

PRACTICAL 4

To determine maximum temperature using Hadoop MapReduce

Map Reduce : works by breaking the processing into two phases: the map phase and the reduce phase. Each phase has key-value pairs as input and output, the types of which may be chosen by the programmer. The programmer also specifies two functions: the map function and the reduce function.

Hadoop MaxTemperature operation occurs in 3 stages –

Mapper Phase

Reducer Phase

Driver code

Dataset – temperature.txt

```
0067011990999991950051507004+68750+023550FM-
12+038299999V0203301N00671220001CN9999999N9+00001+9999999999
0043011990999991950051512004+68750+023550FM-
12+038299999V0203201N00671220001CN9999999N9+00221+9999999999
0043011990999991950051518004+68750+023550FM-
12+038299999V0203201N00261220001CN9999999N9-00111+9999999999
0043012650999991949032412004+62300+010750FM-
12+048599999V0202701N00461220001CN0500001N9+01111+9999999999
0043012650999991949032418004+62300+010750FM-
12+048599999V0202701N00461220001CN0500001N9+00781+9999999999
```

These lines are presented to the map function as the key-value pairs

```
(0, 0067011990999991950051507004...9999999N9+00001+9999999999...) (106,
0043011990999991950051512004...9999999N9+00221+9999999999...) (212,
0043011990999991950051518004...9999999N9-00111+9999999999...) (318,
0043012650999991949032412004...0500001N9+01111+9999999999...) (424,
0043012650999991949032418004...0500001N9+00781+9999999999...)
```

The keys are the line offsets within the file, which we ignore in our map function. The map function merely extracts the year and the air temperature (indicated in bold text), and emits them as its output (the temperature values have been interpreted as integers):

```
(1950, 0)
```

(1950, 22)

The output from the map function is processed by the MapReduce framework before being sent to the reduce function. This processing sorts and groups the key-value pairs by key. So, continuing the example, our reduce function sees the following input:

(1949, [111, 78])

(1950, [0, 22, -11])

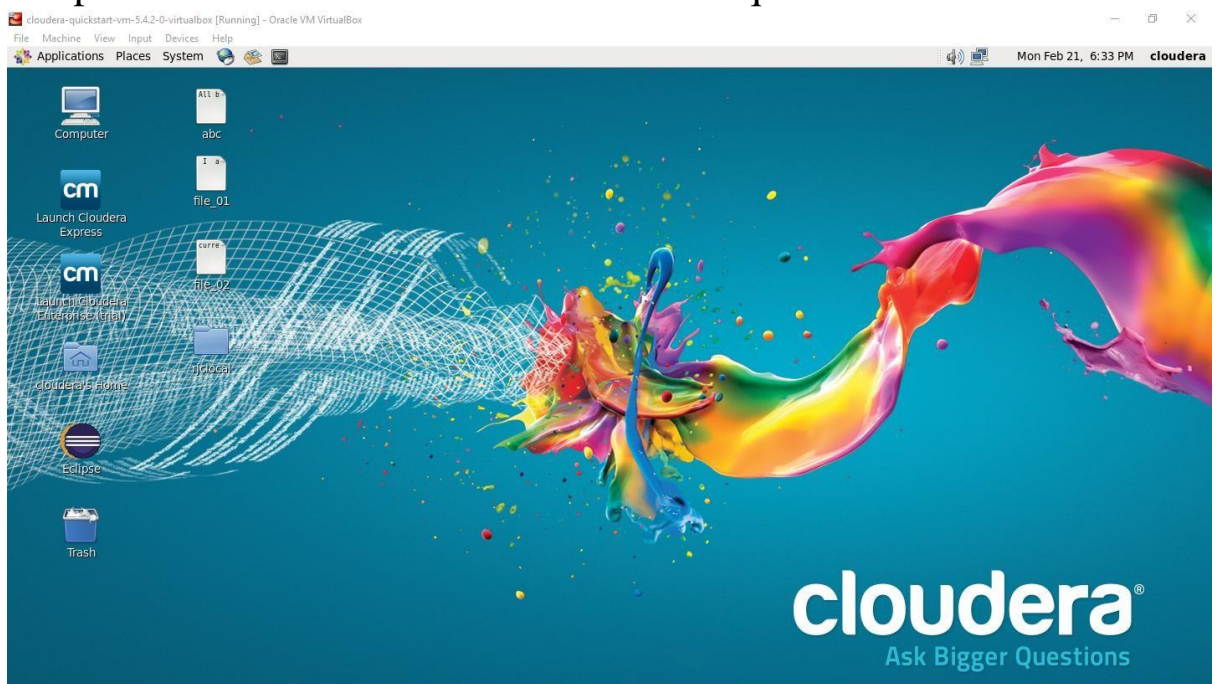
Each year appears with a list of all its air temperature readings. All the reduce function has to do now is iterate through the list and pick up the maximum reading:

(1949, 111)

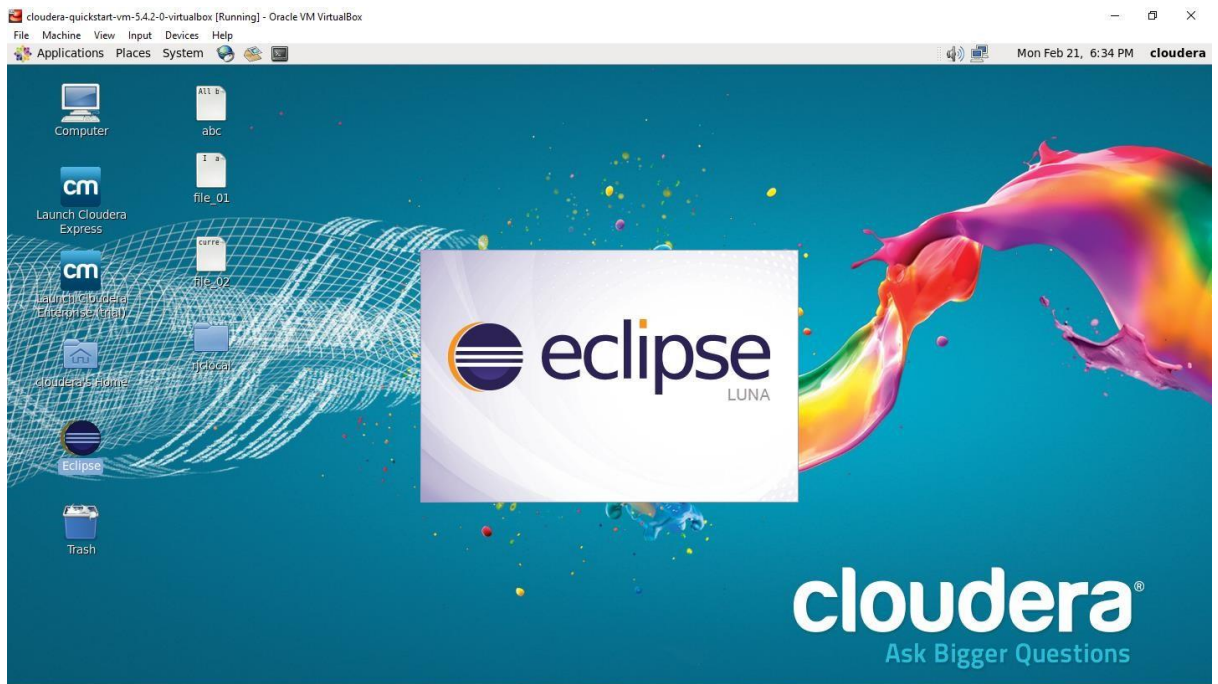
(1950, 22)

- Steps to determine maximum temperature using Hadoop MapReduce in Cloudera:

