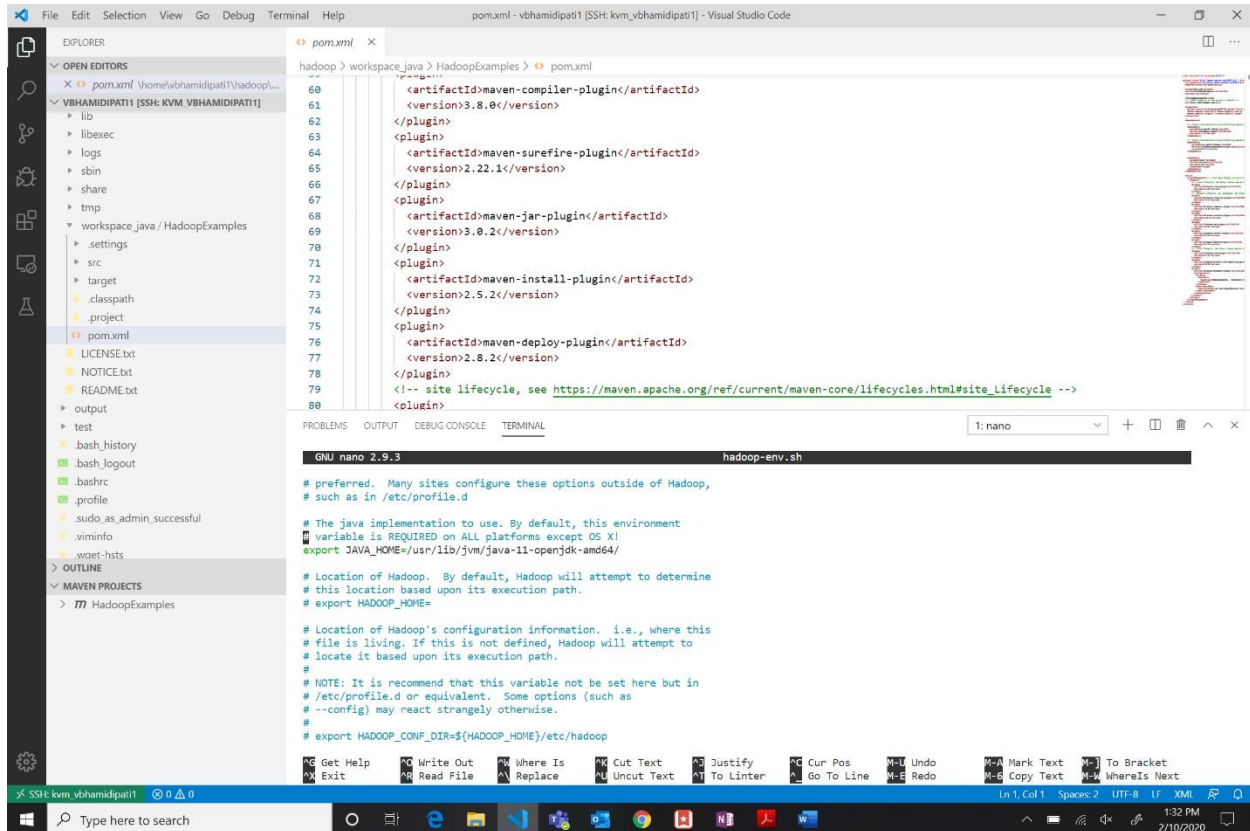


Setting up Hadoop on Ubuntu:

1. Install Java Using
 - a. Sudo apt install default-jdk
2. Set the root of the java installation
 - a. vim hadoop-env.sh
 - b. Update the line port JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64



