**PRINCIPLES OF BIG DATA MANAGEMENT**

**PROJECT PHASE 1**

**Submitted by,**

**Sneha Lagandula**

**Abhiram Ampabathina**

**Dinesh Reddy Paduru**

**(Team-5)**

**Code for Tweet Collection:**

```python
from tweepy.streaming import StreamListener
from tweepy import OAuthHandler
from tweepy import Stream
import json, time, sys , tweepy
import re
import unicodedata
import codecs
#customer keys
consumer_key = 'yGRxlXyDG39iHr1akrN0zRU9F'
consumer_secret = '4LiY6VeolUtPPaV9dmEsh6jJvvXCAcGpHXIrbQZaxHpcc0ANMt'
access_key = '3700831395-31MrA3aG8vPRWlKzA2eu8dONZ7GhGfso2AuUe9g'
access_secret = 'ft3LSqsjk8hEe7AABYG5uga9320mLoOTZIA0yPGI6iwY1'
#auth keys
auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_key, access_secret)
api = tweepy.API(auth)
#stream listener
class StdOutListener(StreamListener):
    def on_status(self, status):
        twit_id = status.id
        text = status.text
        created = status.created_at
        time_zone = status.user.time_zone
        location = status.user.location
        language = status.lang
        hashtags = re.findall(r"#(\w+)" , text)


     for tag in hashtags:
```
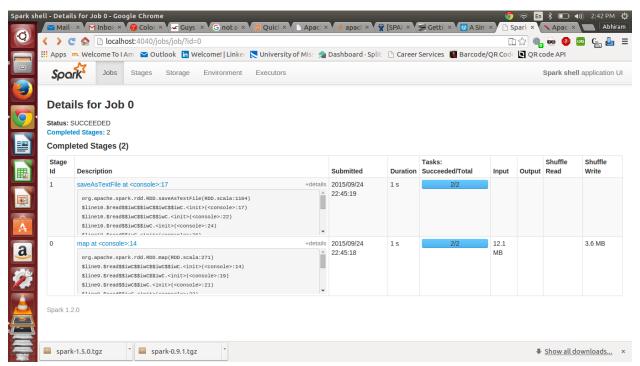
```python
        record = unicode(twit_id) + '\t' +  unicode(created) + '\t' + unicode(language) + '\t' +
        unicode(location)  + '\t' + unicode(time_zone) + '\t' +  unicode(tag) + '\t' + unicode(text) +
        '\t'

        tweetfile = codecs.open('sampletweets4.txt' , 'a' , 'utf8')

         tweetfile.write(unicode(record))

        #new line for every tweet

         tweetfile.write('\n')

         #close tweet file

         tweetfile.close()


         return True


    def on_error(self, status):

        print 'Error on status', status


    def on_limit(self, status):

        print 'Limit threshold exceeded', status


    def on_timeout(self, status):

        print 'Stream disconnected; continuing...'


stream = Stream(auth, StdOutListener())

stream.filter(track = ['#NFL', '#football', '#baseball', '#royals', '#seattlemariners', '@tamba', '#Orioles',
'#Yankees', '#marlins', '#redsox', '#MLB', '@StarSports', '#FMRedHawks', '@USABaseball',
'@ForrestWhitley', '@MickeyMoniak', '@colestobbe', '#USAbaseball', '@baseball', 'Bryce Harper', 'Albert
Pujols', 'David Ortiz', 'Alex Rodriguez', 'Derek Jeter', 'Clayton Kershaw', 'Andrew McCutchen', 'Yasiel
Puig', 'Miguel Cabrera', 'Mike Trout', 'Johnny Bench', '#sports', '#kcroyals', 'pirates', 'rockies', 'sports',
'#indians', '#twins', '#gaints', '#padres', '#marlins', '#phillies', '#royals', '#mariners', '#cardinals', '#brewers',
'#mets', '#reds', '#pirates', '#rockies', '#rangers', '#athletics', '#whitesox', '#yankees', '#orioles', '#nationals',
'#dodgers', '#diamondbacks', '#rays', '#redsox', '#mariners', '#angels', '#LaMarcus Aldridge', '#Carmelo
Anthony', '#Harrison Barnes', '#Boston Celtics', '#brooklyn nets', '#houston rockets', '#new york knicks'])
```
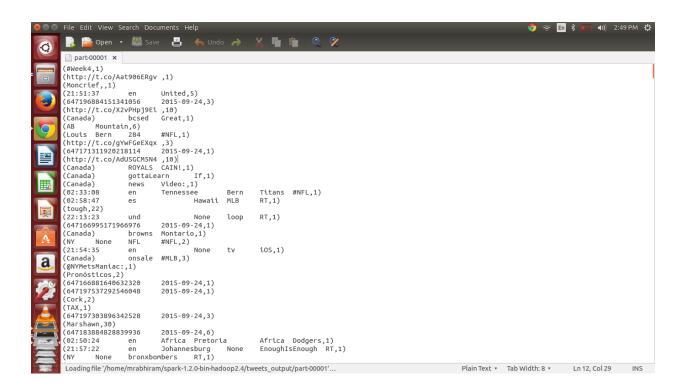
**Code for Spark Word Count:**

```
val textFile = spark.textFile("sample_tweets.txt")
val counts = textFile.flatMap(line => line.split(" "))
          .map(word => (word, 1))
          .reduceByKey(_ + _)
counts.saveAsTextFile("tweets_output")
```

