# Batch-19(IT-A)……
# Installing Hadoop for Windows

## Choosing a Hadoop distribution to install

In this chapter I will do an installation of Apache Hadoop for Windows. This is a version of Hadoop that is totally free, and is the basis of all Hadoop distributions. To test whether Hadoop can be installed on pretty much any Windows PC, I will do a local installation on Windows 8.
While the screenshots I show are from Windows 8, the process on Windows 10 and Windows Server is similar, and you shouldn't have any trouble making the necessary adjustments. It's important that you don't feel you need the latest shiny, new PC to run Hadoop, though later when we look at multi-node installations, we will use multiple Windows Server 2016 machines. The PC in this exercise meets the requirements shown in Table 1, with 8 GB of RAM, a Quad Core 2.4 GHz AMD processor, and solid-state drives.

## Apache Hadoop installation prerequisites

- **JAVA 8 or later**: You can download the 64-bit Windows .jdk file **jdk-8u381-windows-x64.exe** from https://www.oracle.com/java/technologies/downloads/#java8-windows. It is important to state why a prerequisite is required, so the nature of the dependency on the prerequisite is understood. Hadoop is a Java-based application that creates various dependencies on Java. For example, in a single-node Hadoop installation, there is a single Java process running all Hadoop functions. Without Java, all those functions would be unavailable. It is essential to have the right version and architecture of Java, and a 64-bit JDK higher than 1.6 should always be chosen to install Hadoop for Windows.

- **Hadoop 2.0 or later**: You can download the Hadoop binary file **hadoop-3.2.4.tar.gz** from https://hadoop.apache.org/releases.html.

- **Microsoft Windows**: Windows 7, 8, 10, and Windows Server 2008 and above.

- **Additional prerequisites**: You'll also need a text editor, such as Notepad or Notepad ++, for writing short amounts of code, and Winutils 3.1, which you can download from https://github.com/killerangaswamy56/hadoop_config.git

