

11. *What is the primary purpose of the following statement?*

Clean the dataset by

rices & Clean the non-number fields):

15  
16  
17  
18  
19

A screenshot of a Linux desktop environment. The top bar includes a menu with 'Edit', 'View', 'Window', and 'Help' options, and a system status bar showing battery level (100%), signal strength, and temperature (37°). The desktop has a dark theme with a dock containing icons for file, terminal, and system settings. A terminal window is open with the following content:

```
drn — hassan@hassan-vm: ~/Desktop/Lab4_Project — ssh hassan@172.16.186.132 - 124x35
[hassan@hassan-vm:~/Desktop/Lab4_Project$ python3 DataSetLoad.py
Prices Saved to prices_data.txt!
```

34.99  
28.91  
17.49  
18.66  
29.12  
97.68  
12.99  
28.49

38.49  
18.16  
84.61  
33.92  
14.99  
4.99  
34.39  
12.88  
117.26  
9.39  
17.85  
27.5  
10.99  
159.99  
12.63  
11.88  
26.99  
34.27  
18.7  
14.95  
6.91  
5.99  
36.37

## Task 3: Writing MapReduce Program (Mapper, Reducer, Driver, Combiner)

### 1. Mapper Class:

```
1 package edu.hassan;
2 import java.io.IOException;
3 import org.apache.hadoop.io.DoubleWritable;
4 import org.apache.hadoop.io.LongWritable;
5 import org.apache.hadoop.io.Text;
6 import org.apache.hadoop.mapreduce.Mapper;
7
8 public class MapperClass extends Mapper<LongWritable, Text, Text, DoubleWritable> {
9     private final static Text word = new Text("price");
10
11     @Override
12     protected void map(LongWritable key, Text value, Context context) throws IOException,
13     InterruptedException {
13         double price = Double.parseDouble(value.toString().trim());
14         context.write(word, new DoubleWritable(price));
15     }
16 }
17 }
```

### 2. Reducer Class:

```
1 package edu.hassan;
2 import java.io.IOException;
3 import org.apache.hadoop.io.DoubleWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Reducer;
6
7 public class ReducerClass extends Reducer<Text, DoubleWritable, Text, DoubleWritable> {
8     private double sum = 0, min = Double.MAX_VALUE, max = Double.MIN_VALUE, sumOfSquares =
9 0; private int count = 0;
10
11     @Override
12     protected void reduce(Text key, Iterable<DoubleWritable> values, Context context){
13
14         for (DoubleWritable val : values) {
15             double value = val.get();
16
17             switch (key.toString()) {
18                 case "Count":
19                     count += (int) value;
20                     break;
21                 case "Sum":
22                     sum += value;
23                     break;
24                 case "Min":
25                     min = Math.min(min, value);
26                     break;
27                 case "Max":
28                     max = Math.max(max, value);
29                     break;
30                 case "SumOfSquares":
31                     sumOfSquares += value;
32                     break;
33             }
34         }
35     }
36     @Override
37     protected void cleanup(Context context) throws IOException, InterruptedException {
38         double mean = sum / count;
39         double variance = (sumOfSquares / count) - (mean * mean);
40         double stdDev = Math.sqrt(variance);
41
42         context.write(new Text("Count"), new DoubleWritable(count));
43         context.write(new Text("Sum"), new DoubleWritable(sum));
44         context.write(new Text("Min"), new DoubleWritable(min));
45         context.write(new Text("Max"), new DoubleWritable(max));
46         context.write(new Text("Mean"), new DoubleWritable(mean));
47         context.write(new Text("Standard Deviation"), new DoubleWritable(stdDev));
48     }
49 }
50
51 }
```