**Spark Job:**

**Word Count Output:**

**Word Count Execution:**

```
root@mrabhiram-ThinkPad-Edge-E530: /home/mrabhiram/spark-1.2.0-bin-hadoop2.4

scala> val textFile = sc.textFile("sampletweets.txt")
15/09/24 22:41:49 INFO MemoryStore: ensureFreeSpace(163705) called with curMem=0
, maxMem=278302556
15/09/24 22:41:49 INFO MemoryStore: Block broadcast_0 stored as values in memory
 (estimated size 159.9 KB, free 265.3 MB)
15/09/24 22:41:49 INFO MemoryStore: ensureFreeSpace(22692) called with curMem=16
3705, maxMem=278302556
15/09/24 22:41:49 INFO MemoryStore: Block broadcast_0_piece0 stored as bytes in
memory (estimated size 22.2 KB, free 265.2 MB)
15/09/24 22:41:49 INFO BlockManagerInfo: Added broadcast_0_piece0 in memory on l
ocalhost:42889 (size: 22.2 KB, free: 265.4 MB)
15/09/24 22:41:49 INFO BlockManagerMaster: Updated info of block broadcast_0_pie
ce0
15/09/24 22:41:49 INFO SparkContext: Created broadcast 0 from textFile at <conso
le>:12
textFile: org.apache.spark.rdd.RDD[String] = sampletweets.txt MappedRDD[1] at te
xtFile at <console>:12

scala> counts = text_file.flatMap(lambda line: line.split(" ")) \.map(lambda wor
d: (word, 1)) \.reduceByKey(lambda a, b: a + b)
<console>:1: error: ')' expected but '(' found.
       counts = text_file.flatMap(lambda line: line.split(" ")) \.map(lambda wor
d: (word, 1)) \.reduceByKey(lambda a, b: a + b)
                                                                    ^
```

```
root@mrabhiram-ThinkPad-Edge-E530: /home/mrabhiram/spark-1.2.0-bin-hadoop2.4

textFile: org.apache.spark.rdd.RDD[String] = sampletweets.txt MappedRDD[1] at te
xtFile at <console>:12

scala> counts = text_file.flatMap(lambda line: line.split(" ")) \.map(lambda wor
d: (word, 1)) \.reduceByKey(lambda a, b: a + b)
<console>:1: error: ')' expected but '(' found.
       counts = text_file.flatMap(lambda line: line.split(" ")) \.map(lambda wor
d: (word, 1)) \.reduceByKey(lambda a, b: a + b)
                                               ^
<console>:1: error: ';' expected but ')' found.
       counts = text_file.flatMap(lambda line: line.split(" ")) \.map(lambda wor
d: (word, 1)) \.reduceByKey(lambda a, b: a + b)
                                                                    ^
scala> val counts = textFile.flatMap(line => line.split(" ")).map(word => (word,
 1)).reduceByKey(_ + _)
15/09/24 22:44:15 INFO FileInputFormat: Total input paths to process : 1
counts: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[4] at reduceByKey
at <console>:14

scala> counts.saveAsTextFile("tweets_output")
15/09/24 22:45:17 INFO deprecation: mapred.tip.id is deprecated. Instead, use ma
preduce.task.id
15/09/24 22:45:17 INFO deprecation: mapred.task.id is deprecated. Instead, use m
```

```
scala> counts.saveAsTextFile("tweets_output")
15/09/24 22:45:17 INFO deprecation: mapred.tip.id is deprecated. Instead, use ma
preduce.task.id
15/09/24 22:45:17 INFO deprecation: mapred.task.id is deprecated. Instead, use m
apreduce.task.attempt.id
15/09/24 22:45:17 INFO deprecation: mapred.task.is.map is deprecated. Instead, u
se mapreduce.task.ismap
15/09/24 22:45:17 INFO deprecation: mapred.task.partition is deprecated. Instead
, use mapreduce.task.partition
15/09/24 22:45:17 INFO deprecation: mapred.job.id is deprecated. Instead, use ma
preduce.job.id
15/09/24 22:45:18 INFO SparkContext: Starting job: saveAsTextFile at <console>:1
7
15/09/24 22:45:18 INFO DAGScheduler: Registering RDD 3 (map at <console>:14)
15/09/24 22:45:18 INFO DAGScheduler: Got job 0 (saveAsTextFile at <console>:17)
with 2 output partitions (allowLocal=false)
15/09/24 22:45:18 INFO DAGScheduler: Final stage: Stage 1(saveAsTextFile at <con
sole>:17)
15/09/24 22:45:18 INFO DAGScheduler: Parents of final stage: List(Stage 0)
15/09/24 22:45:18 INFO DAGScheduler: Missing parents: List(Stage 0)
15/09/24 22:45:18 INFO DAGScheduler: Submitting Stage 0 (MappedRDD[3] at map at
<console>:14), which has no missing parents
15/09/24 22:45:18 INFO MemoryStore: ensureFreeSpace(3544) called with curMem=186
```