

Write mapper and reducer program and then compile a jar of it and use hadoop command to run them.

Run the jar using following command

```
hadoop jar <jar_name> <class_name> <input_path_of_file> <output_path_HDFS_location>
```

### **Task 1:**

Driver class(TVSalesInvalidDataJob.java)

```
package mapreduce.assignment.Task1;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

public class TVSalesInvalidDataJob {
    @SuppressWarnings("deprecation")
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Remove Invalid Data");
        job.setJarByClass(TVSalesInvalidDataJob.class);

        // Set number of reducer to 0
        job.setNumReduceTasks(0);

        // Set output key and Value classes for Mapper
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(Text.class);

        // Set Mapper class
        job.setMapperClass(TVSalesInvalidDataMapper.class);

        // Set Input and Output format class
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);

        // Provide input and output path
        FileInputFormat.addInputPath(job, new Path(args[1]));
        FileOutputFormat.setOutputPath(job, new Path(args[2]));

        // execute the job
        job.waitForCompletion(true);
    }
}
```

### Mapper class(TVSalesInvalidDataMapper.java)

```
package mapreduce.assignment.Task1;

import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

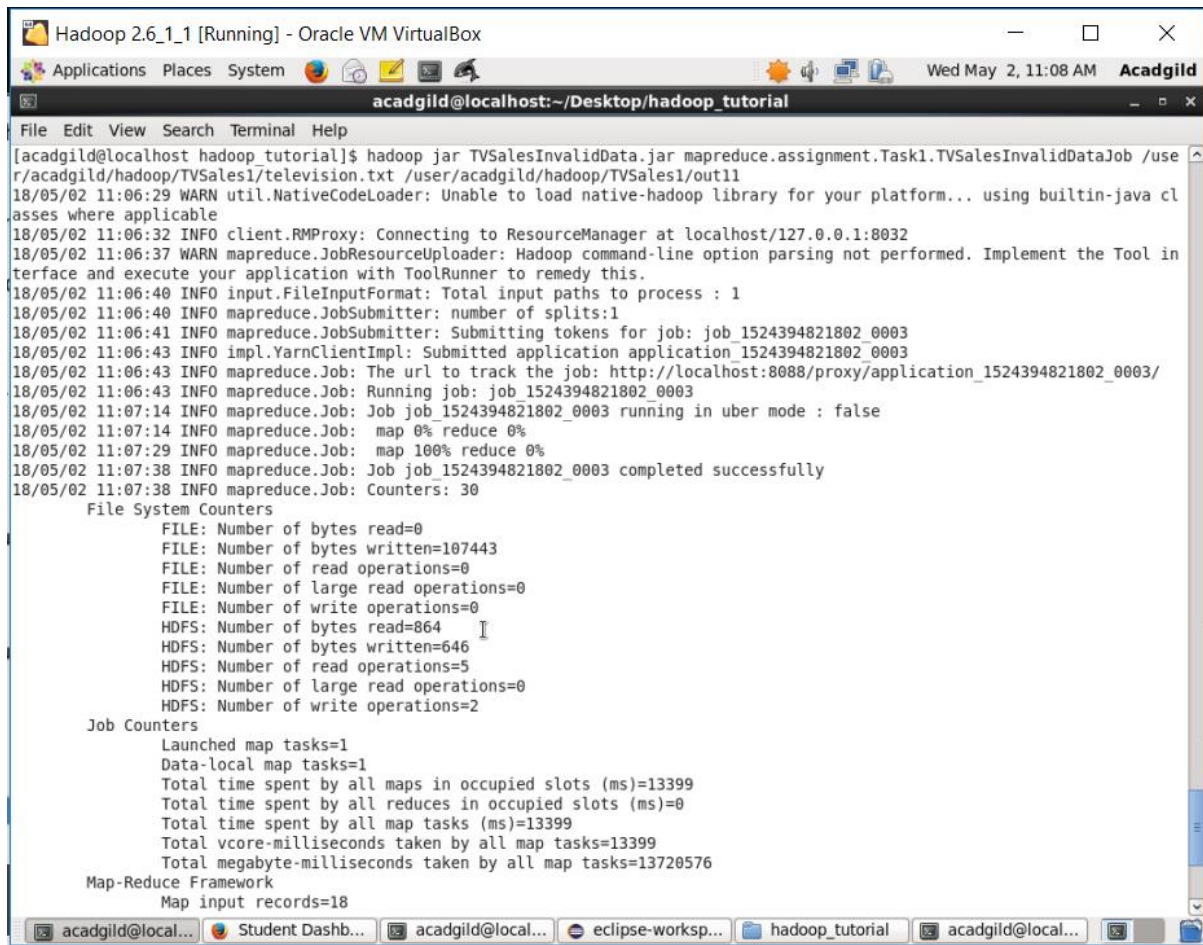
public class TVSalesInvalidDataMapper extends Mapper<LongWritable,
Text, Text, Text> {
    public void map(LongWritable key, Text value, Context context)
throws IOException, InterruptedException {

        Text word = new Text();
        String wholeLine = value.toString();
        int count = 0;
        StringTokenizer tokenWords = new
StringTokenizer(wholeLine, "|");

        // loop till tokens are available
        while(tokenWords.hasMoreTokens()) {
            word.set(tokenWords.nextToken());
            // finding if a string contains NA
            if(word.toString().equalsIgnoreCase("NA")) {
                count++;
            }
        }
        // excluding record that had invalid data
        if(count == 0) {
            Text newLine = new Text(wholeLine);
            context.write(newLine, null);
        }

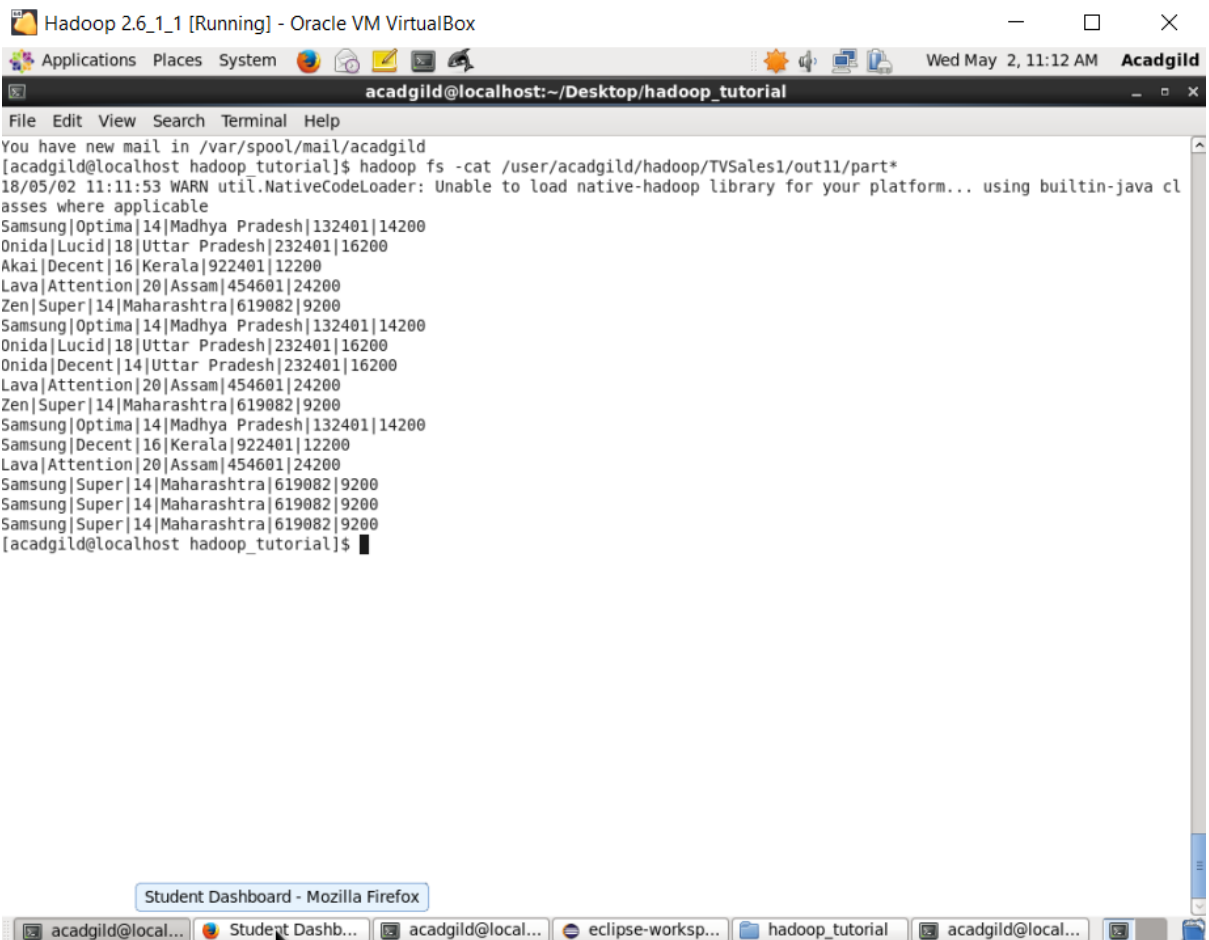
    }
}
```

