# Spark Programming

# Access Cluster at AHPCC

# Access Cluster

- Open a web browser, on the address bar type in
  - hpc-portal2.hpc.uark.edu
  - On the popup window, enter your uark ID and password
- When you are on the dashboard page
  - Click on Clusters --> _Karpinski Shell Access to open the shell
  - Click on Files --> /karpinski/uarkID to open the FTP for uploading/downloading files
- Off-campus access to cluster
  - Install and run GlobalProtect vpn first
  - https://its.uark.edu/network-access/vpn/index.php

# Allocate Cluster Node

- In the shell, type in <u>spark-node4me.sh</u> to allocate a cluster node for 2 hours to program.

- Once you finish using the cluster node, you can type in exit to release the node.

- The shell supports "Ctrl-V" for pasting. In the shell, when you select some texts, they are automatically copied.

# Load Module

- Load modules
  - module load spark/2.3.0

- (Optional) Check available modules
  - module avail [java/python/spark/hadoop]

# Submitting Python Applications

# Submitting Python Applications

- Write your Python code and save as **.py

- Upload the Python code and input files to the cluster

- In shell type in spark-submit **.py [path_to_inputfile]

`spark-submit pywordcount.py pg100.txt output`
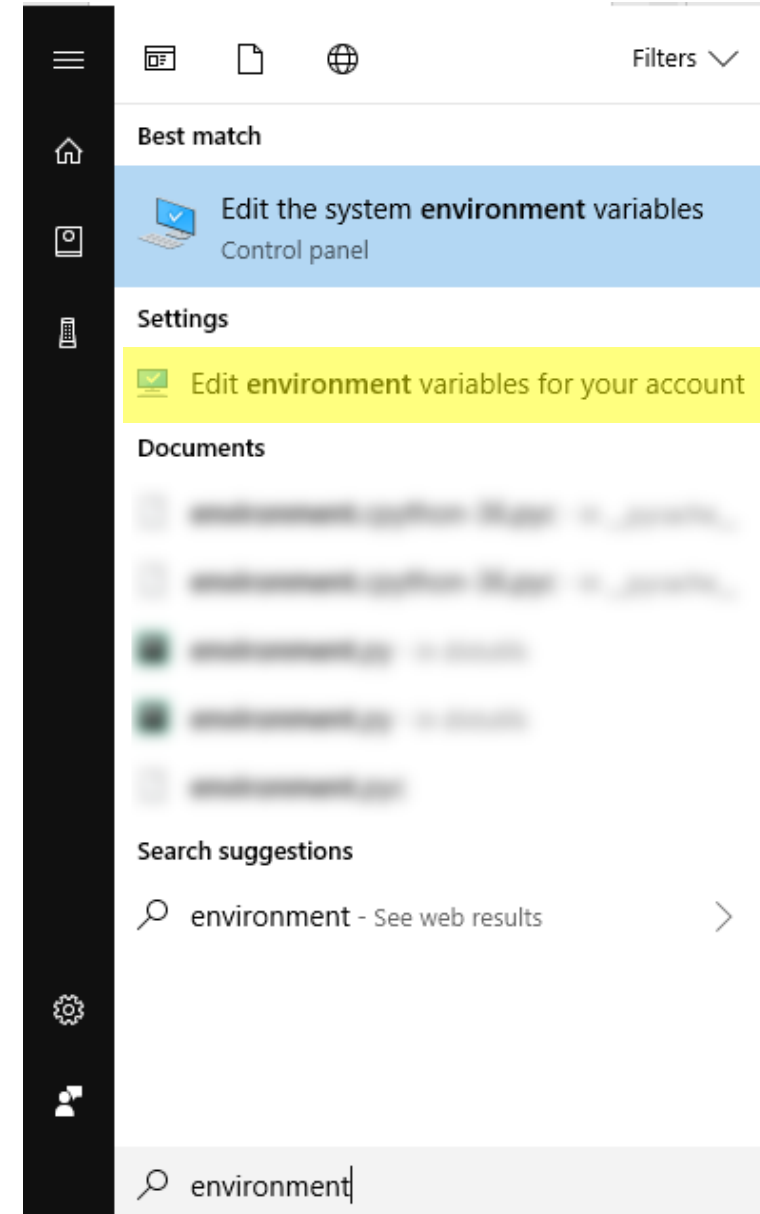
# Submitting Java Applications