The screenshot shows the Visual Studio Code interface. The Explorer pane on the left shows the project structure with files like pom.xml, LICENSE.txt, NOTICE.txt, README.txt, and output. The main editor shows the pom.xml file with Maven plugins for maven-compiler-plugin, maven-surefire-plugin, maven-jar-plugin, maven-install-plugin, and maven-deploy-plugin. The bottom terminal window shows the hadoop-env.sh file being edited in nano. The file contains comments and export statements for JAVA_HOME, HADOOP_HOME, and HADOOP_CONF_DIR.

```
hadoop > workspace.java > HadoopExamples > pom.xml
60 <artifactId>maven-compiler-plugin</artifactId>
61 <version>3.8.0</version>
62 </plugin>
63 <plugin>
64 <artifactId>maven-surefire-plugin</artifactId>
65 <version>2.22.1</version>
66 </plugin>
67 <plugin>
68 <artifactId>maven-jar-plugin</artifactId>
69 <version>3.0.2</version>
70 </plugin>
71 <plugin>
72 <artifactId>maven-install-plugin</artifactId>
73 <version>2.5.2</version>
74 </plugin>
75 <plugin>
76 <artifactId>maven-deploy-plugin</artifactId>
77 <version>2.8.2</version>
78 </plugin>
79 <!-- site lifecycle, see https://maven.apache.org/ref/current/maven-core/lifecycles.html#site_lifecycle -->
80 </plugin>
```

```
GNU nano 2.9.3 hadoop-env.sh
# preferred. Many sites configure these options outside of Hadoop,
# such as in /etc/profile.d

# The java implementation to use. By default, this environment
# variable is REQUIRED on ALL platforms except OS X!
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64/

# Location of Hadoop. By default, Hadoop will attempt to determine
# this location based upon its execution path.
# export HADOOP_HOME=

# Location of Hadoop's configuration information. i.e., where this
# file is living. If this is not defined, Hadoop will attempt to
# locate it based upon its execution path.
#
# NOTE: It is recommend that this variable not be set here but in
# /etc/profile.d or equivalent. Some options (such as
# --config) may react strangely otherwise.
# export HADOOP_CONF_DIR=${HADOOP_HOME}/etc/hadoop
```

3. Core-site.xml
 - a. This is used to specify the default file system and defaults to the local file system that's why it needs to be set to a HDFS address. This is important for client configuration as well so your local configuration file should include this element.
 - b. When client request for the datanode information to the namenode then the namenode give them, the information about the datanode i.e hdfs details and then client write the dataset into the datanode.
 - c. Add the following configurations
 - i. <configuration>
 - ii. <property>
 - iii. <name>hadoop.tmp.dir</name>
 - iv. <value>/home/vbhamidipati1/hadoop/tmp</value>
 - v. </property>
 - vi. <property>
 - vii. <name>fs.defaultFS</name>
 - viii. <value>hdfs://localhost:9000</value>

- ix. `<!-- <value>file:///</value> -->`
- x. `</property>`
- xi. `</configuration>`
- d. Here fs denotes file system and default.name denotes namenode
- e. Here 9000 denotes port on which datanode will send heartbeat to namenode. And full address is the machine name that is converted to hostname.
- f. The fs.defaultFS makes HDFS a file abstraction over a cluster, so that its root is not the same as the local system's. You need to change the value in order to create the distributed file system. The fs.defaultFS in core-site.xml gives the datanode address of namenode. The datanode looks here for the namenode address and tries to contact it using RPC.