### 3. Driver Class:

```
1 package edu.hassan;
2 import org.apache.hadoop.conf.Configuration;
3 import org.apache.hadoop.fs.Path;
4 import org.apache.hadoop.io.DoubleWritable;
5 import org.apache.hadoop.io.Text;
6 import org.apache.hadoop.mapreduce.Job;
7 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
8 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
9
10 public class DriverClass {
11
12     public static void main(String[] args) throws Exception {
13         if (args.length != 2) {
14             System.err.println("Usage: StatisticalMeasures <input path> <output path>");
15             System.exit(-1);
16         }
17         Configuration conf = new Configuration();
18         Job job = Job.getInstance(conf, "Statistical Measures Calculation");
19
20         job.setJarByClass(DriverClass.class);
21         job.setMapperClass(MapperClass.class);
22         job.setCombinerClass(CombinerClass.class);
23         job.setReducerClass(ReducerClass.class);
24
25         job.setOutputKeyClass(Text.class);
26         job.setOutputValueClass(DoubleWritable.class);
27
28         FileInputFormat.addInputPath(job, new Path(args[0]));
29         FileOutputFormat.setOutputPath(job, new Path(args[1]));
30
31         System.exit(job.waitForCompletion(true) ? 0 : 1);
32     }
33 }
34
```

### 4. Combiner Class (Optional):

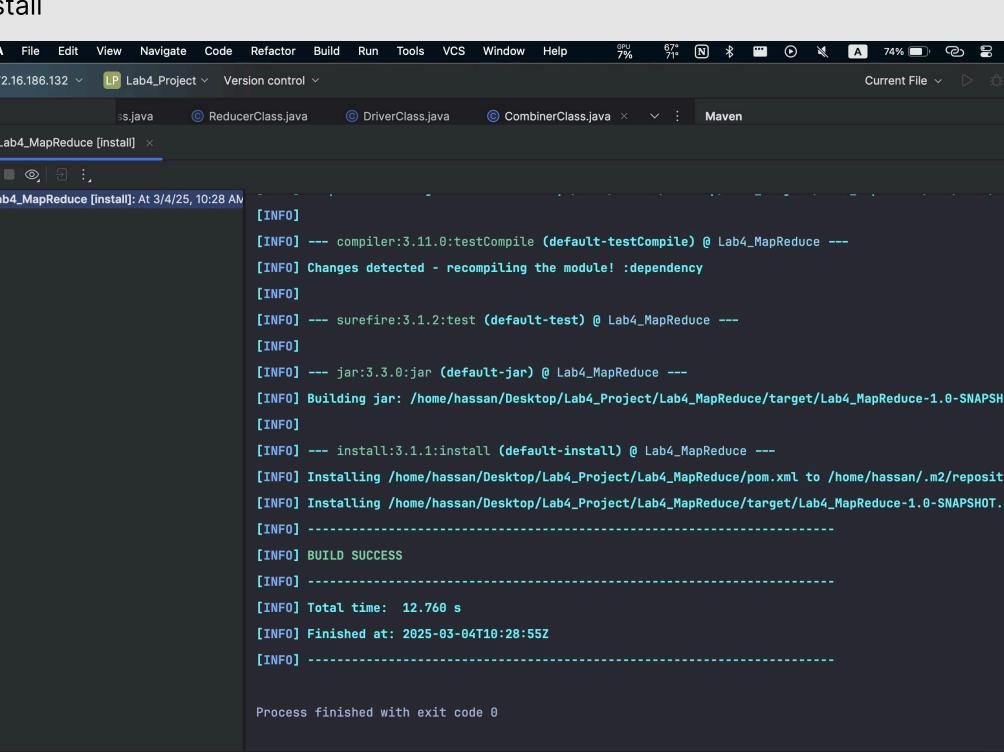
```
1 package edu.hassan;
2 import java.io.IOException;
3 import org.apache.hadoop.io.DoubleWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Reducer;
6
7 public class CombinerClass extends Reducer<Text, DoubleWritable, Text, DoubleWritable> {
8     @Override
9     protected void reduce(Text key, Iterable<DoubleWritable> values, Context context)
10         throws IOException, InterruptedException {
11
12         double sum = 0, min = Double.MAX_VALUE, max = Double.MIN_VALUE, sumOfSquares =
13 0; int count = 0;
14
15         for (DoubleWritable val : values) {
16             double value = val.get();
17
18             switch (key.toString()) {
19                 case "Count":
20                     count += (int) value;
21                     break;
22                 case "Sum":
23                     sum += value;
24                     break;
25                 case "Min":
26                     min = Math.min(min, value);
27                     break;
28                 case "Max":
29                     max = Math.max(max, value);
30                     break;
31                 case "SumOfSquares":
32                     sumOfSquares += value;
33                     break;
34             }
35         }
36     }
37     @Override
38     protected void cleanup(Context context) throws IOException, InterruptedException {
39
40         double mean = sum / count;
41         double variance = (sumOfSquares / count) - (mean * mean);
42         double stdDev = Math.sqrt(variance);
43
44         context.write(new Text("Count"), new DoubleWritable(count));
45         context.write(new Text("Sum"), new DoubleWritable(sum));
46         context.write(new Text("Min"), new DoubleWritable(min));
47         context.write(new Text("Max"), new DoubleWritable(max));
48         context.write(new Text("Mean"), new DoubleWritable(mean));
49         context.write(new Text("Standard Deviation"), new DoubleWritable(stdDev));
50     }
51 }
```

```
26     context.write(new Text("Min"), new DoubleWritable(min));
27     context.write(new Text("Max"), new DoubleWritable(max));
28     context.write(new Text("SumOfSquares"), new DoubleWritable(sumOfSquares));
29 }
30 }
```

