

Database Management Systems

CS-236

Project Report

Team Members:

Ankit Mane [862393644] [amane011@ucr.edu]
Daiyaan Ahmed Shaik [862393960] [dshai005@ucr.edu]

1 Project Description

1.1 Experimental Setup

Installed Hadoop 3.2.1 and the compatible version of Java to create the MapReduce jobs.

1.2 Combining the datasets / Logic for join(s)

The first dataset, 'hotel-bookings' contains data on reservations made between 2015 and 2016. The second dataset 'customer-reservations' offers details on reservations made between 2017 and 2018.

1.2.1 Mapper

Mapper handles the merging and initial data processing for both the hotel booking datasets. It transforms input records into a set of intermediate key-value pairs. We used a separate set of variables to store the year, month and price for each of the datasets. The Mapper works as follows:

- The map function initially checks if the current line contains **avg_price_per_room** which is in the header row in both the datasets to ensure only data rows are processed.
- We then check if the current line is from dataset 1 or 2 by confirming the number of fields of the current line (**13 for dataset 1, 10 for dataset 2**). Either way, the price per room is parsed from the correct field and multiplied by the number of weekend and weekday nights (**2nd and 3rd fields in dataset 1, 8th and 9th field in dataset 2**).
- The calculated total price is then written to the context with a key combining the year and month and the value being the total price.

After all the lines are processed, the context will contain key-value pairs where the key is a combination of year and month and the value is the total revenue for that booking. This data is then passed to the reducer for further

processing, where revenues for each month will be summed up across all years and then sorted to find the most profitable month.

1.2.2 Reducer

The Reducer takes the intermediate key-value pairs from the Mapper and reduces the values for each key to a single output. The reducer functions as follows:

- The reduce function takes in a key and an iterable of values. The key represents the combination of year and month from the Mapper output, and the values are a set of numbers that represent the total revenues for each booking corresponding to that key.
- We then iterate over each value from the values, which contains all the total revenues for bookings of that month across all years. It adds up these revenues to a total revenue variable.
- Lastly, it writes the key (year and month combination) and the total revenue for that key to the context. Therefore, the output will be a list of year and month combinations and the corresponding total revenue for that month across all years.

The output from the Reducer is then written to the final output file on HDFS. The output key-value pairs are sorted by key, so the revenues for each month are in ascending order of year and month.

1.3 Computing the most profitable month

1.3.1 Mapper

The output from the previous reducer is taken as input for the Mapper function. Each line from the input is split into fields based on tab (`\t`). The first field (`fields[0]`) contains the **month and year combination** and the second field (`fields[1]`) contains the **average revenue for that month**. The Mapper is just preparing the data for the next Reducer.

1.3.2 Reducer

This reducer takes the output from the previous Mapper and determines which month in each year had the highest revenue. It works as follows:

- The Reducer has two HashMaps - where one keeps track of the maximum revenue encountered for each year(**yearMaxRevenue**) and the other keeps track of which month that maximum revenue occurred in (**yearMaxMonth**).
- It takes in a key (the month and year) and a set of values (the revenues for that month). For each key-value group, it calculates the total revenue

for that month and year. If this is the first time it has seen this year or if the total revenue is greater than the current maximum for that year, it updates **yearMaxRevenue** and **yearMaxMonth**

- The cleanup method is called at the end of the reduce phase. In this method, the reducer writes the most profitable month and corresponding revenue for each year to the context.

The reducer will finally output one record for each year in the dataset, with the record containing the most profitable month and the corresponding revenue for that year.

