# OBULI SAI NAREN

**18CSR125**

**Installation of Apache Hadoop**

**TUTORIAL - V**

*Tutorial 5*

*Date: 09.05.2021*

# *Apache Hadoop Installation & Execution*

## Installation of Apache Hadoop on Windows 10

### ✓ *Prerequisites: -*
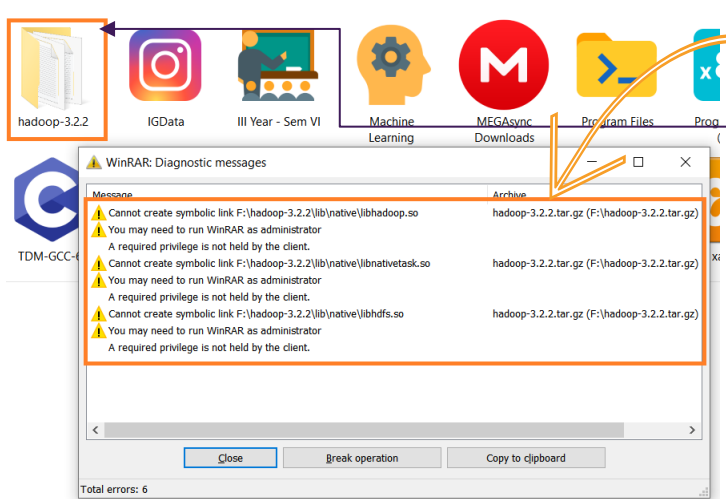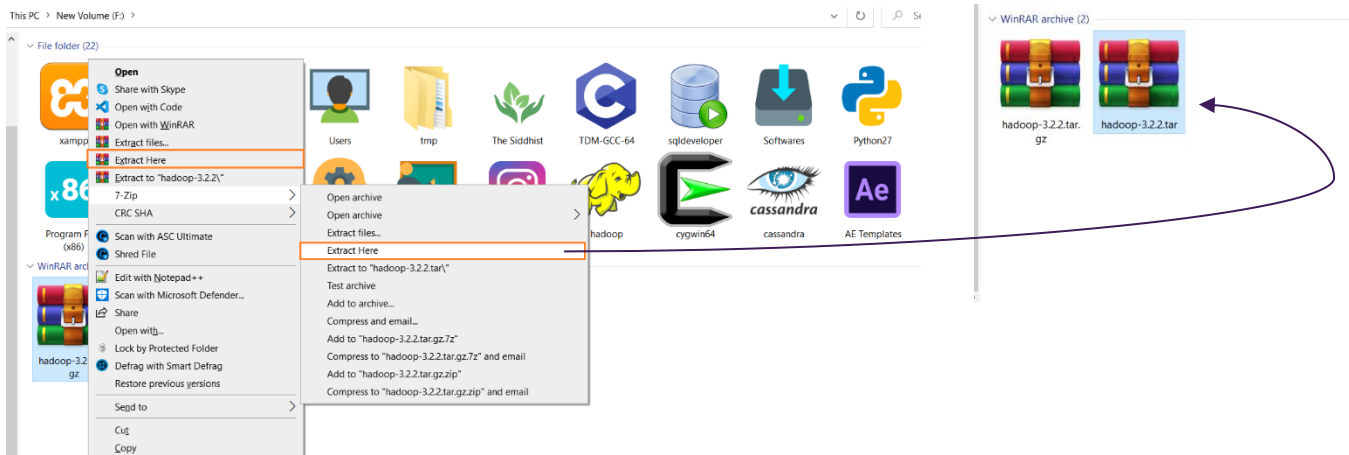
- ♦ Java 8 runtime environment (JRE)
- ♦ Preferably Java Development Kit (JDK) v1.8_***. (seems to work on all JDK)
- ♦ Any Zip Extractors [WinRAR/ 7-Zip]

### ✓ *Hadoop Download: -*

- ♦ Follow the link here ☞ Hadoop Releases .
- ♦ Download the one version older than the latest one. (for stability)
- ♦ Here I have used hadoop-3.2.2 version.
  - ▪ Click here hadoop-3.2.2.tar.gz for direct download.