10

## Task 4: Compile & Build JAR File

→ mvn clean  
→ mvn install

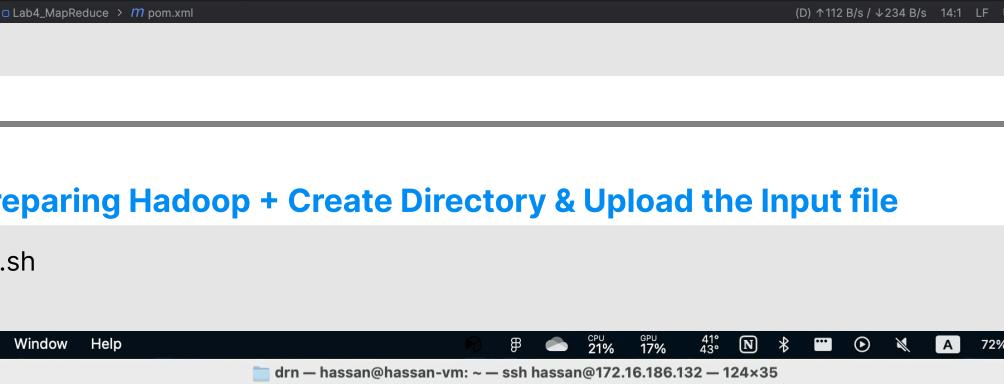


```
[INFO]
[INFO] --- compiler:3.11.0:testCompile (default-testCompile) @ Lab4_MapReduce ---
[INFO] Changes detected - recompiling the module! :dependency
[INFO]
[INFO] --- surefire:3.1.2:test (default-test) @ Lab4_MapReduce ---
[INFO]
[INFO] --- jar:3.3.0:jar (default-jar) @ Lab4_MapReduce ---
[INFO] Building jar: /home/hassan/Desktop/Lab4_Project/Lab4_MapReduce/target/Lab4_MapReduce-1.0-SNAPSHOT.jar
[INFO]
[INFO] --- install:3.1.1:install (default-install) @ Lab4_MapReduce ---
[INFO] Installing /home/hassan/Desktop/Lab4_Project/Lab4_MapReduce/pom.xml to /home/hassan/.m2/repository/edu/
[INFO] Installing /home/hassan/Desktop/Lab4_Project/Lab4_MapReduce/target/Lab4_MapReduce-1.0-SNAPSHOT.jar to /
[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 12.760 s
[INFO] Finished at: 2025-03-04T10:28:55Z
[INFO] -----
```