Output:



```
[acadgild@localhost hadoop tutorial]$ hadoop jar TVSalesInvalidData.jar mapreduce.assignment.Task1.TVSalesInvalidDataJob /user/acadgild/hadoop/TVSales1/television.txt /user/acadgild/hadoop/TVSales1/out11
18/05/02 11:06:29 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
18/05/02 11:06:32 INFO client.RMPProxy: Connecting to ResourceManager at localhost/127.0.0.1:8032
18/05/02 11:06:37 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
18/05/02 11:06:40 INFO input.FileInputFormat: Total input paths to process : 1
18/05/02 11:06:40 INFO mapreduce.JobSubmitter: number of splits:1
18/05/02 11:06:41 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524394821802_0003
18/05/02 11:06:43 INFO impl.YarnClientImpl: Submitted application application_1524394821802_0003
18/05/02 11:06:43 INFO mapreduce.Job: The url to track the job: http://localhost:8088/proxy/application_1524394821802_0003/
18/05/02 11:06:43 INFO mapreduce.Job: Running job: job_1524394821802_0003
18/05/02 11:07:14 INFO mapreduce.Job: Job job_1524394821802_0003 running in uber mode : false
18/05/02 11:07:14 INFO mapreduce.Job: map 0% reduce 0%
18/05/02 11:07:29 INFO mapreduce.Job: map 100% reduce 0%
18/05/02 11:07:38 INFO mapreduce.Job: Job job_1524394821802_0003 completed successfully
18/05/02 11:07:38 INFO mapreduce.Job: Counters: 30
  File System Counters
    FILE: Number of bytes read=0
    FILE: Number of bytes written=107443
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=864
    HDFS: Number of bytes written=646
    HDFS: Number of read operations=5
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=13399
    Total time spent by all reduces in occupied slots (ms)=0
    Total time spent by all map tasks (ms)=13399
    Total vcore-milliseconds taken by all map tasks=13399
    Total megabyte-milliseconds taken by all map tasks=13720576
  Map-Reduce Framework
    Map input records=18
```

```
Hadoop 2.6_1_1 [Running] - Oracle VM VirtualBox
Applications Places System
acadgild@localhost: ~/Desktop/hadoop_tutorial
File Edit View Search Terminal Help
18/05/02 11:07:38 INFO mapreduce.Job: Job job_1524394821802_0003 completed successfully
18/05/02 11:07:38 INFO mapreduce.Job: Counters: 30
  File System Counters
    FILE: Number of bytes read=0
    FILE: Number of bytes written=107443
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=864
    HDFS: Number of bytes written=646
    HDFS: Number of read operations=5
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=13399
    Total time spent by all reduces in occupied slots (ms)=0
    Total time spent by all map tasks (ms)=13399
    Total vcore-milliseconds taken by all map tasks=13399
    Total megabyte-milliseconds taken by all map tasks=13720576
  Map-Reduce Framework
    Map input records=18
    Map output records=16
    Input split bytes=131
    Spilled Records=0
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=94
    CPU time spent (ms)=800
    Physical memory (bytes) snapshot=88141824
    Virtual memory (bytes) snapshot=2056757248
    Total committed heap usage (bytes)=32571392
  File Input Format Counters
    Bytes Read=733
  File Output Format Counters
    Bytes Written=646
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost hadoop_tutorial]$
```



## **Task 2:**

Driver class(TVSalesCountByCompanyJob.java)

```
package mapreduce.assignment.Task2;

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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

public class TVSalesCountByCompanyJob {
    @SuppressWarnings("deprecation")
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Total Unit sold");
        job.setJarByClass(TVSalesCountByCompanyJob.class);

        //set mapper output key and value
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(IntWritable.class);

        //set reducer output key and value
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        //set the mapper and reducer class
        job.setMapperClass(TVSalesCountByCompanyMapper.class);
        job.setReducerClass(TVSalesCountByCompanyReducer.class);

        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);

        //input file and the output path
        FileInputFormat.addInputPath(job, new Path(args[1]));
        FileOutputFormat.setOutputPath(job, new Path(args[2]));

        job.waitForCompletion(true);
    }
}
```

