

```

1  /* Edited 04.04.2017 08h00 */
2
3
4          /*****      MOVIE RATINGS REVIEW Design 2 - MapReduce      *****/
5
6  /* Import classes necessary for the movieReview job */
7  import java.io.*;
8  import java.io.BufferedReader; // Class which allows buffering of characters for efficient reading of characters
9  import java.io.FileReader;
10 import java.io.IOException;
11 import java.util.Collections;
12 import java.util.HashSet;
13 import java.util.HashMap;
14 import java.util.Map;
15 import java.util.Set;
16 import java.util.StringTokenizer;
17 import java.util.Scanner;
18 import org.apache.hadoop.conf.Configured;
19 import org.apache.hadoop.filecache.DistributedCache;
20 import org.apache.hadoop.fs.Path;
21 import org.apache.hadoop.io.Text;
22 import org.apache.hadoop.io.IntWritable;
23 import org.apache.hadoop.mapreduce.Job;
24 import org.apache.hadoop.mapreduce.Mapper;
25 import org.apache.hadoop.mapreduce.Reducer;
26 import org.apache.hadoop.mapreduce.Reducer.Context;
27 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
28 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
29 import org.apache.hadoop.util.Tool;
30 import org.apache.hadoop.util.ToolRunner;
31
32
33 /* The main movieReview class consisting of 'reviewMapper', 'reviewCombiner', 'reviewReducer', 'main' and 'run' methods*/
34 public class movieReview extends Configured implements Tool // A tool interface to manage custom arguments provided in the
command line (3 arguments are expected)
35 {
36     /* Data is read from the file into the reviewMapper and emitted to the context as ( Key : Value ) -> ( "reviewRating" :
"reviewWord,1" ) -> ( "9" , "Amazing,1" ) */
37     public static class reviewMapper extends Mapper < Object, Text, Text, Text > // reviewMapper input <Object, Text> and
output as <Text, Text>
38     {
39
40         /* Declare/initialise variables for the reviewMapper */
41         private BufferedReader myReader; // Declare myReader to use BufferedReader
42         private Set<String> excludedWords = new HashSet<String>(); // Declare an empty HashSet in which to load the files or
punctuation characters to be excluded from counts
43
44
45         /* Use setup method to incorporate a third argument to point to the file location of words to be excluded */
46         @Override
47         protected void setup(Context context) throws IOException, InterruptedException

```

```

48     {
49         try
50         {
51             Path[] excludedWordFilePath = new Path[0]; // Declare excludedWordFilePath as a path
52             excludedWordFilePath = context.getLocalCacheFiles(); // Load the cache file location from the context
53             // into variable excludedWordFilePath
54             if (excludedWordFilePath != null && excludedWordFilePath.length > 0) // Ensure path exists (not null)
55                 // and the path length is a positive integer
56             {
57                 for (Path excludedWordPath : excludedWordFilePath)
58                 {
59                     readStopWordFile(excludedWordPath);
60                 }
61             }
62             catch (IOException e)
63             {
64                 System.err.println("Error reading the location of the excluded_word file: " + e);
65             }
66         }
67         /* Method to load file contents (with words to be excluded) into HashSet (excludedWords) */
68         private void readStopWordFile(Path excludedWordPath)
69         {
70             try
71             {
72                 myReader = new BufferedReader(new FileReader(excludedWordPath.toString())); // Read the file provided
73                 // containing the words to be excluded, line by line
74                 String excludedWordLine = null;
75                 while ((excludedWordLine = myReader.readLine()) != null) // Continue reading line by line until there
76                     // are no more lines in the file
77                 {
78                     String wordsTobeExcluded[] = excludedWordLine.split(","); // Split the line being read, by the
79                     // (,) seperating the words to be excluded
80                     for (String excludedWord: wordsTobeExcluded) // Iterate over the comma seperated exclusion words
81                         // and punctuation characters
82                     {
83                         excludedWords.add(excludedWord); // For each excludedWord found, add to the HashSet
84                         excludedWords
85                     }
86                 }
87             }
88             catch (IOException ioe)
89             {
90                 System.err.println("Unable to read contents of the file containing words to be excluded '" +
91                     excludedWordPath + "' : " + ioe.toString());
92             }
93         }
94         /* This map function takes a line of text from the file at a time, splits into a review string and rating

