# Generic Explanation:

This assignment required multiple things such as file handling, string parsing etc. I am using file handling for *positive.txt and negative.txt.* Whereas the input file (*Comments.txt*) is being read by the mapper function itself. I am converting the whole line into a string and then parsing the string by storing the string into a different variable. For this I hardcoded the parser code as I saw the input file and before every comment, there is always a "T".

Another implementation was that in order to see if the reducer has ended, I am using a built-in function Cleaner (). This function is basically called when the reducer has finished its tasks.

Furthermore, for this assignment I am using HDFS as for some reason when I installed HDFS, my Hadoop-standalone stopped working.

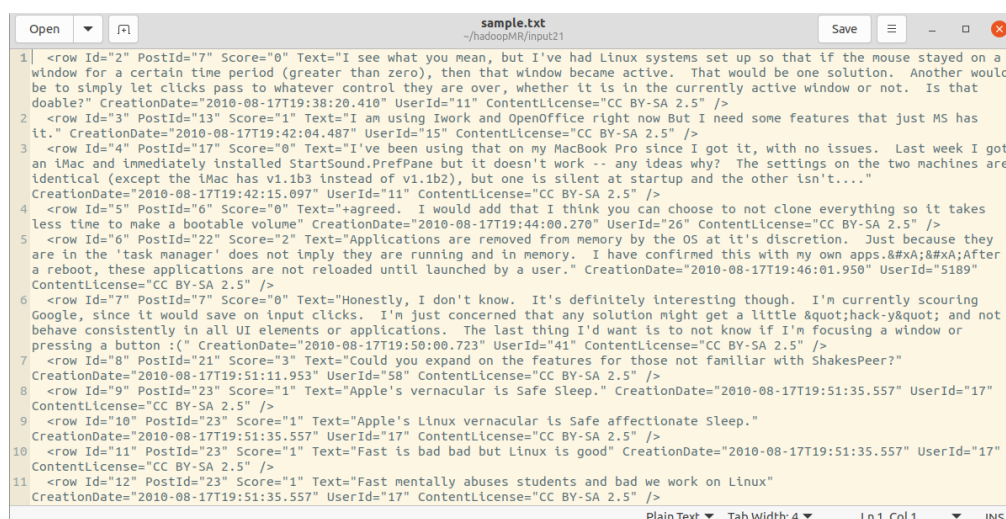Rest of the implementation is requirement specific explained below:

## Keyword Sentiment Analysis:

For this task I first read both positive.txt and negative.txt into respective strings and then parsed them into an array of strings such that each index will contain a whole word.