# Overview

- Install/download Java, Spark, Eclipse, and Maven on your PC

1. Create Maven project in Eclipse
2. Write your java code in the project
3. Build Maven package
4. Upload the Maven package and input files to the cluster
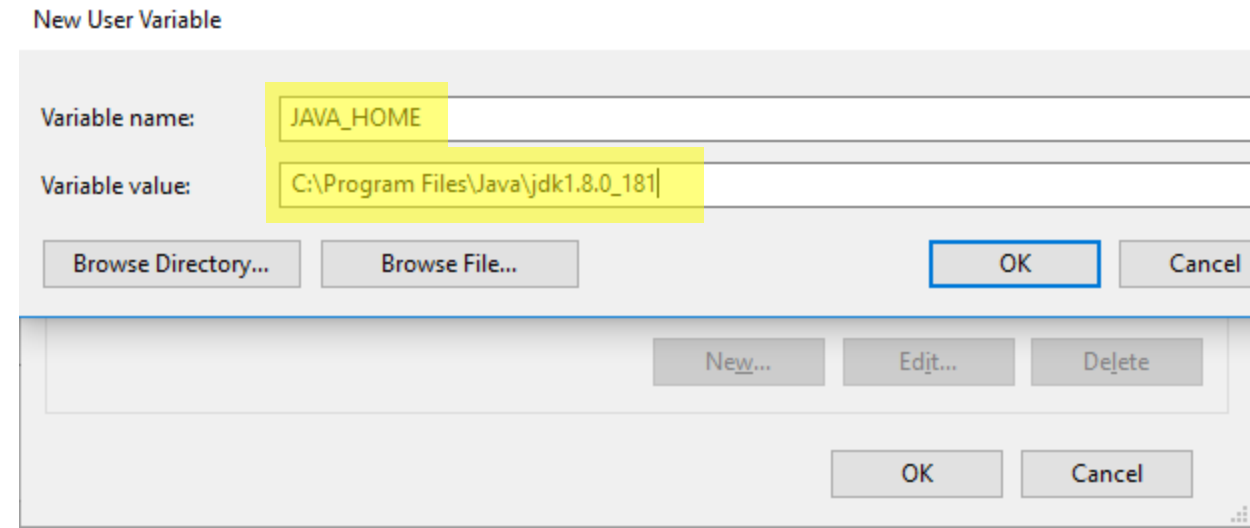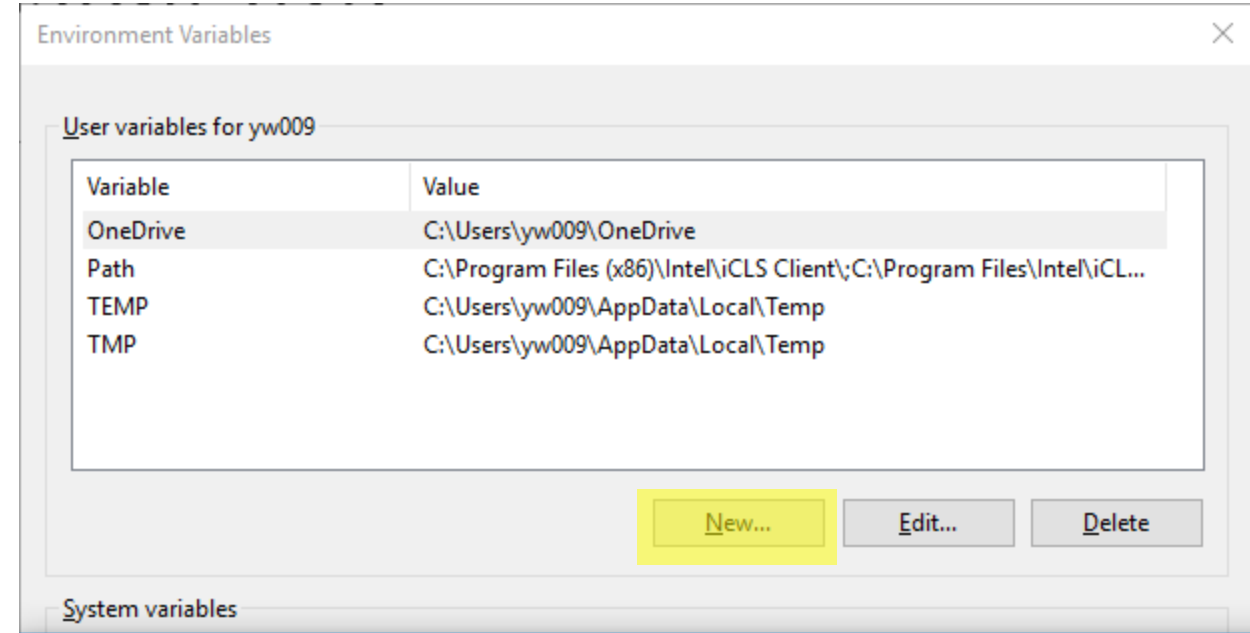5. Submit the application to spark in the shell