```

```

string, and emits a key and a word associated with that key */
91 @Override
92 public void map(Object key, Text value, Context context) throws IOException, InterruptedException
93 {
94     String[] ratingAndReview = value.toString().toLowerCase().split("\t");// Line is of the format
    <"[Review]" [Rating {1-9}]]> --> split this line by the tab (\t) into a string with the reviewRating and
    the movie rating

95     // Get the current line and remove unexpected characters

96
97
98     String rating = ratingAndReview[0]; // Save the review words from the first part of the split, into
    variable 'rating'
99     rating = rating.replaceAll("'", ""); // Remove single quotes from the rating string
100    rating = rating.replaceAll("[^a-zA-Z]", " "); // Remove other unnecessary punctuation from the rating
    string
101    String reviewRating = ratingAndReview[1]; // Save the review rating from the second part of the split,
    into variable 'reviewRating'
102    Scanner itr = new Scanner(rating); // Use scanner to run over the words in the 'rating' string
103    while (itr.hasNext()) // Iterate over words until all the review words have been taken into account
104    {
105        String reviewWord = itr.next(); // Call next word in the sequence
106        if (!excludedWords.contains(reviewWord) && reviewWord.length() > 1 ) // Check if the word being
        iterated over is not in the file containing the words to be excluded
107        {
108            String compoundValue = reviewWord + ",1"; // Concatenate the word and a "1" to form a string
            "word,1" as the value emitted from the mapper in the key-value pair
109            context.write(new Text(reviewRating), new Text(compoundValue) ); // Emit all the
            reviewRatings as the keys and values as "word,1" (for example)
110        }
111    }
112 }
113 }
114
115 /* The combiner class aims to aggregate data locally from each reviewMapper before sending the data to the reviewReducer
    */
116 public static class reviewCombiner extends Reducer < Text, Text ,Text, Text > // reviewCombiner input <Text, Text> and
    output as <Text, Text>
117 {
118     /* Declare/initialise variables for the reviewCombiner */
119     private HashMap <String, Integer> localCombinerCount = new HashMap <String, Integer>(); // Declare
    'localCombinerCount' HashMap to store the word (for a given review) and highest frequency found, nine HashMaps are
    produced - one for each reviewRating

120
121     public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException
122     {
123         for (Text value : values) // Iterate over the values read in from the context
124         {
125             StringTokenizer itr = new StringTokenizer(value.toString(),","); // Seperate the "word,1"
126             String aWord = itr.nextToken(); // Pull out the word and store as a string variable aWord
127             Integer count = Integer.parseInt(itr.nextToken()); // Parse the string, pull out the "1" and convert
            the "1" to an integer value

```

```

128         if (localCombinerCount.containsKey(aWord))
129         {
130             localCombinerCount.put(aWord, localCombinerCount.get(aWord) + 1); // Add 1 to the count/frequency
131             found in the HashMap if the given word is in the HashMap
132         }
133         else
134         {
135             localCombinerCount.put(aWord, count); // If not in the HashMap localCombinerCount, add the word
136             and give it a count of 1
137         }
138     }
139     for (Map.Entry<String, Integer> entry : localCombinerCount.entrySet()) // Iterate over <String,Integer>
140     values in the HashMap entries
141     {
142         String aKeyWord = entry.getKey(); // Saves the key from the HashMap (which is a word) as aKeyWord
143         variable
144         Integer aKeyWordCount = entry.getValue(); // Saves the count associated with the key in the HashMap
145         as aKeyWordCount variable
146         String compoundKeyWithCount = aKeyWord + "," + Integer.toString(aKeyWordCount); // Converts the count
147         to a string and concatenates aKeyWord and aKeyWordCount
148         context.write(key,new Text(compoundKeyWithCount)); // Emits the reviewRating number (key) as "2" and
149         the compound (value) "word,count"
150     }
151     localCombinerCount.clear(); // Clear the HashMap
152 }
153
154 /* Data is read into the reviewReducer, finally aggregated/collated and written out to a file as the final result */
155 public static class reviewReducer extends Reducer < Text, Text ,Text, Text > // reviewReducer receives <Text, Text> and
156 sends <Text, Text>
157 {
158     /* Declare variables for the reviewReducer */
159     Text frequentWord_Frequency = new Text(); // Declare variable frequentWord_Frequency to represent a word/s in the
160     review for which we are counting its frequency
161     HashMap <String, Integer> reducerCount = new HashMap <String, Integer>(); // Declare 'reducerCount' HashMap to store
162     the word (for a given review) and highest frequency found, nine HashMaps are produced - one for each reviewRating
163
164     public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException
165     {
166         StringBuffer wordWithHighestFrequency = new StringBuffer(); // Declare 'wordWithHighestFrequency'
167         variable to store word/s from each final HashMap with highest occurrence
168
169         for (Text value : values) // Iterate over the values read in from the context
170         {
171             StringTokenizer itr = new StringTokenizer(value.toString(),","); // Separate the compoundValue (which
172             contains "word,12")