Similarly, I am reading the keyword from "keyword.txt". Then I am searching that keyword in every comment. If a comment has that particular keyword, then all of its contents will be parsed into an array of strings, each index containing a word of that comment. Then each word is compared with both the files and depending upon the number of instances in each file, respective values are updated. In the end, the values are compared and if the comment has more instances in positive, then 1 is added to "overall" (a variable that stores the overall sentiment of the word) and vice versa. In the end, the variable overall is compared and if it is :
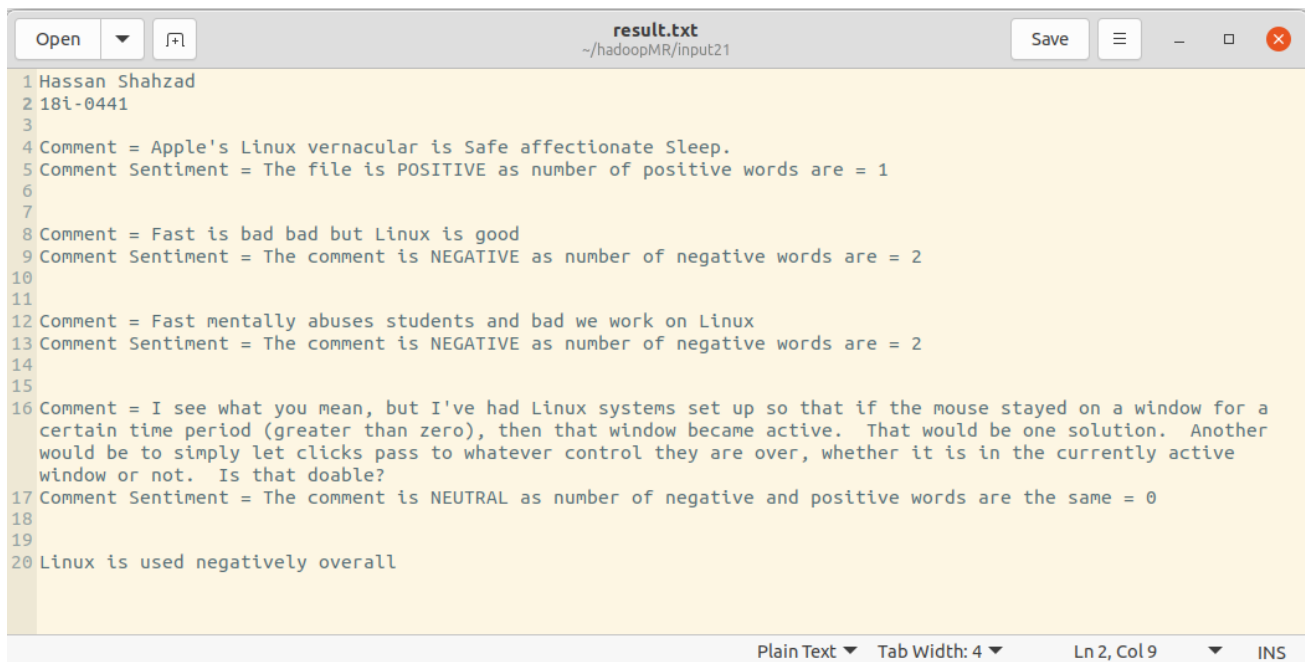
- >0 := This means that the number of positive comments were more and overall sentiment is positive.
- <0 := This means that the number of negative comments were more and overall sentiment is negative.
- =0 := This means that the number of positive and negative comments are equal and overall sentiment is neutral.

Hence, my implementation prints sentiment of each comment along with the overall sentiment. I ran it on a sample dataset to test it out first. For testing, the keyword used was "Linux". The screenshots of sample dataset and its outputs are given as follows:



*Fig 3.1: Test Dataset*

1

*Fig 3.2: Test output (Customized)*



*Fig 3.3: Test Output*

After this, the same code was run on the original dataset. The code that was run and the screenshots of its outputs are given as follows:



*Fig 3.4: Keyword File*

```java
1  //-------------------------//
2  // Hassan Shahzad
3  // 18i-0441
4  // PDC A3
5  //-------------------------//
6
7  import java.io.IOException;
8  import org.apache.hadoop.io.IntWritable;
9  import org.apache.hadoop.io.LongWritable;
10 import org.apache.hadoop.io.Text;
11 import org.apache.hadoop.mapreduce.Mapper;
12
13 public class IdentifyCommentsMapper
14     extends Mapper<LongWritable, Text, Text, IntWritable> {
15 @Override
16   public void map(LongWritable key, Text value, Context context)
17     throws IOException, InterruptedException {
18
19     String line = value.toString();        // Reads a complete line
20     int r_index = 0;                       // Stores index of the character "=" of the word "Row id ="
21     String id = "";                        // Stores the id
22     int c_index = 0;                       // Stores index of the character "T" of the word "Text"
23     String comment = "";                   // Stores the comment
24
25     for (int i=0; i< line.length(); i++){  // Getting the starting index for row id
26         char c = line.charAt(i);
27         if (c == '='){
28             r_index = i+2;
29             break;
30         }
31     }
32     for (int i = r_index; i< line.length(); i++){   // Getting the ending index for row id
33         char c = line.charAt(i);
34         if (c == '"')
35             break;
36         else
37             id = id + "" + c;    // Adding the id
38     }
39
40     for (int i=0; i< line.length(); i++){   // Getting the starting index for comment
41         char c = line.charAt(i);
42         if (c == 'T'){
43             c_index = i+6;
44             break;
45         }
46     }
47     for (int i = c_index; i< line.length(); i++){    // Getting the ending index for comment
48         char c = line.charAt(i);
49         if (c == '"')
50             break;
51         else
52             comment = comment + c ;                         // Adding the comment
53     }
54     int id1 = Integer.parseInt(id);    // Converting the row id into integer
55
56
57     context.write(new Text(comment), new IntWritable(id1));
58
59 } }
```

