

Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

Big Data Processing- Coursework1 – Twitter Analysis

PART A - CONTENT ANALYSIS

TWEET MAPPER CODE:

```
1 import java.io.IOException;
2 import java.util.StringTokenizer;
3 import org.apache.hadoop.io.IntWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Mapper;
6 import org.apache.commons.lang.StringUtils;
7
8 public class TwitterMapper extends Mapper<Object, Text, IntWritable, IntWritable> {
9     //constant variable stores a value of one
10    private final IntWritable one = new IntWritable(1);
11    //mapper method takes takes as input the text file
12    public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
13        // Format per tweet is id;date;hashtags;tweet;
14
15        String dump = value.toString();
16        if(StringUtils.ordinalIndexOf(dump,";",4)>-1){
17            int startIndex = StringUtils.ordinalIndexOf(dump,";",3) + 1;
18            String tweet = dump.substring(startIndex,dump.lastIndexOf(';'));
19            //store tweet length in a variable
20            IntWritable tweetSize=new IntWritable(tweet.length());
21            //send tweet length and its occurrence to reducer
22            context.write(tweetSize,one);
23
24            }//END if statement
25
26    }//END map method
27
28 }//END TwitterMapper
```

TWEET REDUCER CODE:

```
import java.io.IOException;
import java.util.Iterator;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class TweetReducer extends Reducer<IntWritable, IntWritable, IntWritable, IntWritable> {

    private IntWritable result = new IntWritable();
    //reducer method takes key from mapper and value from mapper as input
    public void reduce(IntWritable key, Iterable<IntWritable> occurrences, Context context)
        throws IOException, InterruptedException {

        int sum = 0;
        //each value is iterated over to sum up each occurrence of X length
        for (IntWritable occurrence : occurrences) {

            sum+=occurrence.get();

        }

        //variable stores the above sum of occurrences of X length
        result.set(sum);
        //the results i.e: key: length and total occurrences of length X, Y, Z... are emitted
        context.write(key,result);

    }

}
```

Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

TWEET COUNT CODE:

```
public class TweetCount {  
    public static void runJob(String[] input, String output) throws Exception {  
        Configuration conf = new Configuration();  
  
        Job job = new Job(conf);  
        //sets current class as the class which generates and enables executability of a jar file  
        job.setJarByClass(TweetCount.class);  
        //mapper class is set  
        job.setMapperClass(TwitterMapper.class);  
        //reducer class is set  
        job.setReducerClass(TweetReducer.class);  
        //sets output key's class type  
        job.setMapOutputKeyClass(IntWritable.class);  
        //sets output value's class type  
        job.setMapOutputValueClass(IntWritable.class);  
        //produces three tasks which output three text files with its corresponding results  
        job.setNumReduceTasks(3);  
        Path outputPath = new Path(output);  
        FileInputFormat.setInputPaths(job, StringUtils.join(input, ","));  
        FileOutputFormat.setOutputPath(job, outputPath);  
        outputPath.getFileSystem(conf).delete(outputPath, true);  
        job.waitForCompletion(true);  
    }  
  
    public static void main(String[] args) throws Exception {  
        runJob(Arrays.copyOfRange(args, 0, args.length-1), args[args.length-1]);  
    }  
}
```

EXPLANATION:

The Map Reduce approach with the first part of **(A) Content Analysis** is where a counter is created to store 1 as a value which keeps track of a particular tweet length within the mapper and thus is sent (emitted) to the reducer to total up and count how many times a certain tweet length occurred within the data set (data/olympictweets). Henceforth, a map method's would store the tweet length of every single tweet in a variable and emits it directly to the reducer along with its corresponding occurrence within the data set i.e.: length=300, occurrence=200. Which means for length 300 it would occur in the data set 200x.

