# BIG DATA WEEK 11 TASK REPORT
# ( JOBSHEET MAPREDUCE 3 )

**By:**

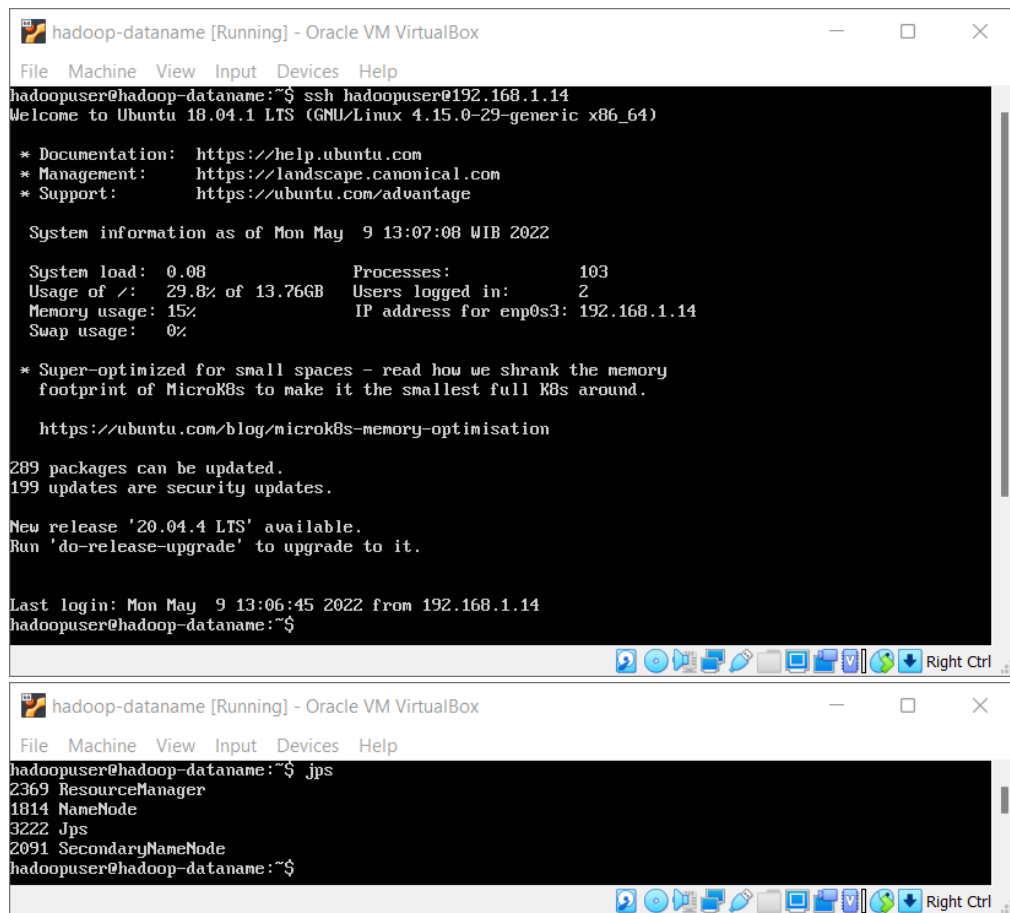**Rajendra Rakha Arya Prabaswara**

**1941720080**

**( TI-3H / 20 )**

**INFORMATICS ENGINEERING STUDY PROGRAM MAJORING INFORMATION TECHNOLOGY**
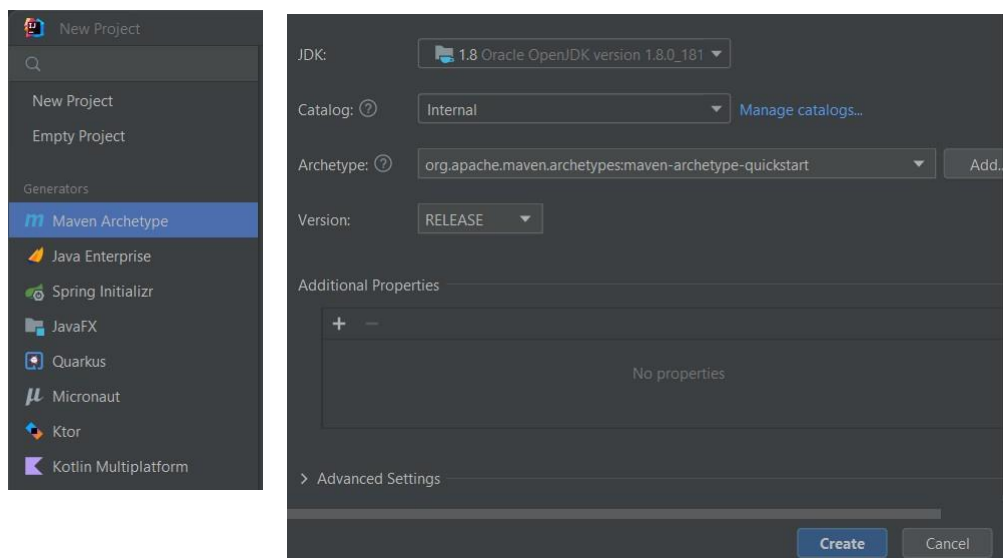
**STATE POLYTECHNIC MALANG**

