MScIT Semester 2
Big Data Analytics
Hadoop Installation

Mumbai University

# Install, configure and run Hadoop and HDFS and explore HDFS on Windows

Steps to Install Hadoop

1. Install Java JDK 1.8
2. Download Hadoop and extract and place under C drive
3. Set Path in Environment Variables
4. Config files under Hadoop directory
5. Create folder datanode and namenode under data directory
6. Edit HDFS and YARN files
7. Set Java Home environment in Hadoop environment
8. Setup Complete. Test by executing start-all.cmd

There are two ways to install Hadoop, i.e.

9.      Single node

10.     Multi node

Here, we use multi node cluster.

1.      Install Java

11.     – Java JDK Link to download

https://www.oracle.com/java/technologies/javase-jdk8-downloads.html

12.     – extract and install Java in C:\Java

13.     – open cmd and type -> javac -version
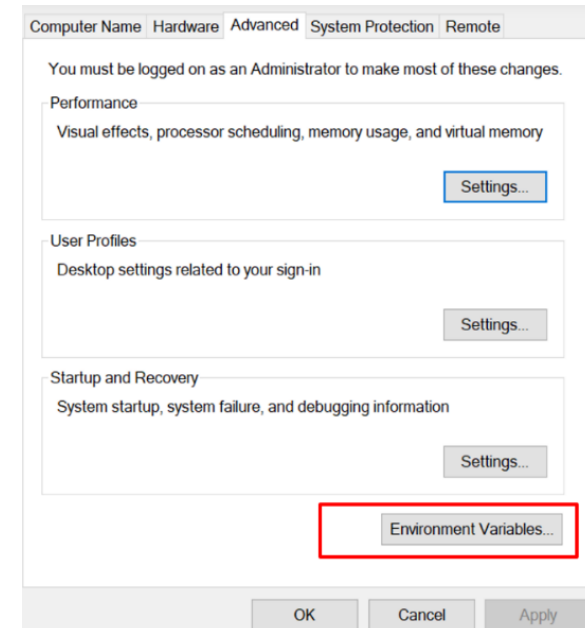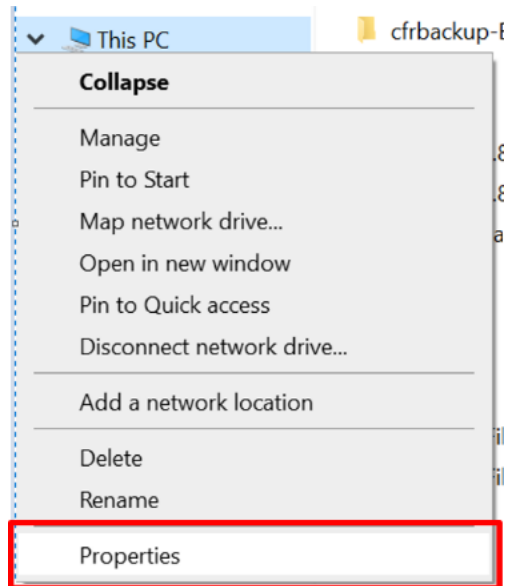
```
C:\>javac -version
javac 20
```