# Steps of Set Environment Variables

- Open Environment Variables
  - Open Start and type "environment",
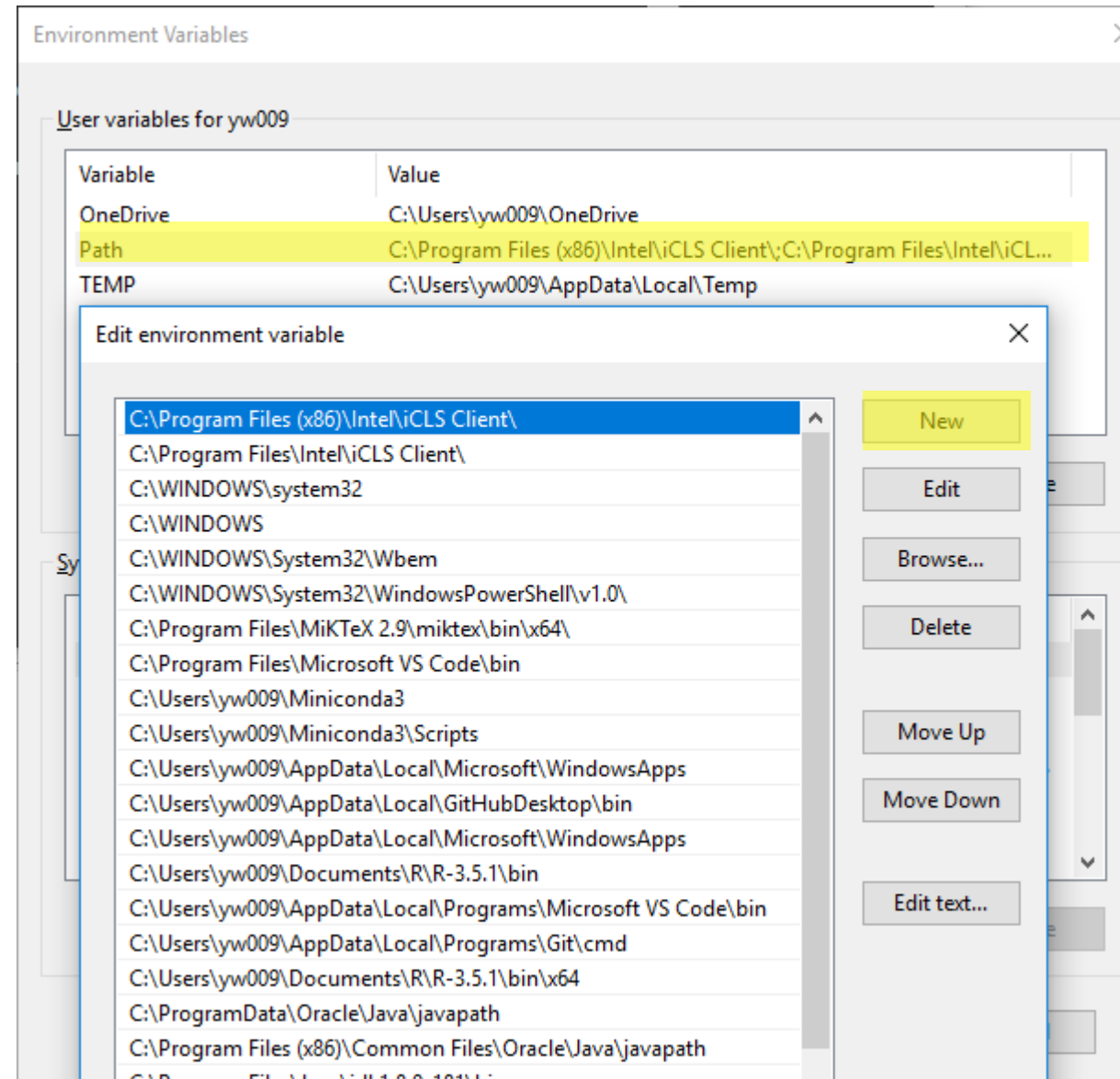  - Select "Edit environment variables for your account". (see highlight of the figure)