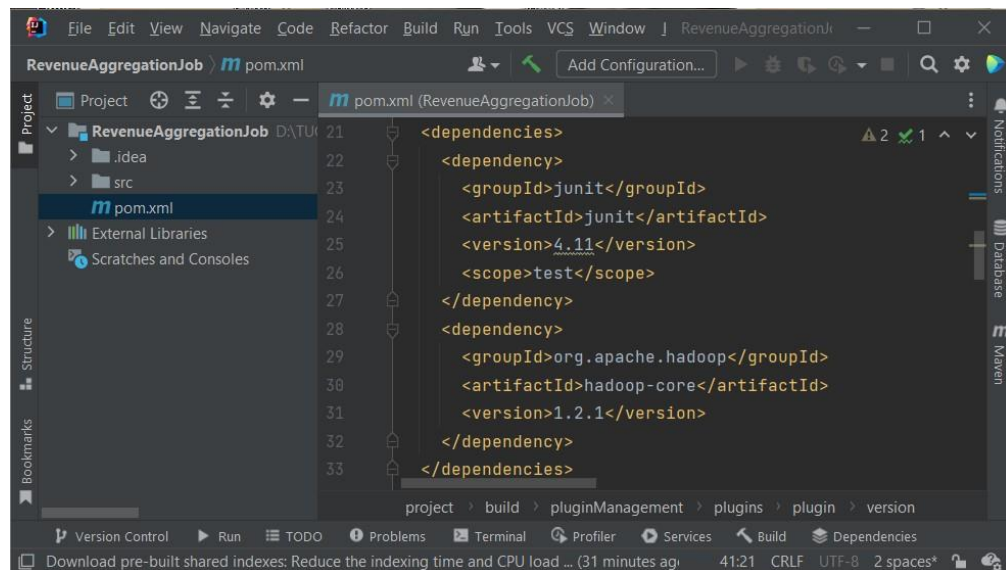**Part I : Creating a MapReduce Job with JetBrains IntelliJ IDEA**

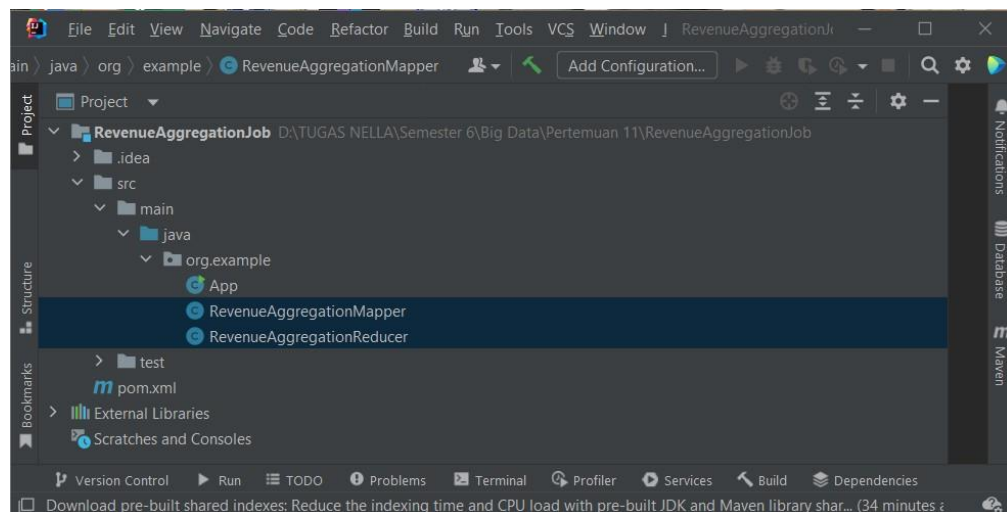1. Install and will be able to access the Hadoop cluster correctly .