2 Code Snippets

2.1 Mapper for combining datasets

```
J Mapper.java X
J Mapper.java > Mapper > map(LongWritable, Text, Context)
1 package org.example;
2
3 import org.apache.hadoop.io.LongWritable;
4 import org.apache.hadoop.io.Text;
5
6 import java.io.IOException;
7
8 public class Mapper extends org.apache.hadoop.mapreduce.Mapper<LongWritable, Text, Text, LongWritable> {
9     public static int d1Year = 4;
10    public static int d1Month = 5;
11    public static int d1Price = 8;
12
13    public static int d2Year = 3;
14    public static int d2Month = 4;
15    public static int d2Price = 11;
16    public static int d1 = 10;
17    public static int d2 = 13;
18
19    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
20        String line = value.toString();
21        String[] fields = line.split(regex:",");
22        if (line.contains(s:"avg_price_per_room")) {
23            return;
24        }
25        if (fields.length == d1) {
26            double price = Double.parseDouble(fields[d1Price]);
27            long total = (long) (price * (Integer.parseInt(fields[1]) + Integer.parseInt(fields[2]]));
28            Text outputKey = new Text(fields[d1Year] + "-" + fields[d1Month]);
29            context.write(outputKey, new LongWritable(total));
30        } else if (fields.length == d2) {
31            double price = Double.parseDouble(fields[d2Price]);
32            long total = (long) (price * (Integer.parseInt(fields[7]) + Integer.parseInt(fields[8]]));
33            Text outputKey = new Text(fields[d2Year] + "-" + fields[d2Month]);
34            context.write(outputKey, new LongWritable(total));
35        }
36    }
37 }
38
```

2.2 Reducer for combining datasets

```
J Reducer.java > ...
1 package org.example;
2
3 import org.apache.hadoop.io.LongWritable;
4 import org.apache.hadoop.io.Text;
5
6 import java.io.IOException;
7
8 public class Reducer extends org.apache.hadoop.mapreduce.Reducer<Text, LongWritable, Text, LongWritable> {
9
10     @Override
11     public void reduce(Text key, Iterable<LongWritable> values, Context context) throws IOException, InterruptedException {
12         long totalRevenue = 0;
13         for (LongWritable value : values) {
14             totalRevenue += value.get();
15         }
16         context.write(key, new LongWritable(totalRevenue));
17     }
18 }
```

2.3 Driver code for combining datasets

```
J Driver.java 1 x
J Driver.java > {} org.example
1 package org.example;
2
3 import org.apache.hadoop.conf.Configuration;
4 import org.apache.hadoop.fs.Path;
5 import org.apache.hadoop.io.LongWritable;
6 import org.apache.hadoop.io.Text;
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 Driver {
12     Run | Debug
13     public static void main(String[] args) throws Exception {
14         if (args.length != 2) {
15             System.err.println(x:"Usage: App <input path> <output path>");
16             System.exit(-1);
17         }
18
19         Configuration conf = new Configuration();
20         Job job = Job.getInstance(conf, "Hotel Revenue");
21         job.setJarByClass(Driver.class);
22         job.setMapperClass(Mapper.class);
23         job.setReducerClass(Reducer.class);
24
25         // Set the types of output key and value for Mapper & Reducer
26         job.setMapOutputKeyClass(Text.class);
27         job.setMapOutputValueClass(LongWritable.class); // change here
28         job.setOutputKeyClass(Text.class);
29         job.setOutputValueClass(LongWritable.class); // and here
30
31         FileInputFormat.addInputPath(job, new Path(args[0]));
32         FileOutputFormat.setOutputPath(job, new Path(args[1]));
33
34         System.exit(job.waitForCompletion(true) ? 0 : 1);
35     }
36 }
```

2.4 Mapper for finding most profitable month

```
src > main > java > org > example > J TopMonthMapper.java > ...
1 package org.example;
2 import org.apache.hadoop.io.LongWritable;
3 import org.apache.hadoop.io.Text;
4 import org.apache.hadoop.mapreduce.Mapper;
5 import java.io.IOException;
6 import org.apache.hadoop.mapreduce.Mapper;
7
8 public class TopMonthMapper extends Mapper<LongWritable, Text, Text, LongWritable> {
9     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
10         String[] fields = value.toString().split(regex:"\\t");
11         String monthYear = fields[0];
12         long avgRevenue = Long.parseLong(fields[1]);
13
14         context.write(new Text(monthYear), new LongWritable(avgRevenue));
15     }
16 }
```