*Fig 3.5: Mapper Class*

```java
1  //-------------------------//
2  // Hassan Shahzad
3  // 18i-0441
4  // PDC A3
5  //-------------------------//
6
7  import java.io.IOException;
8  import java.io.File;
9  import java.io.FileWriter;
10 import java.io.FileNotFoundException;
11 import java.util.Scanner;
12
13 import org.apache.hadoop.io.IntWritable;
14 import org.apache.hadoop.io.Text;
15 import org.apache.hadoop.mapreduce.Reducer;
16
17 public class IdentifyCommentsReducer
18     extends Reducer<Text, IntWritable, Text, IntWritable> {
19
20     int row_id = 0;            // Will be used for break condition
21     String keyword;           // The keyword entered
22     int overall = 0;          // If overall > 0 then positive , else negative
23     @Override
24     public void cleanup(Context context) throws IOException, InterruptedException{    // Last line of file
25         Text key = new Text();
26         if (overall > 0){
27             String stro = keyword + " is used positively overall";
28             FileWriter myWriter3 = new FileWriter("/home/hxn/hadoopMR/input/result.txt", true);
29             myWriter3.write(stro);
30             myWriter3.close();
31             key.set(keyword + " is used positively overall");
32             context.write(key, new IntWritable(overall));
33         }
34         else if (overall < 0){
35             String stro = keyword + " is used negatively overall";
36             FileWriter myWriter3 = new FileWriter("/home/hxn/hadoopMR/input/result.txt", true);
37             myWriter3.write(stro);
38             myWriter3.close();
39             key.set(keyword + " is used negatively overall");
40             context.write(key, new IntWritable(overall));
41         }
42         else if (overall == 0 ){
43             String stro = keyword + " is used neutrally overall";
44             FileWriter myWriter3 = new FileWriter("/home/hxn/hadoopMR/input/result.txt", true);
45             myWriter3.write(stro);
46             myWriter3.close();
47             key.set(keyword + " is used neutrally overall");
48             context.write(key, new IntWritable(overall));
49         }
50     }
51     @Override
52     public void reduce(Text key, Iterable<IntWritable> values, Context context)
53         throws IOException, InterruptedException {
54         row_id ++;
55         try {
56             File kw = new File ("/home/hxn/hadoopMR/input/keyword.txt");     // Getting the keyword from the file
57             Scanner myReader = new Scanner(kw);
58             while (myReader.hasNextLine()) {
59                 keyword = myReader.nextLine();
60             }
61             myReader.close();
62         }catch(FileNotFoundException e) {     // If file is not found
63             System.out.println("File not found");
64             e.printStackTrace();
65         }
66
67         String str = key.toString();        // Converting the comment into string
68
69         if (str.contains(keyword)) {        // Checking if the comment contains that keyword
70             int positive = 0;               // Number of positive words
71             int negative = 0;               // Number of negative words
72             String[] words = str.split("\\W+"); // Splitting the string into array of words
73             String file = "";
74             String file1 = "";
75
```