# Steps of Set Environment Variables

- Set Java_HOME
  1. Click on the "New…" button
  2. In the Pop-up window, Type "JAVA_HOME" and your jdk path

# Steps of Set Environment Variables

- Add to path

- Steps
    1. Double-click on the Path
    2. In the Pop-up window, click on the "New" button, type "%JAVA_HOME%\bin"

# Install Java 8

- Download Java 8 from the link: https://www.java.com/download/ie_manual.jsp and install it.
- Set environmental variables:
  - User variable:
    - Variable: JAVA_HOME;
    - Value: C:\Program Files\Java\jdk1.8.0_341
  - System variable:
    - Variable: PATH
    - Value: %JAVA_HOME%\bin
- Check on cmd, see below:

```
C:\>java -version
java version "1.8.0_341"
Java(TM) SE Runtime Environment (build 1.8.0_341-b10)
Java HotSpot(TM) Client VM (build 25.341-b10, mixed mode, sharing)
```

# Download Spark 2.3.0

- Download it from the following link: https://archive.apache.org/dist/spark/ and extract it into a folder such as C:\spark

- Set environmental variables:
  - User variable:
    - Variable: SPARK_HOME;
    - Value: C:\spark\spark-2.3.0-bin-hadoop2.7
  - System variable:
    - Variable: PATH
    - Value: %SPARK_HOME%\bin