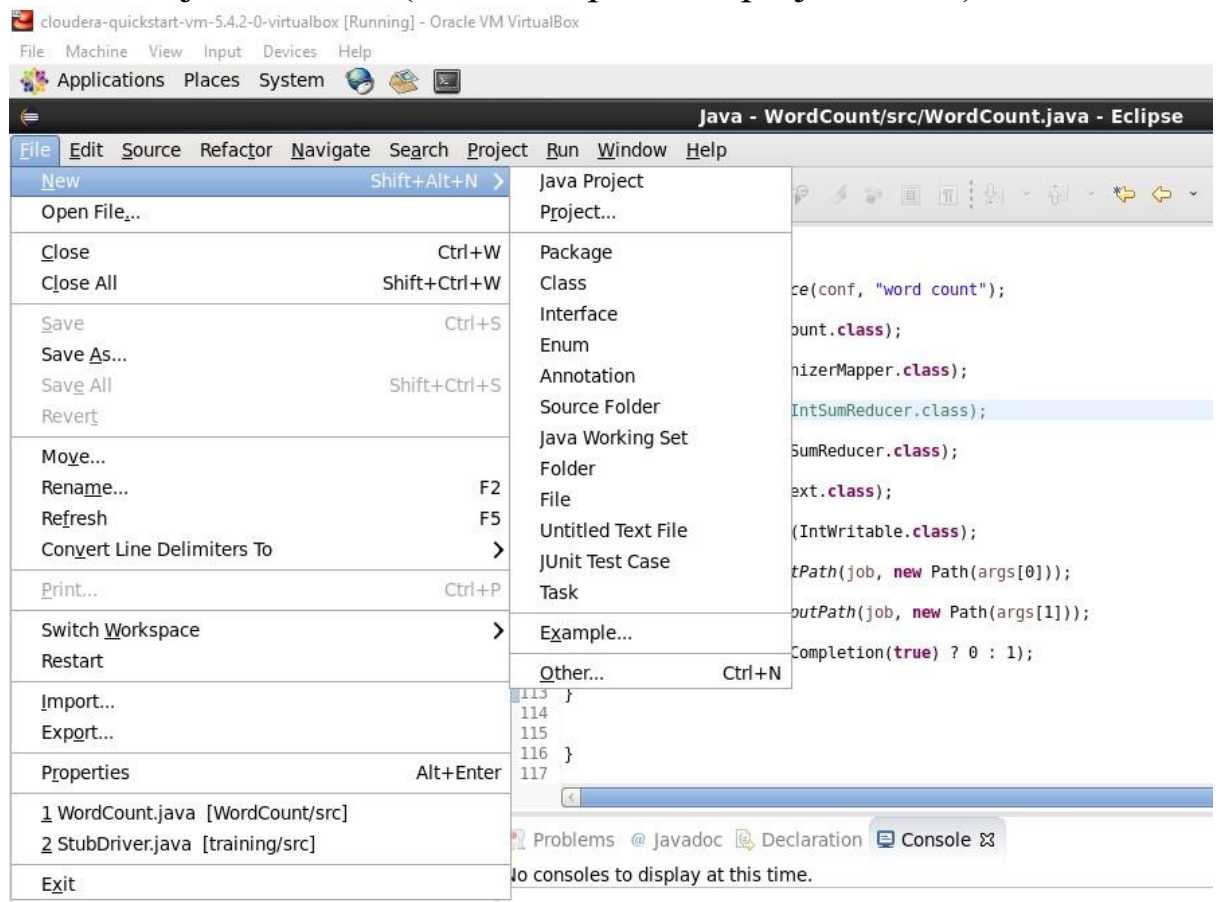
1) Open virtual box and then start cloudera quickstart

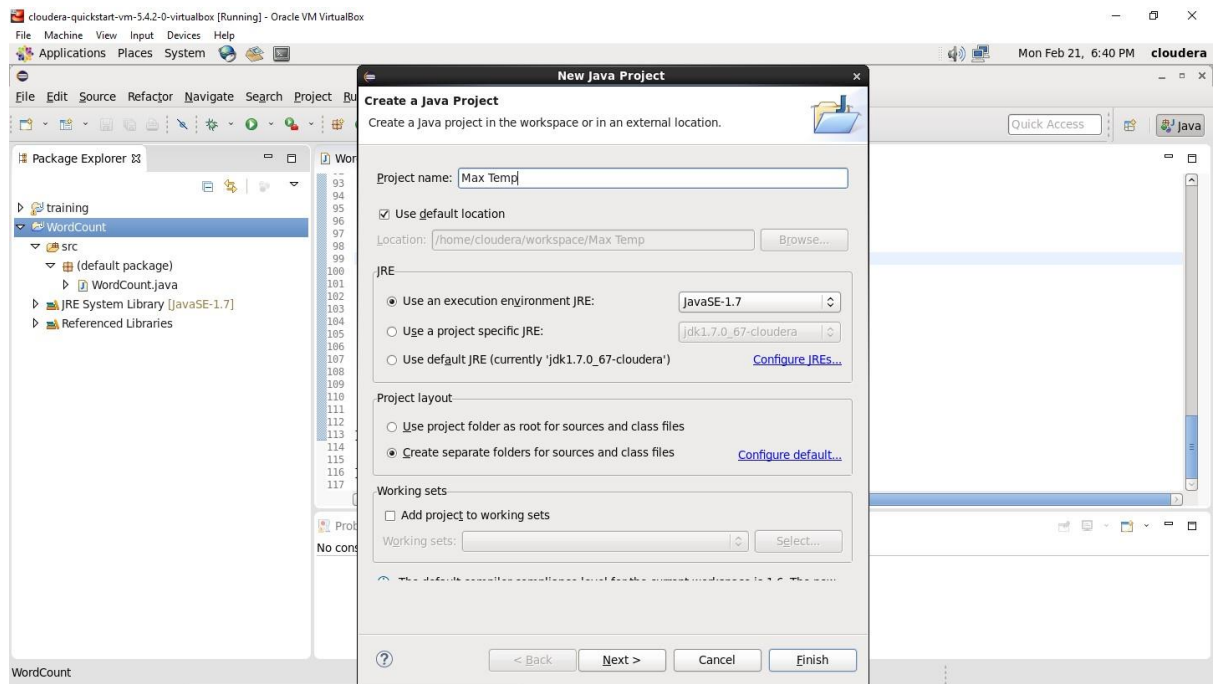


2) Open Eclipse present on the cloudera desktop



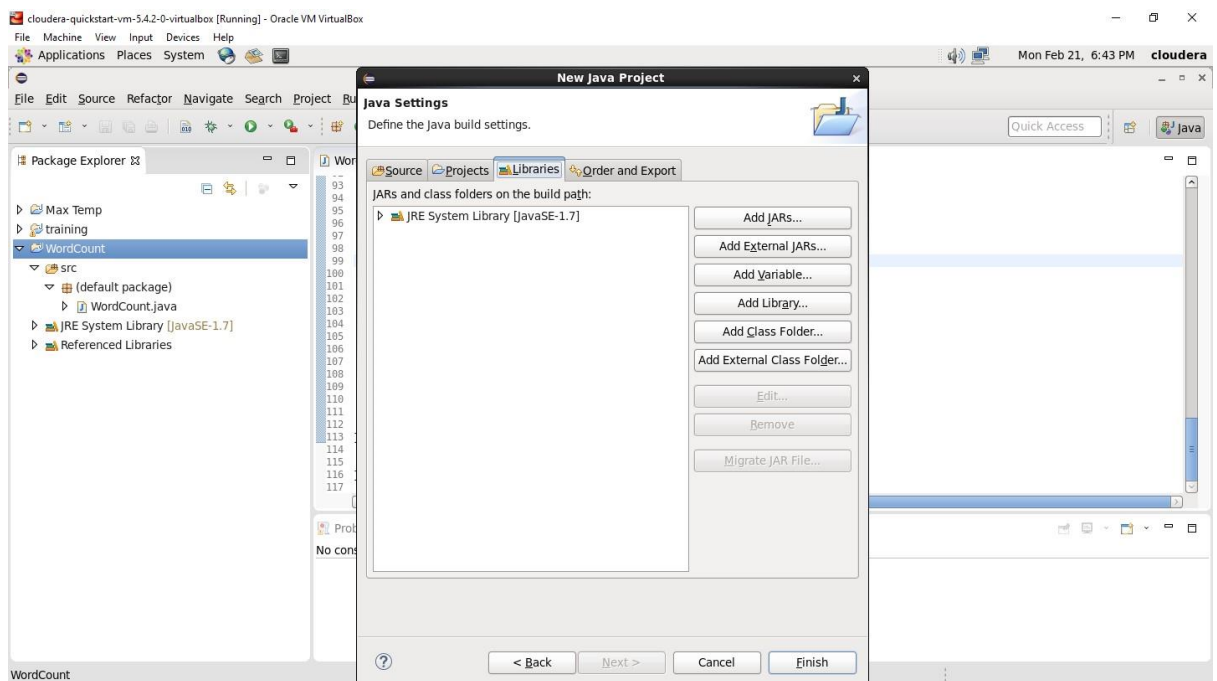
- 3) Create a new Java project clicking: File -> New -> Project -> Java Project -> Next (“MaxTemp” is the project name).