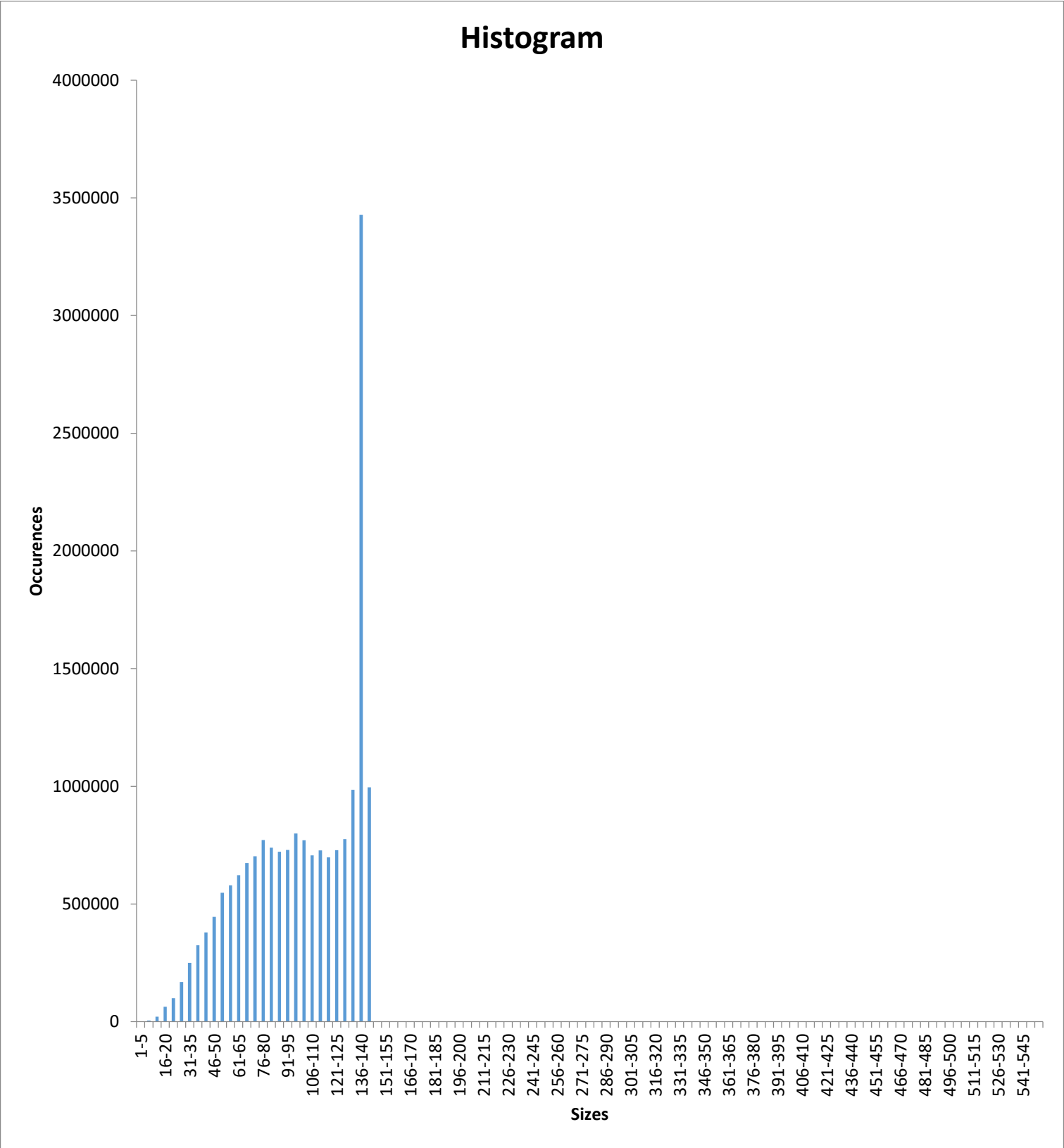
In addition, a reducer method would take as input the key being a length and its occurrence being a value. It would contain an iteration which iterates through each occurrence and sums up the whole occurrence total for X length:

i.e.:

length=12, occurrence=1

length=12, occurrence=1

As a result, length=12, would result in two occurrences; this total of occurrences is set as an output value.



Furthermore, this enabled me to produce the histogram displayed above in Excel, which clearly denotes that the range of lengths 136-140 had the most lengths and that certain ranges had no lengths reported within in the dataset at all.

Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

(ii) AVERAGE LENGTH

TWEET AVERAGE MAPPER CODE:

```
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Mapper;
6 import org.apache.commons.lang.StringUtils;
7
8 public class TweetAverageMapper extends Mapper<Object, Text, Text, IntWritable> {
9     //constant variable stores a value of one
10     private final IntWritable one = new IntWritable(1);
11     //create a text variable
12     Text tweetStr=new Text("");
13     //mapper method takes takes as input the text file
14     public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
15         // Format per tweet is id;date;hashtags;tweet;
16
17         String dump = value.toString();
18         if(StringUtils.ordinalIndexOf(dump,";",4)>-1){
19             int startIndex = StringUtils.ordinalIndexOf(dump,";",3) + 1;
20             String tweet = dump.substring(startIndex,dump.lastIndexOf(';'));
21             //use tweet as a key
22             String tweettweet="tweet";
23             //store tweet length in a variable
24             IntWritable tweetLengths=new IntWritable(tweet.length());
25             //text variable stores a string which will be used as a key
26             tweetStr.set(tweettweet);
27             //send text variable and its contents as a key and tweet length as a value to reducer
28             context.write(tweetStr,tweetLengths);
29
30             }//END if statement
31
32     }//END map method
33
34 }//END TweetAverageMapper
35
```