## Java installation for Hadoop for Windows

Run the downloaded Java installer, following the onscreen instructions to complete the installation. Ensure that you right-click on the Java installation file and choose **Run as administrator** from the menu. You will see a User Account Control message asking you to allow the application to make changes to your computer, to which you answer **Yes**. Follow the onscreen prompts to install Java, including the following screen, where you can accept the default installation path.
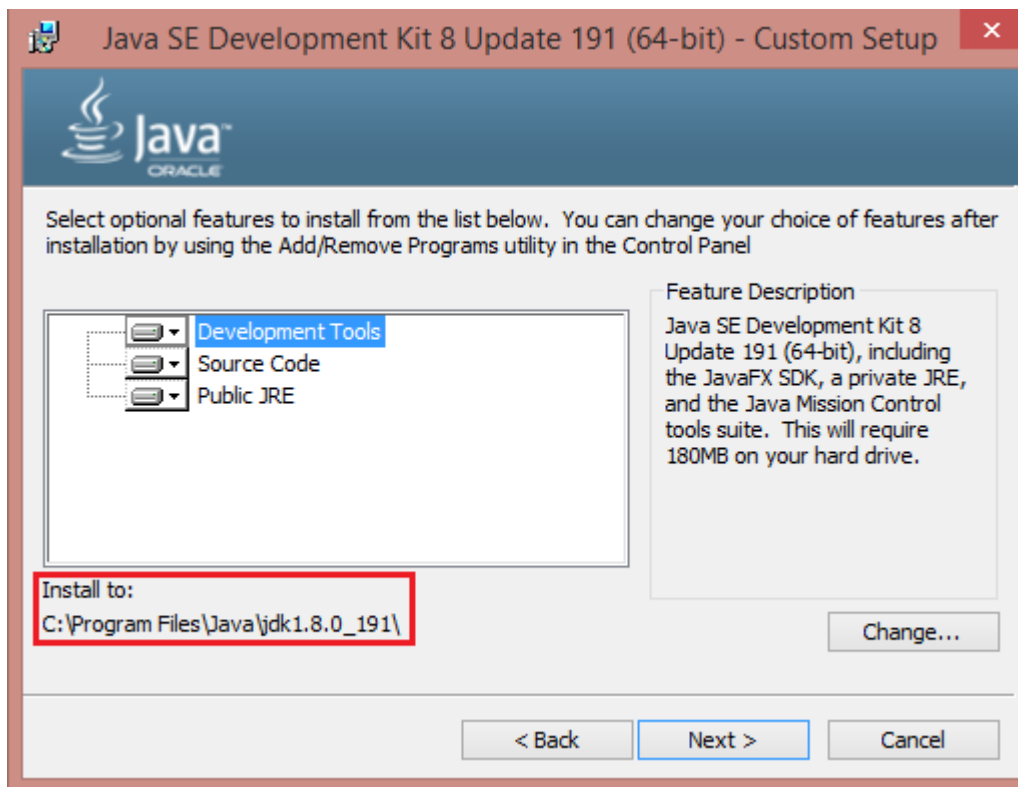
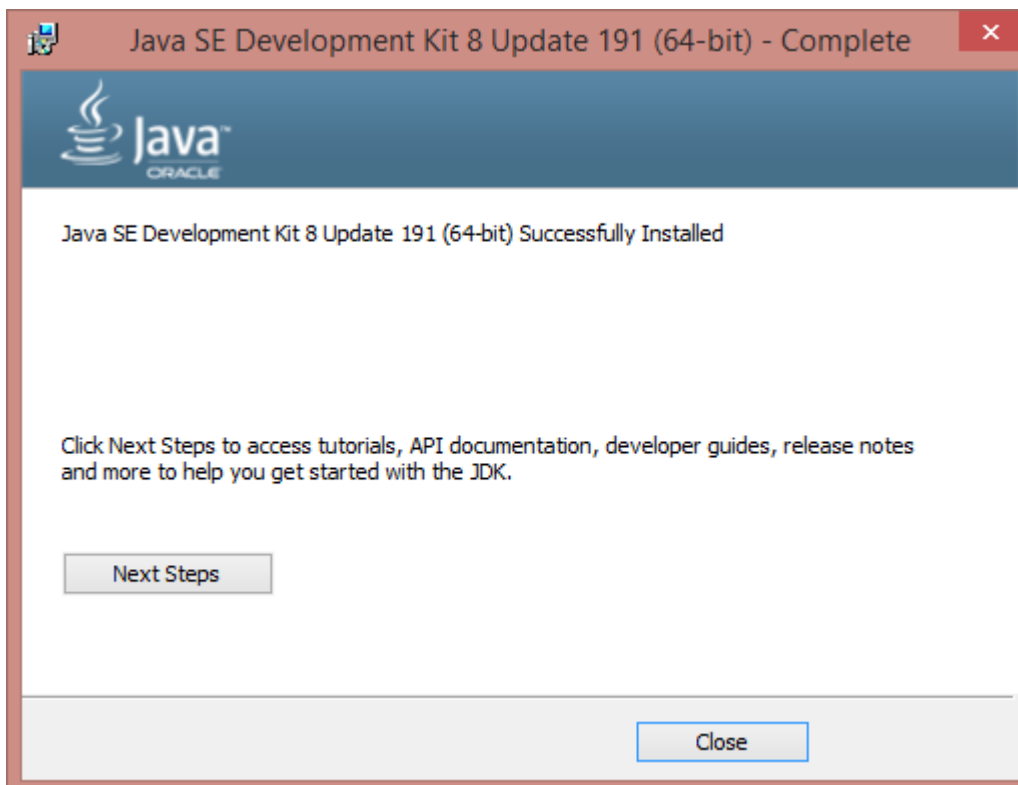*Figure 9: Default Java installation path*



*Figure 10: Successful Java installation*

Ensure that you see the screen informing you that Java has been successfully installed.

Go to **Control Panel** > **System and Security** > **System** and click **Advanced System Settings**, and then click the **System Environment Variables** button. Whether creating a new environment variable for **JAVA_HOME** or editing an existing one, you must alter the Program Files text to text that Hadoop can interpret. On Windows 8, to create a Hadoop-compatible **JAVA_HOME** file instead of entering Program Files, insert **Progra~1** when entering the Java location in the **Variable value** field. On Windows 10 and Windows Server, avoid folder names with blank spaces.



*Figure 11: Hadoop compatible Java Home*

Please ensure that you add the **JAVA_HOME** to the **Path** variable in System Variables. In this instance, it is done by adding **%JAVA_HOME%\bin** between semi colons in the Path **Variable value** field. Use the **java -version** command from a command prompt to ensure that Java is installed and running correctly.



*Figure 12: Adding Java Home to the Path Variable*

## Apache Hadoop installation

1. Create a folder called **C:\hadoop** on your hard drive.

2. Using an application such as 7-Zip File Manager, extract the Hadoop binary file **hadoop-3.2.4.tar.gz** from the Apche hadoop website to a directory of your choice, or directly to **C:\hadoop\hadoop-3.2.4**. If you choose to extract the files to a directory of your choice, then you first have to copy the extracted files to **C:\hadoop**. You may find it more convenient to extract them directly to **C:\hadoop**, which will then have an extracted folder in it called hadoop-3.2.4, so you end up with the C:\hadoop\hadoop-3.2.4 folders.

