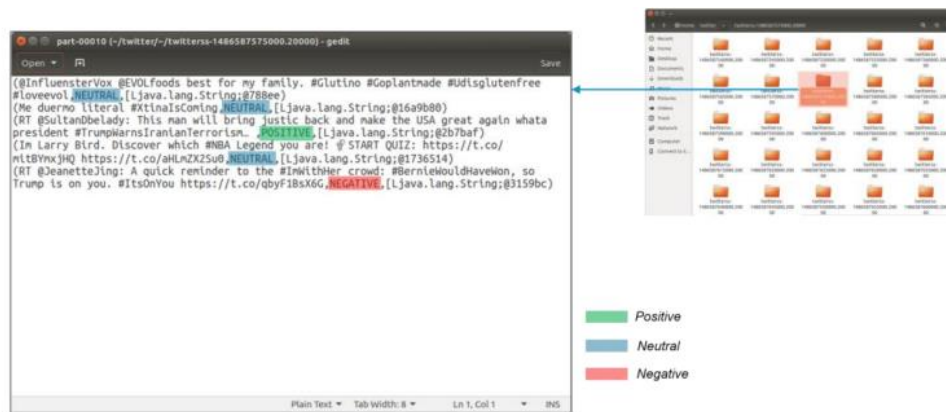
## OUTPUT:



- As we can see from the screenshot, all tweets are categorised as Positive, Neutral and Negative according to the feelings of the contents of the tweets.
- The output of the Sentiments of the Tweets is stored in folders and files, depending on the time they were created. This output can be stored on the local file system or HDFS as appropriate. The output directory looks like this one:
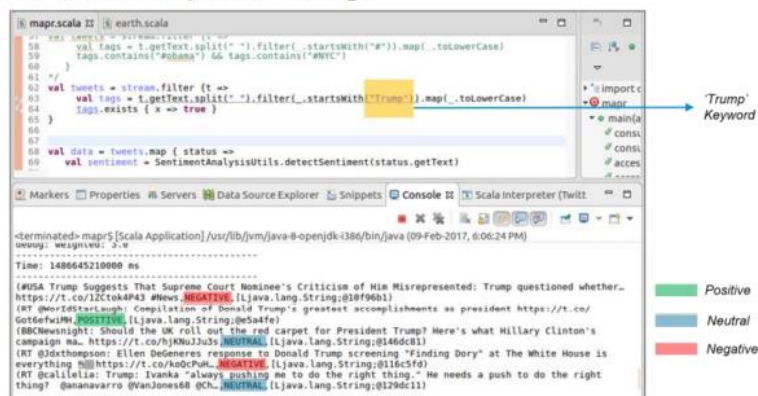


- Here, within the twitter folders, we will find the usernames of the Twitter users and the timestamp for. tweet as seen below:



- Now that we have Twitter usernames and timestamps, let 's look at the Sentiments and Tweets saved in the main directory. Here, every tweet is accompanied by a feeling of emotion. This Sentiment, which is processed, is further used to examine a broad variety of business observations.

- **Tweaking Code:** Now, let's modify our code a little to get feelings about specific hashtags (topics). Donald Trump, the President of the United States, is now travelling through television outlets and online social media. Let us look at the emotions associated with the keyword 'Trump.'



## CONCLUSION:

In this practical, we learnt about spark and installed it and configured it. We explored the spark shell as well.