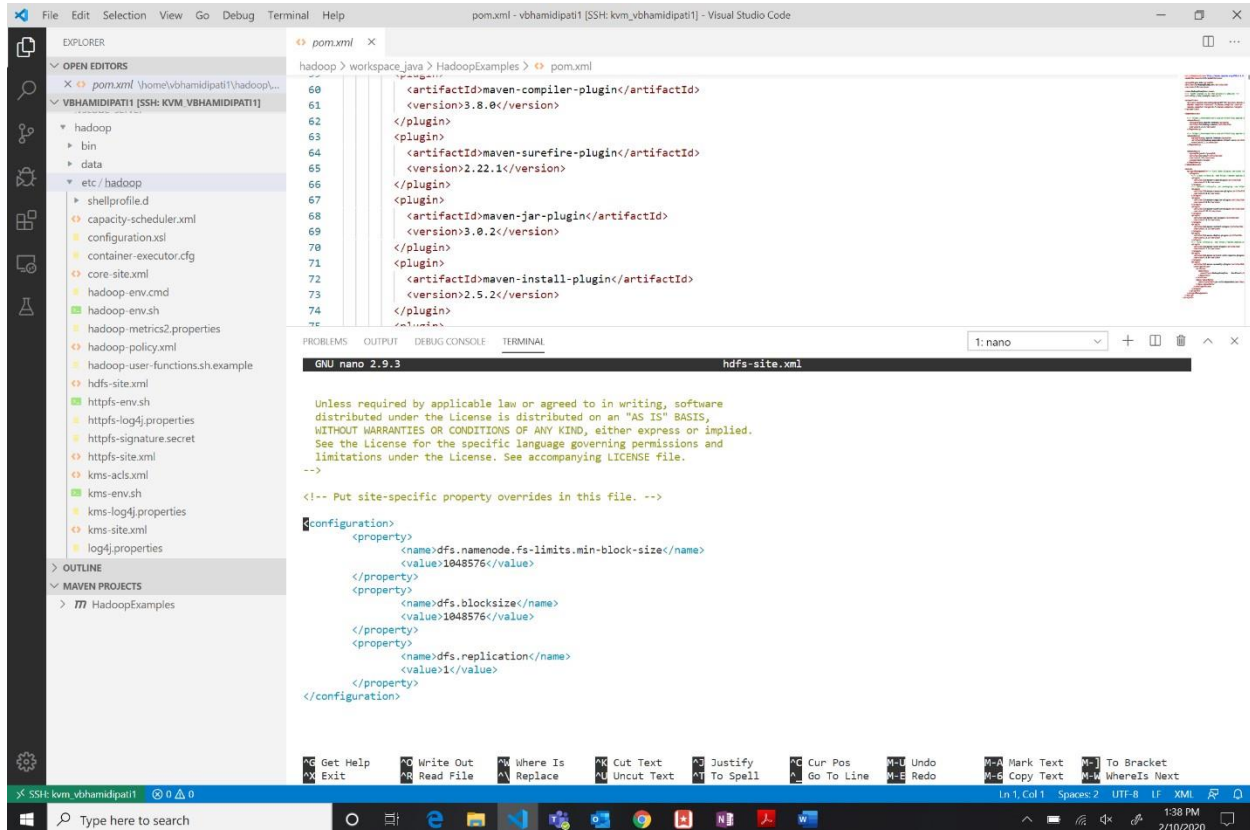
The screenshot shows a Visual Studio Code window with two main panes. The top pane displays the `pom.xml` file, which contains Maven plugin configurations for `maven-compiler-plugin`, `maven-surefire-plugin`, `maven-jar-plugin`, `maven-install-plugin`, and `maven-deploy-plugin`. The bottom pane shows a terminal window running the `nano` editor, editing the `core-site.xml` file. The `core-site.xml` file contains a `<configuration>` block with `<property>` tags for `hadoop.tmp.dir`, `fs.defaultFS`, and a commented-out `hdfs` property. The `fs.defaultFS` property is set to `hdfs://localhost:9000`. The terminal window also shows the license text for GNU nano 2.9.3.

4. hdfs-site.xml

- a. Here we set the following:
 - i. Minimum block size
 - ii. Default block size
 - iii. Replication value
- b. In Hadoop, HDFS splits huge files into small chunks known as data blocks. HDFS Data blocks are the smallest unit of data in a filesystem. We (client and admin) do not have any control over the data block like block location. Namenode decides all such things.
- c. HDFS stores each file as a data block. However, the data block size in HDFS is very large. The default size of the HDFS block is 128MB which you can configure as per your