3.  You can now create a **HADOOP_HOME** similar to how we created one previously, by going back to **Control Panel** > **System and Security** > **System**, clicking **Advanced System Settings**, and then clicking the **Environment Variables** button. Create the Hadoop home by adding the system variable name **HADOOP_HOME**, with the system variable value being the folder that we extracted the Hadoop binary to, which was **C:\hadoop\hadoop-3.2.4**.
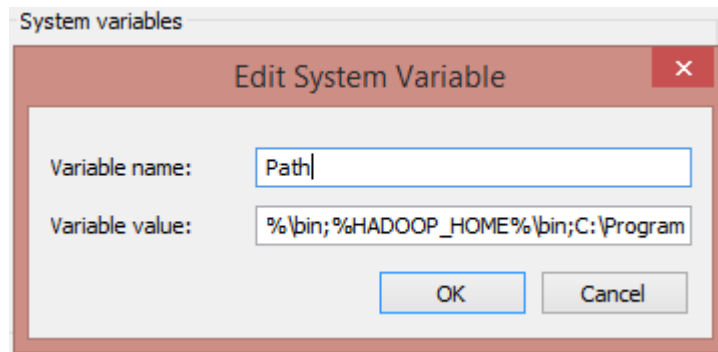


*Figure 13: Creating a Hadoop Home*



*Figure 14: Adding Hadoop Home to the Path variable*

We must add the **HADOOP_HOME** file to the **Path** variable in System variables. In this instance, it is done by adding **%HADOOP_HOME%\bin** between semi colons in the **Variable value** field.

In addition, we must add a second **HADOOP_HOME** to the **Path** variable for the folder in Hadoop called **sbin**. This is done by adding **%HADOOP_HOME%\sbin** between semi colons in the **Variable value** field. You should now have Hadoop and Java homes, and two Hadoop path variables.

*Figure 15: Java and Hadoop homes*

The resource page I mentioned previously is an official Apache resource that will assist us in finishing the installation. The area of the site we need first is "Section 3.1. Example HDFS Configuration," which states:

> "Before you can start the Hadoop Daemons you will need to make a few edits to configuration files. The configuration file templates will all be found in c:\deploy\etc\hadoop, assuming your installation directory is c:\deploy."

Since our installation is at C:\hadoop\hadoop-3.2.4, our configuration file templates will be located at C:\hadoop\hadoop-3.2.4\etc\hadoop\. The first file we need to edit is the **core-site.xml** file. The following code listing shows the format of the core-site.xml file, which is the style that we need to adapt. You will need your code editor at this point (I am using Notepad++).

*Code Listing 1: The core-site.xml file format*

```xml
<configuration>
<property>
<name>fs.default.name</name>
<value>hdfs://0.0.0.0:19000</value>
</property>
</configuration>
```

We need to substitute the name and value elements shown on the core-site.xml file on the Apache Wiki page for values in the installation we are carrying out. The values we require are contained in the following code listing and reflect our current Hadoop installation.

*Code Listing 2: Editing the core-site.xml file*

```xml
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl"href="configuration.xsl"?>

<configuration>
<property>
<name>fs.defaultFS</name>
<value>hdfs://localhost:9000</value>
</property>
</configuration>
```

We need to do the same for the **hdfs-site.xml** file template, and the new values we require are in the following code listing.

*Code Listing 3: Editing the hdfs-site.xml template*

```xml
<configuration>
<property>
<name>dfs.replication</name>
<value>1</value>
</property>
<property>
<name>dfs.namenode.name.dir</name>
<value>file:///C:/Hadoop/hadoop-3.2.4/namenode</value>
</property>
<property>
<name>dfs.datanode.data.dir</name>
<value>file:///C:/Hadoop/hadoop-3.2.4/datanode</value>
</property>

</configuration>
```

**Note: You must create two folders in the C:\Hadoop\hadoop-3.2.4\ folder in Windows Explorer to reflect the** namenode **and** datanode **directories mentioned in**

***Code Listing 3. Note that the Hadoop configuration files use forward slashes instead of backward slashes in file paths, even on Windows systems.***

Next, we need to edit the **mapred-site.xml** configuration file; the values required are shown in the following code listing.

*Code Listing 4: Editing the mapred-site.xml configuration file*

```
<configuration>
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>

</configuration>
```

We also need to edit the **yarn-site.xml** configuration file; the values required are provided in the following code listing.

*Code Listing 5: Editing the yarn-site.xml configuration file*

```
<configuration>


<property>
<name>yarn.nodemanager.aux-services</name>
<value>mapreduce_shuffle</value>
</property>
<property>
<name>yarn.nodemanager.auxservices.mapreduce.shuffle.class</name>
<value>org.apache.hadoop.mapred.ShuffleHandler</value>
</property>

</configuration>
```

*  *Now open etc/Hadoop.env file as shown below.*

```
set JAVA_HOME=C:\Java\jdk-1.8
```

Next, follow these steps:

1. Replace the bin folder at **C:\hadoop\hadoop-3.2.4\bin** with a bin folder extracted from https://github.com/killerangaswamy56/hadoop_config.git

2. Extract the bin folder from the **apache-hadoop-3.1.0-winutils-master** file downloaded from https://github.com/killerangaswamy56/hadoop_config.git.