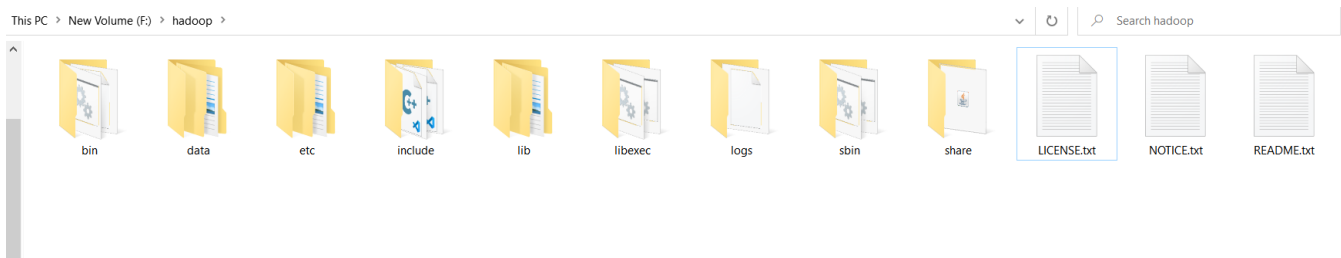### ✓ *Hadoop Extraction & Configuring: -*

- ♦ Extract the downloaded binary [hadoop-3.2.2.tar.gz] file.
  - ▪ You will directly get the **Hadoop-3.2.2** folder if you used WinRAR.
  - ▪ If you have used 7-Zip you will get " *hadoop-3.2.2.tar* " file which you should extract once again.
- ♦ See Image below for Extract Options...!
- ♦ Extract Location: "F:\"
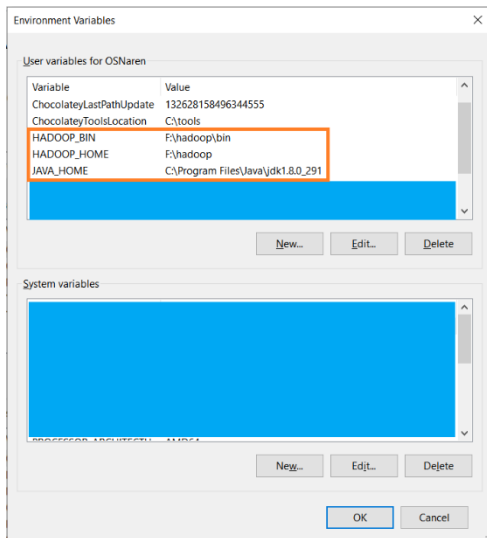
Ignore these warnings that show up at the end.

You will get a folder as in the image after extraction. Rename folder "hadoop-3.2.2"   ⟶   "hadoop"



🔖 Folder format after extraction. *(Ignore 'data' folder)*
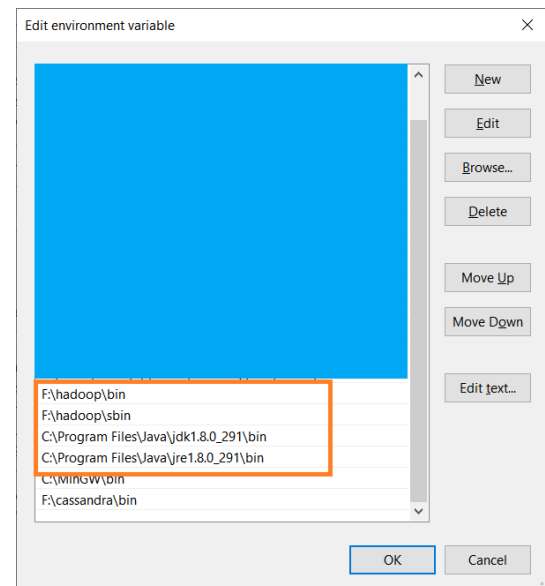


❖ *Setting up environment variables: -*

- JAVA_HOME: JDK installation folder path      :

  C:\Program Files\Java\jdk1.8.0_291

- HADOOP_HOME: Hadoop installation folder path :      F:\hadoop

- HADOOP_BIN: Hadoop Bin path                    :      F:\hadoop\bin

＋ Create 3 New variables as shown.