# Install Eclipse

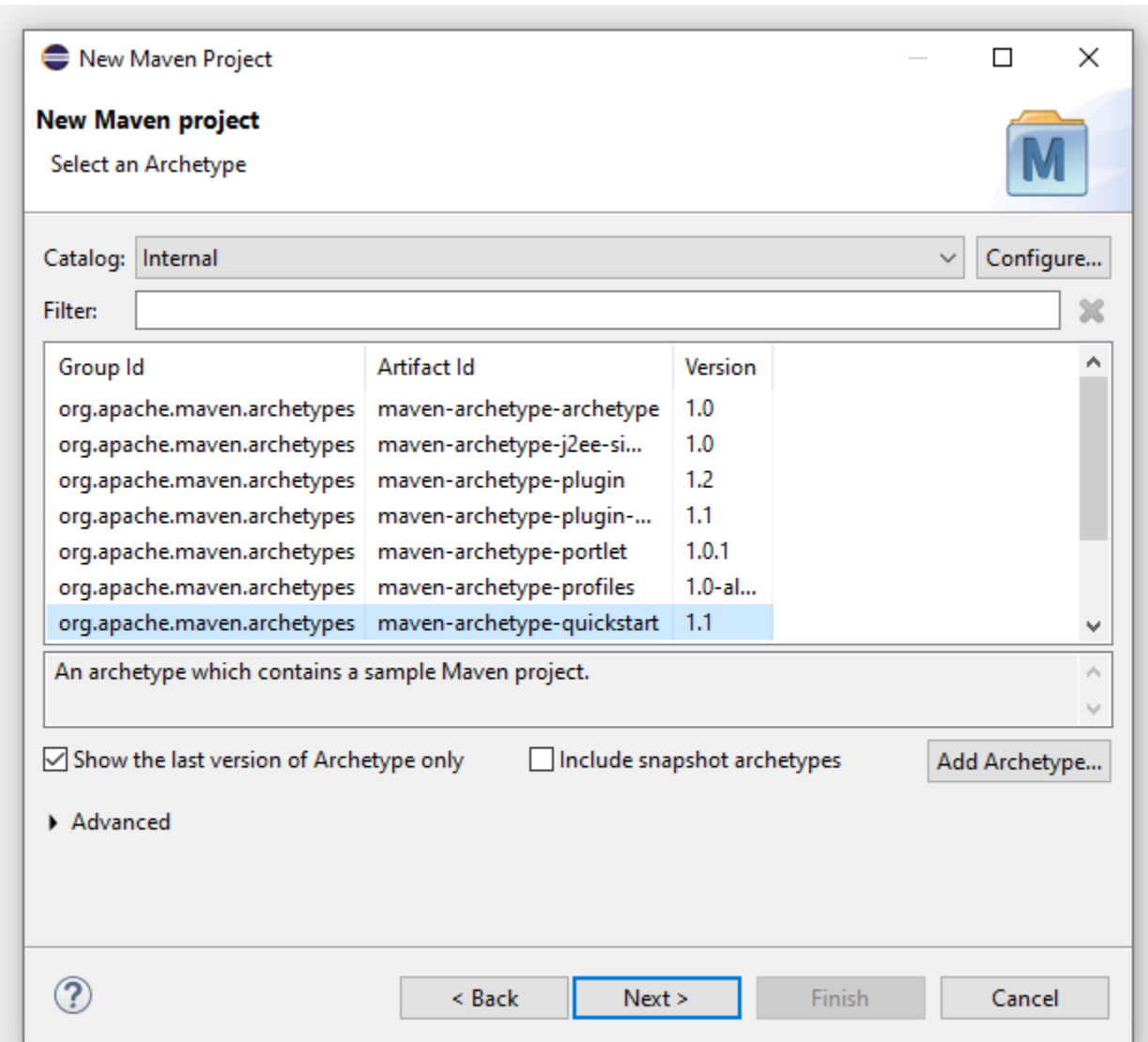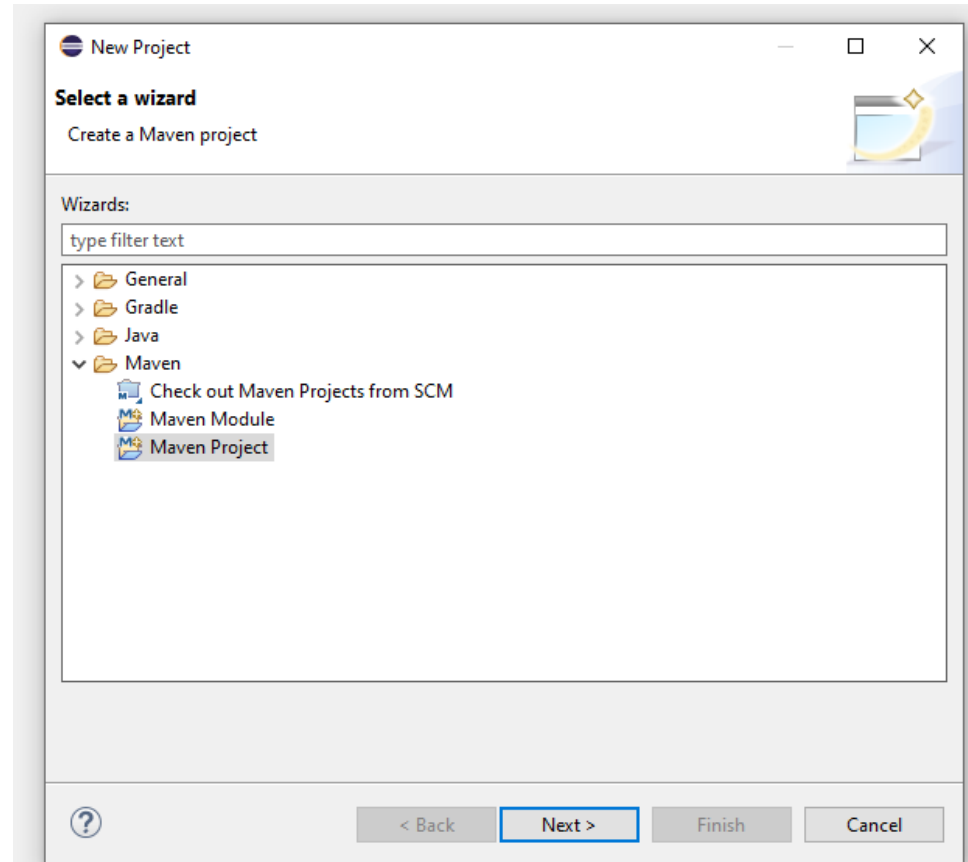- Download it from the link: https://eclipse.org/downloads/ and install it for Java.

