

UNIVERSITY PARTNER



## **DISTRIBUTE AND CLOUD SYSTEMS PROGRAMMING (5CS022)**

### **WEEK 9 WORKSHOP THE APACHE SHARK FRAMEWORK**

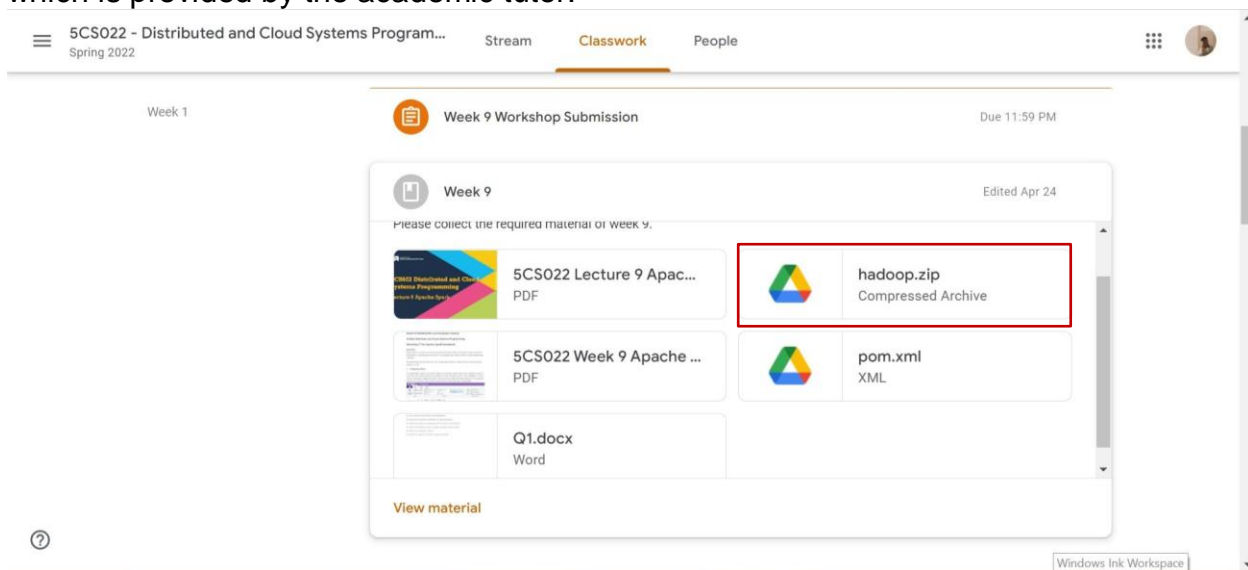
<b>Student Id</b>	: 2065697
<b>Student Name</b>	: Dhiraj Kumar Sah Kanu
<b>Group</b>	: L5CG12
<b>Submitted on</b>	: May 06, 2022

## Overview

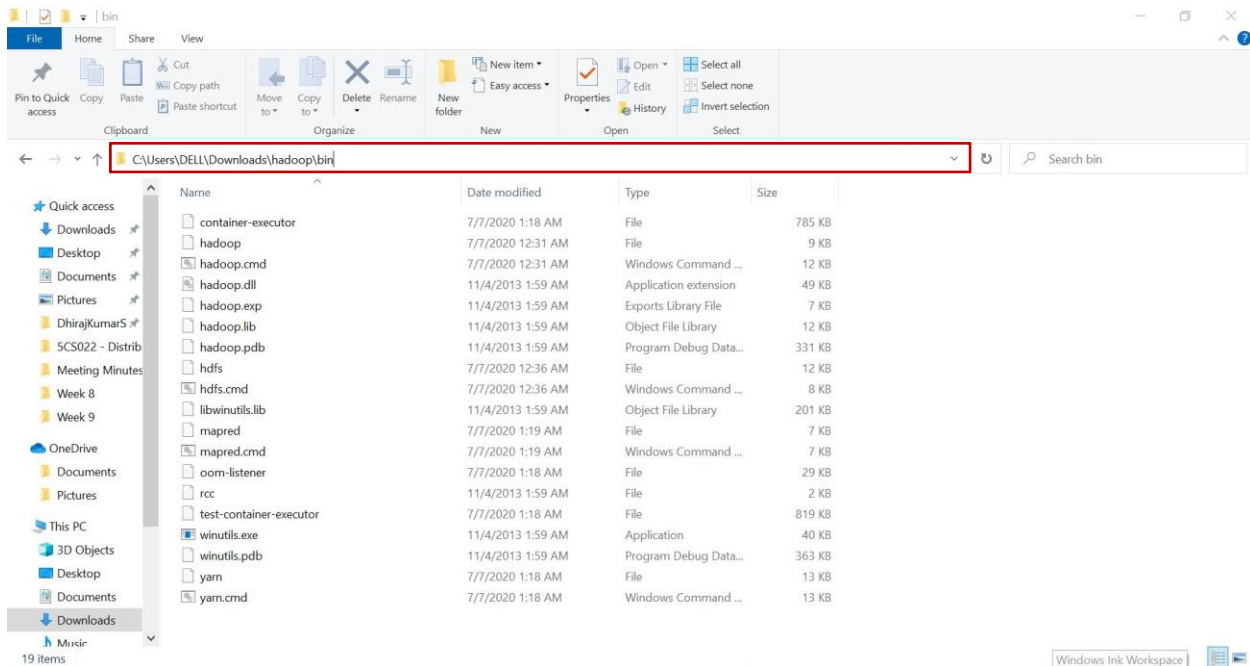
This assignment will help you how to create a small Spark project in Eclipse and run a word and letter counting application in Java. Here, you will set the environment variable to instruct Apache Spark where to find Hadoop. You will start by making a new project in Eclipse. You will edit and replace the code and will run the file. Then you will create the WordCount Program, run the program and see the resultant output. Then you will be creating a Spark program to count letters instead of words.

## Step 1: Downloading Hadoop

To begin, download the file 'Hadoop.zip' from your Week 9 folder in Google Classroom which is provided by the academic tutor.

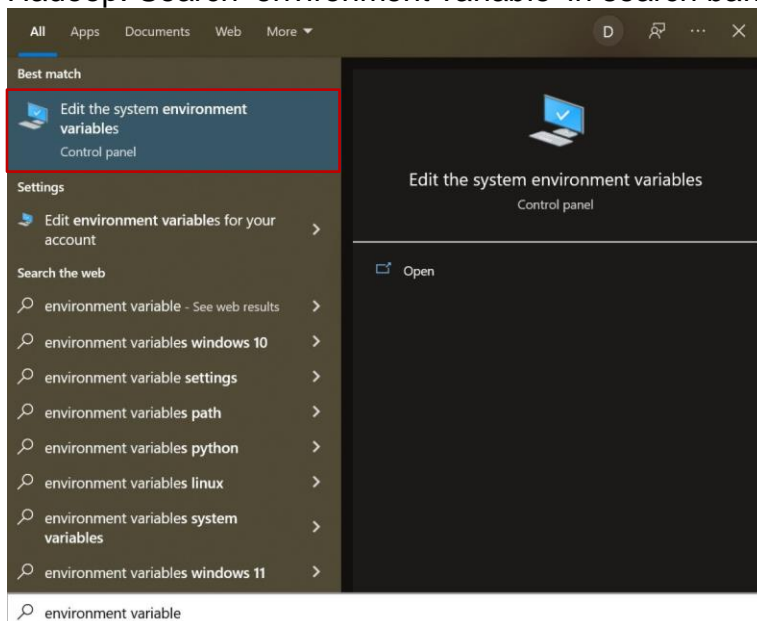


I have downloaded the file and have saved the hadoop file in Downloads, and copy the file location.