2. Open the IntelliJ IDEA IDE IDE and create a new project with the Java Language and maven framework build. And select the archetype "maven quick-start".
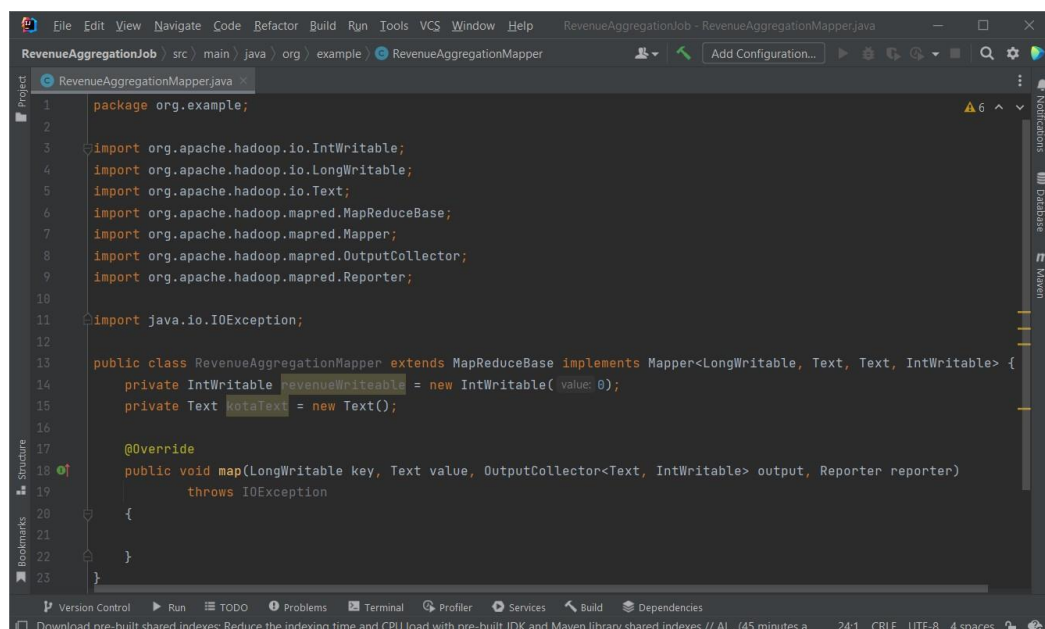
3. After the project is completed. In pom.xml added dependencies in order to create a MapReduce Job. Reload the project by clicking on the 'm' icon.
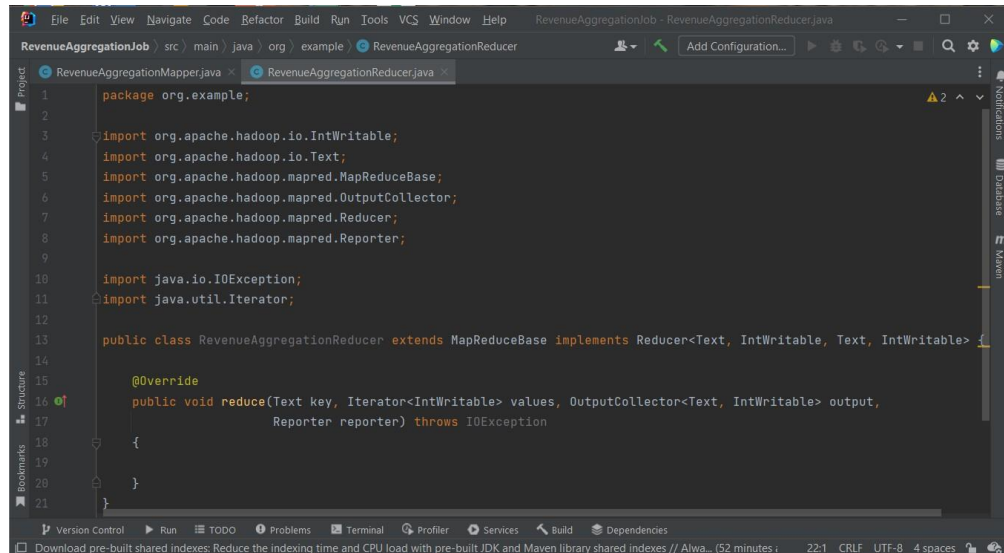