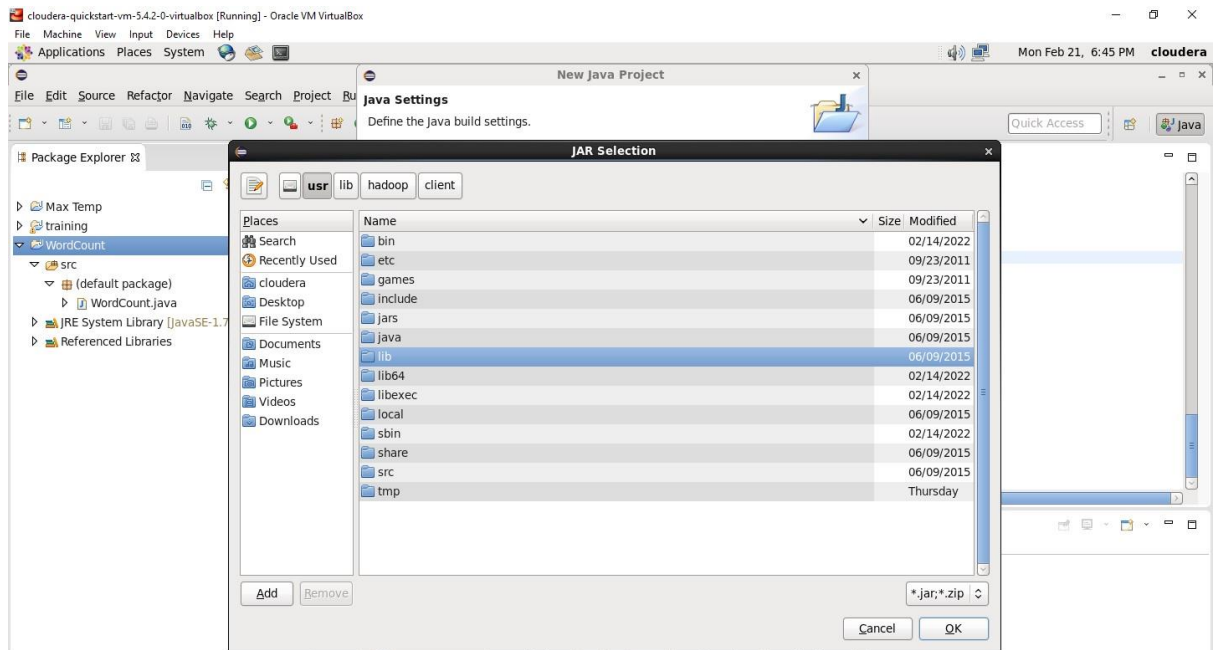
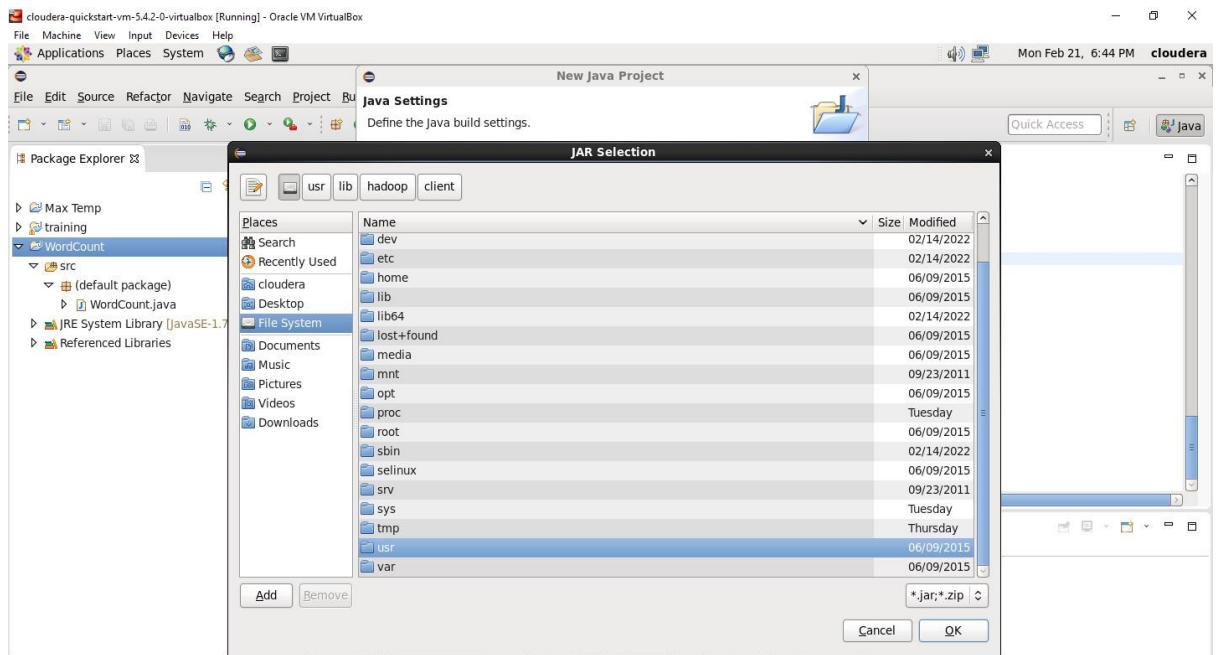