### Mapper class(TVSalesCountByCompanyMapper.java)

```
package mapreduce.assignment.Task2;

import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class TVSalesCountByCompanyMapper extends Mapper<LongWritable,
Text, Text, IntWritable>{
    public void map(LongWritable key, Text value, Context context)
throws IOException, InterruptedException {

        //split the string
        String[] lineArray = value.toString().split("\\|");
        //first element of array is the company name
        Text companyName = new Text(lineArray[0]);

        //add 1 for an occurrence of the company name
        context.write(companyName, new IntWritable(1));

    }
}
```

## Reducer class(TVSalesCountByCompanyReducer.java)

```
package mapreduce.assignment.Task2;

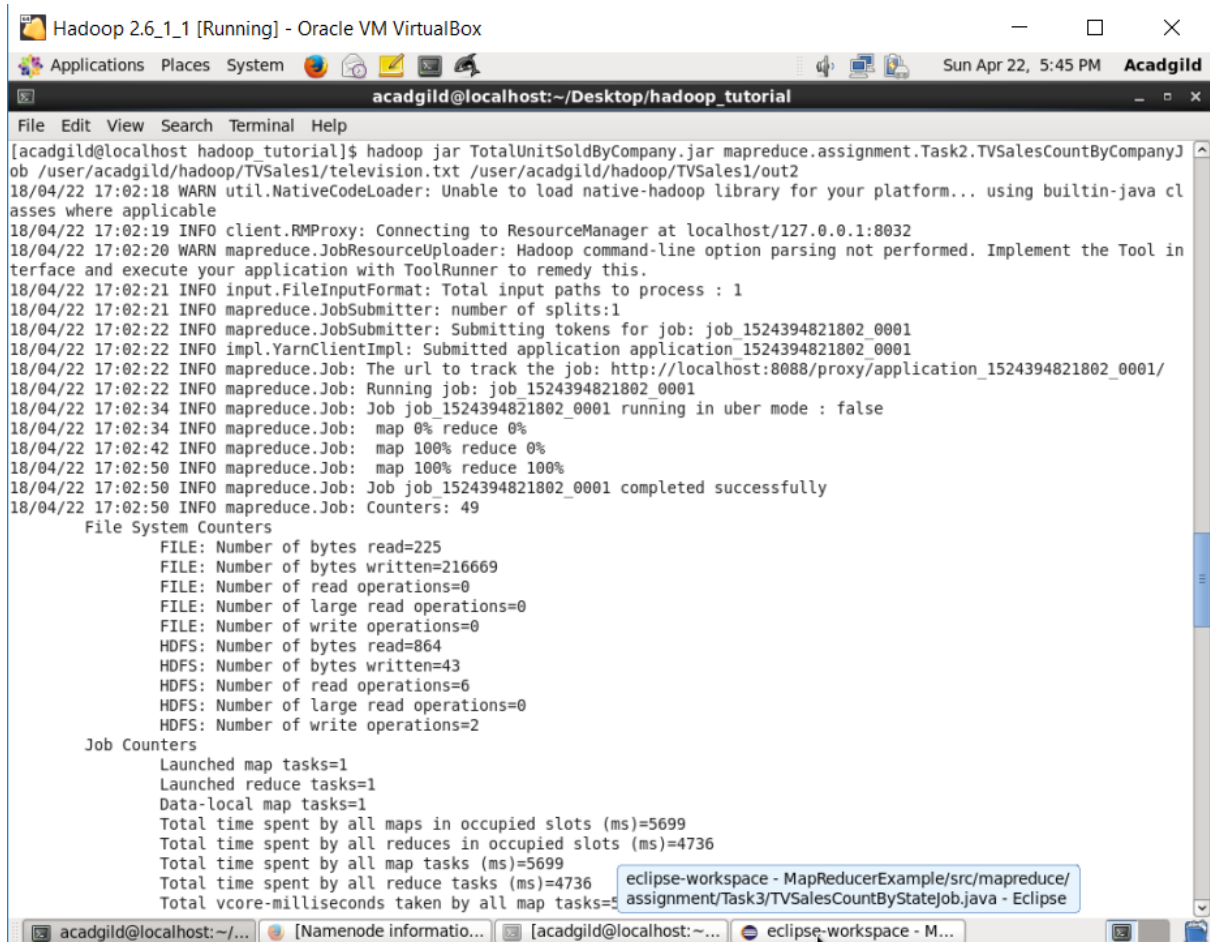
import java.io.IOException;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class TVSalesCountByCompanyReducer extends Reducer<Text,
IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable> values,
Context context) throws IOException, InterruptedException {
        int count = 0;
        for(IntWritable value : values) {
            //increase the count
            count += value.get();
        }
        context.write(key, new IntWritable(count));
    }
}
```