4. Create 2 new classes on the project under the names RevenueAggregationMapper and RevenueAggregationReducer.



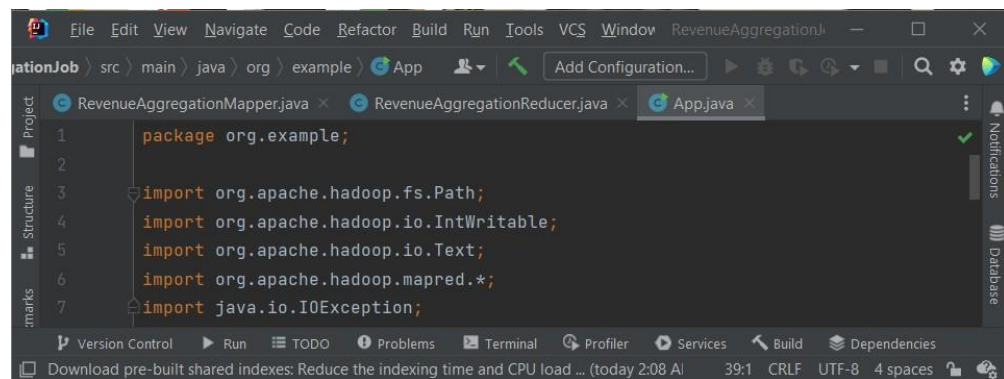5. Modify the Class RevenueAggregationMapper.
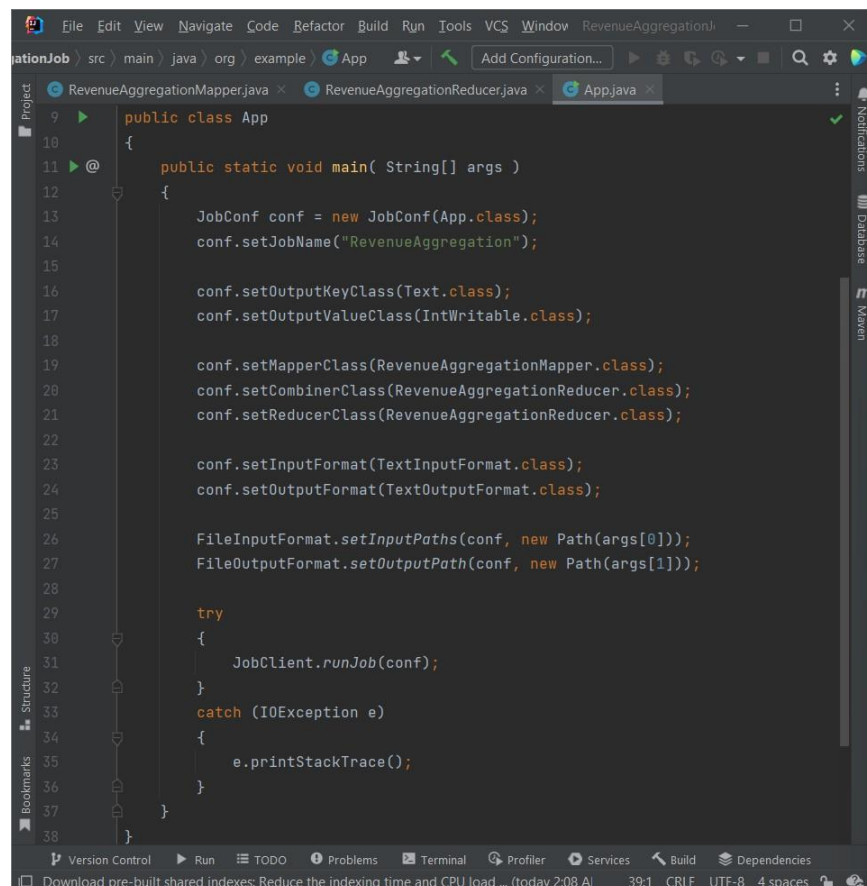
6. Modify the Class RevenueAggregationReducer.



7. Add imports blocks to the App class that is the project's built-in class.



8. Then, add the code to the same class body .

9. Compile the program using the help of Maven by opening the Maven panel located to the right of the IDE. In the pane expand the folder "Lifecycle", select "Clean" and "Installl". Then click Run.
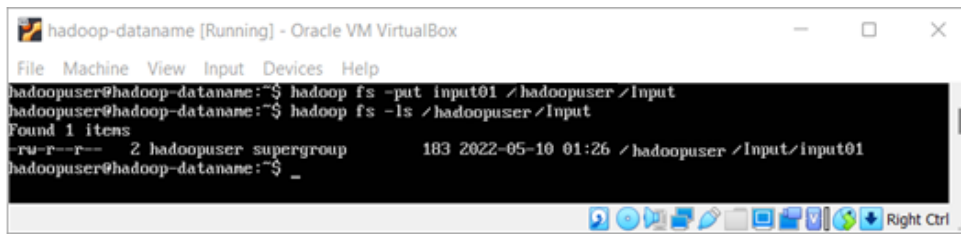


10. Create a plain text file named input01.



Upload the file to HDFS in the /YourName/Input folder. Also create 1 other folder in HDFS under the name /NamAnda/Output.