Process finished with exit code 0

## Task 5: Preparing Hadoop + Create Directory & Upload the Input file

```
$ start-all.sh
$ jps
```



```
drn - hassan@hassan-vm: ~ - ssh hassan@172.16.186.132 - 124x35
[hassan@hassan-vm: ~]$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hassan in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
localhost: namenode is running as process 12361. Stop it first and ensure /tmp/hadoop-hassan-namenode.pid file is empty before retry.
Starting datanodes
localhost: datanode is running as process 12500. Stop it first and ensure /tmp/hadoop-hassan-datanode.pid file is empty before retry.
Starting secondary namenodes [hassan-vm]
hassan-vm: secondarynamenode is running as process 12681. Stop it first and ensure /tmp/hadoop-hassan-secondarynamenode.pid
```

```
file is empty
Starting resource
resourcemanager
before retry.
Starting nodeman
localhost: node
ty before retry.
[hassan@hassan-va
128017 RemoteMa
```

```
136617 RemoteMavenServer30
131206 MavenServerIndexerMain
12500 DataNode
128692 Main
12952 ResourceManager
12681 SecondaryNameNode
12361 NameNode
13083 NodeManager
141375 Jps
hassan@hassan-vm:~$ █

$ hadoop fs -mkdir -p /user/hassan/Lab4/input
$ hadoop fs -ls /user/hassan
$ hadoop fs -put ~/Desktop/Lab4_Project/prices_data.txt /user/hassan/Lab4/input/
$ hadoop fs -ls /user/hassan/Lab4/input
$ hadoop fs -cat /user/hassan/Lab4/input/prices_data.txt

Edit View Window Help
drn — hassan@hassan-vm: ~ — ssh hassan@172.16.186.132 — 124x35
hassan@hassan-vm:~$ hadoop fs -mkdir -p /user/hassan/Lab4/input
hassan@hassan-vm:~$ hadoop fs -ls /user/hassan
Found 3 items
drwxr-xr-x  - hassan  supergroup          0 2025-02-04 08:45 /user/hassan/Folder1
drwxr-xr-x  - hassan  supergroup          0 2025-02-17 07:43 /user/hassan/Lab3
drwxr-xr-x  - hassan  supergroup          0 2025-03-04 10:36 /user/hassan/Lab4
hassan@hassan-vm:~$ hadoop fs -put ~/Desktop/Lab4_Project/prices_data.txt /user/hassan/Lab4/input/
hassan@hassan-vm:~$ hadoop fs -ls /user/hassan/Lab4/input
Found 1 items
-rw-r--r--  1 hassan  supergroup  54510 2025-03-04 10:37 /user/hassan/Lab4/input/prices_data.txt
hassan@hassan-vm:~$ hadoop fs -cat /user/hassan/Lab4/input/prices_data.txt
237.68
99.95
34.99
28.91
17.49
18.66
29.12
97.68
12.99
38.49
18.16
84.61
33.92
14.99
4.99
34.39
12.88
117.26
9.39
17.85
27.5
10.99
159.99
12.63

rowsing HDFS
localhost:9870/explorer.html#/user/hassan/Lab4/input
Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities
File information - prices_data.txt
Download Head the file (first 32K) Tail the file (last 32K)
```

## Browse Directo



Showing 1 to 1 of 1 entries

Generation Stamp: 1014

Size: 54510

Availability:

- hassan-vm

File contents

```
237.68
99.95
34.99
28.91
17.49
18.66
29.12
97.68
```

Size Name

IB prices\_data.txt

Previous 1 Next

## Task 6: Execute the MapReduce Program

```
$ hadoop jar Desktop/Lab4_Project/Lab4_MapReduce/target/Lab4_MapReduce-1.0-SNAPSHOT.jar \
> edu.hassan.DriverClass \
> /user/hassan/Lab4/input/prices_data.txt /user/hassan/Lab4/output
```

(1)