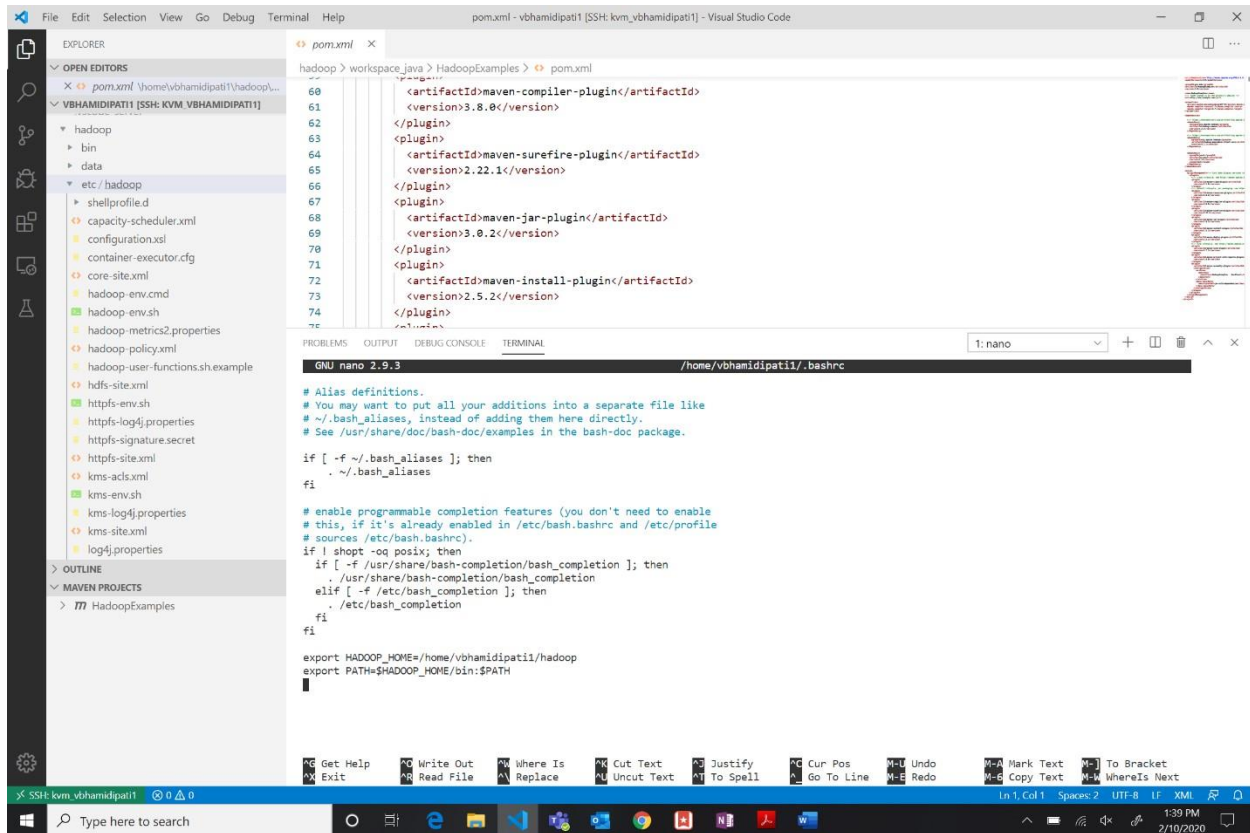
requirement. All blocks of the file are the same size except the last block, which can be either the same size or smaller. The files are split into 128 MB blocks and then stored into the Hadoop file system. The Hadoop application is responsible for distributing the data block across multiple nodes. Here in our example, the block size is 1 MB.

- d. Replication factor dictates how many copies of a block should be kept in your cluster. The replication factor is 3 by default (there would be one original block and two replicas) and hence any file you create in HDFS will have a replication factor of 3 and each block from the file will be copied to 3 different nodes in your cluster.