＋ Add these 4 Path variables as shown.



❖ *Configure Hadoop XMLs: -*

There are five files we should configure in the path below…

- Path: "F:\hadoop\etc\hadoop\"

  o core-site.xml

  o hdfs-site.xml

  o mapred-site.xml

  o yarn-site.xml

  o hadoop-env.cmd

➕ Just download & extract the files from <u>Config Zip</u> and replace the existing files in the path mentioned.

➢ hdfs-site.xml     :        Check the path of Namenode & Datanode.

```xml
<configuration>
    <property>
        <name>dfs.replication</name>
        <value>1</value>
    </property>
    <property>
        <name>dfs.namenode.name.dir</name>
        <value>/F:/hadoop/data/namenode</value>
    </property>
    <property>
        <name>dfs.datanode.data.dir</name>
        <value>/F:/hadoop/data/datanode</value>
    </property>
</configuration>
```

➢ Open "hadoop-env.cmd" & check if the paths are set correct...

```
set JAVA_HOME="C:\Program Files\Java\jdk1.8.0_291"
set HADOOP_PREFIX=%HADOOP_HOME%
set HADOOP_CONF_DIR=%HADOOP_PREFIX%\etc\hadoop
set YARN_CONF_DIR=%HADOOP_CONF_DIR%
set PATH=%PATH%;%HADOOP_PREFIX%\bin
```

➕ No changes in other 3 files...!

➕ Now download & extract the files from <u>Bin Zip</u> and replace the existing files in the path mentioned.
  ▪ Path: "F:\hadoop\bin"

➕ If you are installing some other versions download respective Bin from here ☞ <u>Winutils</u>.

## ✓ *Running Hadoop: -*

➕ Open an Administrator Command Prompt.

➕ Run the command hadoop version and expect an output like below.

```
C:\Users\66nar>hadoop version
Hadoop 3.2.2
Source code repository Unknown -r 7a3bc90b05f257c8ace2f76d74264906f0f7a932
Compiled by hexiaoqiao on 2021-01-03T09:26Z
Compiled with protoc 2.5.0
From source with checksum 5a8f564f46624254b27f6a33126ff4
This command was run using /F:/hadoop/share/hadoop/common/hadoop-common-3.2.2.jar

C:\Users\66nar>
```

➕ Now, to format the Namenode → Run the command hdfs namenode -format

▪ 🚨 Note – Run this command only once.

▪ After long log messages… To the 6th-7th line from the bottom has a message: /F:/hadoop/data/Datanode has been successfully formatted.

➕ We must start 4 process of Hadoop…

○ Namenode

○ Datanode

○ Resource manager

○ Node manager

➢ Deprecated method → Run the command start-all.cmd

```
C:\Users\66nar>start-all.cmd
This script is Deprecated. Instead use start-dfs.cmd and start-yarn.cmd
starting yarn daemons
```

▪ This opens the 4 processes each in a separate cmd window as below.

> Or you can → Run the commands:
>   ▪ start-dfs.cmd - Starts namenode and datanode.
>   ▪ start-yarn.cmd - Starts node manager and resource manager.
> To check whether all the process has started and running properly:
>   ▪ Run the command → *jps*

> ➢ To stop the processes: Use respective commands.

  ▪ stop-all.cmd

  ▪ stop-dfs.cmd

  ▪ stop-yarn.cmd

```
C:\Users\66nar>stop-all.cmd
This script is Deprecated. Instead use stop-dfs.cmd and stop-yarn.cmd
SUCCESS: Sent termination signal to the process with PID 22156.
SUCCESS: Sent termination signal to the process with PID 8116.
stopping yarn daemons
SUCCESS: Sent termination signal to the process with PID 20700.
SUCCESS: Sent termination signal to the process with PID 15488.

INFO: No tasks running with the specified criteria.
```

🞧 Hadoop Web UI: -

> ➢ These 3-localhost address works if hadoop is running without and errors.