2.      Download Hadoop

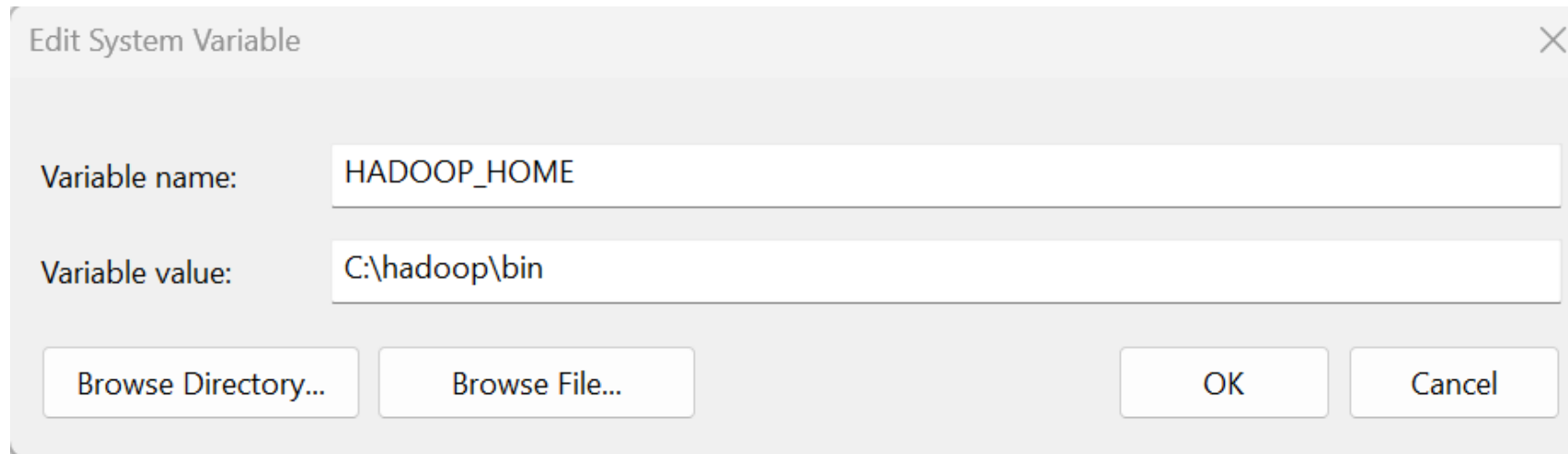https://www.apache.org/dyn/closer.cgi/hadoop/common/hadoop-3.3.0/hadoop-3.3.0.tar.gz


- right click .rar.gz file -> show more options -> 7-zip->and extract to C:\Hadoop-3.3.0\

3.    Set the path JAVA_HOME Environment variable

4.    Set the path HADOOP_HOME Environment variable

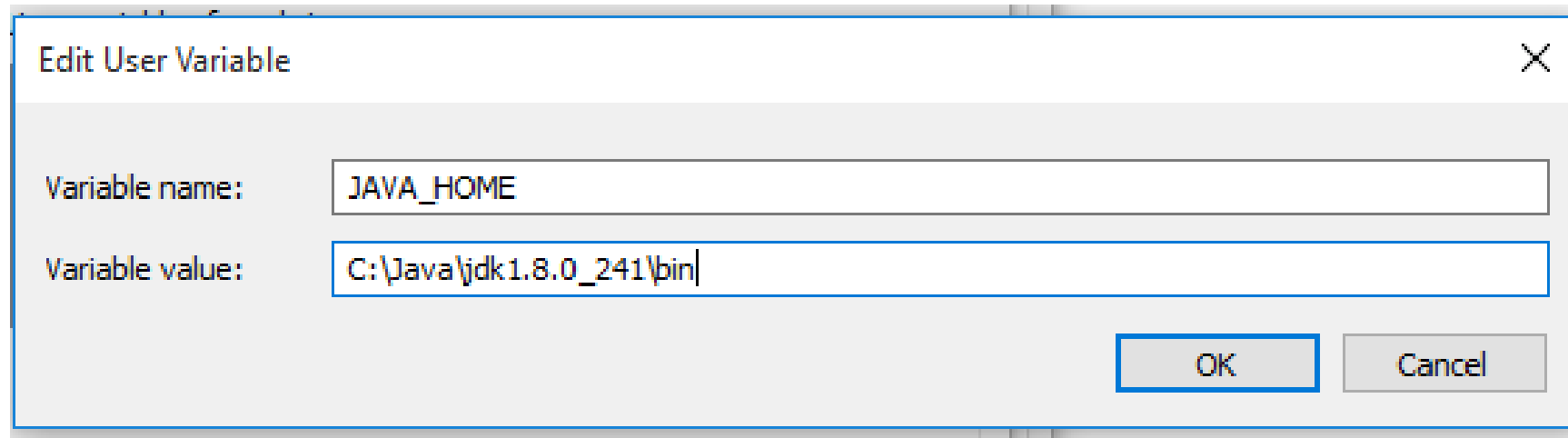# Click on New to both user variables and system variables.