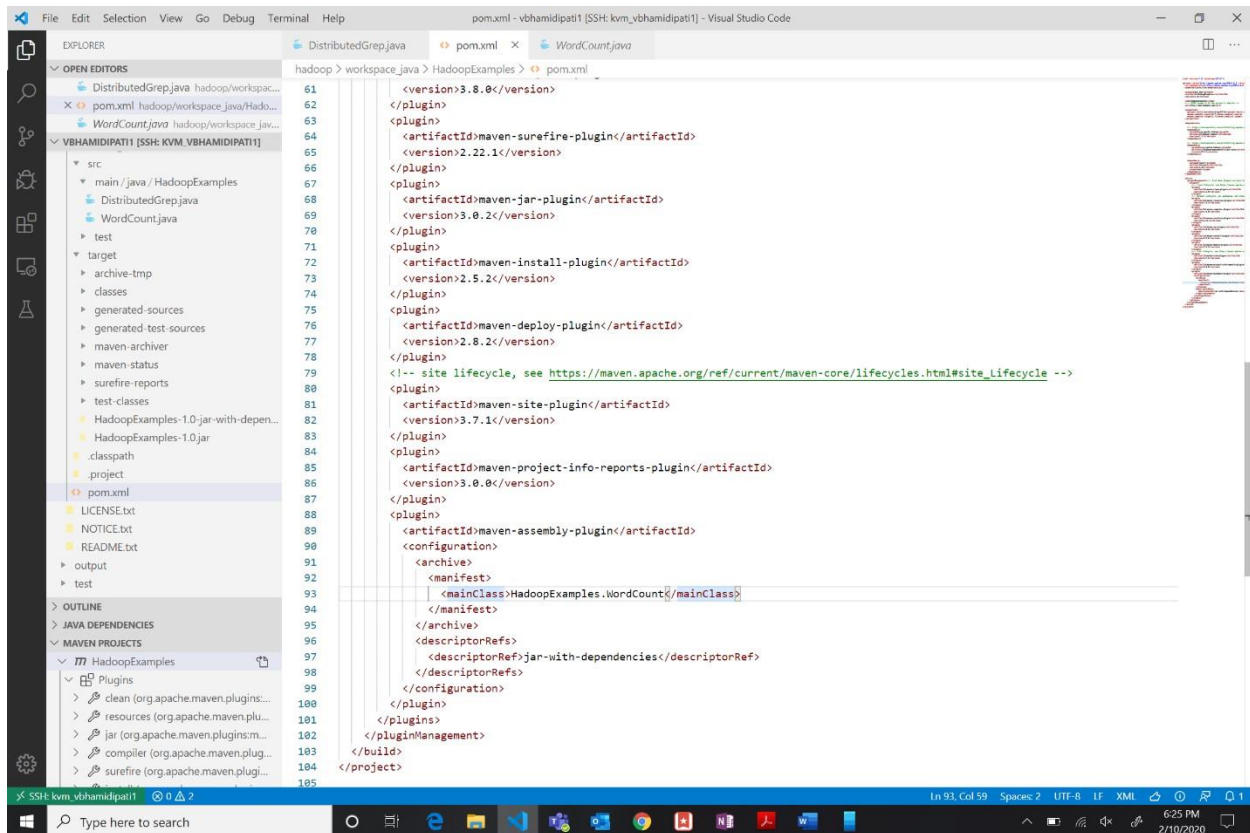