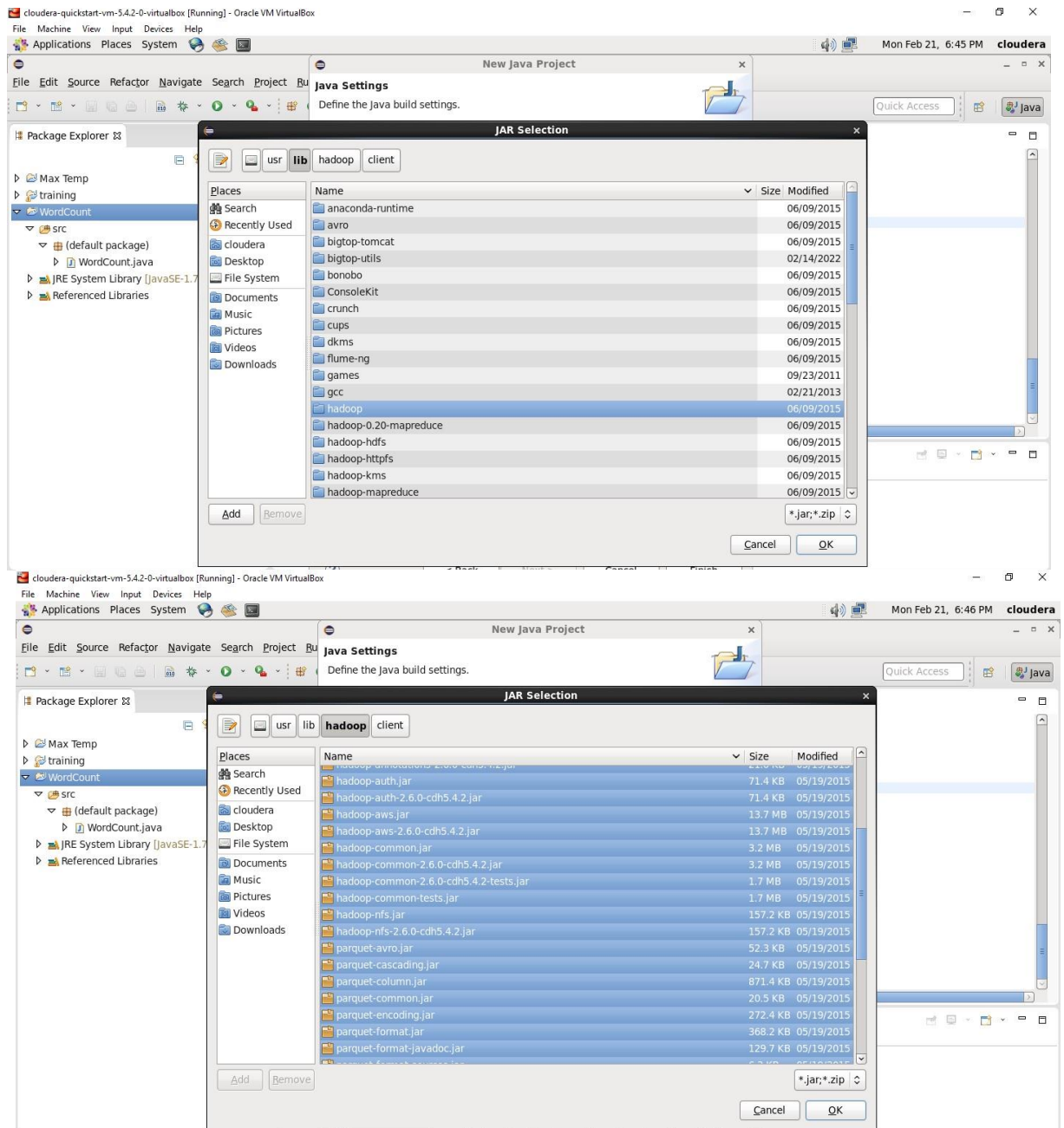


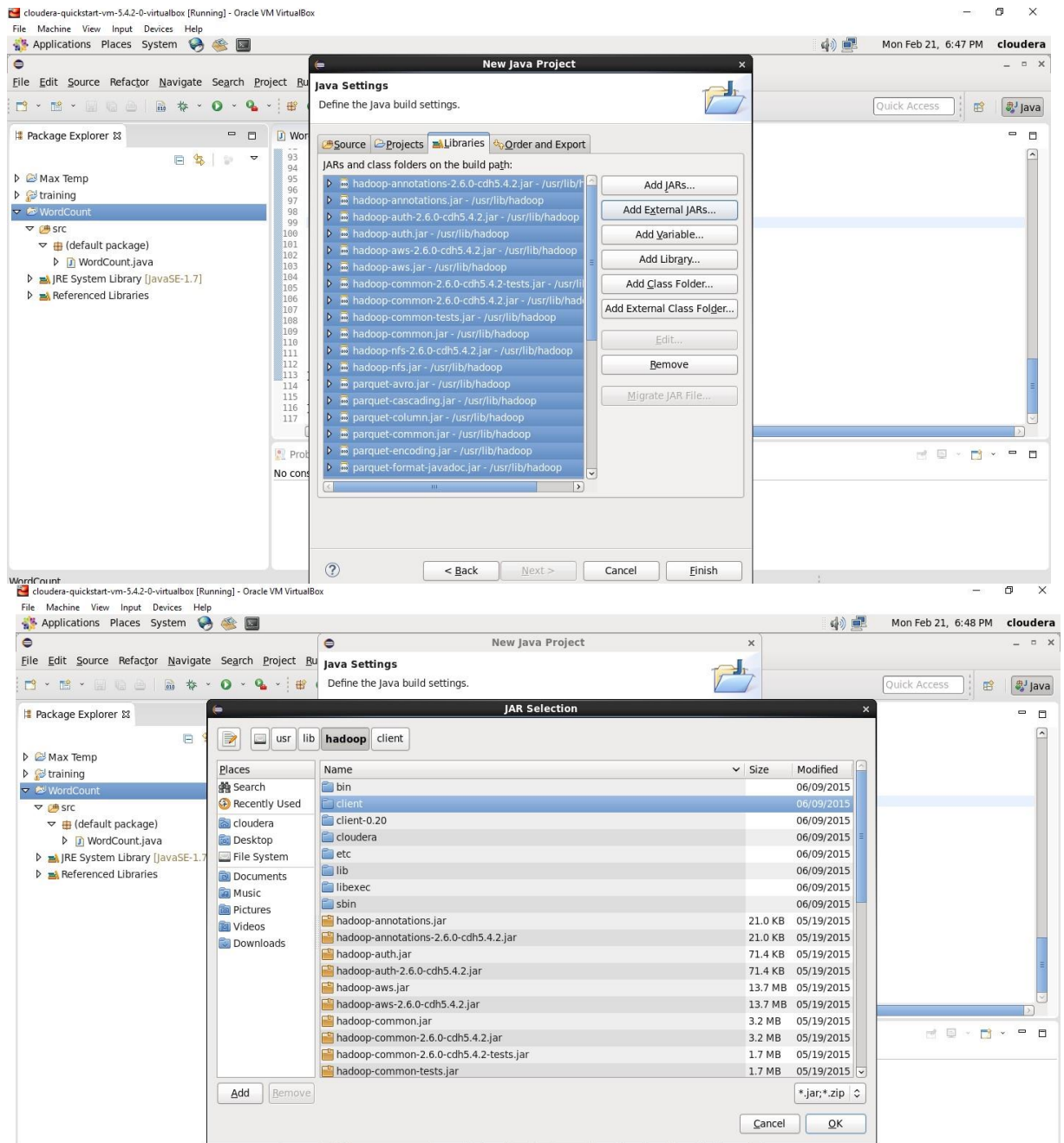
4) Adding the Hadoop libraries to the project Click on Libraries -

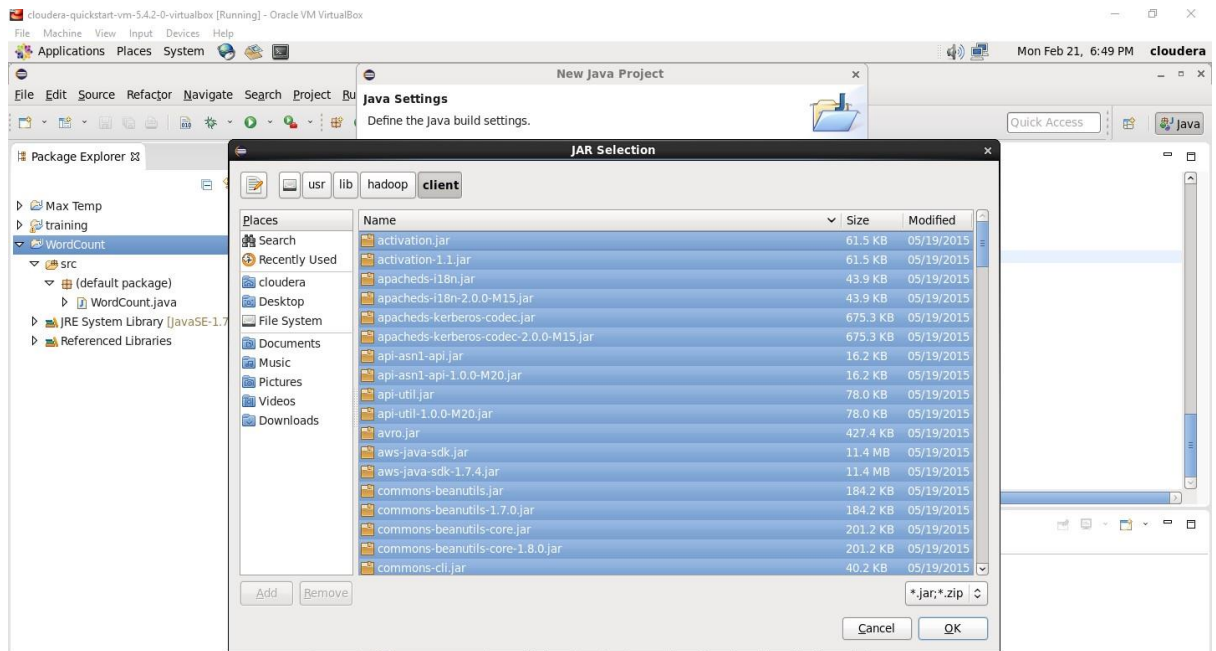
>
Add External JARs Click on File System -> usr -> lib -> hadoop
Select all the libraries (JAR Files) -> click OK Click on Add
External jars, -> client -> select all jar files -> ok -> Finish