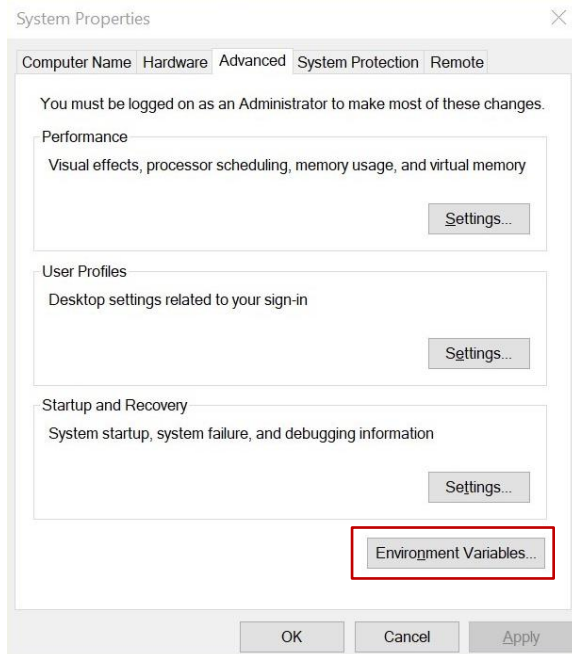


## Step 2: Configuring Hadoop

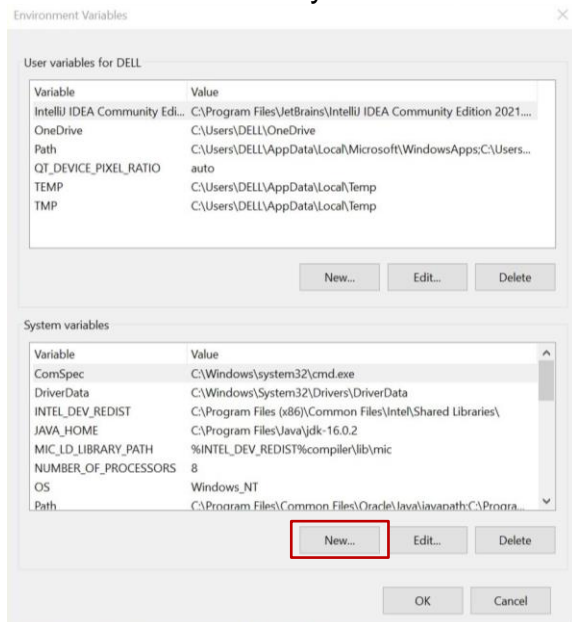
Now you need to set the environment variable to instruct Apache Spark where to find Hadoop. Search 'environment variable' in search bar. Click the first option.



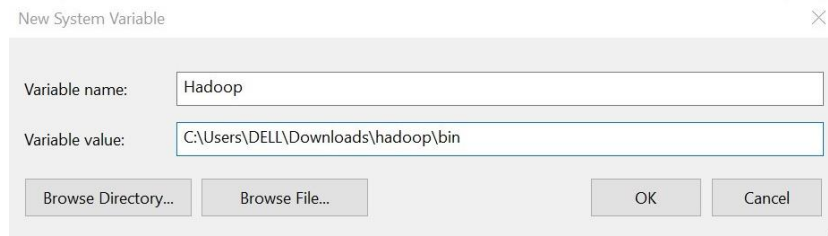
You will see the box of system properties, here click on 'Environment Variables'.



Click on New under system variables section.



Next, you must configure the Environment Variable for this directory as follows, and paste the directory which you have copied previously.



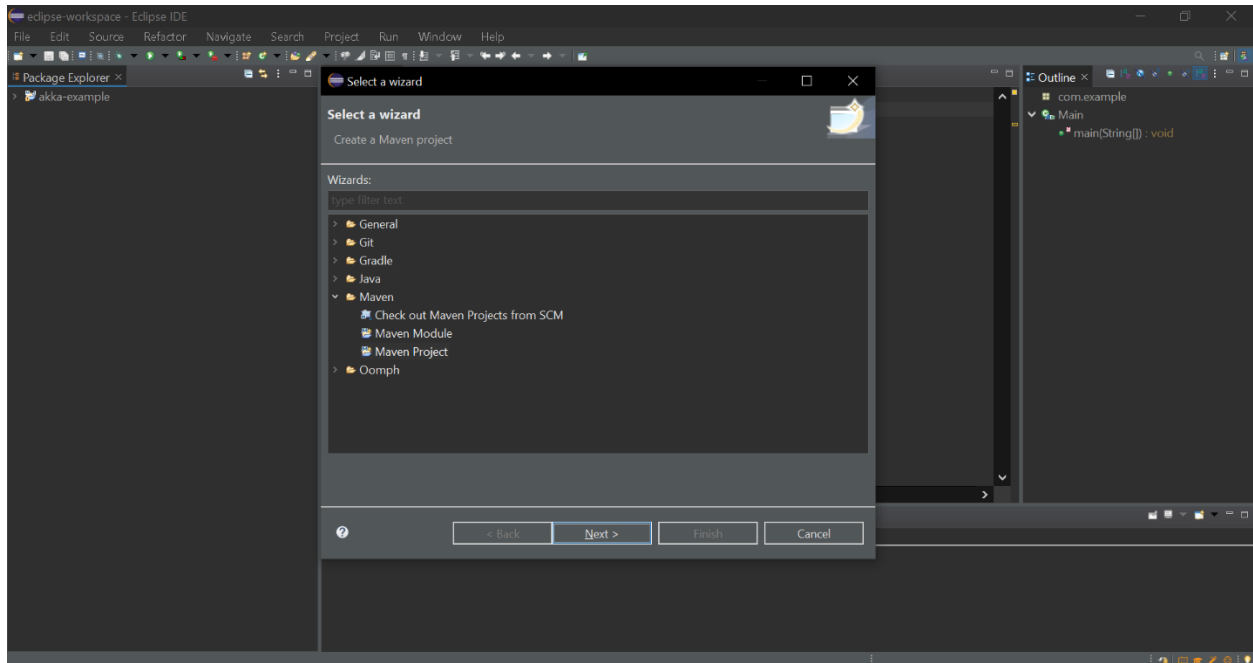
New System Variable

Variable name:

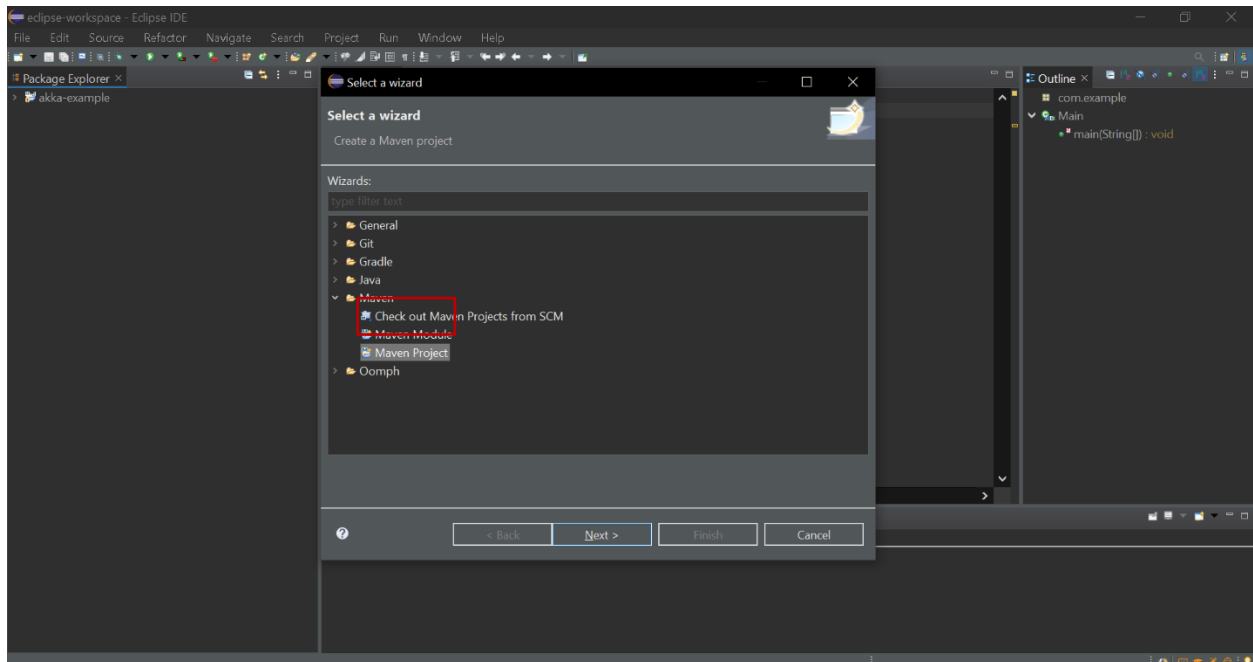
Variable value:

### Step 3: Creating project in Eclipse

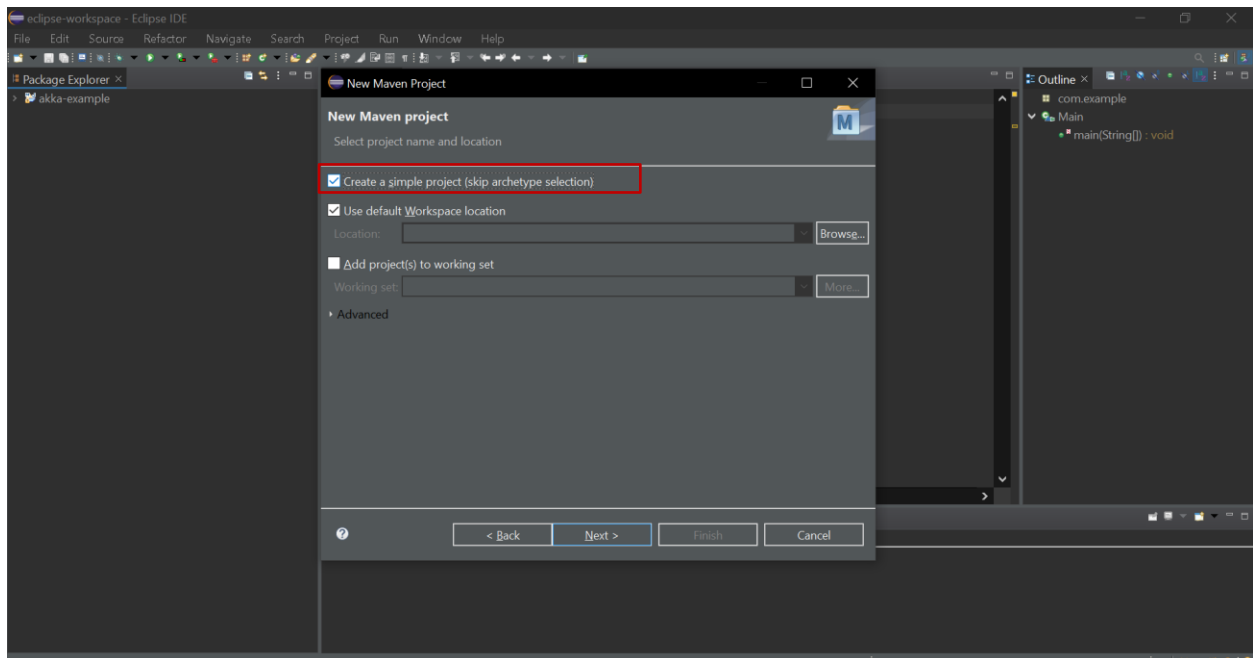
Here you have to start by making a new project in Eclipse. Launch Eclipse and then select New Project.



There should be a list of the many sorts of projects you may build; simply select Maven Project and click next.

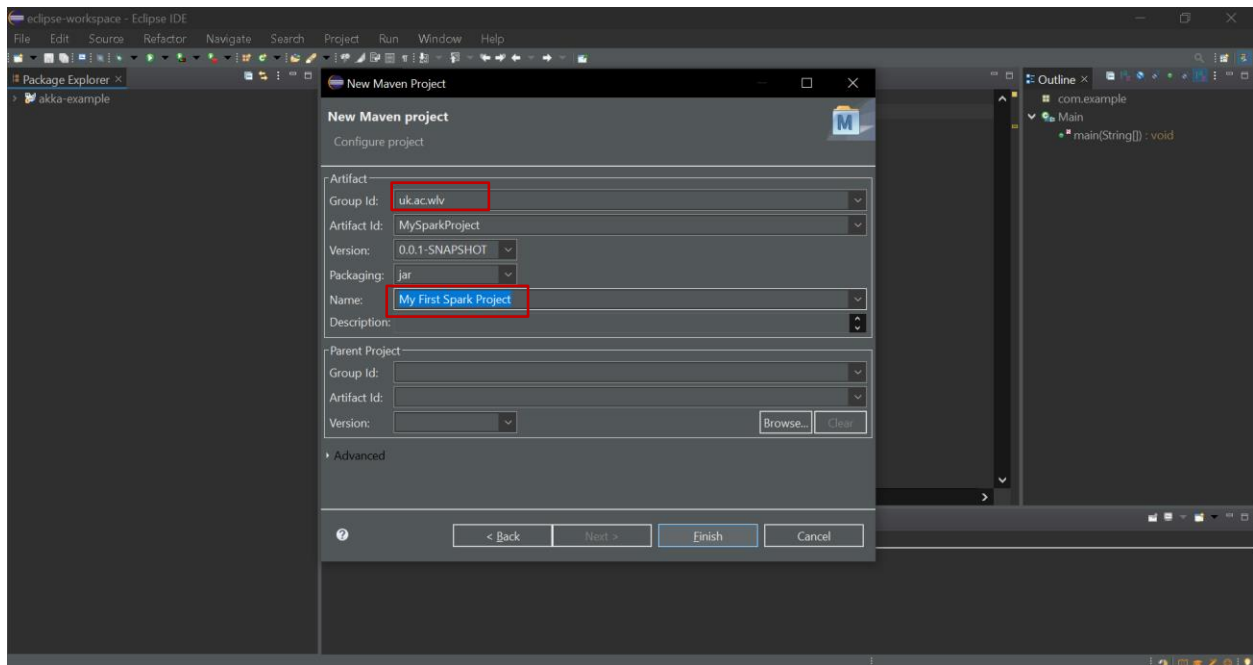


Then, after clicking next. Make sure you choose "simple project" and then click the next button.

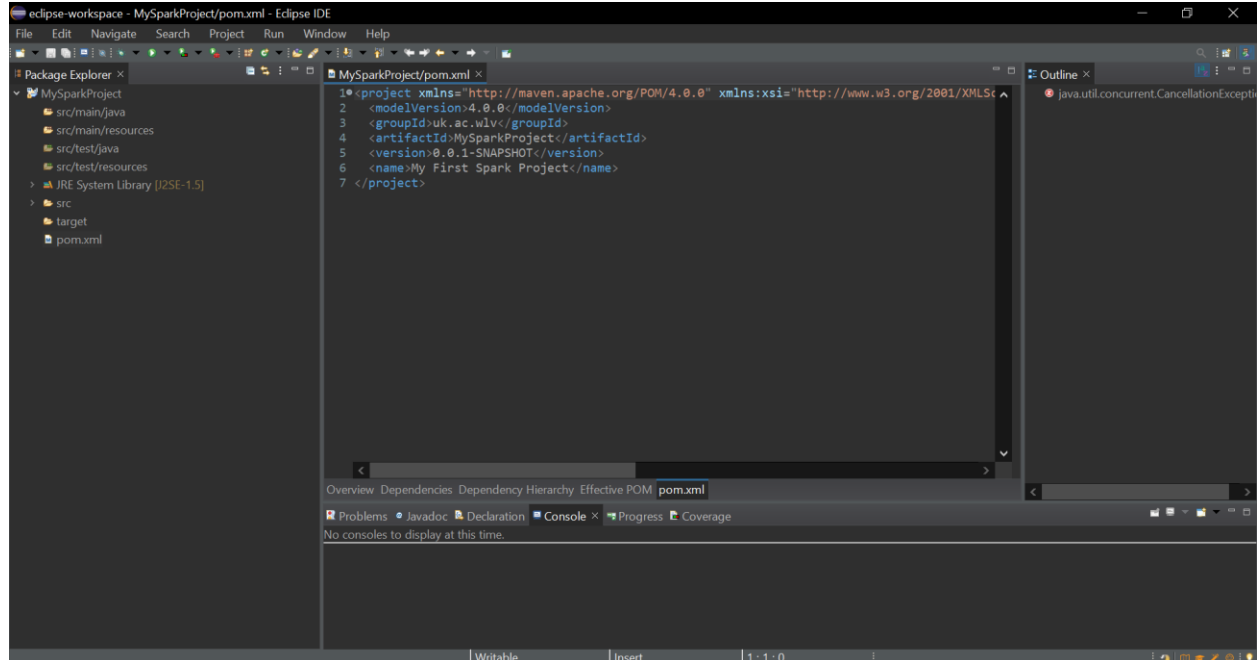


You need to configure. So fill the form as below, name your file as My First Spark Project and click finish.