11. Upload jar files to the Name Node by using SCP or other methods. Execute the JAR using the command: Hadoop jar <nama_file>.jar <package_identifier> <folder_input_di_hdfs> <folder_output_di_hdfs>. JAR files uploaded to the namenode aredibe rinama revenue.jar.



12. When the execution process is complete, the folder /Your Name/Output/output/result in which there is a part-00000 file that when you



read, the contents are still empty. [Question] Why is the output file empty/has no content?

## Bagian II : Mapper & Reducer

13. Modify the RevenueAggregationMapper class to find out how many times Mapper is run.



```java
public class RevenueAggregationMapper extends MapReduceBase implements Mapper<LongWritable, Text, Text, Int
    private IntWritable revenueWriteable = new IntWritable( value: 0);
    private Text kotaText = new Text();


    @Override
    public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter)
            throws IOException
    {
        String line = value.toString();

        System.out.println("Ini adalah isi dari line");
        System.out.println(line);
    }
}
```
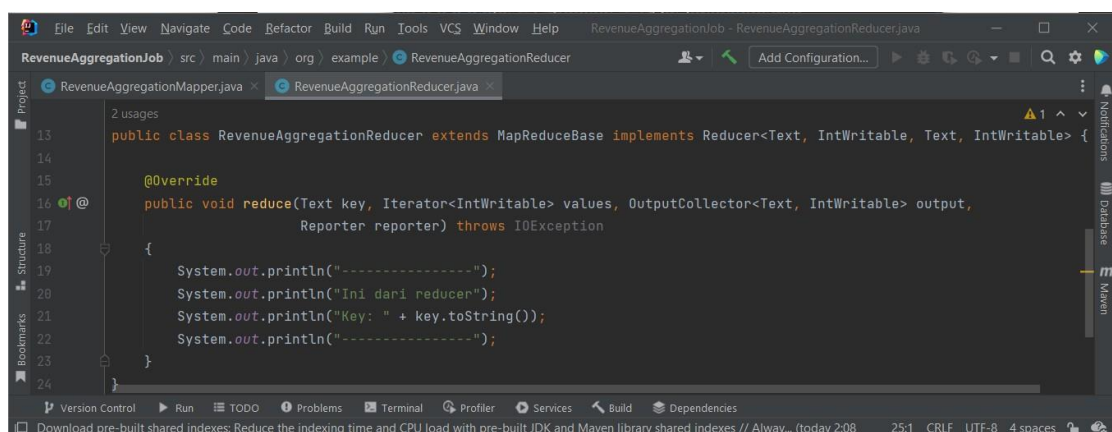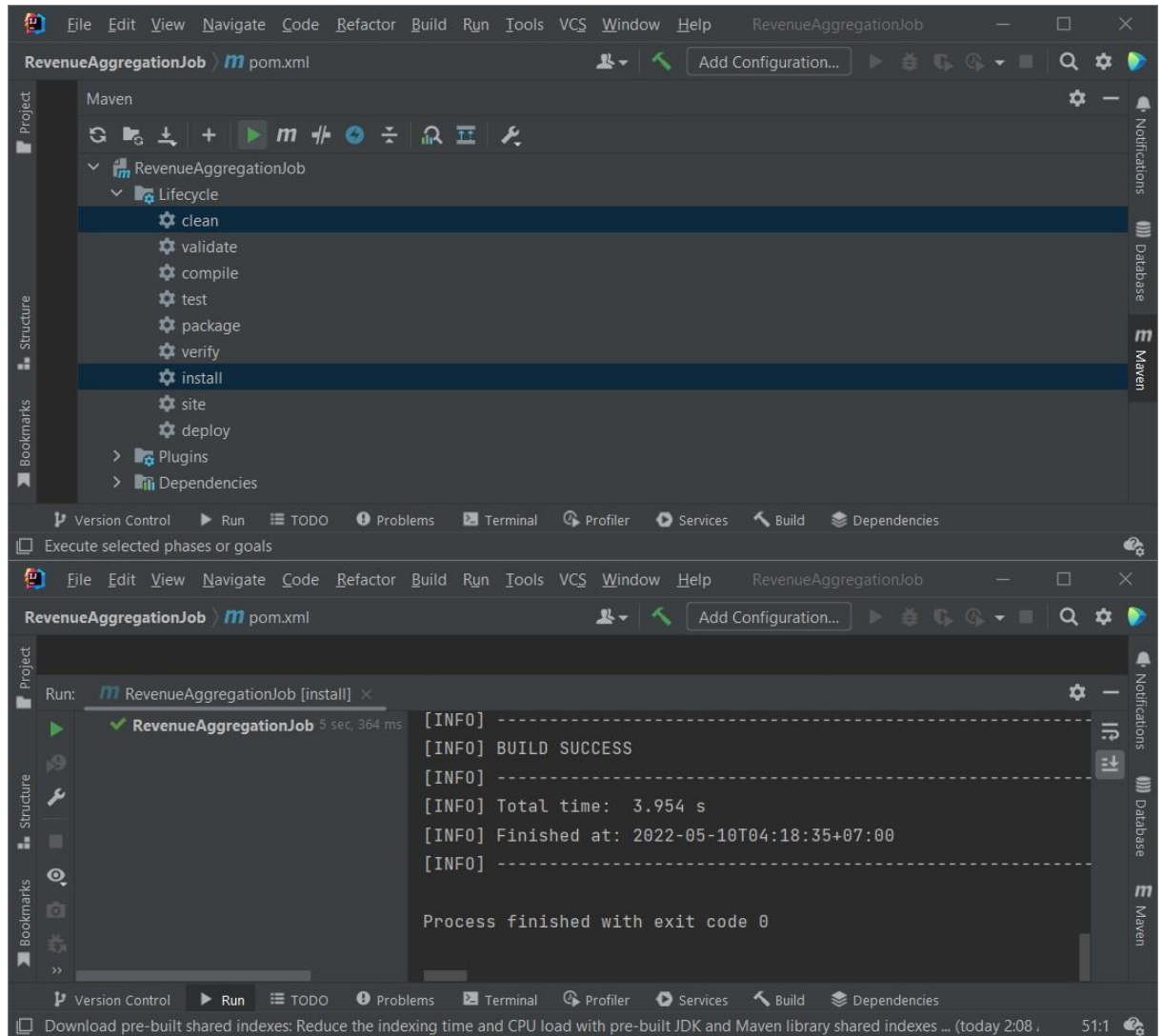
14. Modify the Class RevenueAggregationReducer.



```java
public class RevenueAggregationReducer extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {

    @Override
    public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output,
                       Reporter reporter) throws IOException
    {
        System.out.println("---------------");
        System.out.println("Ini dari reducer");
        System.out.println("Key: " + key.toString());
        System.out.println("---------------");
    }
}
```
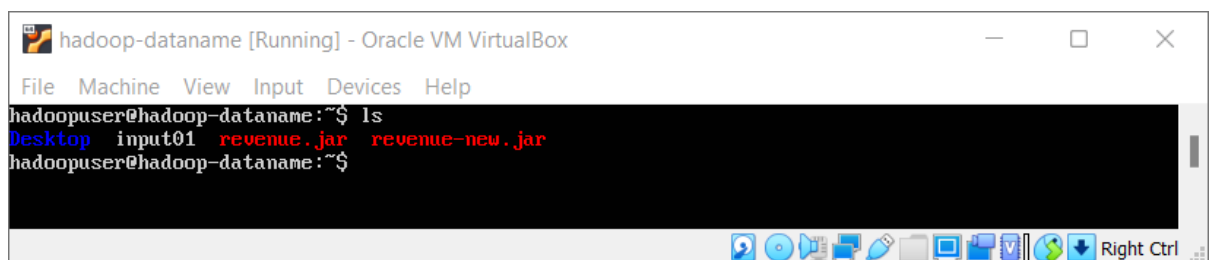
15. Rebuild jar again by using Maven.