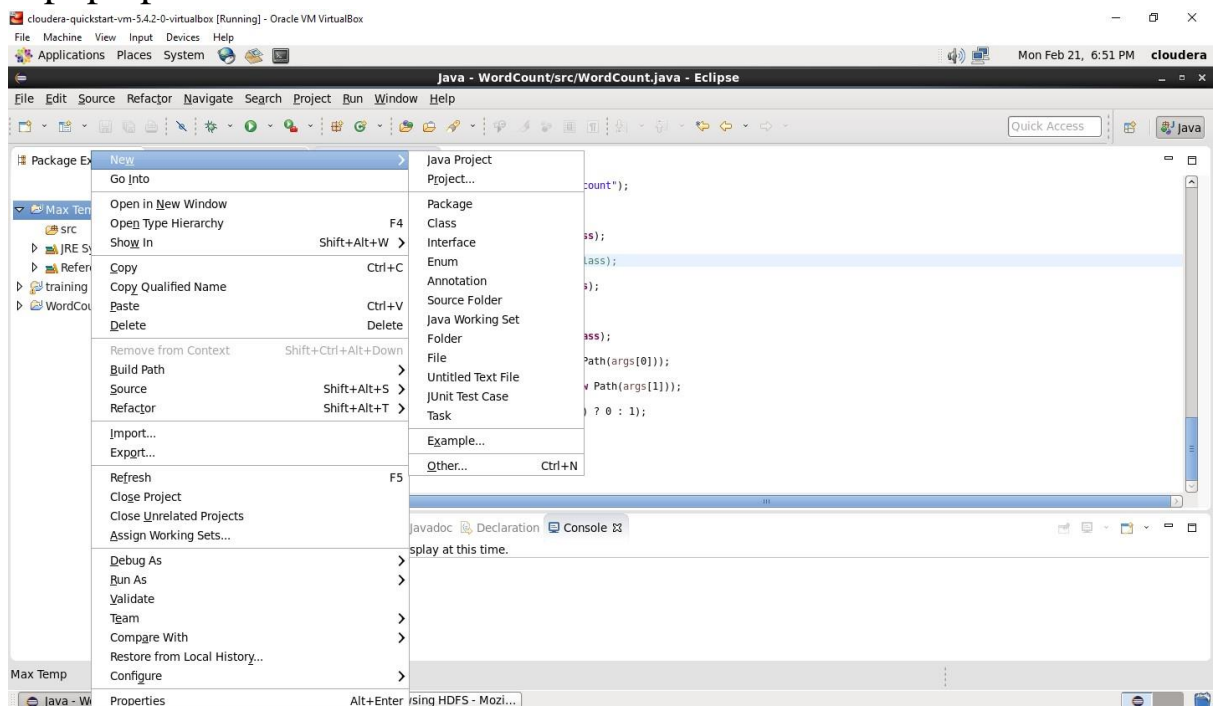


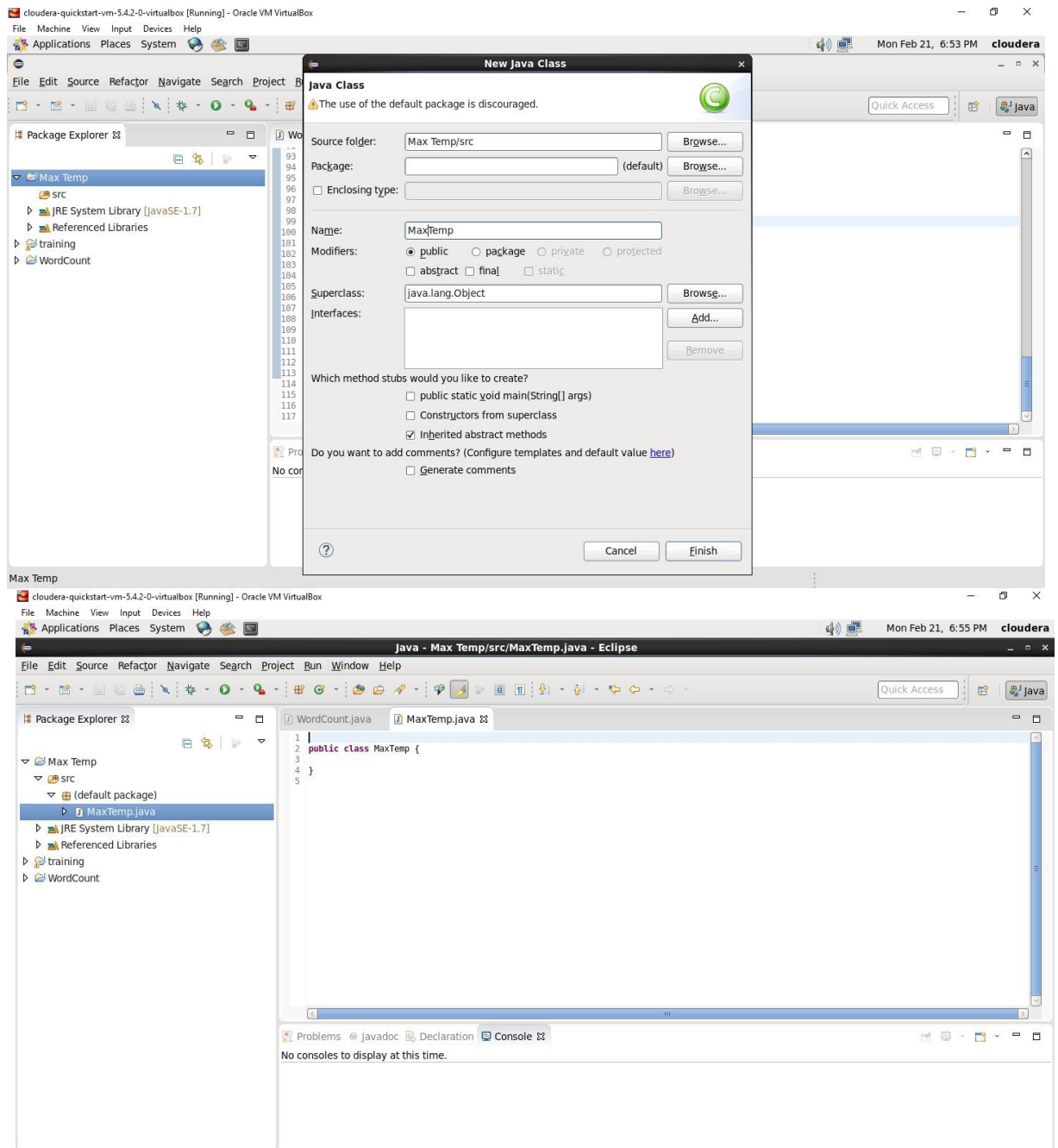


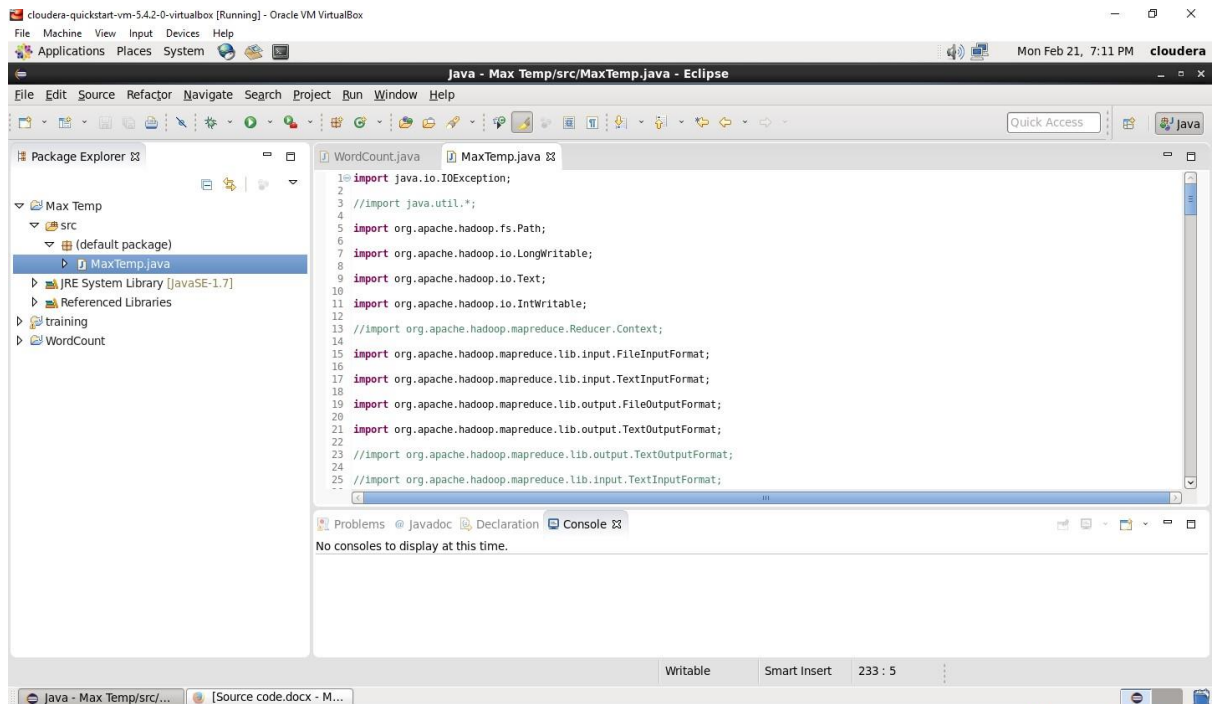




- 5) Right Click on the name of Project “MaxTemp” -> New -> class “MaxTemp” -> Finish Then MaxTemp.java window will pop up







Source code:

```
import java.io.IOException;
//import java.util.*;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.IntWritable;
//import org.apache.hadoop.mapreduce.Reducer.Context;
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;
//import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
//import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.conf.Configuration;
//import MaximumTemp.MaxTemperatureMapper;
//import MaxTemp.MaxTemperatureReducer;
public class MaxTemp {
```

```
// Mapper
/*MaxTemperatureMapper class is static
 * And extends Mapper abstract class
 * having four Hadoop generics type
 * Long Writable, Text, Text, Text.
 */
public static class MaxTemperatureMapper extends
Mapper<LongWritable, Text, Text, IntWritable> {
// the data in our data set with
// this value is inconsistent data
//public static final int MISSING = 9999;

    public void map(LongWritable key, Text value, Context context)
throws IOException, InterruptedException {
        String line=value.toString();
        String year=line.substring(15, 19);
        int airtemp;
        if(line.charAt(87)=='+')
        {
            airtemp=Integer.parseInt(line.substring(88,92));
        }
        else
            airtemp=Integer.parseInt(line.substring(87,92));
        String q=line.substring(92,93);
        if(airtemp!=9999 && q.matches("[01459]"))
        {

            context.write(new Text(year), new IntWritable(airtemp));
        }
    }
}

// Reducer
/*MaxTemperatureReducer class is static
and extends Reducer abstract class
```

having four Hadoop generics type

Text, Text, Text, Text.

*/

public static class MaxTemperatureReducer extends
Reducer<Text, IntWritable, Text, IntWritable> {

/**

* @method reduce

* This method takes the input as key and

* list of values pair from the mapper,

* it does aggregation based on keys and

* produces the final context.

*/

public void reduce(Text key, Iterable<IntWritable> values, Context
context)

throws IOException, InterruptedException {

int maxvalue= Integer.MIN_VALUE;

for (IntWritable value : values) {

maxvalue=Math.max(maxvalue, value.get());

}

context.write(key, new IntWritable(maxvalue));

}

}

/**

* @method main

* This method is used for setting

* all the configuration properties.

* It acts as a driver for map-reduce

* code.