  ▪ http://localhost:8088      -      *Yarn page*



  ▪ http://localhost:9870      -      *Name node page*



### Overview 'localhost:9000' (active)

| | |
|---|---|
| Started: | Mon May 10 12:56:07 +0530 2021 |
| Version: | 3.2.2, r7a3bc90b05f257c8ace2f76d74264906f0f7a932 |
| Compiled: | Sun Jan 03 14:56:00 +0530 2021 by hexiaoqiao from branch-3.2.2 |
| Cluster ID: | CID-67c61df3-bb09-4f98-8b36-2c380271bed8 |
| Block Pool ID: | BP-1940875482-192.168.45.1-1620294652996 |

- [http://localhost:9864](http://localhost:9864)    -    *Data node page*

Hadoop    Overview    Utilities ▾

DataNode on ASUS-Naren:9866

| Cluster ID: | CID-67c61df3-bb09-4f98-8b36-2c380271bed8 |
| Version: | 3.2.2, r7a3bc90b05f257c8ace2f76d74264906f0f7a932 |

## Block Pools

| Namenode Address | Block Pool ID | Actor State | Last Heartbeat | Last Block Report | Last Block Report Size (Max Size) |
| --- | --- | --- | --- | --- | --- |
| localhost:9000 | BP-1940875482-192.168.45.1-1620294652996 | RUNNING | 0s | 24 minutes | 0 B (64 MB) |

# ✓ Handling Errors in Installation Process: –

1. **JAVA_HOME is incorrectly set :–**
   a. Use "Progra~1" instead of "Program Files"
   b. Use "Progra~2" instead of "Program Files(x86)"
      - Replace this wherever you have given JAVA path.
         o hadoop-env.cmd – Path variables – JAVA_HOME

2. **Datanode shutting down / Not starting / Exception in datanode cmd :–**
   a. Copy the Cluster ID from the namenode VERSION file in the directory → *"F:\hadoop\data\namenode\current"*.
   b. Paste the Cluster ID that you copied to the datanode VERSION file in the directory → *"F:\hadoop\data\datanode\current"*.
      - Do not give *hdfs namenode -format* more than once, this resets the Cluster ID in the namenode. So, you have to copy & paste each and every time.

3. *Exception while formatting Namenode / formatting failure :–*

    a. Run the ***hdfs namenode -format*** command in administrator command prompt.

    b. Do not use hadoop-3.2.1

    c. If you are using hadoop-3.2.1…

        ▪ Download hadoop-hdfs-3.2.1.jar file from the link.

        ▪ Rename the file name *hadoop-hdfs-3.2.1.jar* to *hadoop-hdfs-3.2.1.bak* in folder %HADOOP_HOME%\share\hadoop\hdfs

        ▪ Copy the downloaded hadoop-hdfs-3.2.1.jar to folder %HADOOP_HOME%\share\hadoop\hdfs

4. *Exception in Resource manager / Not running :–*

    a. Check whether the file *"hadoop-yarn-server-resourcemanager-3.2.2.jar"* is present in the locations…

        ▪ F:\hadoop\share\hadoop\yarn

        ▪ F:\hadoop\share\hadoop\yarn\sources

    b. If not copy & paste from the directory that has the file to the directory that does not have the file.

# Hadoop MapReduce

✔ *MapReduce Program Source Code*

Hadoop.java.html          Choices.java.html          MapMR.java.html

✔ *MapReduce Sample Data*

Sample Data.html

✔ Hadoop using IntelliJ – Follow this link to run your Hadoop MapReduce program in IntelliJ.

✔ Hadoop using Eclipse – Follow this link to run your Hadoop MapReduce program in Eclipse.