# Index:

| Practical No | Title |
|---|---|
| 1. | Installation of Hadoop on Windows |
| 2. | Installation of Scala and Apache Spark |
| 3. | Spark GraphX |
| 4. | Working with PySpark |
| 5. | Installation of HBase |
| 6. | |
| 7. | |

# Practical 1: Installation of Hadoop on Windows.

1. Download the latest Hadoop Binary version from [here](#).



2. Download Java version 11.0.23 from [here](#).

3. Run the installer and follow the steps to install Java.



4: Make sure to set the Java Environment Variables by creating a new variable and also adding it to the path variable.
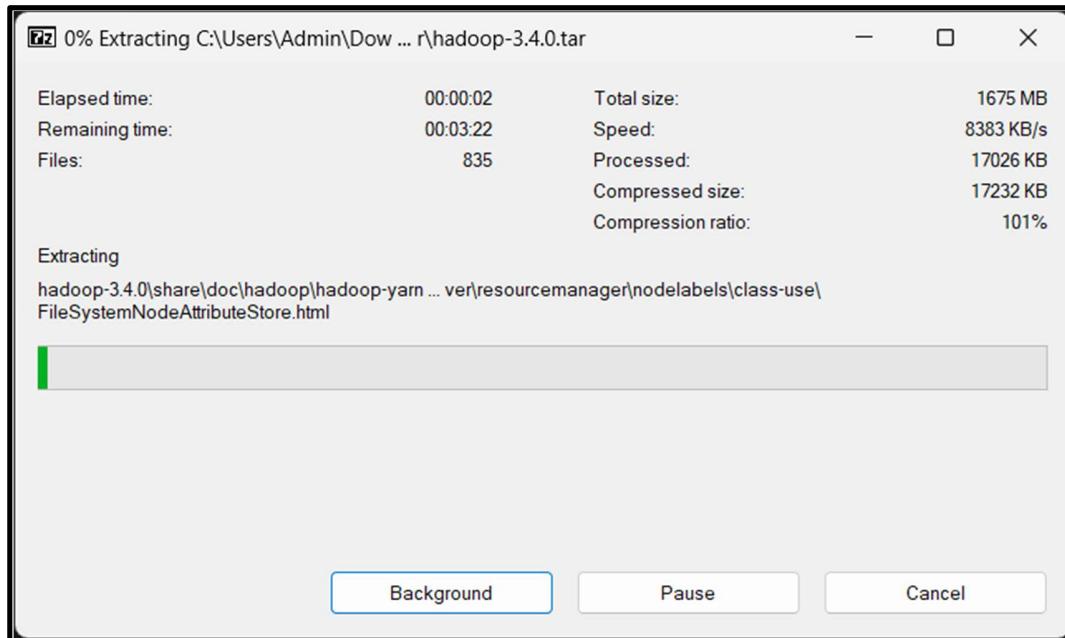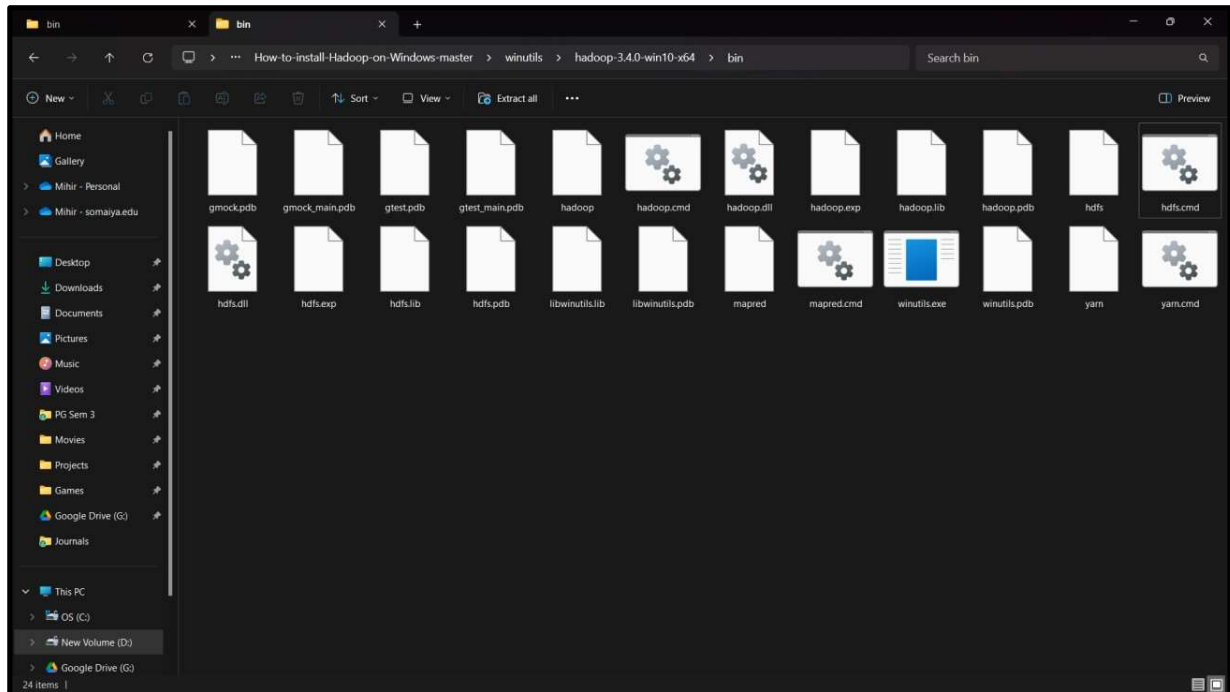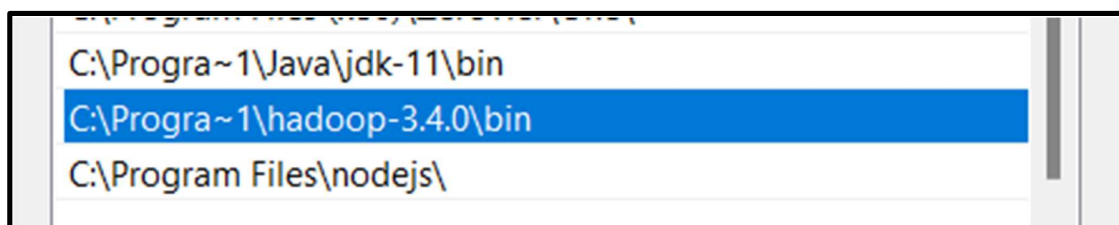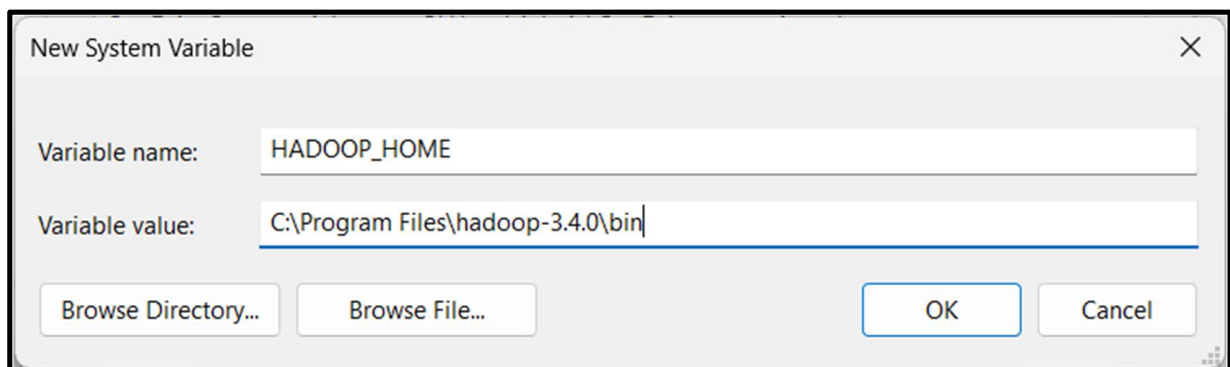
5. Extract the previously downloaded Hadoop archive and copy the contents to your desired location.