# Download Maven 3.3

- Download Apache-Maven-3.3.9 from the link: [http://apache.mivzakim.net/maven/maven-3/3.3.9/binaries/apache-maven-3.3.9-bin.zip](http://apache.mivzakim.net/maven/maven-3/3.3.9/binaries/apache-maven-3.3.9-bin.zip) and extract it into a folder such as C:\apache-maven-3.3.9
- Set environmental variables:
  - User variable:
    - Variable: MAVEN_HOME;
    - Value: C:\apache-maven-3.3.9
  - System variable:
    - Variable: PATH
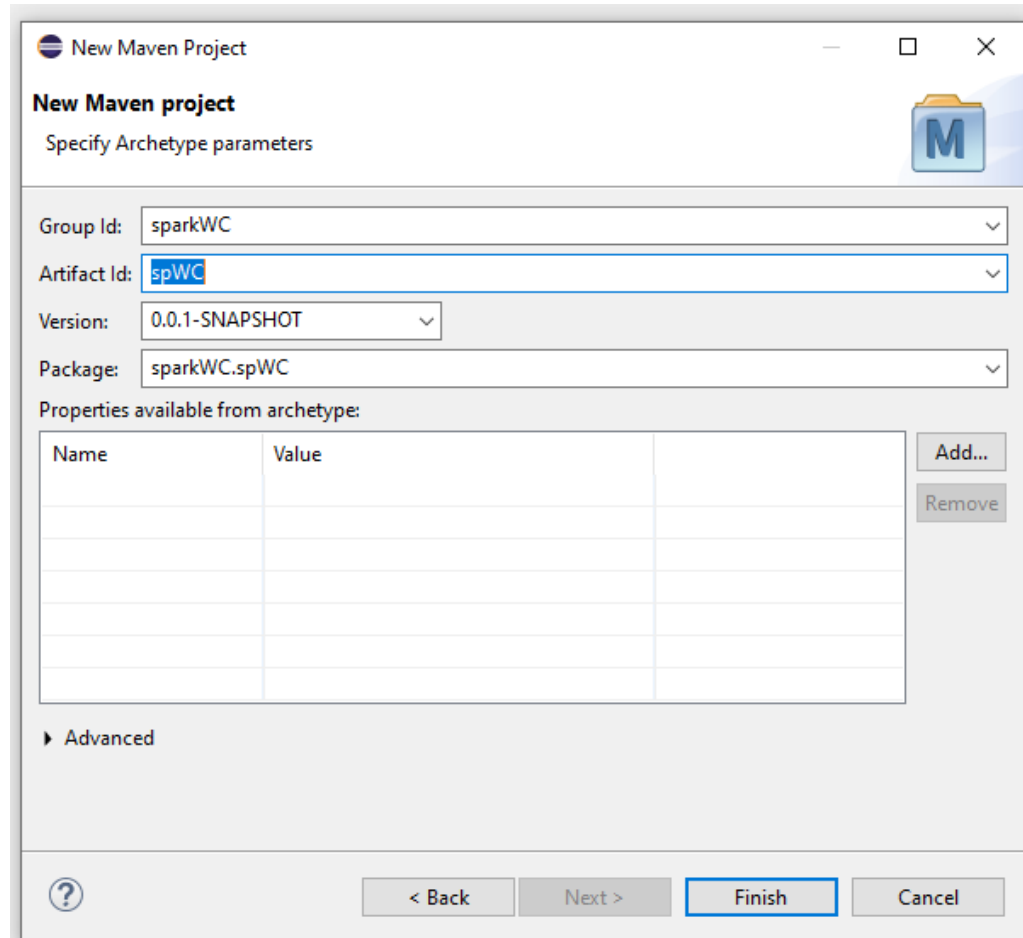    - Value: %MAVEN_HOME%\bin