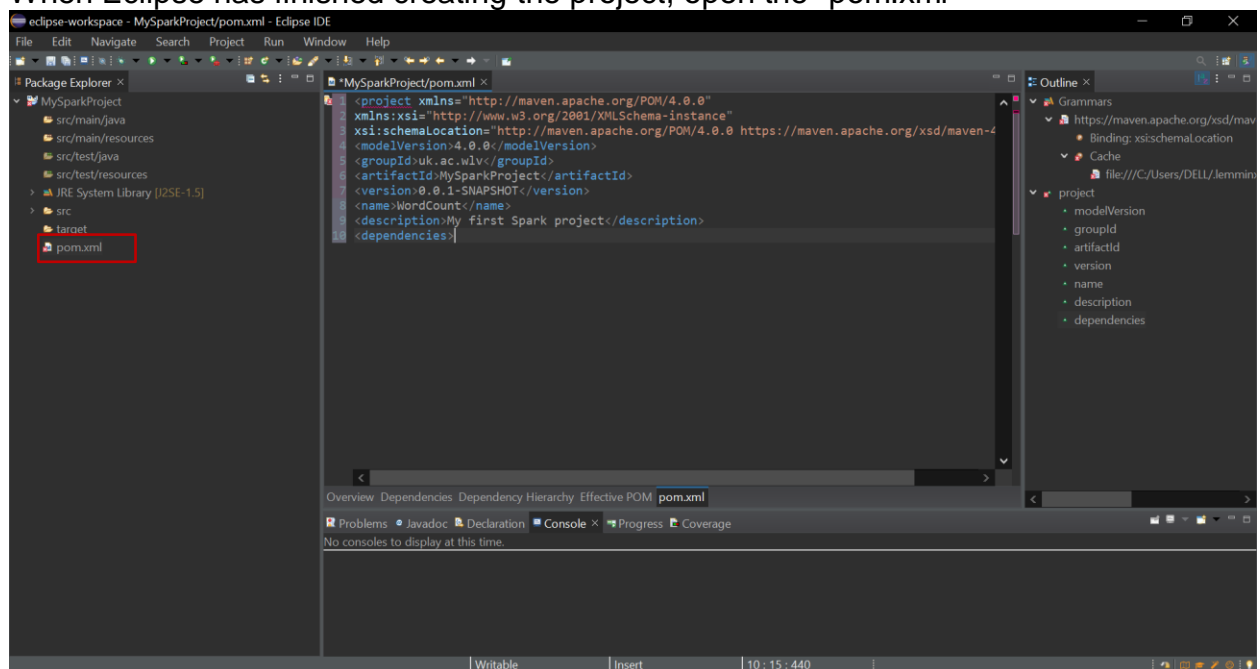
Note: Your group ID must be named as uk.ac.wlv



After clicking on finish, you will see the text editor and package explorer in the left section.



When Eclipse has finished creating the project, open the "pom.xml"





When you open the pom.xml file, you will see a sample java code in the text editor section. Simply replace all of the following code with the following.

```
<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
https://maven.apache.org/xsd/maven-4.0.0.xsd">
  <modelVersion>4.0.0</modelVersion>
  <groupId>uk.ac.wlv</groupId>
  <artifactId>MySparkProject</artifactId>
  <version>0.0.1-SNAPSHOT</version>
  <name>WordCount</name>
  <description>My first Spark project</description>

  <dependencies>
    <dependency>
      <groupId>org.apache.spark</groupId>
      <artifactId>spark-core_2.11</artifactId>
      <version>1.4.0</version>
    </dependency>
  </dependencies>

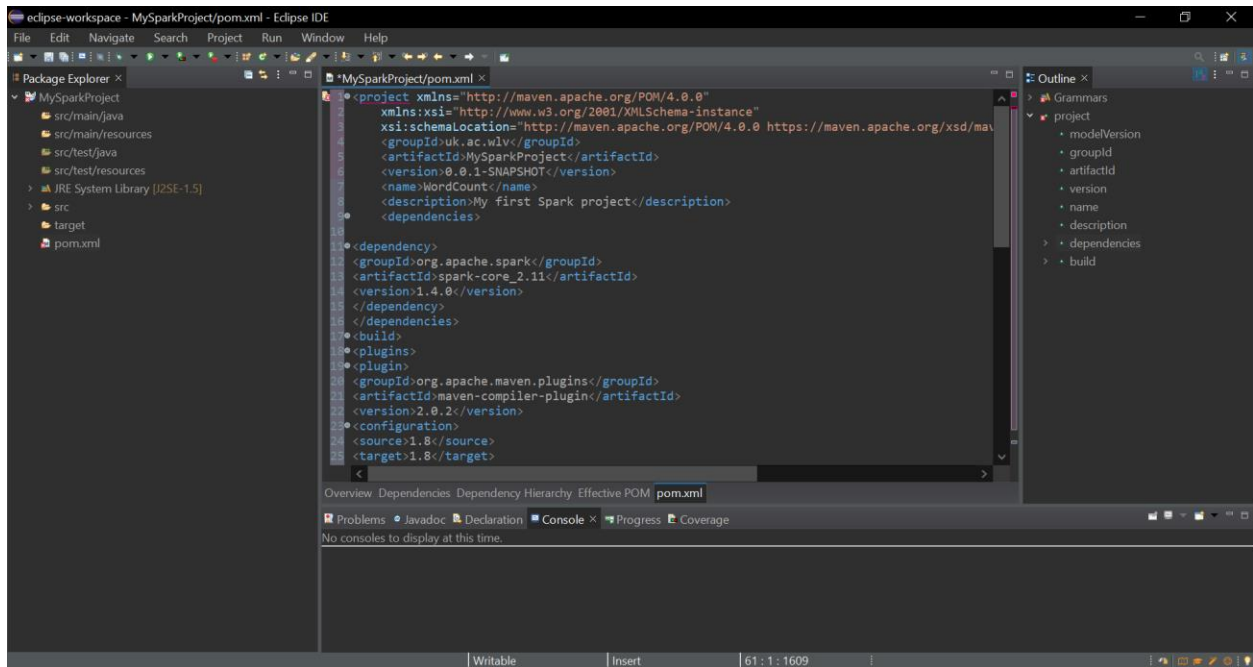
  <build>
    <plugins>
      <plugin>
        <groupId>org.apache.maven.plugins</groupId>
        <artifactId>maven-compiler-plugin</artifactId>
        <version>2.0.2</version>
        <configuration>
          <source>1.8</source>
          <target>1.8</target>
        </configuration>
      </plugin>
      <plugin>
        <groupId>org.apache.maven.plugins</groupId>
        <artifactId>maven-jar-plugin</artifactId>
        <configuration>
          <archive>
            <manifest>
              <addClasspath>true</addClasspath>
              <classpathPrefix>lib/</classpathPrefix>
              <mainClass>uk.ac.wlv.WordCount</mainClass>
            </manifest>
          </archive>
        </configuration>
      </plugin>
    </plugins>
  </build>
</project>
```

```

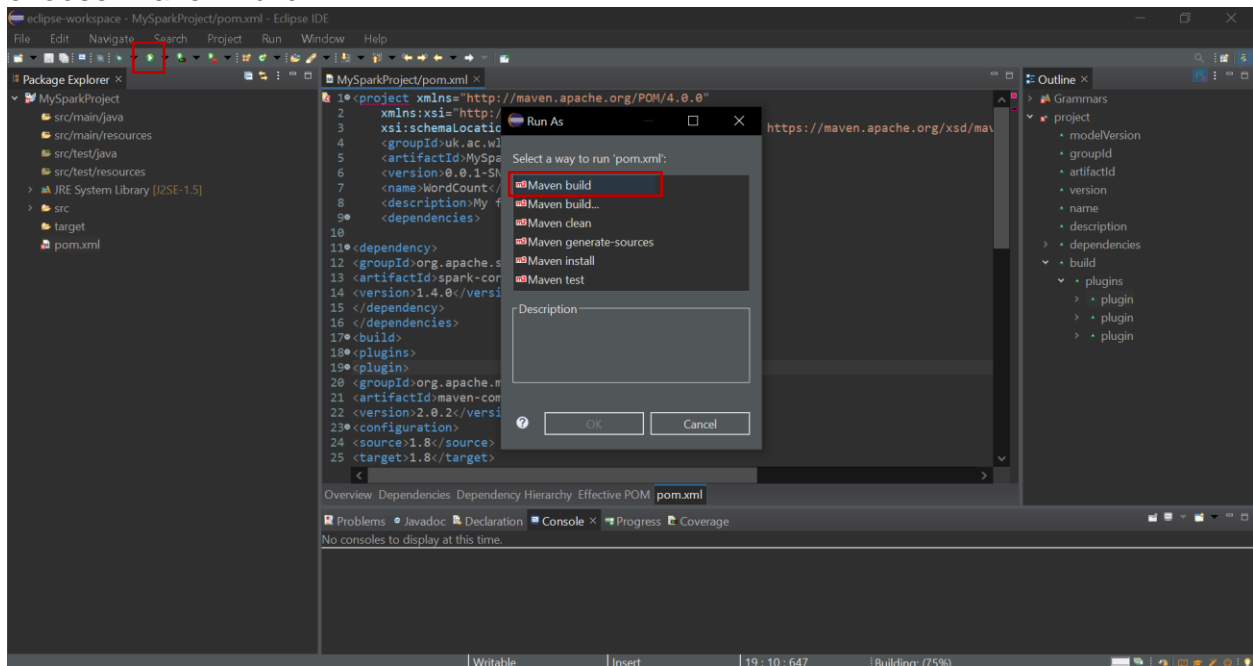
</plugin>
<plugin>
  <groupId>org.apache.maven.plugins</groupId>
  <artifactId>maven-dependency-plugin</artifactId>
  <executions>
    <execution>
      <id>copy</id>
      <phase>install</phase>
      <goals>
        <goal>copy-dependencies</goal>
      </goals>
      <configuration>
        <outputDirectory>${project.build.directory}/lib</outp
utDirectory>
      </configuration>
    </execution>
  </executions>
</plugin>
</plugins>
</build>

</project>

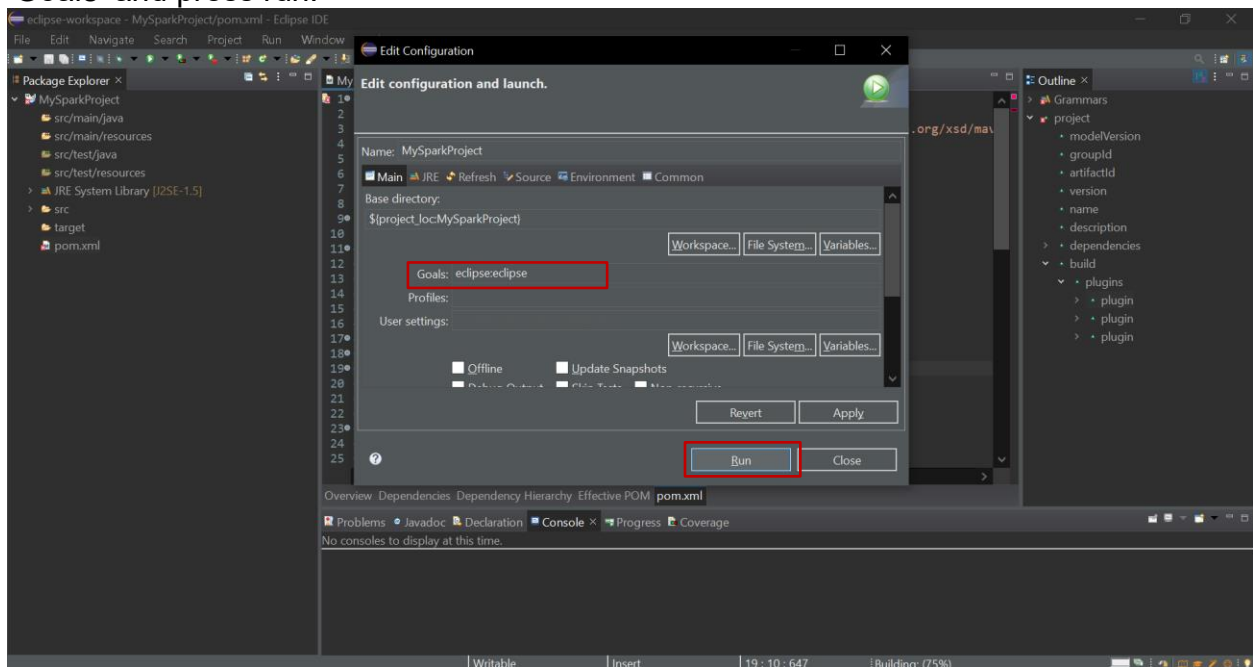
```