```

```

167         String aWord = itr.nextTokn(); // Pull the word out as a token called aWord
168         Integer count = Integer.parseInt(itr.nextTokn()); // Parse over the token string to pull out the
169         integer count
170
171         if (reducerCount.containsKey(aWord))
172         {
173             reducerCount.put(aWord, reducerCount.get(aWord) + count); // Add the parsed integer count to the
174             count/frequency found in the HashMap if the given word is in the HashMap
175         }
176         else
177         {
178             reducerCount.put(aWord, count); // If not in the HashMap reducerCount, add the word and add its
179             count so far
180         }
181     }
182
183     // Iterate through the reducerCount HashMaps and return the key associated with the highest count
184     Integer maximumValue = Collections.max(reducerCount.values()); // Declare 'MaximumValue' as the variable
185     to hold the highest frequency
186     for (Map.Entry<String, Integer> entry : reducerCount.entrySet()) // Iterate over values in the HashMap
187     {
188         if (entry.getValue() == maximumValue) // Compare the frequency of the values frequency being
189         considered, to the highest value stored so far in maximumValue
190         {
191             wordWithHighestFrequency.append(entry.getKey()+ "|"); // For the highest frequency found in the
192             HashMap, retrieve the associated word and add a '|'
193         }
194     }
195
196     wordWithHighestFrequency.append(maximumValue); // Append/add the frequency for the word associated with
197     the highest frequency, to the end of the string
198     frequentWord_Frequency.set(wordWithHighestFrequency.toString()); // Store this result in the variable
199     'frequentWord_Frequency'
200
201     context.write(key, frequentWord_Frequency); // Write the final result to context
202     reducerCount.clear();
203 }
204
205 /* main() program initiates the program */
206 public static void main(String[] args) throws Exception
207 {
208     int exitCode = ToolRunner.run(new movieReview(), args);
209     System.exit(exitCode);
210 }
211
212 /* Configures and defines the jobs to be run in the main() program */
213 public int run(String[] args) throws Exception
214 {
215     if (args.length != 3) // Ensure three arguments are provided to the command line

```

```

210     {
211         System.err.printf(" %s needs three arguments: <movie_data.txt> <results> <wordsTobeExcluded.txt> \n",
            getClass().getSimpleName()); // Note the results directory must not exist, the program will create a new
            directory
            return -1;
212     }
213 }
214
215 //Initialize the Hadoop job and set the jar as well as the name of the Job
216 Job job = new Job(); // Create a new instance of the Job object
217 job.setJarByClass(movieReview.class); // Set the jar class to use
218 job.setJobName("reviewMR"); // Allocate the job a name for logging/tracking purposes
219
220 //Add input and output file paths to job based on the arguments passed
221 FileInputFormat.addInputPath(job, new Path(args[0])); // Assign the first argument provided as the input path
    for the data to be considered in the MapReduce program
222 FileOutputFormat.setOutputPath(job, new Path(args[1])); // Assign the second argument provided as the output
    path for the results from the MapReduce program to be written
223
224 // Use text objects to output the key (ideally the rating) and value (ideally the most frequent word
    associated with the key)
225 job.setOutputKeyClass(Text.class); // Use a text object for output of the key
226 job.setOutputValueClass(Text.class); // Similarly use a text object for output of the value
227
228 //Set the reviewMapper and reviewReducer in the job
229 job.setMapperClass(reviewMapper.class); // Set reviewMapper as the map class for the job
230 job.setCombinerClass(reviewCombiner.class);
231 job.setReducerClass(reviewReducer.class); // Set reviewReducer as the reduce class for the job
232 //job.setNumReduceTasks(1); // Set the number of reducerReducers to be called - this number dictates how many
    result files are created
233
234 // Ensure the program excludes the words to be ommitted from the word count
235 DistributedCache.addCacheFile(new Path(args[2]).toUri(), job.getConfiguration()); // Assign the second
    argument provided as the location of the file containing words/punctuation to be excluded
236
237 //Wait for the job to complete and print if the job was successful or not
238 Integer returnValue = job.waitForCompletion(true) ? 0:1; // Unix notation to describe if a job is successful
    (0) or unsuccessful (1)
239 if(job.isSuccessful())
240 {
241     System.out.println("The movieReview MapReduce job was successful.");
242 }
243 else if(!job.isSuccessful())
244 {
245     System.out.println("The movieReview MapReduce job was unsuccessful...please review the program and
        address errors.");
246 }
247
248 return returnValue; // Return 0 or 1 depending on job success
249 }
250
251

```