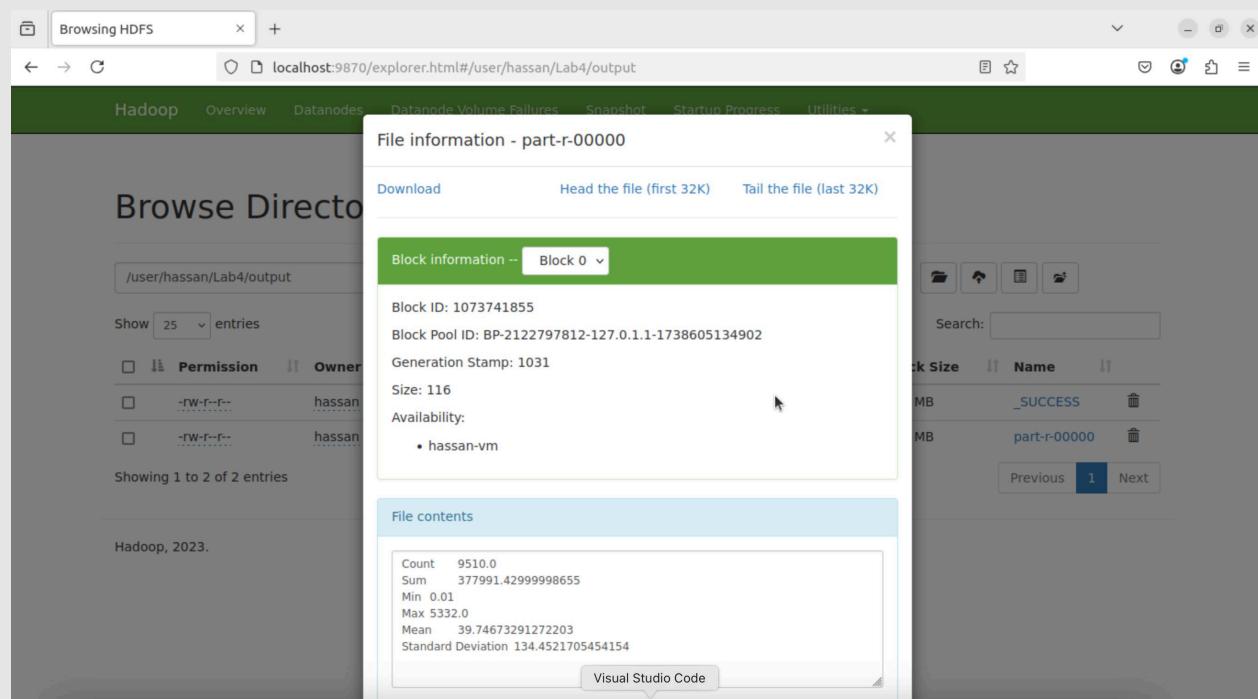
```
Edit View Window Help drn — hassan@hassan-vm: ~ — ssh hassan@172.16.186.132 — 140x46
hassan@hassan-vm:~$ hadoop jar Desktop/Lab4_Project/Lab4_MapReduce/target/Lab4_MapReduce-1.0-SNAPSHOT.jar \
> edu.hassan.DriverClass \
> /user/hassan/Lab4/input/prices_data.txt /user/hassan/Lab4/output
2025-03-04 11:52:55,101 INFO client.DefaultNoHARMfailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2025-03-04 11:52:55,992 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2025-03-04 11:52:56,030 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/hassan/.staging/job_1739764183578_0003
2025-03-04 11:52:56,832 INFO input.FileInputFormat: Total input files to process : 1
2025-03-04 11:52:56,944 INFO mapreduce.JobSubmitter: number of splits:1
2025-03-04 11:52:57,237 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1739764183578_0003
2025-03-04 11:52:57,237 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-03-04 11:52:57,397 INFO conf.Configuration: resource-types.xml not found
2025-03-04 11:52:57,397 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2025-03-04 11:52:59,683 INFO impl.YarnClientImpl: Submitted application application_1739764183578_0003
2025-03-04 11:52:59,966 INFO mapreduce.Job: The url to track the job: http://hassan-vm:8088/proxy/application_1739764183578_0003/
2025-03-04 11:52:59,967 INFO mapreduce.Job: Running job: job_1739764183578_0003
2025-03-04 11:53:21,455 INFO mapreduce.Job: Job job_1739764183578_0003 running in uber mode : false
2025-03-04 11:53:21,466 INFO mapreduce.Job: map 0% reduce 0%
2025-03-04 11:53:37,696 INFO mapreduce.Job: map 100% reduce 0%
2025-03-04 11:53:42,982 INFO mapreduce.Job: map 100% reduce 100%
2025-03-04 11:53:44,034 INFO mapreduce.Job: Job job_1739764183578_0003 completed successfully
2025-03-04 11:53:44,232 INFO mapreduce.Job: Counters: 54
  File System Counters
    FILE: Number of bytes read=87
    FILE: Number of bytes written=553049
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=54207
    HDFS: Number of bytes written=116
    HDFS: Number of read operations=8
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
```

```
Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps
    Total time spent by all reducers
    Total time spent by all map tasks
    Total time spent by all reduce tasks
    Total vcore-milliseconds taken by all map tasks
```

```
Total vcore-milliseconds taken by all reduce tasks=2869
Total megabyte-milliseconds taken by all map tasks=14648320
```

```
-rw-r--r-- 1 hassan supergroup 116 2025  
hassan@hassan-vm:~$ hadoop fs -get /user/hassan/  
hassan@hassan-vm:~$ cd Desktop/Lab4_Project/  
hassan@hassan-vm:~/Desktop/Lab4_Project$ ls  
DataSetLoad.py Lab4_MapReduce marketing_sample  
hassan@hassan-vm:~/Desktop/Lab4_Project$ [1]  
Count 9510.0  
Sum 377991.42999998655  
Min 0.01
```

```
Max      5332.0
Mean     39.74673291272203
Standard Deviation 134.4521705454154
hassan@hassan-vm:~/Desktop/Lab4_Project$
```