Now, you need to save the file. After saving, you must run this file by clicking the Run button in the upper left corner of the window. Hover your mouse over the button and choose 'Maven Build.'



After that in edit the configuration of the build, simply insert 'eclipse:eclipse' inside the 'Goals' and press run.

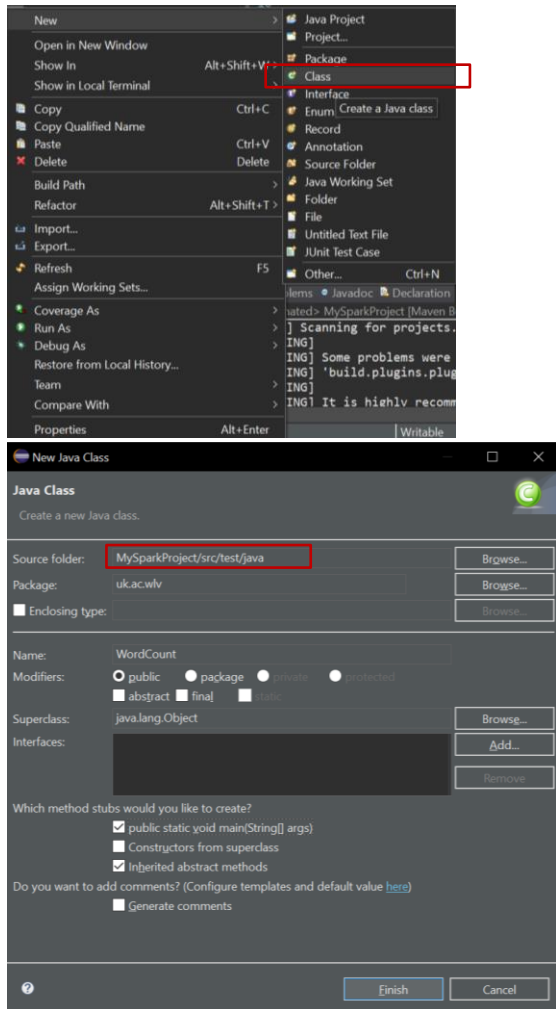


You will see this output in your console.



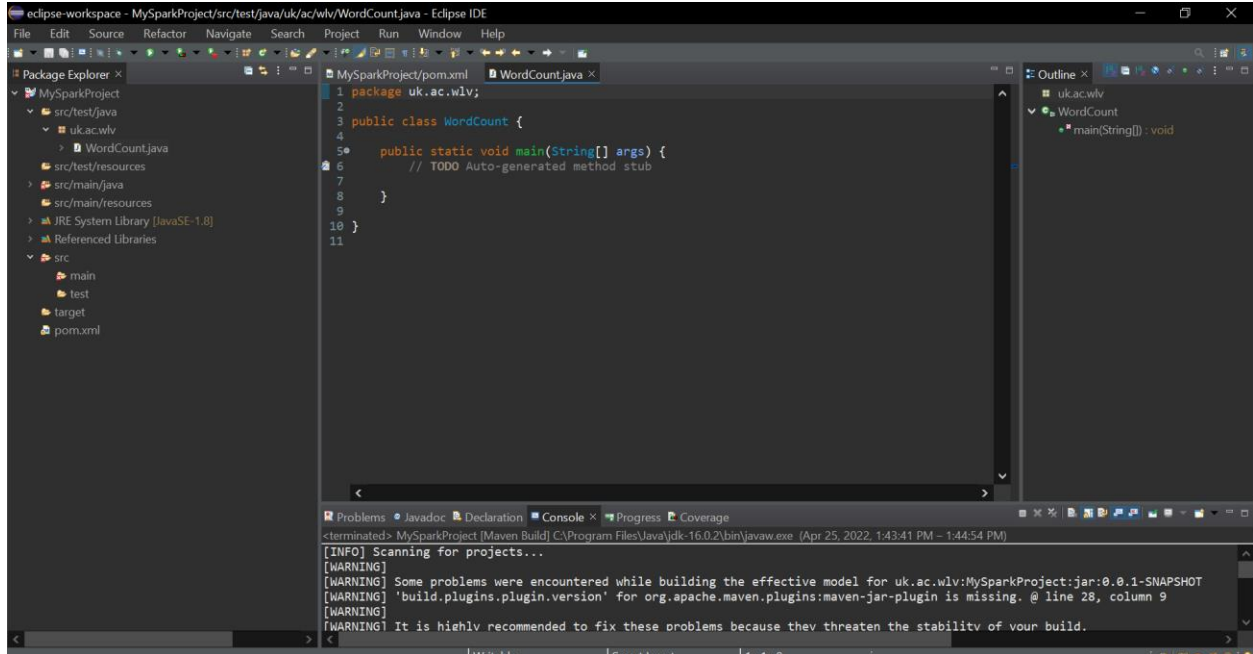
#### Step 4: Creating Word Count program

Now you need to create the WordCount Program. Create a new java class, by right-clicking on src/main/java folder.



Click on finish now.

This generates a new Java class called 'WordCount.java.'