TWEET AVERAGE REDUCER CODE:

```
private IntWritable result = new IntWritable();

//reducer method takes key from mapper and value from mapper as input
public void reduce(Text key, Iterable<IntWritable> tweetLengths, Context context)
    throws IOException, InterruptedException {
    //create a variable to store the count
    int count =0;
    //create a variable to store all the lengths added together
    int totalTweetSize = 0;
    //iterate through the sum of all tweet lengths and store that sum in a variable and count how many lengths exist along
    for (IntWritable tweetLength : tweetLengths) {
        totalTweetSize+=tweetLength.get();
        //increment the counter variable
        count++;
    }
    //store a average of all the lengths totalled up divided by how many lengths which exist in the data set
    int avg=totalTweetSize/count;

    //variable stores the above average calculation of each length
    result.set(avg);

    //the results i.e: key: "tweet" and the average of all lengths are emitted
    context.write(key,result);
}
```

Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

TWEET AVERAGE COUNT CODE:

```
public class TweetAvgCount {  
    public static void runJob(String[] input, String output) throws Exception {  
        Configuration conf = new Configuration();  
  
        Job job = new Job(conf);  
        //sets current class as the class which generates and enables executability of a jar file  
        job.setJarByClass(TweetAvgCount.class);  
        //mapper class is set  
        job.setMapperClass(TweetAverageMapper.class);  
        //reducer class is set  
        job.setReducerClass(TweetAverageReducer.class);  
        //sets output key's class type  
        job.setMapOutputKeyClass(Text.class);  
        //sets output value's class type  
        job.setMapOutputValueClass(IntWritable.class);  
        //produces three tasks which output three text files with its corresponding results  
        job.setNumReduceTasks(3);  
        Path outputPath = new Path(output);  
        FileInputFormat.setInputPaths(job, StringUtils.join(input, ","));  
        FileOutputFormat.setOutputPath(job, outputPath);  
        outputPath.getFileSystem(conf).delete(outputPath, true);  
        job.waitForCompletion(true);  
    }  
  
    public static void main(String[] args) throws Exception {  
        runJob(Arrays.copyOfRange(args, 0, args.length-1), args[args.length-1]);  
    }  
}  
//END TweetAvgCount
```

EXPLANATION:

The Map Reduce approach with the second part of **(A) Content Analysis** enables you to retrieve the total number of tweet lengths and the total size of all tweet lengths totalled up. A key such as a string is used to denote the tweet and the tweet length is set as a value to be emitted to the reducer. Then you can find the average of tweet lengths by adding all the tweet lengths together and dividing them by how many lengths exist within the dataset. The average is then emitted as a value along with "tweet" as a key.

Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

PART B- TIME ANALYSIS

TWEET DATE TIME MAPPER CODE:

```
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.commons.lang.StringUtils;

public class TweetDateTimeMapper extends Mapper<Object, Text, Text, IntWritable> {
    //constant variable stores a value of one
    private final IntWritable one = new IntWritable(1);
    //create a text variable
    Text data=new Text();
    //mapper method takes takes as input the text file
    public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
        // Format per tweet is id;date;hashtags;tweet;

        String dump = value.toString();
        if(StringUtils.ordinalIndexOf(dump,";",4)>-1){
            //set the start index of the dataset to the date of the tweet
            int startIndex = StringUtils.ordinalIndexOf(dump,";",1) + 1;
            //store the date of the tweet in a variable
            String tweet = dump.substring(startIndex,dump.lastIndexOf(';'));
            //delimit each tweet date by removing the commas
            String args[]=tweet.split(",");
            //text variable stores the date of the tweet
            data.set(args[0]);
            //send tweet String and one to reducer
            context.write(data,one);
        }
    }
}

//END map method
//END TweetDateTimeMapper
```

TWEET DATE TIME REDUCER CODE:

```
import java.util.Iterator;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

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

    //reducer method takes key from mapper and value from mapper as input
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
        throws IOException, InterruptedException {
        //produce a count variable which will store the overall total occurrence of X date within the dataset
        int count =0;

        //iterate through each occurrence of X date and total up to find the total number of times X date appeared in the dataset
        for (IntWritable value : values) {
            count+=value.get();
        }

        //set the occurrence of X date as output
        result.set(count);

        //emit the date and its occurrence
        context.write(key,result);
    }
}

//END TweetDateTimeReducer
```