| | | | |
|---|---|---|---|
| 7z 0% Extracting C:\Users\Admin\Dow ... r\hadoop-3.4.0.tar | | — ☐ ✕ | |

| Elapsed time: | 00:00:02 | Total size: | 1675 MB |
|---|---|---|---|
| Remaining time: | 00:03:22 | Speed: | 8383 KB/s |
| Files: | 835 | Processed: | 17026 KB |
| | | Compressed size: | 17232 KB |
| | | Compression ratio: | 101% |

Extracting

hadoop-3.4.0\share\doc\hadoop\hadoop-yarn ... ver\resourcemanager\nodelabels\class-use\
FileSystemNodeAttributeStore.html

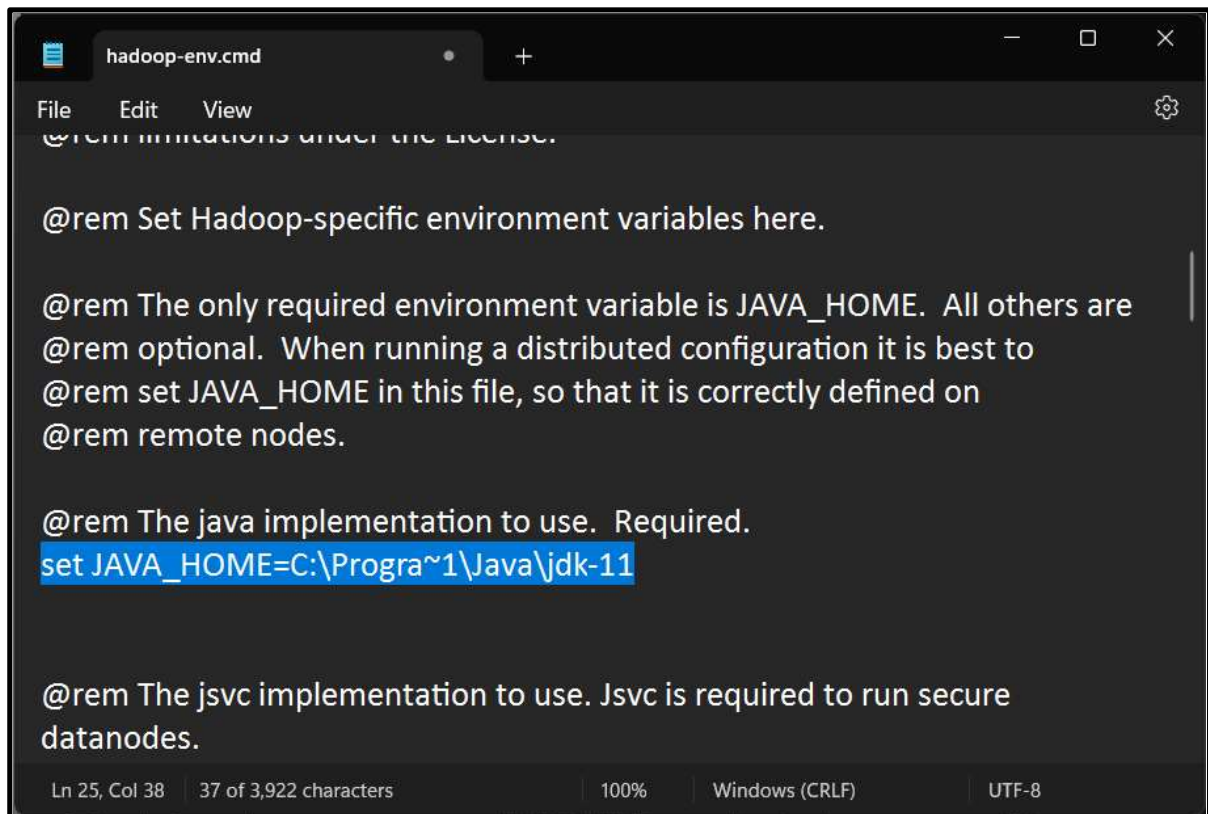| Background | Pause | Cancel |
|---|---|---|

6. Download and Install Hadoop native IO binary from [here](#). To install copy all the contents of the archive to the bin folder of your Hadoop.