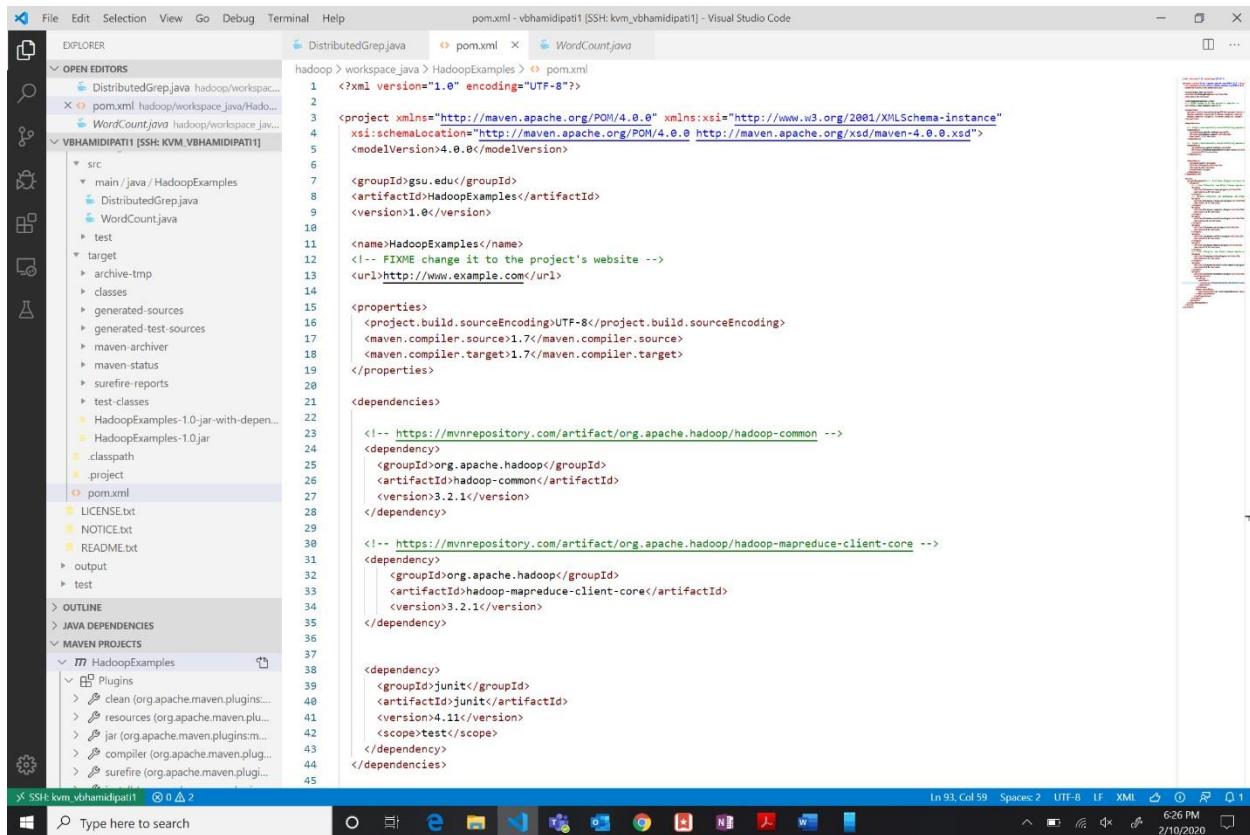
5. bashrc