Re-upload the new JAR to the name node.



Execution.

16. Scroll and/or search (Ctrl + F) on your console, is there a bookmark message from the Reducer class?
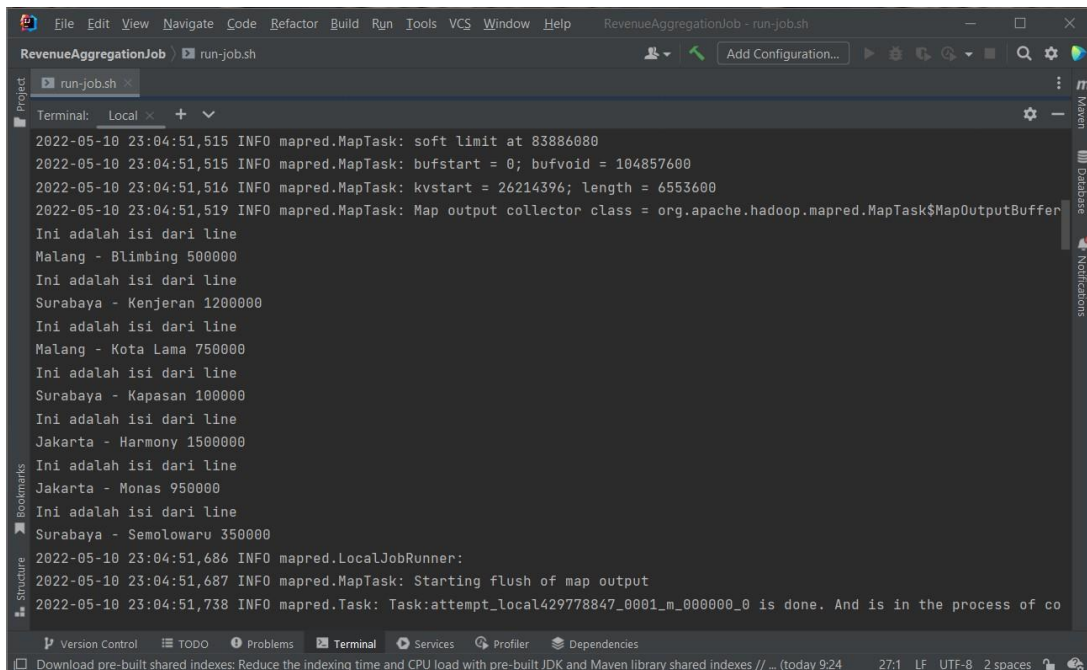