7. Set up the Hadoop Environment Variables.

8. Set the JAVA_HOME in the hadoop_env.cmd file located inside the etc folder.



9. Check the versions of both Java and Hadoop.

10. To verify all the above steps are completed successfully, open the bin folder in terminal and run the command winutils.exe

```
        task kill [TASKNAME]
           Kills task job object

        task processList [TASKNAME]
           Prints to stdout a list of processes in the task
           along with their resource usage. One process per line
           and comma separated info per process
           ProcessId,VirtualMemoryCommitted(bytes),
           WorkingSetSize(bytes),CpuTime(Millisec,Kernel+User)
service         Service operations.

    Usage: service
    Starts the nodemanager Windows Secure Container Executor helper service.
    The service must run as a high privileged account (LocalSystem)
    and is used by the nodemanager WSCE to spawn secure containers on Window
s.


PS C:\Program Files\hadoop-3.4.0\bin>
```
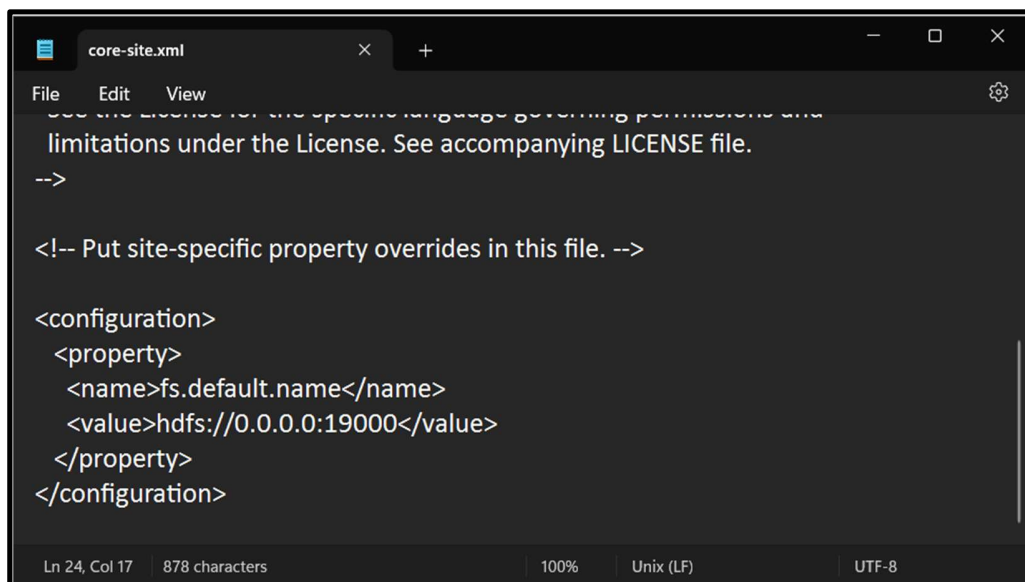
11. Hadoop Configurations. All the files will be located in the etc folder of your Hadoop installation.