*Fig 3.6: Reducer Class (Pt-1)*

```java
      try {
            File pos = new File ("/home/hxn/hadoopMR/input/positive.txt");        // Checking the positive file
            Scanner myReader = new Scanner(pos);
            while (myReader.hasNextLine()) {
                file += myReader.nextLine() + "~";                                // Splitting the positive file
            }
            myReader.close();
      }catch(FileNotFoundException e) {      // If file is not found
            System.out.println("File not found");
            e.printStackTrace();
      }

      try {
            File neg = new File ("/home/hxn/hadoopMR/input/negative.txt");        // Checking the negative file
            Scanner myReader = new Scanner(neg);
            while (myReader.hasNextLine()) {
                file1 += myReader.nextLine() + "~";                              // Splitting the negative file
            }
            myReader.close();
      }catch(FileNotFoundException e) {      // If file is not found
            System.out.println("File not found");
            e.printStackTrace();
      }

      String[] files = file.split("~");       // Splitting the positive string into array of words
      String[] files1 = file1.split("~");     // Splitting the negative string into array of words


      // Checking the instances of the comment in positive file
      for (int i=0; i< words.length; i++) {
          for (int j=0; j<files.length; j++) {
              if (words[i].equals(files[j])) {
                  positive++;                          // Incrementing the positive instances
              }}}

      for (int i=0; i< words.length; i++) {
          for (int j=0; j<files1.length; j++) {
              if (words[i].equals(files1[j])) {
                  negative++;                          // Incrementing the negative instances
              }}}

      FileWriter myWriter = new FileWriter("/home/hxn/hadoopMR/input/result.txt", true);     // Creating a resultant file
      myWriter.write("Comment = " + str + "\n");      // Writing the comment in the file
      myWriter.close();
      key.set("Comment = " + str + "\nOverall = ");
      context.write(key, new IntWritable(overall));
      if (positive > negative){
          String str1 = "Comment Sentiment = The file is POSITIVE as number of positive words are = "+ positive;
          FileWriter myWriter1 = new FileWriter("/home/hxn/hadoopMR/input/result.txt", true);
          myWriter1.write(str1 + "\n\n\n");
          myWriter1.close();
          key.set(str1 + "\nOverall = ");
          context.write(key, new IntWritable(overall));
          overall += 1;
      }
      else if (negative > positive){
          String str1 = "Comment Sentiment = The comment is NEGATIVE as number of negative words are = "+ negative;
          FileWriter myWriter1 = new FileWriter("/home/hxn/hadoopMR/input/result.txt", true);
          myWriter1.write(str1 + "\n\n\n");
          myWriter1.close();
          key.set(str1 + "\nOverall = ");
          context.write(key, new IntWritable(overall));
          overall-=1;
      }
      else if (negative == positive){
          String str1 = "Comment Sentiment = The comment is NEUTRAL as number of negative and positive words are the same = "+ negative;
          FileWriter myWriter1 = new FileWriter("/home/hxn/hadoopMR/input/result.txt", true);
          myWriter1.write(str1+"\n\n\n");
          myWriter1.close();
          key.set(str1 + "\nOverall = ");
          context.write(key, new IntWritable(overall));
          overall+=0;
      }}

   }}
```

*Fig 3.7: Reducer Class (Pt-2)*

```java
//--------------------------//
// Hassan Shahzad
// 18L-0441
// PDC A3
//--------------------------//

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.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class IdentifyComments {
    static String keyword;
  public static void main(String[] args) throws Exception {
    if (args.length != 2) {
      System.err.println("Usage: AvgTemperature <input path> <output path>");
      System.exit(-1);
    }

    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "Identify Comments");
    job.setJarByClass(IdentifyComments.class);
    job.setJobName("Identify Comments");
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path   (args[1]));
    job.setMapperClass(IdentifyCommentsMapper.class);
    job.setReducerClass(IdentifyCommentsReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    System.exit(job.waitForCompletion(true) ? 0 : 1);
  }
}
```

*Fig 3.8: Main Java File*



*Fig 3.9: Execution of the Code*

*Fig 3.10: Input Folder*



*Fig 3.11: Output File (Customized)*



*Fig 3.12: Output File Original*