17. Add a new file. The name of the file: "run-job.sh"

18. Add a script to the file. Adjust the variables in the initial rows of the file to the conditions in the computer and cluster.



19. Execution

**Part III: Logika Aggregation Process**

20. Modify the Class RevenueAggregationMapper.



21. In the RevenueAggregationReducer class, modify the reduce() method.



22. Compile, deploy and rerun the MapReduce program and pay attention to the results. Print marker messages from the Reducer class can find them in the console. Which means, the reduce() method in the class, this time called.

23. Scroll to thebottom of the console, now the results of aggregation should be clearly visible and correct results.



## Part I Questions

1. Why is the output file empty/has no content?

   **Answer:**

   - **Because in the Mapper class and reducer class do not havea command to display the data to be stored, so the process is only run without being stored in the Results directory .**

## Question Part II

1. When and how many folders() in the Mapper class are executed?

   **Answer:**

   - **Method map() is executed when it can be input and the method will be executed as much data from the input.**

2. Why is the reduce() method in the Reducer class never followed?

Answer:

- **Because the mapper result data is not sent into the Reducer class so the reduce() method does not have the data to be executed.**

## Question Part III

1. When and how much is the reduce() method run?

Answer:

- **The method will be run when it can be input from the Mapper class and run 3 times because there are 3 keys, namely Malang, Surabaya, and Jakarta.**

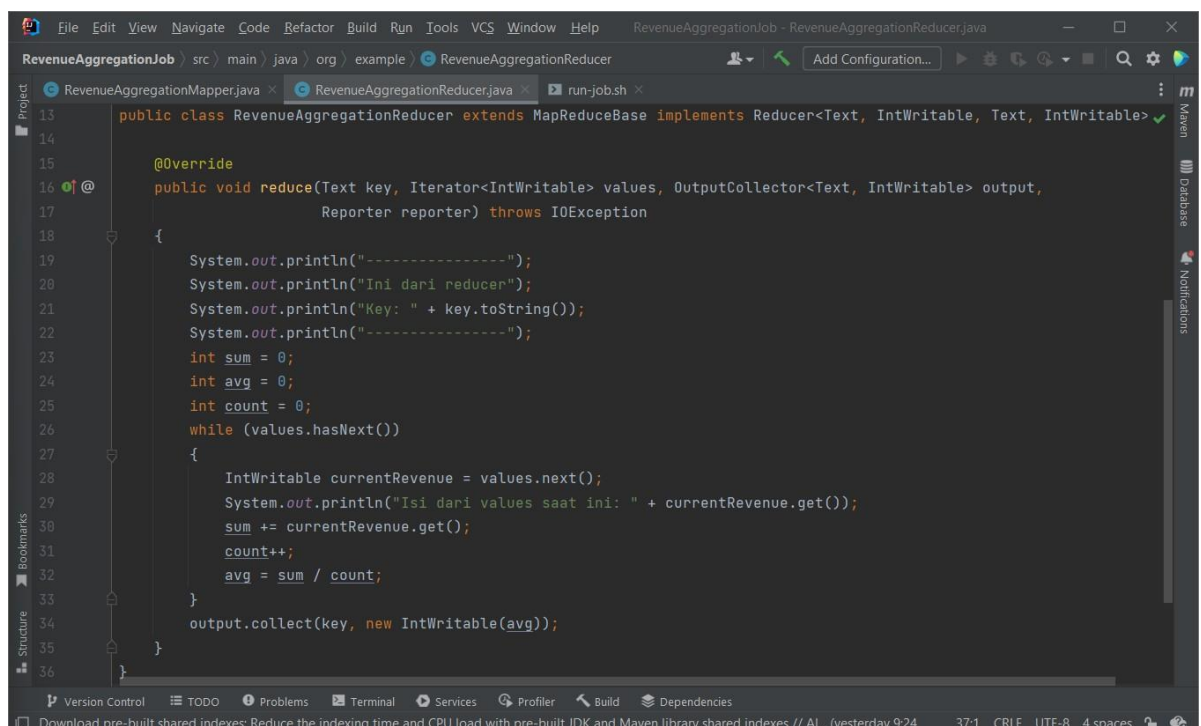2. Why can the console appear the contents of part-0000 files?