*/

public static void main(String[] args) throws Exception {

// reads the default configuration of the

// cluster from the configuration XML files

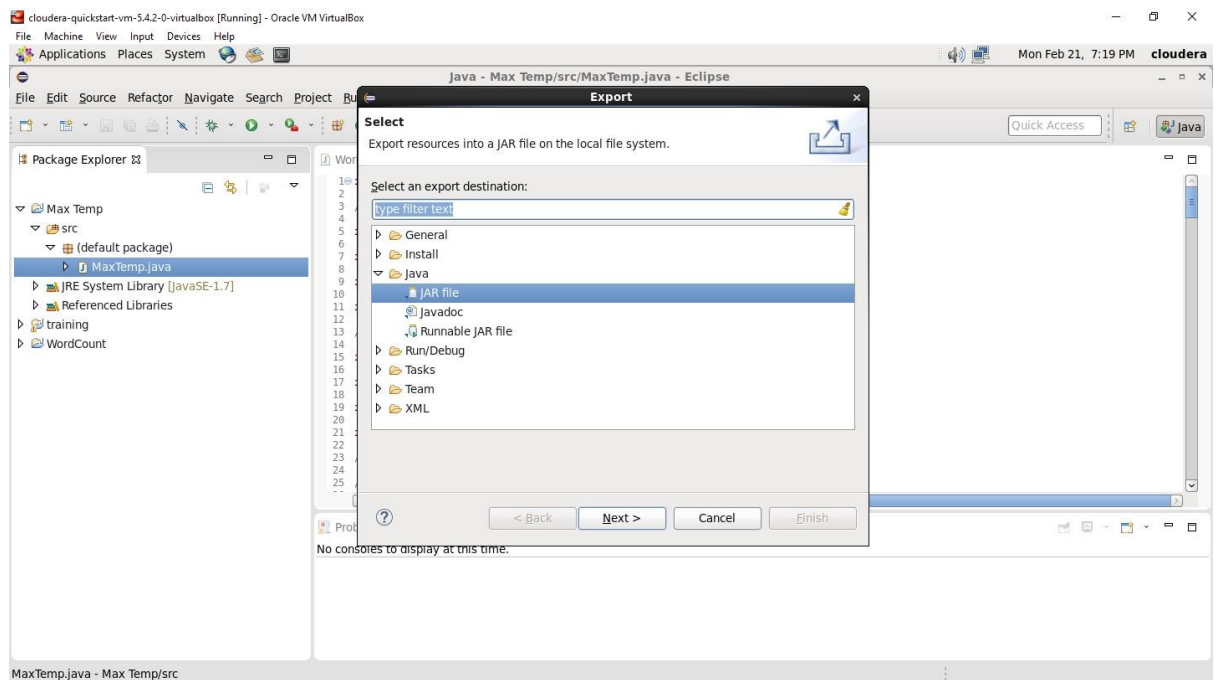
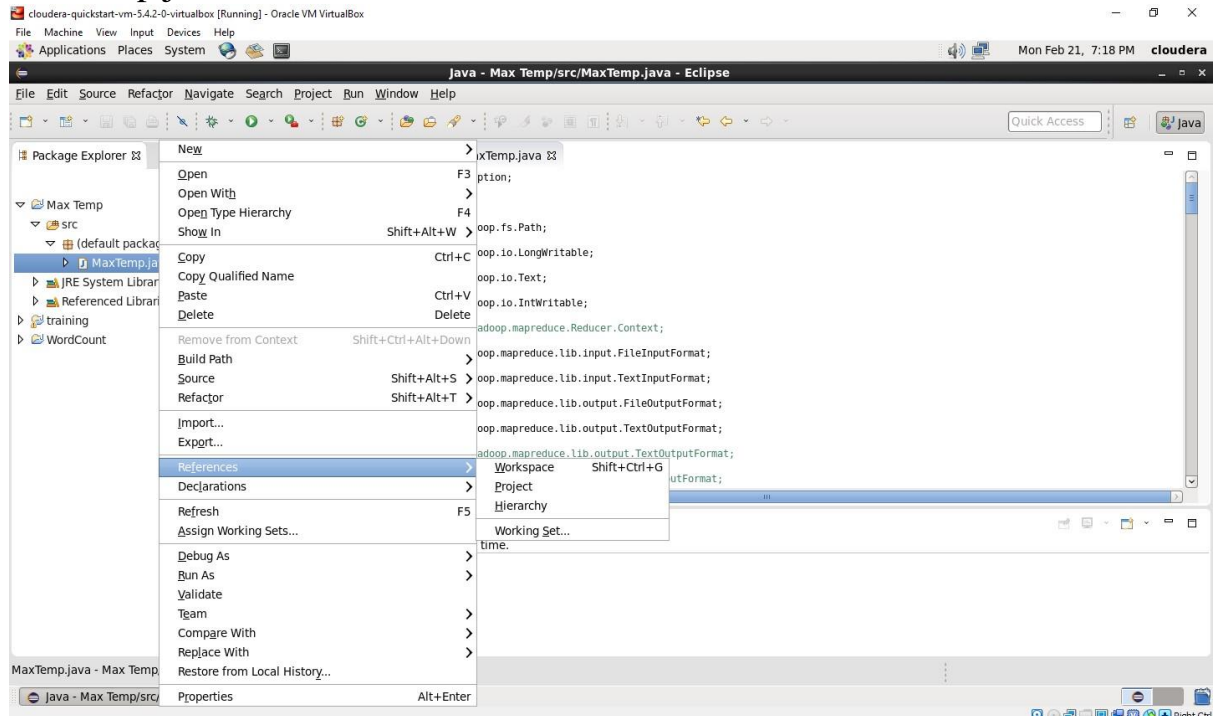
Configuration conf = new Configuration();

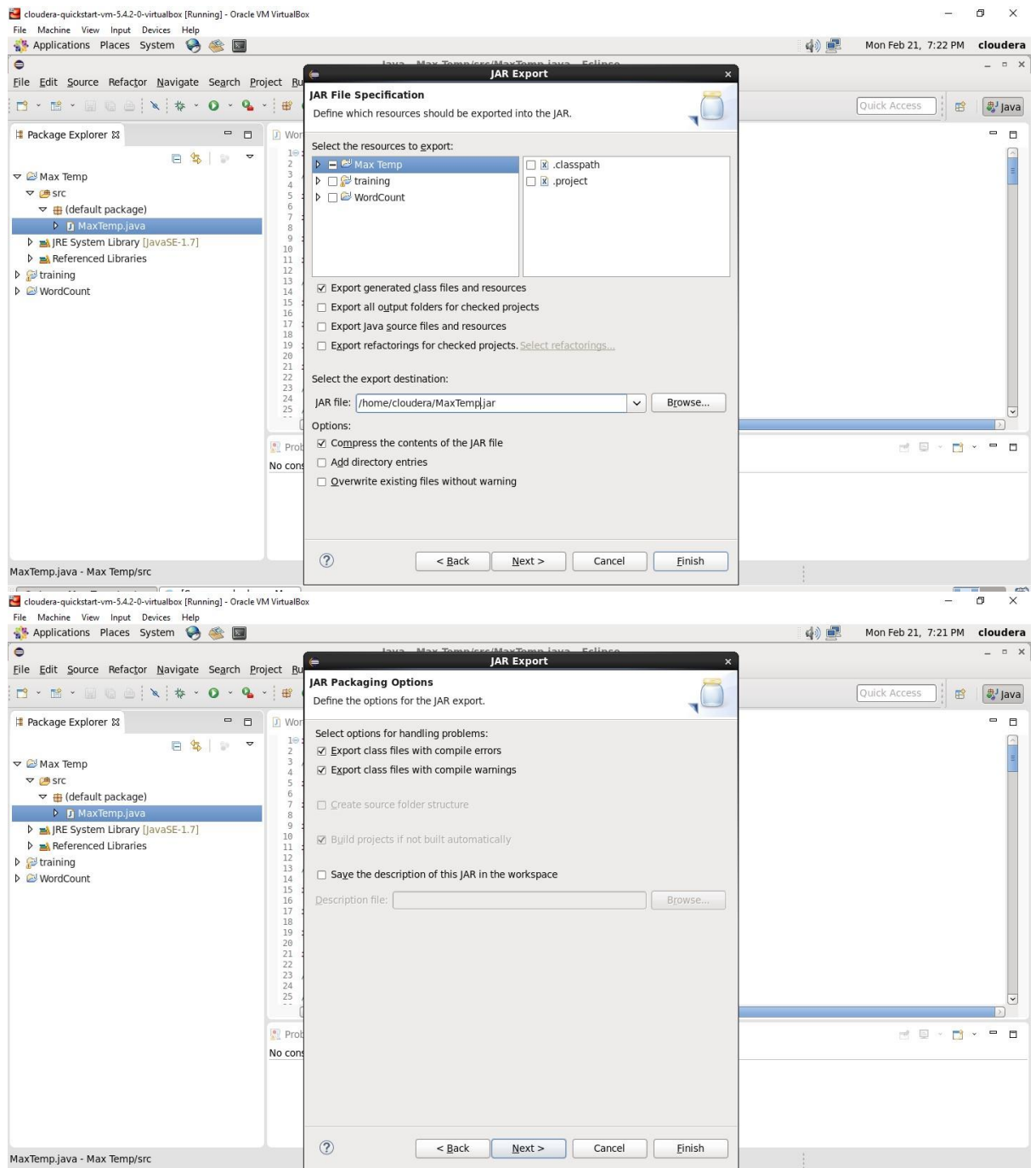
```
// initializing the job with the
// default configuration of the cluster
// Job job = new Job(conf, "weather example");
Job job = Job.getInstance(conf, "weather example");
// Assigning the driver class name
job.setJarByClass(MaxTemp.class);
// Key type coming out of mapper
// job.setMapOutputKeyClass(Text.class);
// value type coming out of mapper
// job.setMapOutputValueClass(Text.class);
// Defining the mapper class name
job.setMapperClass(MaxTemperatureMapper.class);
// Defining the reducer class name
job.setReducerClass(MaxTemperatureReducer.class);
// defining input Format class which is
// responsible to parse the dataset
// into a key value pair
job.setInputFormatClass(TextInputFormat.class);
// Defining output Format class which is
// responsible to parse the dataset
// into a key value pair
job.setOutputFormatClass(TextOutputFormat.class);

job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);
FileInputFormat.addInputPath(job, new Path(args[0]));
FileOutputFormat.setOutputPath(job, new Path(args[1]));
System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}
```

