## WRITE A PARTITIONER

○ **Write a customised Partitioner to separate the output of weblog by years.**

---

## Files and Directories Used in this Exercise

Eclipse project: partitioner

Java files: (Need to be created or copied)
YearPartitioner.java (Partitioner)
ProcessLogs.java (Driver)
LogFileMapper.java (Mapper)
LogFileReducer.java (Reducer)

Test data (HDFS):
Weblog (full web server access log)

Exercise directory: ~/workspace/partitioner

```java
 1  package stubs;
 2  import org.apache.hadoop.mapreduce.Job;
 3  import org.apache.hadoop.fs.Path;
 4  import org.apache.hadoop.io.IntWritable;
 5  import org.apache.hadoop.io.Text;
 6  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
 7  import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
 8
 9  public class ProcessLogs {
10
11    public static void main(String[] args) throws Exception {
12
13      if (args.length != 2) {
14        System.out.printf("Usage: ProcessLogs <input dir> <output dir>\n");
15        System.exit(-1);
16      }
17
18      Job job = new Job();
19      job.setJarByClass(ProcessLogs.class);
20      job.setJobName("Process Logs");
21
22      /*
23       * TODO implement
24       */
25      // Defining the input/output paths.
26      FileInputFormat.setInputPaths(job, new Path(args[0]));
27      FileOutputFormat.setOutputPath(job, new Path(args[1]));
28
29      // Set the Mapper and the Reducer.
30      job.setMapperClass(LogFileMapper.class);
31      job.setReducerClass(LogFileReducer.class);
32      job.setPartitionerClass(YearPartitioner.class);
33      job.setNumReduceTasks(3);
34
35      // Intermediate output key/value produced by Mapper.
36      job.setMapOutputKeyClass(Text.class);
37      job.setMapOutputValueClass(IntWritable.class);
38
39      // Reducer output key/value class
40      job.setOutputKeyClass(Text.class);
41      job.setOutputValueClass(IntWritable.class);
42
43      boolean success = job.waitForCompletion(true);
44      System.exit(success ? 0 : 1);
45    }
```

# Driver Code

4

# MAPPER CODE

```java
1 package stubs;
2 import java.io.IOException;
3
4 import org.apache.hadoop.io.IntWritable;
5 import org.apache.hadoop.io.LongWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Mapper;
8
9 /**
10  * Example input line:
11  * 96.7.4.14 - - [24/Apr/2011:04:20:11 -0400] "GET /cat.jpg HTTP/1.1" 200 12433
12  *
13  */
14 public class LogFileMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
15
16     @Override
17     public void map(LongWritable key, Text value, Context context)
18         throws IOException, InterruptedException {
19
20         String line = value.toString();
21         String[] data = line.trim().split("\\[\\d{2}/\\w{3}/");
22
23         if (data.length > 0) {
24             String part_data = data[1];
25             String year = part_data.substring(0,4);
26             context.write(new Text(year), new IntWritable(1));
27         }
28     }
29 }
```

# YearPartitioner Code

```java
1  package stubs;
2
3  import java.util.HashMap;
4
5  import org.apache.hadoop.io.Text;
6  import org.apache.hadoop.io.IntWritable;
7  import org.apache.hadoop.conf.Configurable;
8  import org.apache.hadoop.conf.Configuration;
9  import org.apache.hadoop.mapreduce.Partitioner;
10
11 public class YearPartitioner<K2, V2> extends Partitioner<Text, IntWritable> implements
12     Configurable {
13
14   private Configuration configuration;
15   HashMap<String, Integer> years = new HashMap<String, Integer>();
16   private String[] YearList = {"2009","2010","2011"};
17   private String tg;
18   private boolean found;
19   private int tgv;
20
21   /**
22    * Set up the months hash map in the setConf method.
23    */
24   @Override
25   public void setConf(Configuration configuration) {
26     /*
27      * Add the months to a HashMap.
28      */
29     for (int y = 0; y < YearList.length; y++) {
30         years.put(YearList[y], y);
31     }
32   }
33
34   /**
35    * Implement the getConf method for the Configurable interface.
36    */
37   @Override
38   public Configuration getConf() {
39     return configuration;
40   }
41
42   /**
43    * You must implement the getPartition method for a partitioner class.
44    * This method receives the three-letter abbreviation for the month
45    * as its value. (It is the output value from the mapper.)
46    * It should return an integer representation of the month.
47    * Note that January is represented as 0 rather than 1.
48    *
49    * For this partitioner to work, the job configuration must have been
50    * set so that there are exactly 12 reducers.
51    */
52   public int getPartition(Text key, IntWritable value, int numReduceTasks) {
53
54       tg = key.toString();
55       found = false;
56       for (int y = 0; y < YearList.length; y++) {
57           if (tg.equals(YearList[y])) {
58               found = true;
59               tgv = y;
60               break;
61           }
62       }
63
64       if (found) {
65           return tgv;
66       }
67       else {
68           return 0;
69       }
70   }
71 }
```

6

## Reducer Code

```java
1  package stubs;
2
3  import java.io.IOException;
4
5  import org.apache.hadoop.io.IntWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Reducer;
8
9
10 public class LogFileReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
11
12   @Override
13     public void reduce(Text key, Iterable<IntWritable> values, Context context)
14             throws IOException, InterruptedException {
15
16       int wordCount = 0;
17
18       for (IntWritable value : values) {
19             wordCount += value.get();
20       }
21
22       context.write(key, new IntWritable(wordCount));
23     }
24 }
```

7

```
[training@192-168-1-109 W6T1]$ ls
_logs   part-r-00000   part-r-00001   part-r-00002   _SUCCESS
[training@192-168-1-109 W6T1]$ cat part-r-00000
2009      50216
[training@192-168-1-109 W6T1]$ cat part-r-00001
2010      1712429
[training@192-168-1-109 W6T1]$ cat part-r-00002
2011      2715198
[training@192-168-1-109 W6T1]$ █
```

# RESULT