

Lab 4 – Total Order Sorting Report

Net ID: tw2770

Name: Tzu-An Wang

1. TotalOrderSortMapper

```
1  import java.io.IOException;
2  import javax.naming.Context;
3  import org.apache.hadoop.io.IntWritable;
4  import org.apache.hadoop.io.LongWritable;
5  import org.apache.hadoop.io.Text;
6  import org.apache.hadoop.mapreduce.Mapper;
7
8
9
10 public class TotalOrderSortMapper extends Mapper<LongWritable, Text, Text, Text> {
11
12     @Override
13     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
14         // Convert the line to a string
15         String line = value.toString();
16
17         // Check if the line is not empty and has at least ten characters
18         if (line != null && line.length() >= 10) {
19             // Extract the first ten bytes as the sort key
20             String sortKey = line.substring(beginIndex:0, endIndex:10);
21
22             // Output the sort key and the original line as the value
23             context.write(new Text(sortKey), new Text(line.substring(beginIndex:10, line.length())));
24         }
25     }
26 }
```

2. TotalOrderSortReducer

```
1  import java.io.IOException;
2  import javax.naming.Context;
3  import org.apache.hadoop.io.Text;
4  import org.apache.hadoop.mapreduce.Reducer;
5
6
7 public class TotalOrderSortReducer extends Reducer<Text, Text, Text, Text> {
8     @Override
9     public void reduce(Text key, Iterable<Text> values, Context context)
10         throws IOException, InterruptedException {
11         // Iterate through all values associated with the key and write them to the context
12         for (Text value : values) {
13             context.write(key, value);
14         }
15     }
16 }
```

3. TotalOrderSort – Job 1

```
15 public class TotalOrderSort {
16     Run | Debug
17     public static void main(String[] args) throws Exception {
18         Path inputPath = new Path(args[0]);
19         Path outputPath = new Path(args[1]);
20         Path stagingPath = new Path(args[1] + "_staging");
21         Path partitionFile = new Path(args[1] + "_partition");
22
23         // Configure job to prepare for sampling
24         Job sampleJob = Job.getInstance();
25         sampleJob.setJarByClass(TotalOrderSort.class);
26
27         // Use the mapper implementation with zero reduce tasks
28         sampleJob.setMapperClass(TotalOrderSortMapper.class);
29         sampleJob.setNumReduceTasks(0);
30
31         sampleJob.setOutputKeyClass(Text.class);
32         sampleJob.setOutputValueClass(Text.class);
33
34         TextInputFormat.setInputPaths(sampleJob, inputPath);
35
36         // Set the output format to a sequence file
37         sampleJob.setOutputFormatClass(SequenceFileOutputFormat.class);
38         SequenceFileOutputFormat.setOutputPath(sampleJob, stagingPath);
39
40         // Submit the job
41         sampleJob.waitForCompletion(true);
42     }
43 }
```

4. TotalOrderSort – Job 2

```
43 Job orderJob = Job.getInstance();
44 orderJob.setJarByClass(TotalOrderSort.class);
45
46 // Here, use the identity mapper to output the key/value pairs in the SequenceFile
47 orderJob.setReducerClass(TotalOrderSortReducer.class);
48
49 // Set the number of reduce tasks to an appropriate number for the amount of data being sorted
50 orderJob.setNumReduceTasks(10);
51
52 // Use Hadoop's TotalOrderPartitioner class
53 orderJob.setPartitionerClass(TotalOrderPartitioner.class);
54
55
56 // Set the partition file
57 TotalOrderPartitioner.setPartitionFile(orderJob.getConfiguration(), partitionFile);
58
59 orderJob.setOutputKeyClass(Text.class);
60 orderJob.setOutputValueClass(Text.class);
61
62 // Set the input to the previous job's output
63 orderJob.setInputFormatClass(SequenceFileInputFormat.class);
64 SequenceFileInputFormat.setInputPaths(orderJob, stagingPath);
65
66 // Set the output path to the command line parameter
67 TextOutputFormat.setOutputPath(orderJob, outputPath);
68
69 // Use the InputSampler to go through the output of the previous job, sample it, and create the partition file
70 InputSampler.writePartitionFile(orderJob, new InputSampler.RandomSampler(.001, 10000));
71
72 // Submit the job
73 orderJob.waitForCompletion(true);
74
75 }
```

5. Program runs successfully

```
tw2770_nyu_edu@nyu-dataproc-m:~$ hadoop fs -ls /user/tw2770_nyu_edu/lab4/output
Found 11 items
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 0 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/_SUCCESS
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 833200 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00000
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 1056000 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00001
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 1525700 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00002
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 838300 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00003
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 729800 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00004
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 1176500 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00005
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 1242600 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00006
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 746300 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00007
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 1143100 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00008
-rw-r--r-- 1 tw2770_nyu_edu tw2770_nyu_edu 708500 2024-03-06 04:49 /user/tw2770_nyu_edu/lab4/output/part-r-00009
```