Output:



The screenshot shows a terminal window titled "Hadoop 2.6\_1\_1 [Running] - Oracle VM VirtualBox". The terminal displays the output of a Hadoop MapReduce job. The command executed is `hadoop jar TotalUnitSoldByCompany.jar mapreduce.assignment.Task2.TVSalesCountByCompanyJob /user/acadgild/hadoop/TVSales1/television.txt /user/acadgild/hadoop/TVSales1/out2`. The output includes various status messages, progress reports, and a final summary of counters.

```
[acadgild@localhost hadoop tutorial]$ hadoop jar TotalUnitSoldByCompany.jar mapreduce.assignment.Task2.TVSalesCountByCompanyJob /user/acadgild/hadoop/TVSales1/television.txt /user/acadgild/hadoop/TVSales1/out2
18/04/22 17:02:18 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
18/04/22 17:02:19 INFO client.RMProxy: Connecting to ResourceManager at localhost/127.0.0.1:8032
18/04/22 17:02:20 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool in interface and execute your application with ToolRunner to remedy this.
18/04/22 17:02:21 INFO input.FileInputFormat: Total input paths to process : 1
18/04/22 17:02:21 INFO mapreduce.JobSubmitter: number of splits:1
18/04/22 17:02:22 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524394821802_0001
18/04/22 17:02:22 INFO impl.YarnClientImpl: Submitted application application_1524394821802_0001
18/04/22 17:02:22 INFO mapreduce.Job: The url to track the job: http://localhost:8088/proxy/application_1524394821802_0001/
18/04/22 17:02:22 INFO mapreduce.Job: Running job: job_1524394821802_0001
18/04/22 17:02:34 INFO mapreduce.Job: Job job_1524394821802_0001 running in uber mode : false
18/04/22 17:02:34 INFO mapreduce.Job: map 0% reduce 0%
18/04/22 17:02:42 INFO mapreduce.Job: map 100% reduce 0%
18/04/22 17:02:50 INFO mapreduce.Job: map 100% reduce 100%
18/04/22 17:02:50 INFO mapreduce.Job: Job job_1524394821802_0001 completed successfully
18/04/22 17:02:50 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=225
    FILE: Number of bytes written=216669
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=864
    HDFS: Number of bytes written=43
    HDFS: Number of read operations=6
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=5699
    Total time spent by all reduces in occupied slots (ms)=4736
    Total time spent by all map tasks (ms)=5699
    Total time spent by all reduce tasks (ms)=4736
    Total vcore-milliseconds taken by all map tasks=5699
```