Now change the codes with the ones shown below.

```
package uk.ac.wlv;
import java.util.Arrays;
import java.util.Iterator;
import org.apache.spark.SparkConf;
import org.apache.spark.api.java.JavaPairRDD;
import org.apache.spark.api.java.JavaRDD;
import org.apache.spark.api.java.JavaSparkContext;
import org.apache.spark.api.java.function.FlatMapFunction;
import org.apache.spark.api.java.function.Function2;
import org.apache.spark.api.java.function.PairFunction;
import scala.Tuple2;
public class WordCount {
    /* This function splits the text file into an arraylist of words */
    static class SplitFunction implements FlatMapFunction < String, String > {
        public Iterable < String > call(String s) {
            return Arrays.asList(s.split(" "));
        }
    }
    public static void main(String[] args) {
        SparkConf sparkConf = new SparkConf();
        sparkConf.setAppName("Spark WordCount example using Java");
        /* Tell Spark that we are running on this computer alone */
        sparkConf.setMaster("local");
        JavaSparkContext sparkContext = new JavaSparkContext(sparkConf);
```

```

/*Reading input file*/

JavaRDD < String > textFile = sparkContext.textFile("input.txt");
/*Creating RDD of words from each line of input file*/
JavaRDD < String > words = textFile.flatMap(new SplitFunction());
/*Generate Pair of Word with count */
JavaPairRDD < String, Integer > pairs = words.mapToPair(new PairFunction
< String, String, Integer > () {
    public Tuple2 < String, Integer > call(String s) {
        return new Tuple2 < String, Integer > (s, 1);
    }
});
/* Aggregate Pairs of Same Words with count */
JavaPairRDD < String, Integer > counts = pairs.reduceByKey(
    new Function2 < Integer, Integer, Integer > () {
        public Integer call(Integer a, Integer b) {
            return a + b;
        }
    });
/*Saving the result file */
counts.saveAsTextFile("output");
sparkContext.stop();
sparkContext.close();
}
}

```

The screenshot shows the Eclipse IDE with the following components:

- Package Explorer:** Shows the project structure for 'MySparkProject', including 'src/test/java/uk.ac.wlv/WordCount.java'.
- Editor:** Displays the 'WordCount.java' file. The code includes imports for Spark and Scala, a 'SplitFunction' class, and a 'main' method that runs the word count logic.
- Outline:** Shows the class structure with 'WordCount' and 'main(String[]): void'.
- Console:** Displays the output of the program, including warnings about Maven build issues and the final result of the word count.

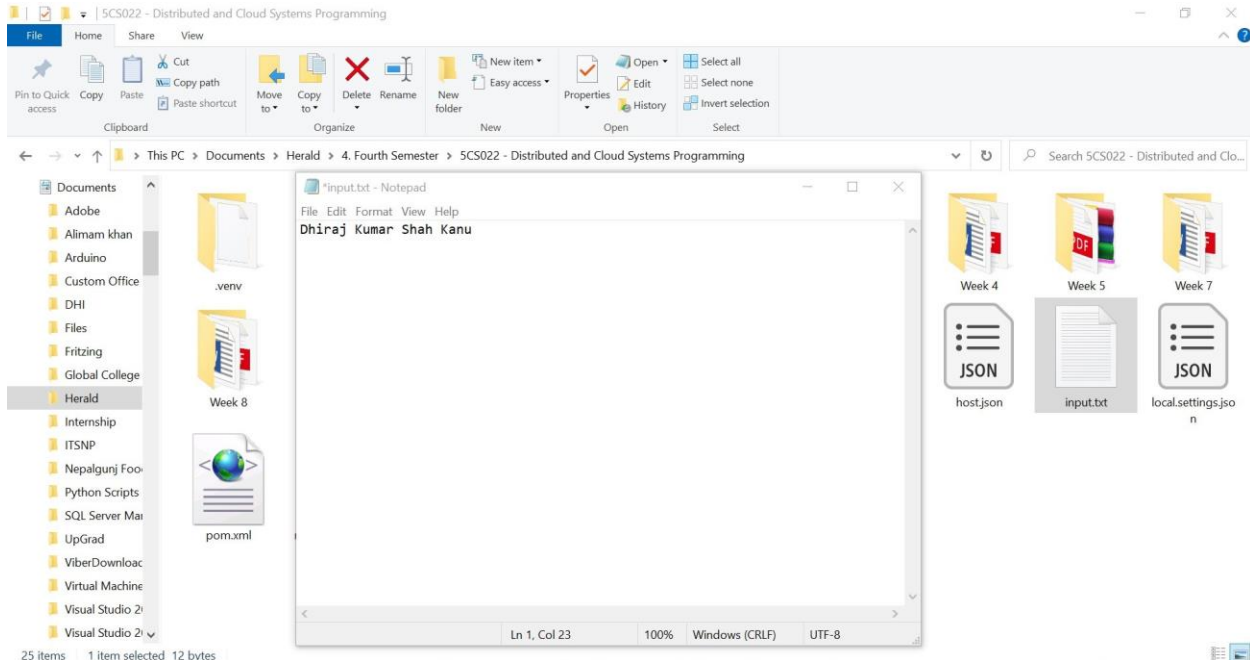
```

terminated: MySparkProject [Maven Build] C:\Program Files\Java\jdk-16.0.2\bin\javaw.exe (Apr 25, 2022, 1:43:41 PM - 1:44:54 PM)
[INFO] Scanning for projects...
[WARNING]
[WARNING] Some problems were encountered while building the effective model for uk.ac.wlv:MySparkProject:jar:0.0.1-SNAPSHOT
[WARNING] 'build.plugins.plugin.version' for org.apache.maven.plugins:maven-jar-plugin is missing. @ line 28, column 9
[WARNING]
[WARNING] It is highly recommended to fix these problems because they threaten the stability of your build.

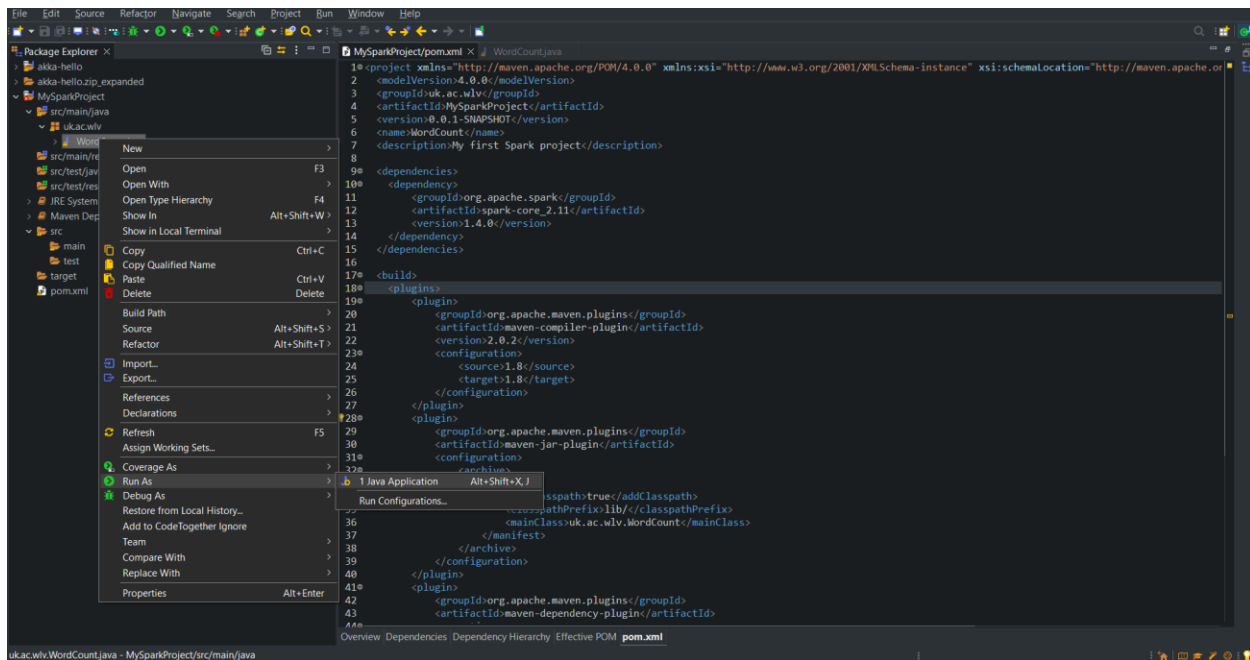
```