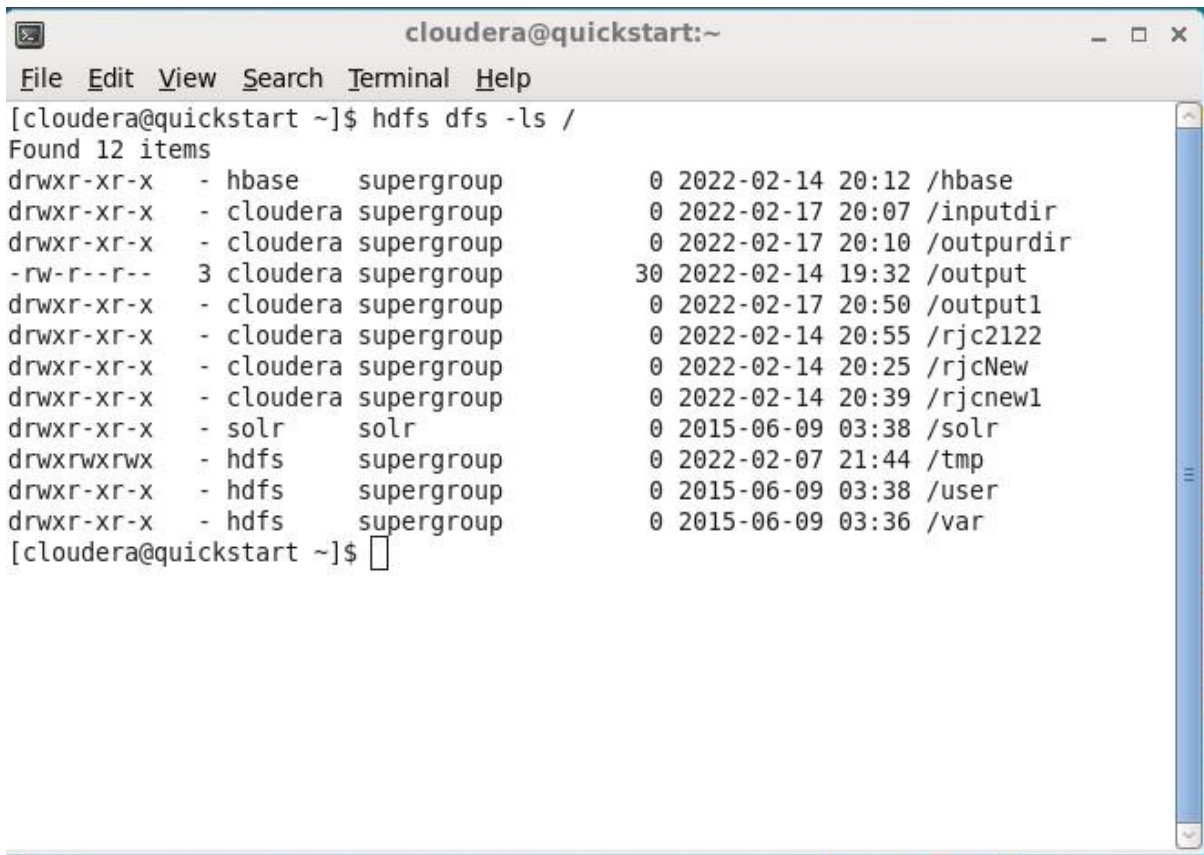

- 6) Right Click on the project name MaxTemp -> Export -> Java -> JAR File -> Next -> for select the export destination for JAR file: browse -> Name :

MaxTemp.jar -> save in folder -> cloudera -> Finish -> OK





- 7) Open terminal and type hdfs dfs -ls/ command Here listing all the directory present in hdfs.

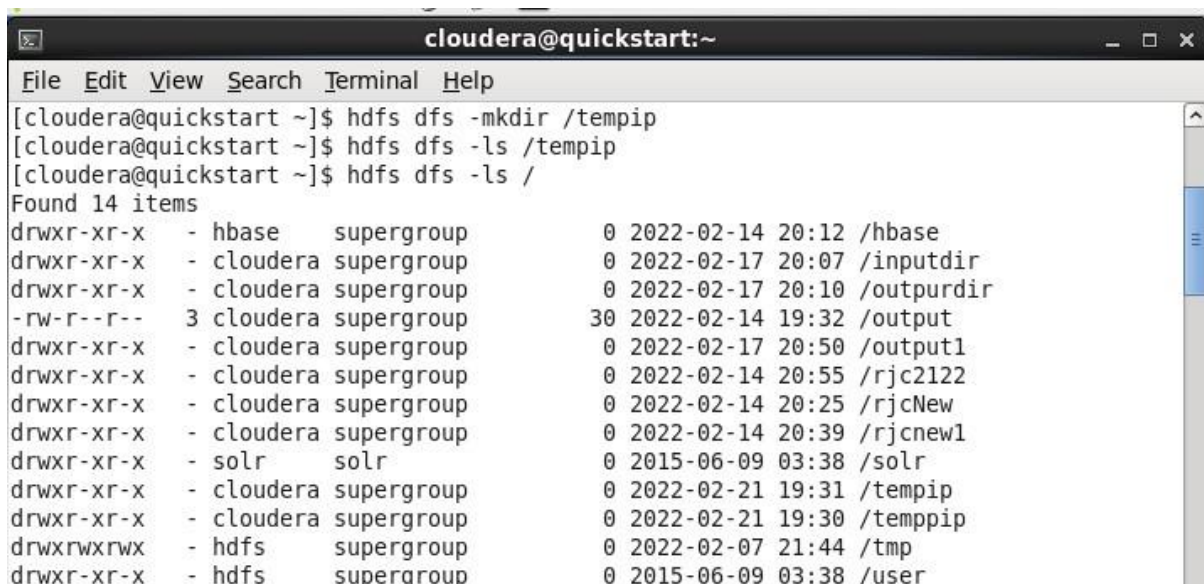


```

cloudera@quickstart:~
File Edit View Search Terminal Help
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 12 items
drwxr-xr-x  - hbase      supergroup      0 2022-02-14 20:12 /hbase
drwxr-xr-x  - cloudera  supergroup      0 2022-02-17 20:07 /inputdir
drwxr-xr-x  - cloudera  supergroup      0 2022-02-17 20:10 /outpurdir
-rw-r--r--  3 cloudera  supergroup     30 2022-02-14 19:32 /output
drwxr-xr-x  - cloudera  supergroup      0 2022-02-17 20:50 /output1
drwxr-xr-x  - cloudera  supergroup      0 2022-02-14 20:55 /rjc2122
drwxr-xr-x  - cloudera  supergroup      0 2022-02-14 20:25 /rjcNew
drwxr-xr-x  - cloudera  supergroup      0 2022-02-14 20:39 /rjcnew1
drwxr-xr-x  - solr      solr            0 2015-06-09 03:38 /solr
drwxrwxrwx  - hdfs      supergroup      0 2022-02-07 21:44 /tmp
drwxr-xr-x  - hdfs      supergroup      0 2015-06-09 03:38 /user
drwxr-xr-x  - hdfs      supergroup      0 2015-06-09 03:36 /var
[cloudera@quickstart ~]$

```

- 8) Input file named as Temp which is present on desktop i.e. in local file system.
- 9) Now we have to move this input file to hdfs. For this we create a direcopy on hdfs using command **hdfs dfs -mkdir /tempip**.



```

cloudera@quickstart:~
File Edit View Search Terminal Help
[cloudera@quickstart ~]$ hdfs dfs -mkdir /tempip
[cloudera@quickstart ~]$ hdfs dfs -ls /tempip
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 14 items
drwxr-xr-x  - hbase      supergroup      0 2022-02-14 20:12 /hbase
drwxr-xr-x  - cloudera  supergroup      0 2022-02-17 20:07 /inputdir
drwxr-xr-x  - cloudera  supergroup      0 2022-02-17 20:10 /outpurdir
-rw-r--r--  3 cloudera  supergroup     30 2022-02-14 19:32 /output
drwxr-xr-x  - cloudera  supergroup      0 2022-02-17 20:50 /output1
drwxr-xr-x  - cloudera  supergroup      0 2022-02-14 20:55 /rjc2122
drwxr-xr-x  - cloudera  supergroup      0 2022-02-14 20:25 /rjcNew
drwxr-xr-x  - cloudera  supergroup      0 2022-02-14 20:39 /rjcnew1
drwxr-xr-x  - solr      solr            0 2015-06-09 03:38 /solr
drwxr-xr-x  - cloudera  supergroup      0 2022-02-21 19:31 /tempip
drwxr-xr-x  - cloudera  supergroup      0 2022-02-21 19:30 /temppip
drwxrwxrwx  - hdfs      supergroup      0 2022-02-07 21:44 /tmp
drwxr-xr-x  - hdfs      supergroup      0 2015-06-09 03:38 /user

```

- 10) Move the input file i.e. temperature to this directory created in hdfs by using either put command or copyFromLocal command.