3. Make a copy of the bin folder at **C:\hadoop\hadoop-3.2.4\bin**, and then delete the folder you made the copy from.

4. Copy the bin folder you extracted from the **apache-hadoop-3.1.0-winutils-master** file to **C:\hadoop\hadoop-3.2.4\**; it replaces the bin folder you deleted.

Now we must follow the instructions in section. Format the FileSystem." This is done by executing the following command (with administrator privileges) from a command shell:

*Code Listing 6: Format of the Filesystem*

```
hdfs namenode -format
```

You should now see the following on your screen.



*Figure 16: Successful formatting of the Filesystem*

You must now copy the **hadoop-yarn-server-timelineservice-3.2.4** file from **C:\hadoop\hadoop-3.2.4\share\hadoop\yarn\timelineservice** to the folder **C:\hadoop\hadoop-3.2.4\share\hadoop\yarn**. We can start Hadoop with the instructions in sections 3.5 and 3.6. of the Hadoop Wiki page, called "3.5. Start HDFS Daemons" and "3.6. Start YARN Daemons and run a YARN job."

You start HDFS daemons by running the following code from the command prompt.

*Code Listing 7: Start HDFS Daemons command*

```
start-dfs
```

You start YARN daemons and run a YARN job by running the following code.

*Code Listing 8: Start YARN daemons and run a YARN job command*

```
start-yarn
```

You should now see the Hadoop **namenode** and **datanode** successfully started.



*Figure 17: Hadoop namenode and datanode started successfully*

In addition, you will see the YARN **resourcemanager** and YARN **nodemanager** successfully started.

*Figure 18: Yarn resourcemanager and Yarn nodemanager successfully started*

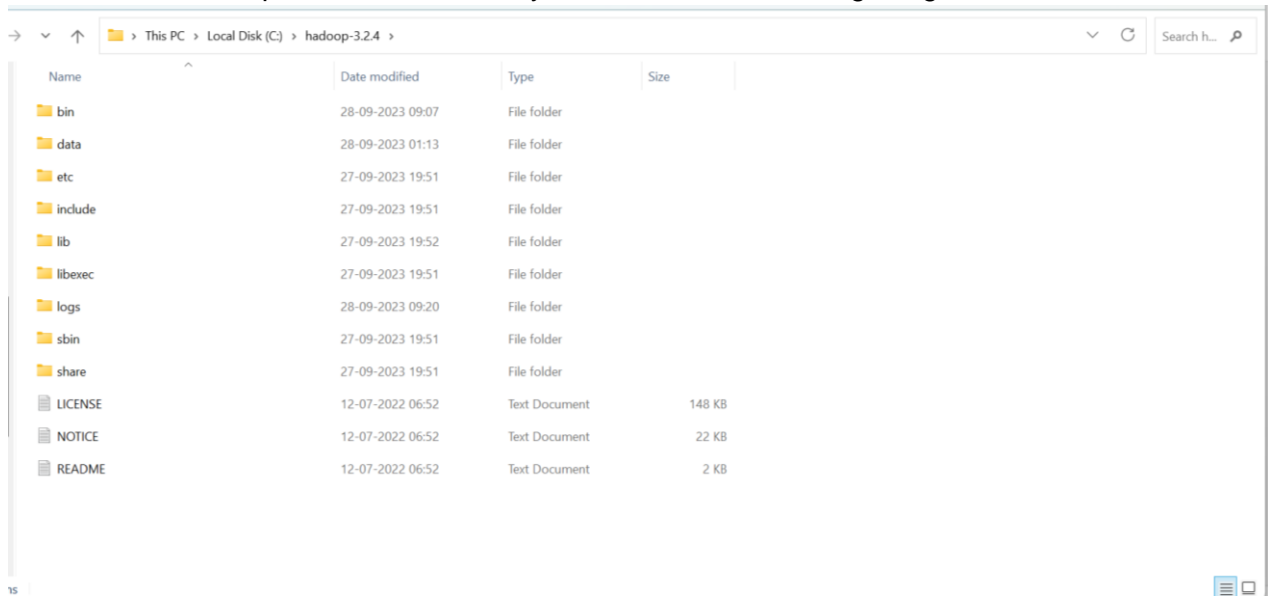The finished Hadoop installation directory is shown in the following image.



*Figure 19: Final Hadoop installation directory*





Now Follow the below steps to view your data in server

Step:1 After data files are Loaded into the HDFS server then type the below url in your default browser. http://localhost:9870

Step:2 Then it will navigates to the below Hadoop page



Step:3 Now Click on the Utilities then click on "Browse the file system"



Step:4 Now Check your data file is loaded into the server or not as shown below