Hadoop 2.6\_1\_1 [Running] - Oracle VM VirtualBox

Applications Places System Sun Apr 22, 5:47 PM Acadgild

acadgild@localhost: ~/Desktop/hadoop\_tutorial

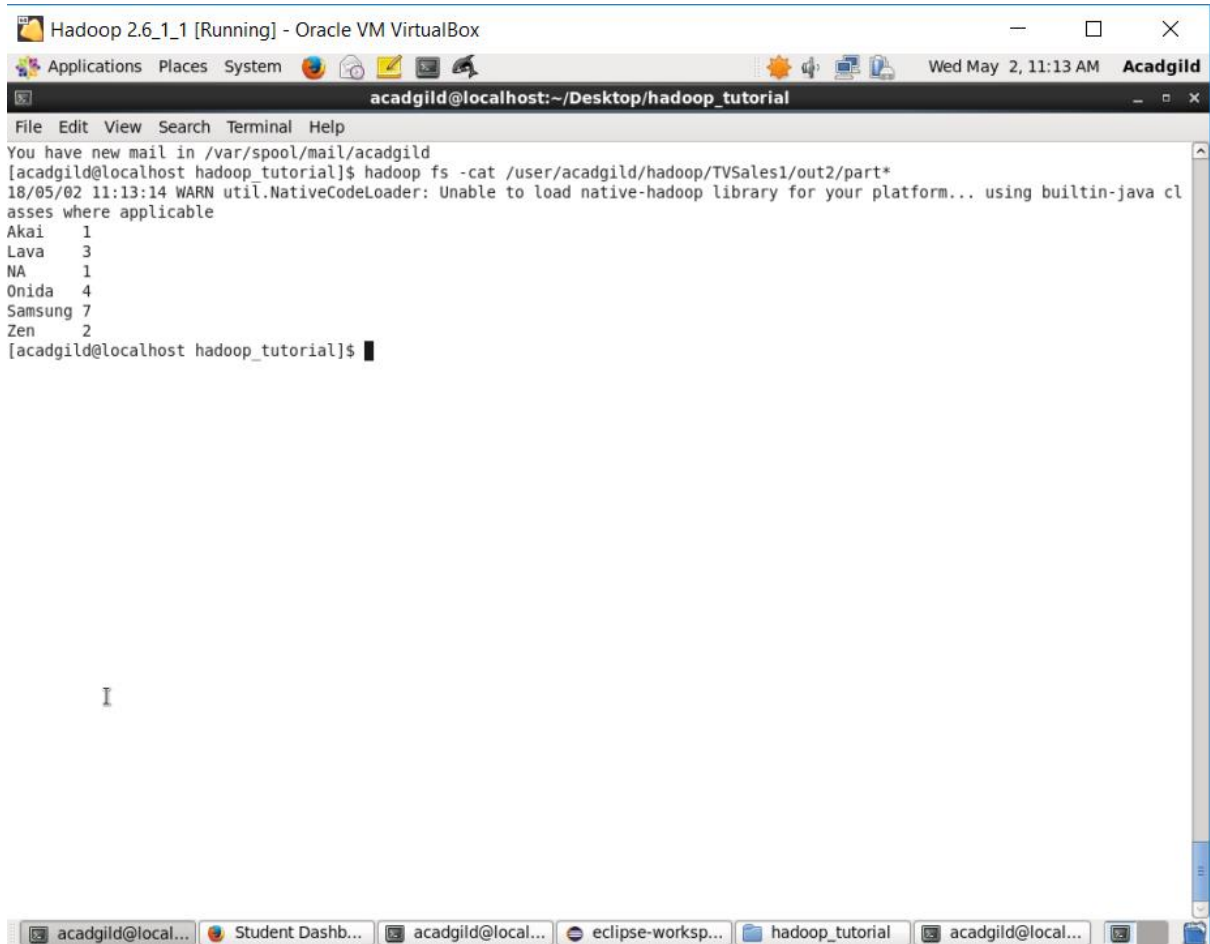
File Edit View Search Terminal Help

```
Total time spent by all map tasks (ms)=5699
Total time spent by all reduce tasks (ms)=4736
Total vcore-milliseconds taken by all map tasks=5699
Total vcore-milliseconds taken by all reduce tasks=4736
Total megabyte-milliseconds taken by all map tasks=5835776
Total megabyte-milliseconds taken by all reduce tasks=4849664
Map-Reduce Framework
  Map input records=18
  Map output records=18
  Map output bytes=183
  Map output materialized bytes=225
  Input split bytes=131
  Combine input records=0
  Combine output records=0
  Reduce input groups=6
  Reduce shuffle bytes=225
  Reduce input records=18
  Reduce output records=6
  Spilled Records=36
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=119
  CPU time spent (ms)=1090
  Physical memory (bytes) snapshot=287375360
  Virtual memory (bytes) snapshot=4118216704
  Total committed heap usage (bytes)=170004480
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=733
File Output Format Counters
  Bytes Written=42
```

You have new mail in /

Namenode information - Mozilla Firefox

acadgild@localhost: ~/... [Namenode informatio... acadgild@localhost: ~/... eclipse-workspace - M...