Create a new .txt file and save it in the same directory as the pom.xml file,



## Step 5: Running the WordCount Spark Program



Make sure that the file WordCount.java is opened as the current file in the Eclipse editor and click the run button on the Eclipse menu toolbar.



```

<terminated> WordCount [Java Application] C:\Program Files\Java\jdk-11.0.14\bin\javaw.exe (Apr 25, 2022, 2:16:26 PM - 2:16:31 PM)
22/04/25 14:16:30 INFO DAGScheduler: ShuffleMapStage 0 (mapPair at WordCount.java:52) finished in 0.132 s
22/04/25 14:16:30 INFO DAGScheduler: looking for newly runnable stages
22/04/25 14:16:30 INFO DAGScheduler: running: Set()
22/04/25 14:16:30 INFO DAGScheduler: waiting: Set(ResultStage 1)
22/04/25 14:16:30 INFO DAGScheduler: failed: Set()
22/04/25 14:16:30 INFO DAGScheduler: Missing parents for ResultStage 1: List()
22/04/25 14:16:30 INFO DAGScheduler: Submitting ResultStage 1 (MapPartitionsRDD[5] at saveAsTextFile at WordCount.java:45), which is now runnable
22/04/25 14:16:30 INFO MemoryStore: ensureFreeSpace(95664) called with curMem=92869, maxMem=2301162946
22/04/25 14:16:30 INFO MemoryStore: Block broadcast_2 stored as values in memory (estimated size 93.4 KB, free 2.1 GB)
22/04/25 14:16:30 INFO MemoryStore: ensureFreeSpace(31305) called with curMem=188533, maxMem=2301162946
22/04/25 14:16:30 INFO MemoryStore: Block broadcast_2_piece0 stored as bytes in memory (estimated size 30.6 KB, free 2.1 GB)
22/04/25 14:16:30 INFO BlockManagerInfo: Added broadcast_2_piece0 in memory on localhost:60648 (size: 30.6 KB, free: 2.1 GB)
22/04/25 14:16:30 INFO SparkContext: Created broadcast 2 from broadcast at DAGScheduler.scala:874
22/04/25 14:16:30 INFO DAGScheduler: Submitting 1 missing tasks from ResultStage 1 (MapPartitionsRDD[5] at saveAsTextFile at WordCount.java:45)
22/04/25 14:16:30 INFO TaskSchedulerImpl: Adding task set 1.0 with 1 tasks
22/04/25 14:16:30 INFO TaskSetManager: Starting task 0.0 in stage 1.0 (TID 1, localhost, PROCESS_LOCAL, 1165 bytes)
22/04/25 14:16:30 INFO Executor: Running task 0.0 in stage 1.0 (TID 1)
22/04/25 14:16:30 INFO deprecation: mapred.output.value.class is deprecated. Instead, use mapreduce.job.output.value.class
22/04/25 14:16:30 INFO deprecation: mapred.output.dir is deprecated. Instead, use mapreduce.output.fileoutputformat.outputdir
22/04/25 14:16:30 INFO deprecation: mapred.working.dir is deprecated. Instead, use mapreduce.job.working.dir
22/04/25 14:16:30 INFO deprecation: mapred.output.key.class is deprecated. Instead, use mapreduce.job.output.key.class
22/04/25 14:16:30 INFO ShuffleBlockFetcherIterator: Getting 1 non-empty blocks out of 1 blocks
22/04/25 14:16:30 INFO ShuffleBlockFetcherIterator: Started 0 remote fetches in 7 ms
22/04/25 14:16:30 INFO BlockManagerInfo: Removed broadcast_1_piece0 on localhost:60648 in memory (size: 2.5 KB, free: 2.1 GB)
22/04/25 14:16:30 INFO FileOutputCommitter: Saved output of task 'attempt_202204251416_0001_m_000000_1' to file:/D:/College/Semester 3/Cloud Computing and Distributed System/week 4/MySparkProject/output/_tempo
22/04/25 14:16:30 INFO SparkHadoopMapRedUtil: attempt_202204251416_0001_m_000000_1: Committed
22/04/25 14:16:30 INFO Executor: Finished task 0.0 in stage 1.0 (TID 1). 886 bytes result sent to driver
22/04/25 14:16:30 INFO DAGScheduler: ResultStage 1 (saveAsTextFile at WordCount.java:45) finished in 0.276 s

```

If it runs successfully, you should see the Spark logging output in Eclipse

**Step 6: View the output results**

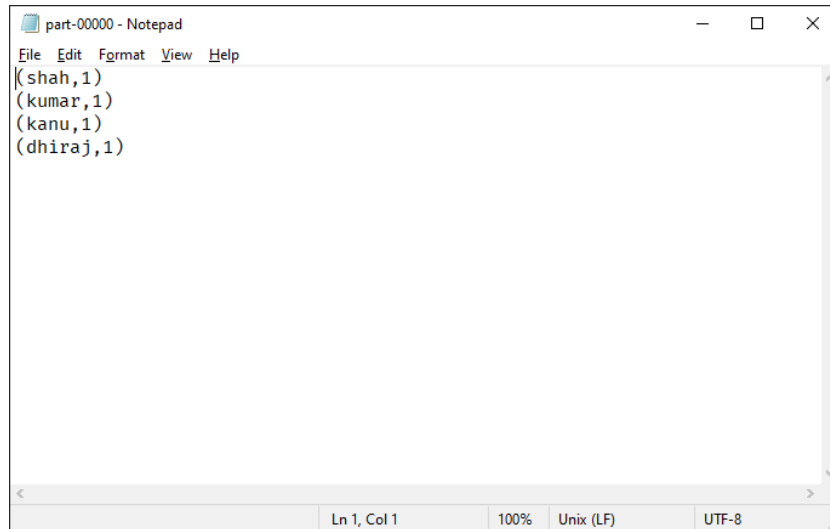
You will see this output if your code runs successfully. Now click on output.

Name	Date modified	Type	Size
.settings	5/6/2022 9:35 PM	File folder	
output	5/6/2022 9:52 PM	File folder	
src	5/6/2022 9:33 PM	File folder	
target	5/6/2022 9:34 PM	File folder	
.classpath	5/6/2022 9:41 PM	CLASSPATH File	14 KB
.project	5/6/2022 9:41 PM	PROJECT File	1 KB
input	5/6/2022 9:52 PM	Text Document	1 KB
pom	5/6/2022 9:37 PM	XML Source File	3 KB

You will see the following files.

Name	Date modified	Type	Size
._SUCCESS.crc	5/6/2022 9:52 PM	CRC File	1 KB
.part-00000.crc	5/6/2022 9:52 PM	CRC File	1 KB
._SUCCESS	5/6/2022 9:52 PM	File	0 KB
part-00000	5/6/2022 9:52 PM	File	1 KB

The results that we are looking for will be in the file "part-0000". Open that file in Notepad, and you should see a list of words and their counts.