# Create First WordCount Project

- Open Eclipse and do File->New->project->Maven Project.

# Create First WordCount Project

- Enter Group id, Artifact id, and click finish.

# Create First WordCount Project

- Edit pom.xml. Change the <dependencies> part to the following.

```xml
<dependencies>
    <dependency>
        <groupId>org.apache.spark</groupId>
        <artifactId>spark-core_2.11</artifactId>
        <version>2.2.0</version>
        <scope>provided</scope>
    </dependency>
    <dependency>
        <groupId>org.apache.spark</groupId>
        <artifactId>spark-sql_2.11</artifactId>
        <version>2.3.0</version>
    </dependency>
    <dependency>
        <groupId>junit</groupId>
        <artifactId>junit</artifactId>
        <version>3.8.1</version>
        <scope>test</scope>
    </dependency>
</dependencies>
```
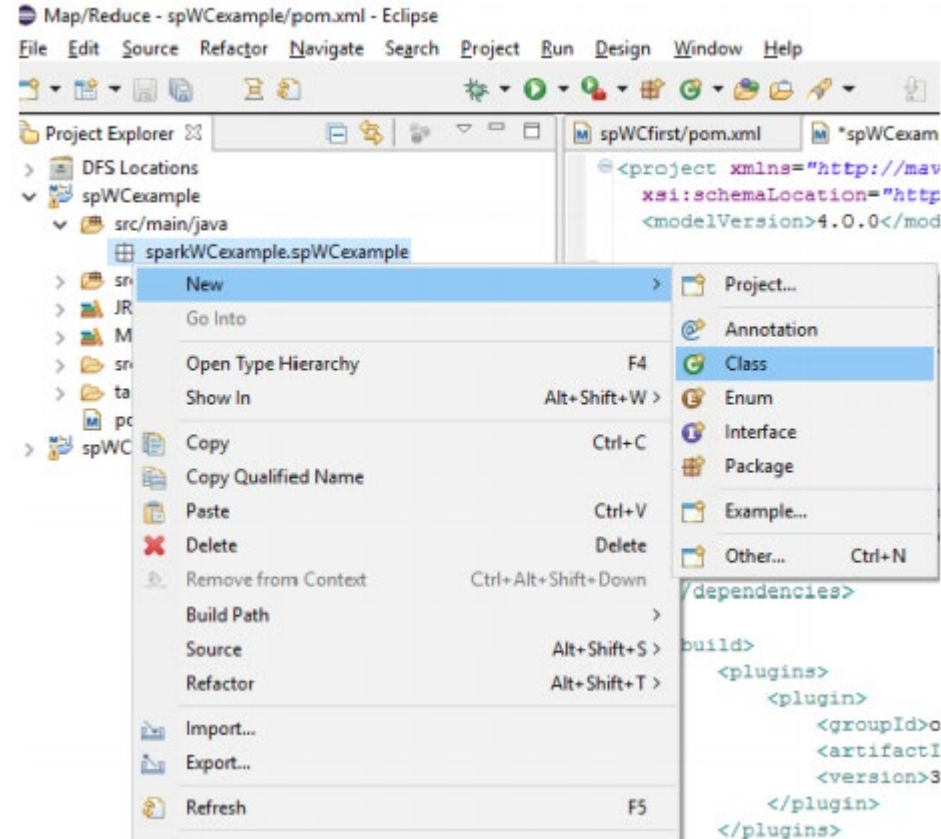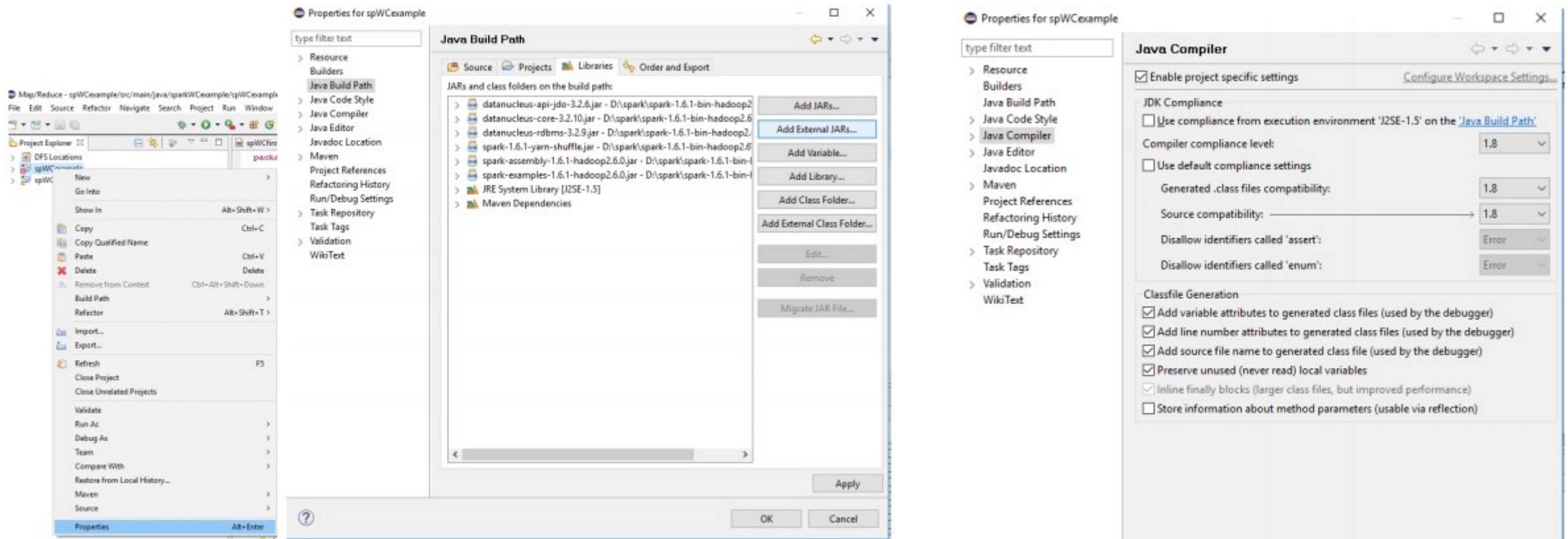
# Create First WordCount Project

- Write your java code in class JavaWordCount

# Create First WordCount Project

- Add external jar from the location C:\spark\spark-2.3.0-bin-hadoop2.7\jar and set Java 8 for compilation; see below.

# Create First WordCount Project

- Build the project: Open cmd, go to the following location (where we stored the project, e.g., C:\Users\eclipse-workspace\spWCexample):
  - Type in <u>mvn package</u>
  - Will build a Maven package, e.g., **spWCexample-0.0.1-SNAPSHOT.jar**
- Execute the project on the cluster:
  - Upload **spWCexample-0.0.1-SNAPSHOT.jar** and input files to the cluster
  - Open the shell, type in

  <u>spark-submit --class sparkWCexample.spWCexample.JavaWordCount --master local[2] ./spWCexample-0.0.1-SNAPSHOT.jar ./input.txt ./output</u>