**Edit System Variable** ✕

Variable name: HADOOP_HOME

Variable value: C:\hadoop\bin

Browse Directory... | Browse File... | OK | Cancel

**Edit User Variable** ✕

Variable name: JAVA_HOME

Variable value: C:\Java\jdk1.8.0_241\bin

OK | Cancel

- Click on user variable -> path -> edit-> add path for Hadoop and java upto 'bin'

5.      Configurations

Edit file C:\hadoop\etc\hadoop\core-site.xml,

paste the xml code in folder and save

```xml
<configuration>

<property>

    <name>fs.defaultFS</name>

    <value>hdfs://localhost:9000</value>

 </property>

</configuration>
```

Rename "mapred-site.xml.template" to "mapred-site.xml" and edit this file C:/hadoop/etc/hadoop/mapred-site.xml, paste xml code and save this file

```xml
<configuration>

   <property>

      <name>mapreduce.framework.name</name>

      <value>yarn</value>

   </property>

</configuration>
```

- Create folder "data" under "C:\Hadoop-3.3.0"
- Create folder "datanode" under "C:\Hadoop-3.3.0\data"
- Create folder "namenode" under "C:\Hadoop-3.3.0\data"

Edit file C:\Hadoop-3.3.0/etc/hadoop/hdfs-site.xml,

paste xml code and save this file.

```xml
<configuration>
<property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>
```

```xml
<property>
    <name>dfs.namenode.name.dir</name>
    <value>/hadoop-3.3.0/data/namenode</value>
  </property>
  <property>
    <name>dfs.datanode.data.dir</name>
    <value>/hadoop-3.3.0/data/datanode</value>
  </property>
  </configuration>
```

Edit file C:/Hadoop-3.3.0/etc/hadoop/yarn-site.xml,

paste xml code and save this file.

```xml
<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>
  <property>
                          <name>yarn.resourcemanager.address</name>
                          <value>127.0.0.1:8032</value>
  </property>
  <property>
                          <name>yarn.resourcemanager.scheduler.address</name>
                          <value>127.0.0.1:8030</value>
   </property>
   <property>
                          <name>yarn.resourcemanager.resource-tracker.address</name>
                          <value>127.0.0.1:8031</value>
    </property>
</configuration>
```

6.      Edit file C:/Hadoop-3.3.0/etc/hadoop/hadoop-env.cmd

Find "JAVA_HOME=%JAVA_HOME%" and replace it as

set JAVA_HOME="C:\java\jdk1.8.0_121"

7.    Download "redistributable" package

Download and run VC_redist.x64.exe

This is a "redistributable" package of the Visual C runtime code for 64-bit applications, from Microsoft. It contains certain shared code that every application written with Visual C expects to have available on the Windows computer it runs on.

8.	Hadoop Configurations

Download bin folder from

https://github.com/s911415/apache-hadoop-3.1.0-winutils

– Copy the bin folder to c:\hadoop-3.3.0. Replace the existing bin folder.

9.	copy "hadoop-yarn-server-timelineservice-3.0.3.jar" from ~\hadoop-3.0.3\share\hadoop\yarn\timelineservice to ~\hadoop-3.0.3\share\hadoop\yarn folder.

Format the NameNode

– Open cmd 'Run as Administrator' and type command "hdfs namenode –format"

. Testing

– Open cmd 'Run as Administrator' and change directory to C:\Hadoop-3.3.0\sbin

– type start-all.cmd

 OR

 - type start-dfs.cmd

– type start-yarn.cmd

```
C:\hadoop-3.3.0\sbin>start-all.cmd
This script is Deprecated. Instead use start-dfs.cmd and start-yarn.cmd
The filename, directory name, or volume label syntax is incorrect.
The filename, directory name, or volume label syntax is incorrect.
starting yarn daemons
The filename, directory name, or volume label syntax is incorrect.
```

# Type start-all.cmd

- You will get 4 more running threads for Datanode, namenode, resouce manager and node manager

Type JPS command to start-all.cmd command prompt, you will get following output.



```
C:\hadoop-3.3.0\sbin>jps
5632 Jps
7572 DataNode
3752 ResourceManager
7992 NameNode
8028 NodeManager
```

- Run http://localhost:9870/ from any browser



Overview 'localhost:9000' (✔active)

| Started: | Wed Mar 15 12:10:54 +0530 2023 |
| --- | --- |
| Version: | 3.3.0, raa96f1871bfd858f9bac59cf2a81ec470da649af |
| Compiled: | Tue Jul 07 00:14:00 +0530 2020 by brahma from branch-3.3.0 |
| Cluster ID: | CID-1986aba8-0ed3-43a2-9db7-42944ec518b2 |
| Block Pool ID: | BP-1049743432-192.168.56.1-1678862097216 |



Browse Directory

Show 25 entries                                           Search:

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | No data available in table | | | | |

Showing 0 to 0 of 0 entries                              Previous  Next