Answer:

- **Because the reducer result was successfully stored on the part-0000 file.**

## Assignment

Try to replace its aggregation operations to be average per city instead of counting the totals!

## Task Result:

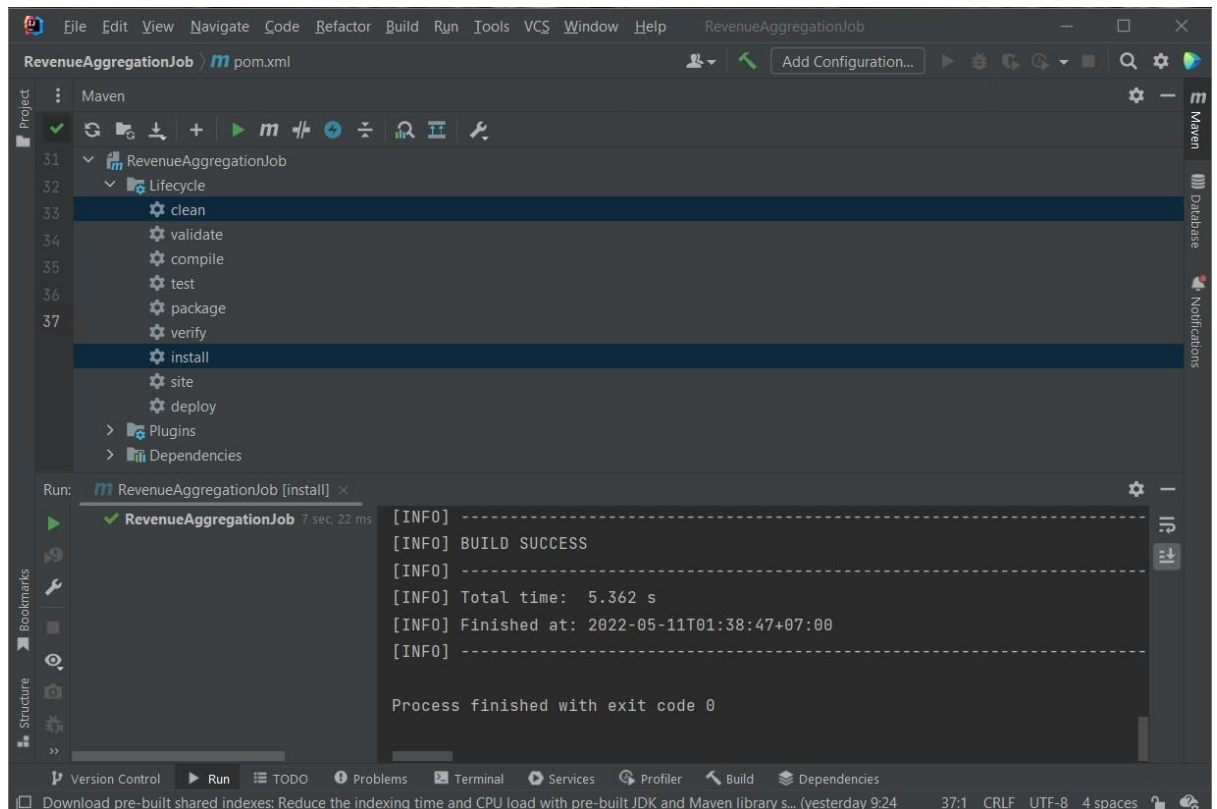- Step 1: Modification of the RevenueAggregationReducer class

- Step 2: Compile and Build Project



- Startl

Isi dari values saat ini: 950000
Isi dari values saat ini: 1500000
----------------
Ini dari reducer
Key: Malang
----------------
Isi dari values saat ini: 750000
Isi dari values saat ini: 500000
----------------
Ini dari reducer
Key: Surabaya
----------------
Isi dari values saat ini: 350000

----------------
Ini dari reducer
Key: Jakarta
----------------
Isi dari values saat ini: 1225000
----------------
Ini dari reducer
Key: Malang
----------------
Isi dari values saat ini: 625000
----------------