### **Task 3:**

Driver class(TVSalesCountByStateJob.java)

```
package mapreduce.assignment.Task3;

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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

public class TVSalesCountByStateJob {
    @SuppressWarnings("deprecation")
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Total Onida Unit sold ");
        job.setJarByClass(TVSalesCountByStateJob.class);

        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(IntWritable.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.setMapperClass(TVSalesCountByStateMapper.class);
        job.setReducerClass(TVSalesCountByStateReducer.class);

        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);

        FileInputFormat.addInputPath(job, new Path(args[1]));
        FileOutputFormat.setOutputPath(job, new Path(args[2]));

        job.waitForCompletion(true);
    }
}
```

### Mapper class(TVSalesCountByStateMapper.java)

```
package mapreduce.assignment.Task3;

import java.io.IOException;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class TVSalesCountByStateMapper extends Mapper<LongWritable,
Text, Text, IntWritable> {
    public void map(LongWritable key, Text value, Context context)
throws IOException, InterruptedException {

        String[] lineArray = value.toString().split("\\\\|");
        if(lineArray[0].equalsIgnoreCase("ONIDA")) {
            Text stateName = new Text(lineArray[3]);
            context.write(stateName, new IntWritable(1));
        }
        else {
            Text stateName = new Text(lineArray[3]);
            context.write(stateName, new IntWritable(0));
        }
    }
}
```

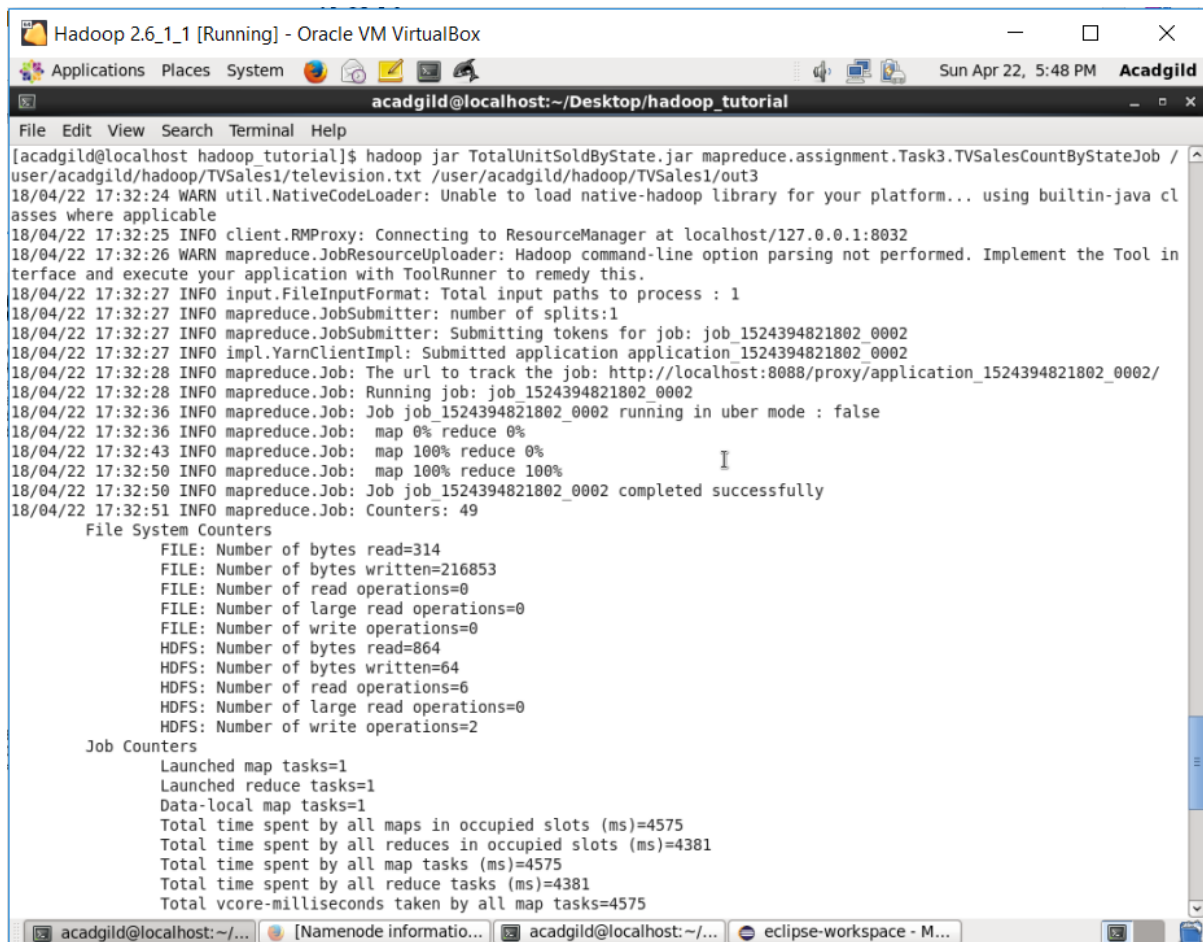
### Reducer class(TVSalesCountByStateReducer.java)

```
package mapreduce.assignment.Task3;

import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class TVSalesCountByStateReducer extends Reducer<Text,
IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable> values,
Context context) throws IOException, InterruptedException {
        int count = 0;
        for(IntWritable value : values) {
            count += value.get();
        }
        context.write(key, new IntWritable(count));
    }
}
```