- a. `.bashrc` is a shell script that Bash runs whenever it is started interactively. It initializes an interactive shell session. You can put any command in that file that you could type at the command prompt.
- b. Here we set the Hadoop path to identify it globally



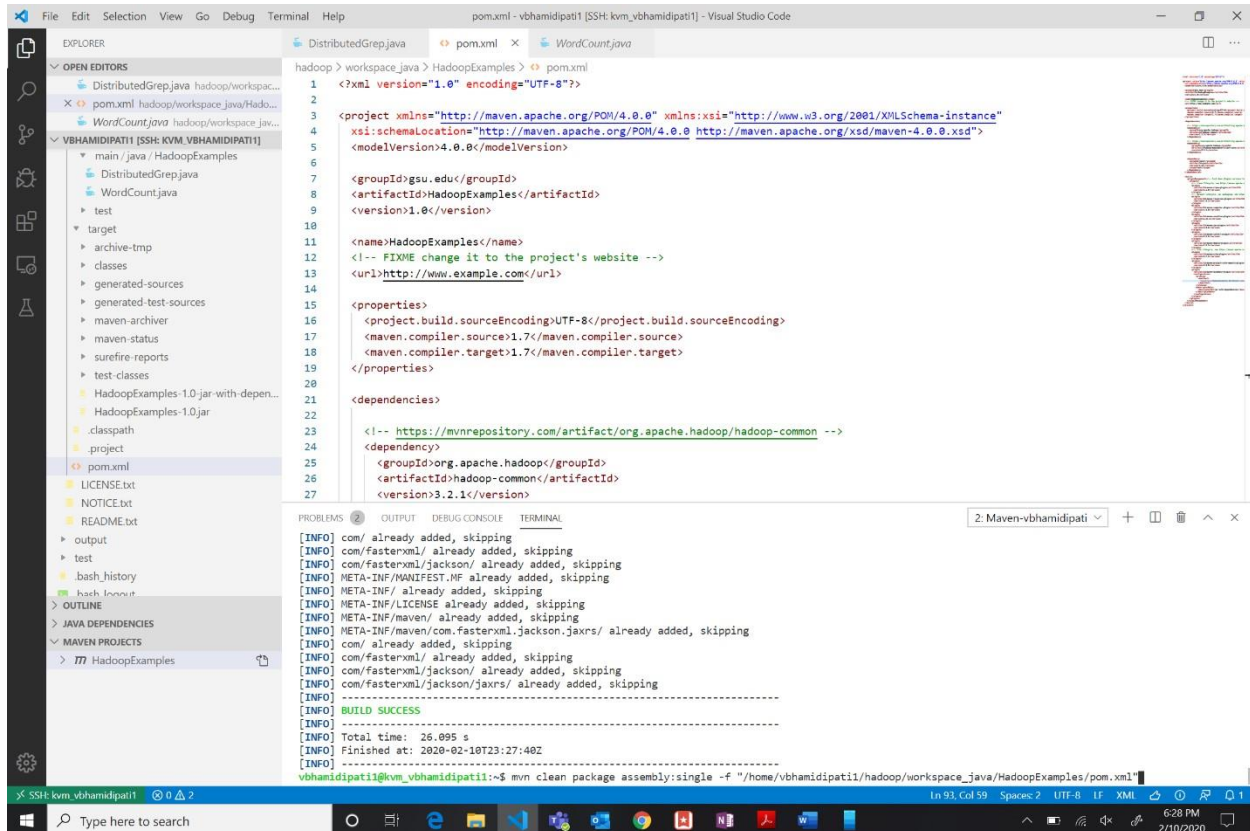
6. Updating the pom.xml file

- Add the following dependencies
 - Apache Hadoop Common 3.2.1
 - Apache Hadoop MapReduce Core 3.2.1
- Add the dependencies by accessing the dependency code from the maven repository to the pom.xml file
- Also specify the main class file for execution as <PackageName.ClassName>
- In our example we used HadoopExamples.WordCount as our main class



7. Maven Assembly pulgin and build the project

- It will create a runnable jar file for the execution of wordcount program
- The jars can be found in the target folder



8. Running the jar file

- We run the jar file using the command

```
hadoop jar /home/vbhamidipati1/hadoop/workspace_java/HadoopExamples/target/HadoopExamples-1.0-jar-with-dependencies.jar file:///home/vbhamidipati1/hadoop/data/peterpan.txt file:///home/vbhamidipati1/hadoop/data/output
```

Explanation:

- The command uses the complete path of jar as stated by "jar /home/vbhamidipati1/hadoop/workspace_java/HadoopExamples/target/HadoopExamples-1.0-jar-with-dependencies.jar"
- It takes in the input text file as <file:///home/vbhamidipati1/hadoop/data/peterpan.txt>. It specifies that the file has to be taken from the local (Ubuntu) file system not the HDFS.
- It outputs text file at file:///home/vbhamidipati1/hadoop/data/output. It specifies that the file is written to local (Ubuntu) file system not the HDFS.