The screenshot shows a web browser window titled "Browsing HDFS" with the URL "localhost:9870/explorer.html#/user/hassan/Lab4/output". The main page displays a list of files in the directory "/user/hassan/Lab4/output", showing 2 entries: "SUCCESS" and "part-r-00000". A modal window titled "File information - part-r-00000" is open, providing detailed information about the file. The modal has tabs for "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". The "Block information" tab is selected, showing Block ID: 1073741855, Block Pool ID: BP-2122797812-127.0.1.1-1738605134902, Generation Stamp: 1031, Size: 116, and Availability: hassan-vm. The "File contents" tab shows statistical data: Count 9510.0, Sum 377991.42999998655, Min 0.01, Max 5332.0, Mean 39.74673291272203, and Standard Deviation 134.4521705454154. A "Visual Studio Code" button is at the bottom of the modal.

## Explanation of the Implementation

- This lab discusses about calculating statistical measures (Count, Sum, Min, Max, Mean, and Standard Deviation) from a price dataset of amazon products using the MapReduce mechanism and the Hadoop Distributed File System.
- First, I created a Java program that consists of four classes: one for Mapping, one for Combining, one for Reducing, and one for Running the job of statistical calculation. Each class

- Mapper Class: Reads the input
- Combiner Class: Performs local reduction

Reducer.

- Reducer Class: Aggregates the final results and calculates all the required statistical measures.
- Driver Class: Configures and manages the MapReduce job execution.

• After developing the program, I built it into a JAR file and created an input file (prices\_data.txt) containing the cleaned price data.

• The input file was uploaded to the Hadoop File System in a new directory "/Lab4/input" inside an existing directory created previously.

• Next, I executed the JAR program with the input file uploaded to HDFS. Hadoop distributed the statistical calculation function and generated an output file containing the calculated Count, Sum, Minimum, Maximum, Mean, and Standard Deviation of the prices.

---

- The System & Tools I'm using:
  - Macbook Air M1
  - VMware Fusion "Virtual Machine Hypervisor"
  - Ubuntu for Servers 22.04 "With ubuntu-desktop package"
  - Hadoop version 3.3.6
  - JDK version 11
  - IntelliJ IDEA with SSH feature (for remote programming)
  - GitHub (to upload the files on a remote repository)

---