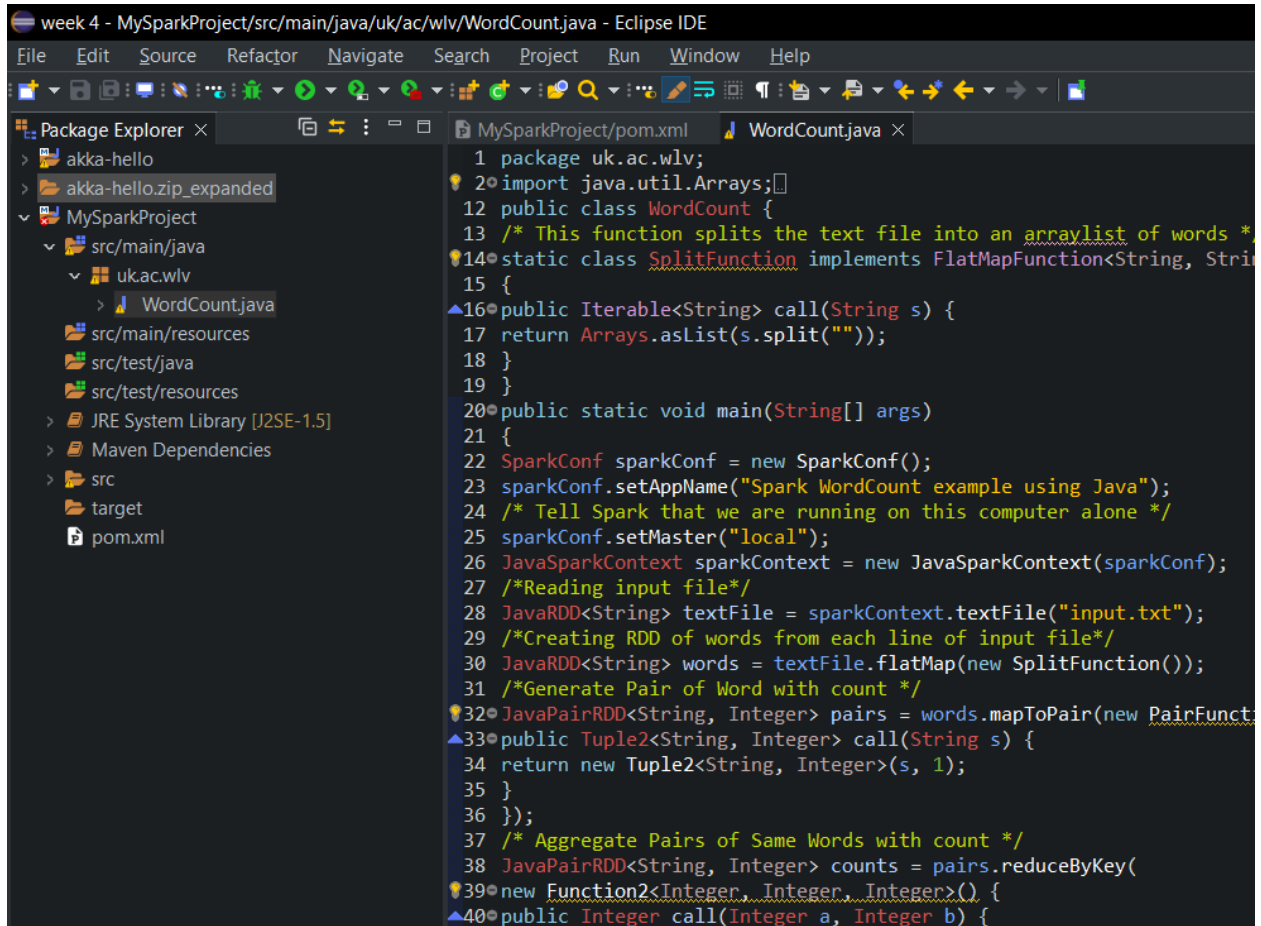


```
part-00000 - Notepad
File Edit Format View Help
(shah,1)
(kumar,1)
(kanu,1)
(dhiraj,1)
Ln 1, Col 1 100% Unix (LF) UTF-8
```

## Task

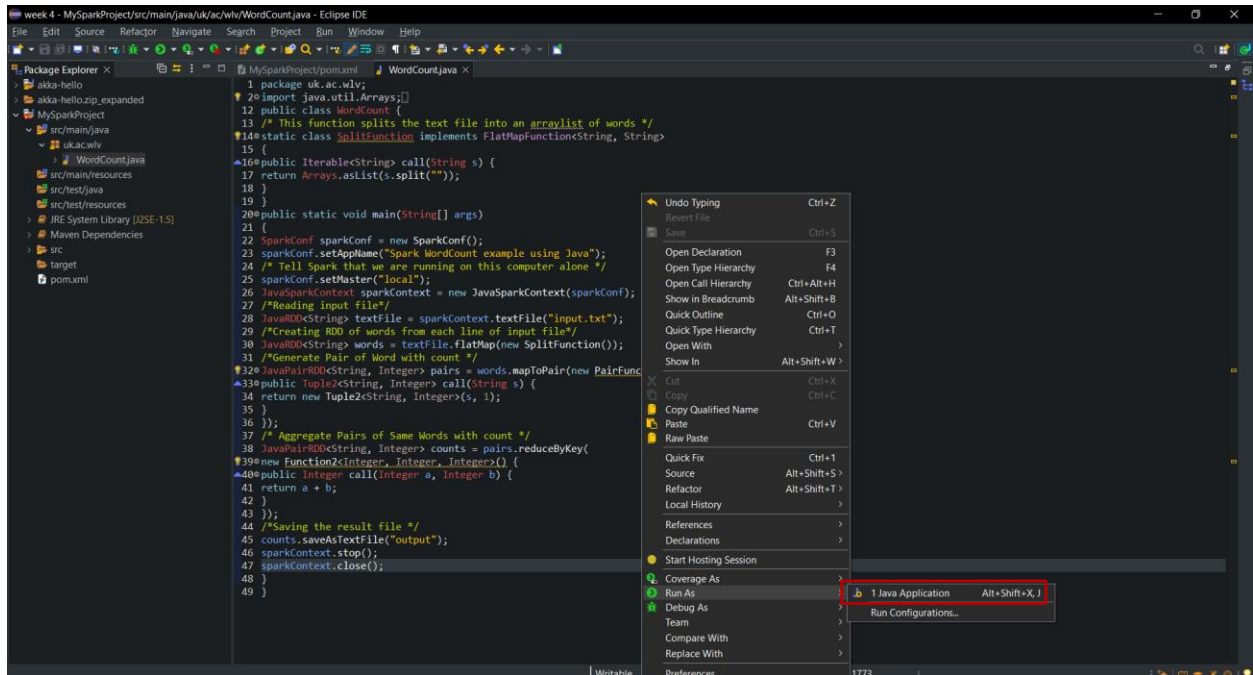
Create a Spark program to count letters instead of words.

Here, in line number 17, remove the space between inverted commas.

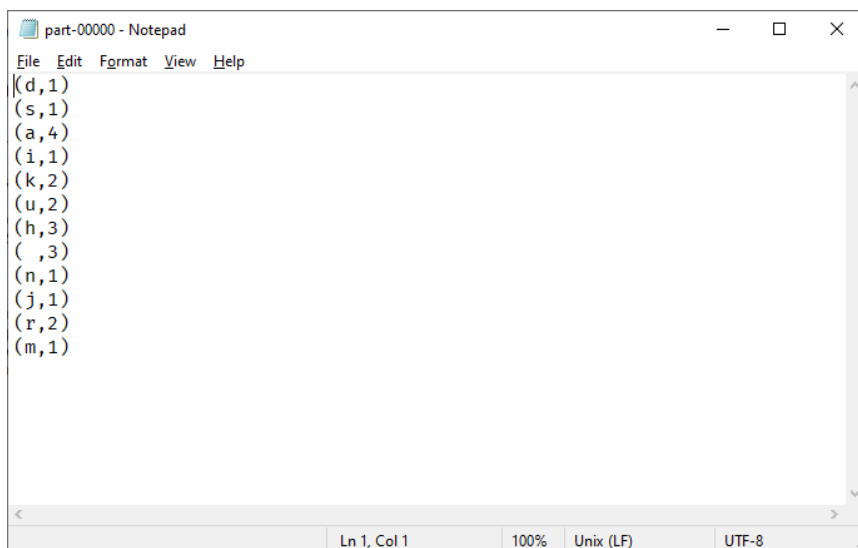


```
1 package uk.ac.wlv;
2 import java.util.Arrays;
3 public class WordCount {
4     /* This function splits the text file into an arraylist of words */
5     static class SplitFunction implements FlatMapFunction<String, String> {
6     {
7         public Iterable<String> call(String s) {
8             return Arrays.asList(s.split(' '));
9         }
10    }
11    public static void main(String[] args)
12    {
13        SparkConf sparkConf = new SparkConf();
14        sparkConf.setAppName("Spark WordCount example using Java");
15        /* Tell Spark that we are running on this computer alone */
16        sparkConf.setMaster("local");
17        JavaSparkContext sparkContext = new JavaSparkContext(sparkConf);
18        /*Reading input file*/
19        JavaRDD<String> textFile = sparkContext.textFile("input.txt");
20        /*Creating RDD of words from each line of input file*/
21        JavaRDD<String> words = textFile.flatMap(new SplitFunction());
22        /*Generate Pair of Word with count */
23        JavaPairRDD<String, Integer> pairs = words.mapToPair(new PairFunction2<String, Integer, Integer>() {
24        public Tuple2<String, Integer> call(String s) {
25            return new Tuple2<String, Integer>(s, 1);
26        }
27        });
28        /* Aggregate Pairs of Same Words with count */
29        JavaPairRDD<String, Integer> counts = pairs.reduceByKey(
30        new Function2<Integer, Integer, Integer>() {
31        public Integer call(Integer a, Integer b) {
```

Right click on the file and run as a Java program



You will see the final output like this, which will be in the file "part-0000". Open that file in Notepad, and you should see a list of words and their counts.



**Summary**

So, first we downloaded the Hadoop file from the google classroom. Then we edited the environment variables, started a new spark project, modify the codes, and develop a functional program in this manner. The program we just created counts the total amount of words entered into the 'input.txt' file. The output is then shown as seen above.