Visual Studio Code interface showing a Maven project named `wordcount` in the `workspace.java` directory. The `pom.xml` file is open, displaying the project configuration for `org.apache.hadoop` with version `4.0.0`. The terminal output shows the execution of `hadoop jar` command, indicating the start of a MapReduce job. The output includes various log messages from the Hadoop framework, such as `impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).` and `mapred.LocalJobRunner: Starting task: attempt_local775807566_0001_m_000000_0`. The job is running on a local host, and the output is being written to the `hadoop/output` directory. The terminal output ends with `mapred.Task: Final Counters for attempt_local775807566_0001_m_000000_0: Counters: 24`.

```
<?xml version="1.0" encoding="UTF-8"?>
<project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org/xsd/maven-4.0.0.xsd">
  <modelVersion>4.0.0</modelVersion>
  <groupId>gsu.edu</groupId>
  <artifactId>wordcount</artifactId>
  <version>1.0</version>
  <packaging>jar</packaging>
  <name>wordcount</name>
  <description>A simple wordcount application</description>
  <dependencies>
    <dependency>
      <groupId>org.apache.hadoop</groupId>
      <artifactId>hadoop-common</artifactId>
      <version>4.0.0</version>
      <scope>compile</scope>
    </dependency>
    <dependency>
      <groupId>org.apache.hadoop</groupId>
      <artifactId>hadoop-mapreduce</artifactId>
      <version>4.0.0</version>
      <scope>compile</scope>
    </dependency>
  </dependencies>
  <build>
    <plugins>
      <plugin>
        <groupId>org.apache.maven.plugins</groupId>
        <artifactId>maven-compiler-plugin</artifactId>
        <version>3.8.1</version>
        <configuration>
          <source>1.8</source>
          <target>1.8</target>
        </configuration>
      </plugin>
    </plugins>
  </build>
  <profiles>
    <profile>
      <id>dev</id>
      <activation>
        <activeByDefault>true</activeByDefault>
      </activation>
      <build>
        <plugins>
          <plugin>
            <groupId>org.apache.maven.plugins</groupId>
            <artifactId>maven-compiler-plugin</artifactId>
            <version>3.8.1</version>
            <configuration>
              <source>1.8</source>
              <target>1.8</target>
            </configuration>
          </plugin>
        </plugins>
      </build>
    </profile>
  </profiles>
</project>
```

```
2020-02-10 23:33:02,733 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2020-02-10 23:33:02,733 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2020-02-10 23:33:02,871 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2020-02-10 23:33:04,162 INFO input.FileInputFormat: Total input files to process : 1
2020-02-10 23:33:04,206 INFO mapreduce.JobSubmitter: number of splits:1
2020-02-10 23:33:04,573 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local775807566_0001
2020-02-10 23:33:04,576 INFO mapreduce.JobSubmitter: Executing with tokens: []
2020-02-10 23:33:04,796 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2020-02-10 23:33:04,797 INFO mapreduce.Job: Running job: job_local775807566_0001
2020-02-10 23:33:04,801 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2020-02-10 23:33:04,814 INFO output.FileOutputCommitter: FileOutputCommitter Algorithm version is 2
2020-02-10 23:33:04,818 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore c
leanup failures: false
2020-02-10 23:33:04,815 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2020-02-10 23:33:04,872 INFO mapred.LocalJobRunner: Waiting for map tasks
2020-02-10 23:33:04,873 INFO mapred.LocalJobRunner: Starting task: attempt_local775807566_0001_m_000000_0
2020-02-10 23:33:04,916 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2020-02-10 23:33:04,918 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore c
leanup failures: false
2020-02-10 23:33:04,961 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2020-02-10 23:33:04,969 INFO mapred.MapTask: Processing split: file:///home/vbhamidipati1/hadoop/data/peterpan.txt:0+291927
2020-02-10 23:33:05,251 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
2020-02-10 23:33:05,252 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2020-02-10 23:33:05,252 INFO mapred.MapTask: soft limit at 83886080
2020-02-10 23:33:05,253 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2020-02-10 23:33:05,253 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2020-02-10 23:33:05,267 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2020-02-10 23:33:05,275 INFO input.LineRecordReader: Found UTF-8 BOM and skipped it
2020-02-10 23:33:05,488 INFO mapred.LocalJobRunner:
2020-02-10 23:33:05,489 INFO mapred.MapTask: Starting flush of map output
2020-02-10 23:33:05,489 INFO mapred.MapTask: Spilling map output
2020-02-10 23:33:05,489 INFO mapred.MapTask: bufstart = 0; bufend = 486076; bufvoid = 104857600
2020-02-10 23:33:05,489 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 26011340(104845360); length = 203057/6553600
2020-02-10 23:33:05,716 INFO mapred.MapTask: Finished spill 0
2020-02-10 23:33:05,737 INFO mapred.Task: Task:attempt_local775807566_0001_m_000000_0 is done. And is in the process of committing
2020-02-10 23:33:05,739 INFO mapred.LocalJobRunner: map
2020-02-10 23:33:05,743 INFO mapred.Task: Task 'attempt_local775807566_0001_m_000000_0' done.
2020-02-10 23:33:05,750 INFO mapred.Task: Final Counters for attempt_local775807566_0001_m_000000_0: Counters: 24
```