- 11) Now checking whether the “Temp” present in /tempip directory of hdfs or not using **hdfs dfs -ls /tempip** command

```

cloudera@quickstart:~
File Edit View Search Terminal Help
[cloudera@quickstart ~]$ hdfs dfs -put /home/cloudera/Desktop/Temp /tempip
/
[cloudera@quickstart ~]$ hdfs dfs -ls /tempip
Found 1 items
-rw-r--r--  1 cloudera supergroup          530 2022-02-21 19:39 /tempip/Temp

```

- 12) As we can see “Temp” file is present in /tempip directory of hdfs. Now we will see the content of this file using **hdfs dfs -cat /tempip/temperature** command

```

cloudera@quickstart:~
File Edit View Search Terminal Help
[cloudera@quickstart ~]$ hdfs dfs -cat /tempip/Temp
0067011990999991950051507004+68750+023550FM-12+038299999V0203301N00671220001CN99999999N9
+00001+999999999999
0043011990999991950051512004+68750+023550FM-12+038299999V0203201N00671220001CN99999999N9
+00221+999999999999
0043011990999991950051518004+68750+023550FM-12+038299999V0203201N00261220001CN99999999N9
-00111+999999999999
0043012650999991949032412004+62300+010750FM-12+048599999V0202701N00461220001CN0500001N9
+01111+999999999999
0043012650999991949032418004+62300+010750FM-12+048599999V0202701N00461220001CN0500001N9
+00781+999999999999

```

- 13) Running Mapreduce Program on Hadoop, syntax is **hadoop jar jarFileName.jar ClassName /InputFileAddress /outputdir**

```

cloudera-quickstart-vm-5.4.2-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System
cloudera@quickstart:~
File Edit View Search Terminal Help
[cloudera@quickstart ~]$ hadoop jar /home/cloudera/MaxTemp.jar MaxTemp /tempip/Temp /tempop1
22/02/21 19:50:21 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
22/02/21 19:50:23 WARN mapreduce.JobSubmitter: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
22/02/21 19:50:23 INFO input.FileInputFormat: Total input paths to process : 1
22/02/21 19:50:23 INFO mapreduce.JobSubmitter: number of splits:1
22/02/21 19:50:24 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1644894610889_0003
22/02/21 19:50:24 INFO impl.YarnClientImpl: Submitted application application_1644894610889_0003
22/02/21 19:50:25 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1644894610889_0003/
22/02/21 19:50:25 INFO mapreduce.Job: Running job: job_1644894610889_0003
22/02/21 19:50:38 INFO mapreduce.Job: Job job_1644894610889_0003 running in uber mode : false
22/02/21 19:50:38 INFO mapreduce.Job: map 0% reduce 0%
22/02/21 19:50:57 INFO mapreduce.Job: map 100% reduce 0%
22/02/21 19:51:10 INFO mapreduce.Job: map 100% reduce 100%
22/02/21 19:51:10 INFO mapreduce.Job: Job job_1644894610889_0003 completed successfully
22/02/21 19:51:10 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=61
    FILE: Number of bytes written=221193
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=638
    HDFS: Number of bytes written=17
    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)=16122
    Total time spent by all reduces in occupied slots (ms)=10551
    Total time spent by all map tasks (ms)=16122
    Total time spent by all reduce tasks (ms)=10551
    Total vcore-seconds taken by all map tasks=16122
    Total vcore-seconds taken by all reduce tasks=10551
    Total megabyte-seconds taken by all map tasks=16508928
    Total megabyte-seconds taken by all reduce tasks=10804224
  Map-Reduce Framework

```

Map-Reduce Framework

```

cloudera-quickstart-vm-5.4.2-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System
cloudera@quickstart:~
File Edit View Search Terminal Help
Total time spent by all maps in occupied slots (ms)=16122
Total time spent by all reduces in occupied slots (ms)=10551
Total time spent by all map tasks (ms)=16122
Total time spent by all reduce tasks (ms)=10551
Total vcore-seconds taken by all map tasks=16122
Total vcore-seconds taken by all reduce tasks=10551
Total megabyte-seconds taken by all map tasks=16908928
Total megabyte-seconds taken by all reduce tasks=10804224
Map-Reduce Framework
  Map input records=5
  Map output records=5
  Map output bytes=45
  Map output materialized bytes=61
  Input split bytes=108
  Combine input records=0
  Combine output records=0
  Reduce input groups=2
  Reduce shuffle bytes=61
  Reduce input records=5
  Reduce output records=2
  Spilled Records=10
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=279
  CPU time spent (ms)=1610
  Physical memory (bytes) snapshot=329367552
  Virtual memory (bytes) snapshot=3007225856
  Total committed heap usage (bytes)=226365440
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=530
File Output Format Counters
  Bytes Written=17
  
```

- 14) Then we can verify the content of tempop1 directory and in that part-r file has the actual output by using the command `Hdfs dfs -cat /tempop1/part-r-00000` This will give us final output. The same file can also be accessed using a browser. For every execution of this program we need to delete the output directory or give a new name to the output directory every time. 1st we are checking whether the tempop1 directory is created in hdfs or not using command **`hdfs dfs -ls /`**
- 15) Now let's check what we have inside this **tempop1** directory using command as **`hdfs dfs -ls /tempop1`**
- 16) Now we want to read the content of the **part-r-00000** file which present inside the **tempop1** using command **`hdfs dfs -cat /tempop1/part-r-00000`**


```
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 15 items
drwxr-xr-x - hbase supergroup 0 2022-02-14 20:12 /hbase
drwxr-xr-x - cloudera supergroup 0 2022-02-17 20:07 /inputdir
drwxr-xr-x - cloudera supergroup 0 2022-02-17 20:10 /outputdir
-rw-r--r-- 3 cloudera supergroup 30 2022-02-14 19:32 /output
drwxr-xr-x - cloudera supergroup 0 2022-02-17 20:50 /output1
drwxr-xr-x - cloudera supergroup 0 2022-02-14 20:55 /rjc2122
drwxr-xr-x - cloudera supergroup 0 2022-02-14 20:25 /rjcNew
drwxr-xr-x - cloudera supergroup 0 2022-02-14 20:39 /rjcnew1
drwxr-xr-x - solr solr 0 2015-06-09 03:38 /solr
drwxr-xr-x - cloudera supergroup 0 2022-02-21 19:39 /tempip
drwxr-xr-x - cloudera supergroup 0 2022-02-21 19:51 /tempop1
drwxr-xr-x - cloudera supergroup 0 2022-02-21 19:30 /tempip1
drwxrwxrwx - hdfs supergroup 0 2022-02-07 21:44 /tmp
drwxr-xr-x - hdfs supergroup 0 2015-06-09 03:38 /user
drwxr-xr-x - hdfs supergroup 0 2015-06-09 03:36 /var
[cloudera@quickstart ~]$ hdfs dfs -ls /tempop1
Found 2 items
-rw-r--r-- 1 cloudera supergroup 0 2022-02-21 19:51 /tempop1/_SUCCESS
-rw-r--r-- 1 cloudera supergroup 17 2022-02-21 19:51 /tempop1/part-r-00000
[cloudera@quickstart ~]$ hdfs dfs -cat /tempop1/part-r-00000
1949 111
1950 22
[cloudera@quickstart ~]$
```

So the maximum temperature for the year 1949 is 111 and for the year 1950 is 22.

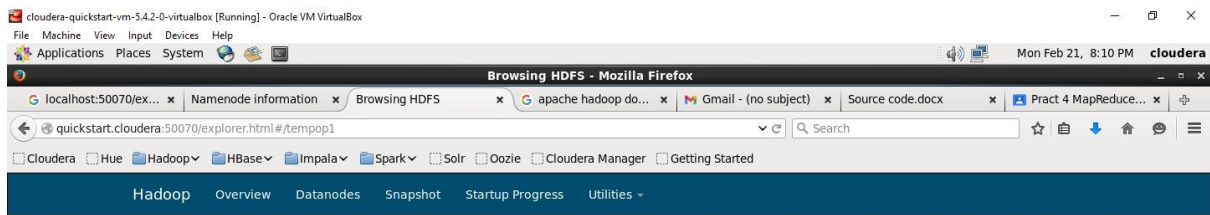
18) Browse the Directory by
Hadoop->HDFS Namenode->Utilities ->Browse the file system

The screenshot shows a web browser window with the title "Browsing HDFS - Mozilla Firefox". The address bar shows "localhost:50070/explorer.html#". The browser displays a table of files and directories in the HDFS file system. The table has columns for Permission, Owner, Group, Size, Replication, Block Size, and Name. The files listed are:

Permission	Owner	Group	Size	Replication	Block Size	Name
drwxr-xr-x	hbase	supergroup	0 B	0	0 B	hbase
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	inputdir
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	outputdir
-rw-r--r--	cloudera	supergroup	30 B	3	128 MB	output
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	output1
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	rjc2122
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	rjcNew
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	rjcnew1
drwxr-xr-x	solr	solr	0 B	0	0 B	solr
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	tempip
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	tempop1
drwxr-xr-x	cloudera	supergroup	0 B	0	0 B	tempip1
drwxrwxrwx	hdfs	supergroup	0 B	0	0 B	tmp
drwxr-xr-x	hdfs	supergroup	0 B	0	0 B	user
drwxr-xr-x	hdfs	supergroup	0 B	0	0 B	var

The file "tempop1" is highlighted, and its details are shown below the table:

Permission	Owner	Group	Size	Replication	Block Size	Name
-rw-r--r--	1 cloudera	supergroup	0	2022-02-21 19:51		/tempop1/_SUCCESS
-rw-r--r--	1 cloudera	supergroup	17	2022-02-21 19:51		/tempop1/part-r-00000



Browse Directory

/tempop1						Go!
Permission	Owner	Group	Size	Replication	Block Size	Name
-rw-r--r--	cloudera	supergroup	0 B	1	128 MB	_SUCCESS
-rw-r--r--	cloudera	supergroup	17 B	1	128 MB	part-r-00000

Hadoop, 2014.

Now downloading the **part-r-00000** file.