Output:



```
Hadoop 2.6.1_1 [Running] - Oracle VM VirtualBox
Applications Places System
acadgild@localhost:~/Desktop/hadoop_tutorial
File Edit View Search Terminal Help
[acadgild@localhost hadoop_tutorial]$ hadoop jar TotalUnitSoldByState.jar mapreduce.assignment.Task3.TVSalesCountByStateJob /
user/acadgild/hadoop/TVSales1/television.txt /user/acadgild/hadoop/TVSales1/out3
18/04/22 17:32:24 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
18/04/22 17:32:25 INFO client.RMProxy: Connecting to ResourceManager at localhost/127.0.0.1:8032
18/04/22 17:32:26 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool in
terface and execute your application with ToolRunner to remedy this.
18/04/22 17:32:27 INFO input.FileInputFormat: Total input paths to process : 1
18/04/22 17:32:27 INFO mapreduce.JobSubmitter: number of splits:1
18/04/22 17:32:27 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1524394821802_0002
18/04/22 17:32:27 INFO impl.YarnClientImpl: Submitted application application_1524394821802_0002
18/04/22 17:32:28 INFO mapreduce.Job: The url to track the job: http://localhost:8088/proxy/application_1524394821802_0002/
18/04/22 17:32:28 INFO mapreduce.Job: Running job: job_1524394821802_0002
18/04/22 17:32:36 INFO mapreduce.Job: Job job_1524394821802_0002 running in uber mode : false
18/04/22 17:32:36 INFO mapreduce.Job: map 0% reduce 0%
18/04/22 17:32:43 INFO mapreduce.Job: map 100% reduce 0%
18/04/22 17:32:50 INFO mapreduce.Job: map 100% reduce 100%
18/04/22 17:32:50 INFO mapreduce.Job: Job job_1524394821802_0002 completed successfully
18/04/22 17:32:51 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=314
    FILE: Number of bytes written=216853
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=864
    HDFS: Number of bytes written=64
    HDFS: Number of read operations=6
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=4575
    Total time spent by all reduces in occupied slots (ms)=4381
    Total time spent by all map tasks (ms)=4575
    Total time spent by all reduce tasks (ms)=4381
    Total vcore-milliseconds taken by all map tasks=4575
```

Hadoop 2.6.1\_1 [Running] - Oracle VM VirtualBox

Applications Places System Sun Apr 22, 5:49 PM Acadgild

acadgild@localhost: ~/Desktop/hadoop\_tutorial

```
File Edit View Search Terminal Help
Total time spent by all map tasks (ms)=4575
Total time spent by all reduce tasks (ms)=4381
Total vcore-milliseconds taken by all map tasks=4575
Total vcore-milliseconds taken by all reduce tasks=4381
Total megabyte-milliseconds taken by all map tasks=4684800
Total megabyte-milliseconds taken by all reduce tasks=4486144
Map-Reduce Framework
  Map input records=18
  Map output records=18
  Map output bytes=272
  Map output materialized bytes=314
  Input split bytes=131
  Combine input records=0
  Combine output records=0
  Reduce input groups=5
  Reduce shuffle bytes=314
  Reduce input records=18
  Reduce output records=5
  Spilled Records=36
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=116
  CPU time spent (ms)=1040
  Physical memory (bytes) snapshot=287027200
  Virtual memory (bytes) snapshot=4117905408
  Total committed heap usage (bytes)=170004480
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=733
File Output Format Counters
  Bytes Written=64
You have new mail in /var/spool/mail/acadgild
acadgild@localhost: ~/eclipse/java-oxygen/eclipse
```

acadgild@localhost: ~/... [NameNode informatio... acadgild@localhost: ~/... eclipse-workspace - M...