11a. Configure core-site (core-site.xml)

```
See the License for the specific language governing permissions and
  limitations under the License. See accompanying LICENSE file.
-->


<!-- Put site-specific property overrides in this file. -->

<configuration>
  <property>
    <name>fs.default.name</name>
    <value>hdfs://0.0.0.0:19000</value>
  </property>
</configuration>
```
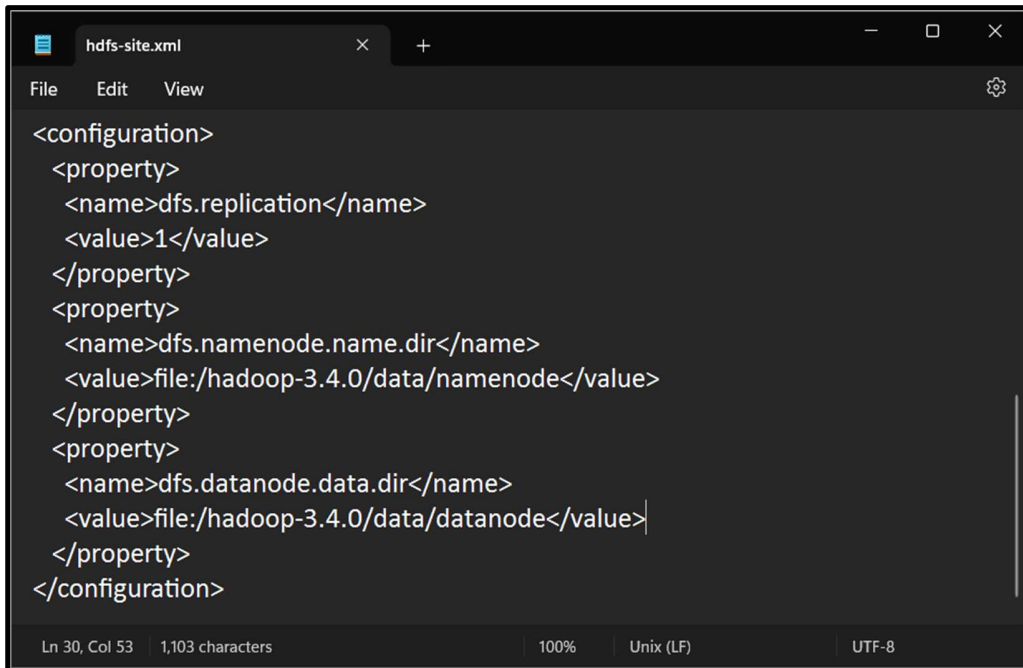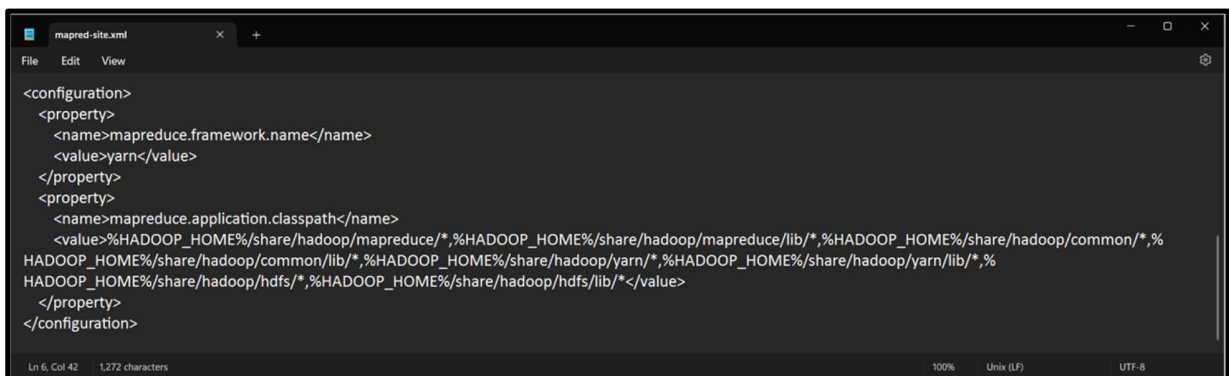
## 11b. Configure HDFS (hdfs-site.xml)