2.5 Reducer for finding most profitable month

```
1 package org.example;
2 import org.apache.hadoop.io.LongWritable;
3 import org.apache.hadoop.io.Text;
4 import org.apache.hadoop.mapreduce.Reducer;
5 import java.io.IOException;
6 import java.util.HashMap;
7 import java.util.Map;
8
9 public class TopMonthReducer extends Reducer<Text, LongWritable, Text, LongWritable> {
10     private Map<String, Long> yearMaxRevenue = new HashMap<>();
11     private Map<String, String> yearMaxMonth = new HashMap<>();
12
13     @Override
14     public void reduce(Text key, Iterable<LongWritable> values, Context context) throws IOException, InterruptedException {
15         String keyStr = key.toString();
16         String year = keyStr.split(regex:"-")[0];
17
18         long revenueSum = 0;
19         for (LongWritable val : values) {
20             revenueSum += val.get();
21         }
22
23         if (!yearMaxRevenue.containsKey(year) || revenueSum > yearMaxRevenue.get(year)) {
24             yearMaxRevenue.put(year, revenueSum);
25             yearMaxMonth.put(year, keyStr);
26         }
27     }
28
29     @Override
30     protected void cleanup(Context context) throws IOException, InterruptedException {
31         for (Map.Entry<String, String> entry : yearMaxMonth.entrySet()) {
32             context.write(new Text(entry.getValue()), new LongWritable(yearMaxRevenue.get(entry.getKey())));
33         }
34     }
35 }
```

2.6 Driver code for finding most profitable month

```
J Driver.java M X
src > main > java > org > example > J Driver.java > ...
1  package org.example;
2
3  import org.apache.hadoop.conf.Configuration;
4  import org.apache.hadoop.fs.Path;
5  import org.apache.hadoop.io.LongWritable;
6  import org.apache.hadoop.io.Text;
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 Driver {
12     Run | Debug
13     public static void main(String[] args) throws Exception {
14         if (args.length != 2) {
15             System.err.println(x:"Usage: App <input path> <output path>");
16             System.exit(-1);
17         }
18
19         Configuration conf = new Configuration();
20         Job job = Job.getInstance(conf, "Top Month");
21         job.setJarByClass(Driver.class);
22         job.setMapperClass(TopMonthMapper.class);
23         job.setReducerClass(TopMonthReducer.class);
24
25         job.setMapOutputKeyClass(Text.class);
26         job.setMapOutputValueClass(LongWritable.class);
27         job.setOutputKeyClass(Text.class);
28         job.setOutputValueClass(LongWritable.class);
29
30         FileInputFormat.addInputPath(job, new Path(args[0]));
31         FileOutputFormat.setOutputPath(job, new Path(args[1]));
32
33         System.exit(job.waitForCompletion(true) ? 0 : 1);
34     }
35 }
```

3 Results

3.1 Estimated time

3.1.1 Time taken for combining the datasets

```
Total time spent by all map tasks (ms)=6383
Total time spent by all reduce tasks (ms)=2660
```

3.1.2 Time taken for computing the most profitable months

```
Total time spent by all map tasks (ms)=2341
Total time spent by all reduce tasks (ms)=2443
```

3.2 Most profitable month

The most profitable months for each year calculated by our MapReduce logic are mentioned below:

- 2015 - August (1698057)
- 2016 - August (2896166)
- 2017 - October (477451)
- 2018 - October (1126826)

4 Team member contribution

4.1 Ankit Mane:

Loaded the datasets into HDFS. Implemented the logic for MapReduce to find the most profitable month of the year for the reservations made by the customers based on revenue from both the given datasets.

4.2 Daiyaan Ahmed Shaik:

Developed the logic for combining both the datasets, modified the dataset for faster execution and made the report.