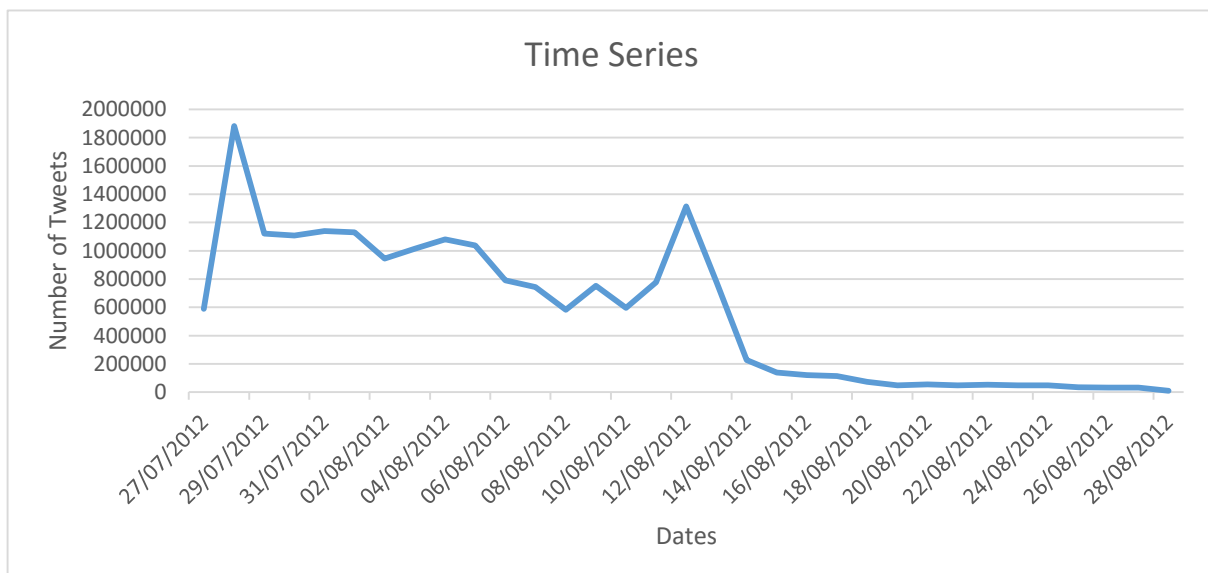
Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

TWEET DATE TIME COUNT CODE:

```
7  import org.apache.hadoop.mapreduce.Job;
8  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
9  import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
10
11  public class TweetDateTimeCount {
12
13      public static void runJob(String[] input, String output) throws Exception {
14
15          Configuration conf = new Configuration();
16
17          Job job = new Job(conf);
18          //sets current class as the class which generates and enables executability of a jar file
19          job.setJarByClass(TweetDateTimeCount.class);
20          //mapper class is set
21          job.setMapperClass(TweetDateTimeMapper.class);
22          //reducer class is set
23          job.setReducerClass(TweetDateTimeReducer.class);
24          //sets output key's class type
25          job.setMapOutputKeyClass(Text.class);
26          //sets output value's class type
27          job.setMapOutputValueClass(IntWritable.class);
28          //produces three tasks which output three text files with its corresponding results1
29          job.setNumReduceTasks(3);
30          Path outputPath = new Path(output);
31          FileInputFormat.setInputPaths(job, StringUtils.join(input, ","));
32          FileOutputFormat.setOutputPath(job, outputPath);
33          outputPath.getFileSystem(conf).delete(outputPath, true);
34          job.waitForCompletion(true);
35      }
36
37      public static void main(String[] args) throws Exception {
38          runJob(Arrays.copyOfRange(args, 0, args.length-1), args[args.length-1]);
39      }
40  }
41  //END TweetDateTimeCount
42
43
```

EXPLANATION:

The Map Reduce approach with the second part of **(B) Time Analysis** enables you to retrieve the date that a tweet was posted and count each date and store an occurrence of each date within the dataset and thus, compare how many tweets were made on a particular date.



As the time series diagram show above, the number of tweets fluctuate varying day to day, sometimes increasing across a number of days and sometimes decreasing. Furthermore, the day which had the greatest number of tweet at a value of just over 180, 000 tweets occurred 27 July 2012 - 29 July 2012.

Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

PART C- HASHTAG ANALYSIS

TWEET CHEER MAPPER COUNT CODE:

```
1 import java.io.IOException;
2 import java.util.StringTokenizer;
3 import org.apache.hadoop.io.IntWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Mapper;
6 import org.apache.commons.lang.StringUtils;
7
8 public class TweetCheerMapper extends Mapper<Object, Text, Text, IntWritable> {
9     //constant variable stores a value of one
10    private final IntWritable one = new IntWritable(1);
11    //create a text variable
12    Text data=new Text();
13    //mapper method takes takes as input the text file
14    public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
15
16        String dump = value.toString();
17        if(StringUtils.ordinalIndexOf(dump,";",4)>=1){
18            int startIndex = StringUtils.ordinalIndexOf(dump,";",3 + 1);
19            String tweet = dump.substring(startIndex,dump.lastIndexOf(";"));
20            //create a matrix which stores the countries from the method below
21            String[][] tweetInformation = getCountryTweets();
22            //iterate through the matrix and check whether the tweet within the dataset contains any key:cheer celebration and va
23            for(int i=0;i<tweetInformation.length;i++){
24                if(tweet.contains(tweetInformation[i][0])){
25                    //sets the contents within the array to lower case and stores that within the text variable which will be sent to the red
26                    data.set(tweetInformation[i][1].toLowerCase());
27                    //send country as a key and its occurrence as a value to the reducer
28                    context.write(data,one);
29                }//END for loop
30            }//END if statement
31
32
33        }//END if statement
34
35
36    }
```

TWEET CHEER REDUCER CODE:

```
1 import java.io.IOException;
2 import java.util.Iterator;
3 import org.apache.hadoop.io.IntWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Reducer;
6
7
8 public class TweetCheerReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
9
10    private IntWritable result = new IntWritable();
11    //reducer method takes key from mapper and value from mapper as input
12    public void reduce(Text key, Iterable<IntWritable> occurrences, Context context)
13        throws IOException, InterruptedException {
14        //use a variable to track the total number of times X country appears within the dataset
15        int total = 0;
16        //iterate through the occurrences of X country and total up how many times X country was cheered
17        for (IntWritable occurrence : occurrences) {
18
19            total=total+occurrence.get();
20
21        }
22        //set the total occurrences as output
23        result.set(total);
24        //emit the country and its total occurrence within the dataset
25        context.write(key,result);
26    }
27
28
29
30
31 }
```


Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

TWEET CHEER COUNT CODE:

```
7 import org.apache.hadoop.mapreduce.Job;
8 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
9 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
10
11 public class TweetCheerCount {
12
13     public static void runJob(String[] input, String output) throws Exception {
14
15         Configuration conf = new Configuration();
16
17         Job job = new Job(conf);
18         //sets current class as the class which generates and enables executability of a jar file
19         job.setJarByClass(TweetCheerCount.class);
20         //mapper class is set
21         job.setMapperClass(TweetCheerMapper.class);
22         //reducer class is set
23         job.setReducerClass(TweetCheerReducer.class);
24         //sets output key's class type
25         job.setMapOutputKeyClass(Text.class);
26         //sets output value's class type
27         job.setMapOutputValueClass(IntWritable.class);
28         //produces three tasks which output three text files with its corresponding results
29         job.setNumReduceTasks(3);
30         Path outputPath = new Path(output);
31         FileInputFormat.setInputPaths(job, StringUtils.join(input, ","));
32         FileOutputFormat.setOutputPath(job, outputPath);
33         outputPath.getFileSystem(conf).delete(outputPath, true);
34         job.waitForCompletion(true);
35     }
36
37     public static void main(String[] args) throws Exception {
38         runJob(Arrays.copyOfRange(args, 0, args.length-1), args[args.length-1]);
39     }
40
41 }
42
```

EXPLANATION:

The Map Reduce approach with the final part of **(C) Hash Tag Analysis** enables you to find the countries that had been cheered on via tweets that were posted, so each tweet would be compared against a matrix to see whether or not it had contained a cheer tweet of country X or country Y and if it had such content i.e: if gocountryX or teamcountryX was found, it would return country X along with the amount of times a tweet had been written with a cheer associated with it.

Student Name: Hussein Ahmed Tejan
Student Id No.: 130108432

Countries	Tweet Count
china	2
jamaica	5
kenya	3
niger	2
nigeria	2
sweden	1
bahamas	1
bulgaria	1
france	2
poland	4
canada	9
hungary	2
indonesia	2

As you can see from the table displayed above, there were not many countries that had been found that had been cheered/celebrated i.e.: goCountryX or teamCountryX. However, of the countries that were found being cheered on in tweets were the countries above with Canada having the most amount of cheers.