Before configuring HDFS, make a folder called as data in your Hadoop installation. Make two subfolders named as namenode and datanode.

```
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file:/hadoop-3.4.0/data/namenode</value>
  </property>
  <property>
    <name>dfs.datanode.data.dir</name>
    <value>file:/hadoop-3.4.0/data/datanode</value>
  </property>
</configuration>
```

## 11c. Configure MapReduce (mapred-site.xml)

```
<configuration>
  <property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
  </property>
  <property>
    <name>mapreduce.application.classpath</name>
    <value>%HADOOP_HOME%/share/hadoop/mapreduce/*,%HADOOP_HOME%/share/hadoop/mapreduce/lib/*,%HADOOP_HOME%/share/hadoop/common/*,%HADOOP_HOME%/share/hadoop/common/lib/*,%HADOOP_HOME%/share/hadoop/yarn/*,%HADOOP_HOME%/share/hadoop/yarn/lib/*,%HADOOP_HOME%/share/hadoop/hdfs/*,%HADOOP_HOME%/share/hadoop/hdfs/lib/*</value>
  </property>
</configuration>
```

## 11d. Configure YARN (yarn-site.xml)



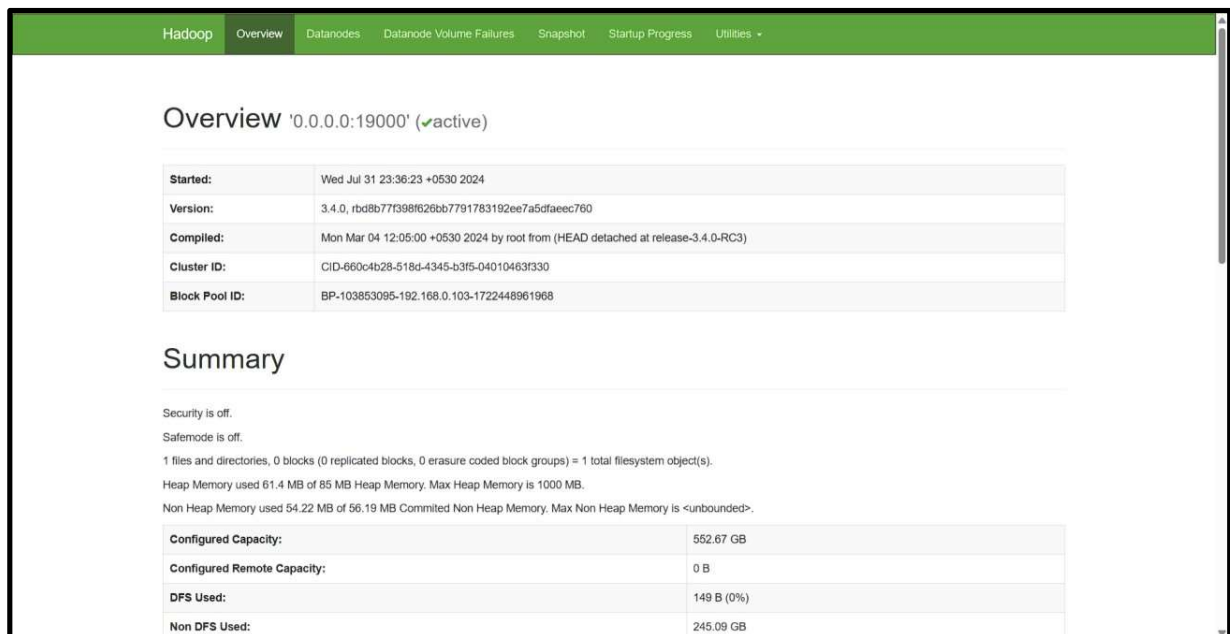## 12. Initialise HDFS using the following command hdfs namenode -format

13. Start HDFS Daemons by going to the sbin folder of Hadoop installation and running the command start-dfs.cmd. Processes will start running.



14. Verify HDFS web portal UI through this link.



15. Start YARN Daemons by using the following command start-yarn.cmd.

16. Verify YARN resource manager through this [link](link).