

DATA SCIENCE AND BIG DATA ANALYTICS

EXPERIMENT-8: Installation of Big data technologies and building a Hadoop Cluster

AIM: Installation of Big data technologies and building a Hadoop cluster.

Description:

Apache Hadoop is an open source framework that is used to efficiently store and process large datasets ranging in size from gigabytes to petabytes of data. Instead of using one large computer to store and process the data, Hadoop allows clustering multiple computers to analyze massive datasets in parallel more quickly.

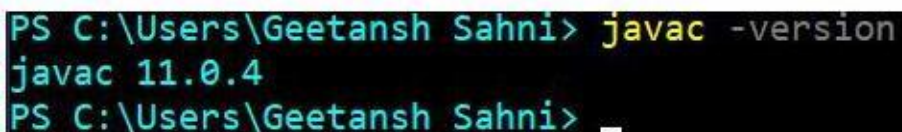
Hadoop consists of four main modules:

- Hadoop Distributed File System (HDFS) – A distributed file system that runs on standard or low-end hardware. HDFS provides better data throughput than traditional file systems, in addition to high fault tolerance and native support of large datasets.
- Yet Another Resource Negotiator (YARN) – Manages and monitors cluster nodes and resource usage. It schedules jobs and tasks.
- MapReduce – A framework that helps programs do the parallel computation on data. The map task takes input data and converts it into a dataset that can be computed in key value pairs. The output of the map task is consumed by reduce tasks to aggregate output and provide the desired result.
- Hadoop Common – Provides common Java libraries that can be used across all modules.

PROCEDURE:

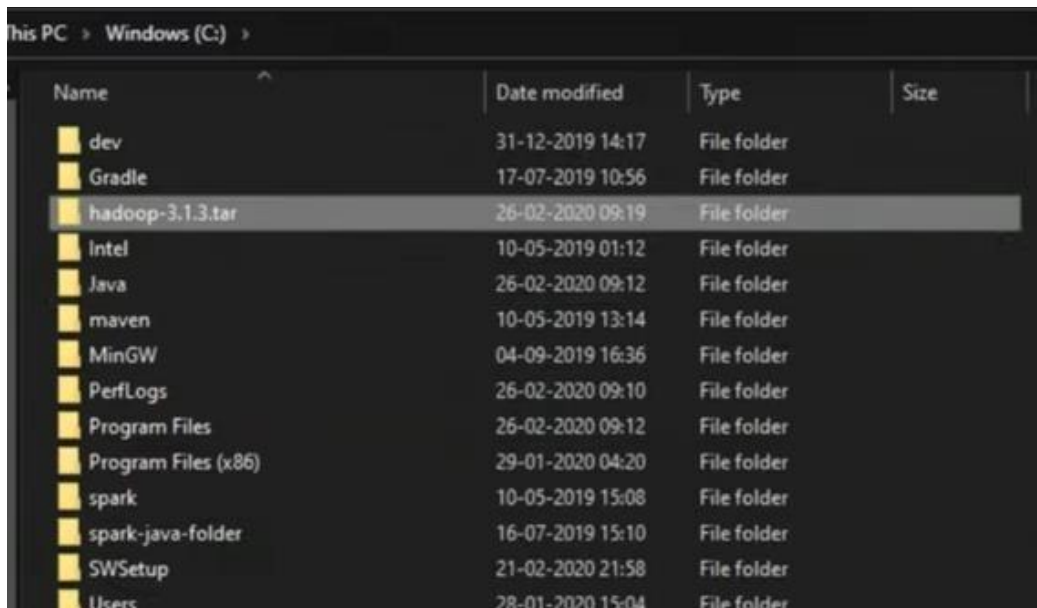
Step 1: Verify the Java installed

javac -version



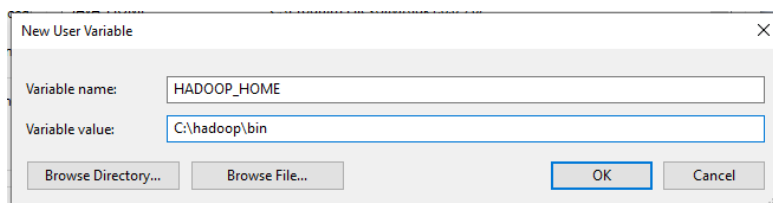
```
PS C:\Users\Geetansh Sahni> javac -version
javac 11.0.4
PS C:\Users\Geetansh Sahni> _
```

Step 2: Extract Hadoop at C:\Hadoop



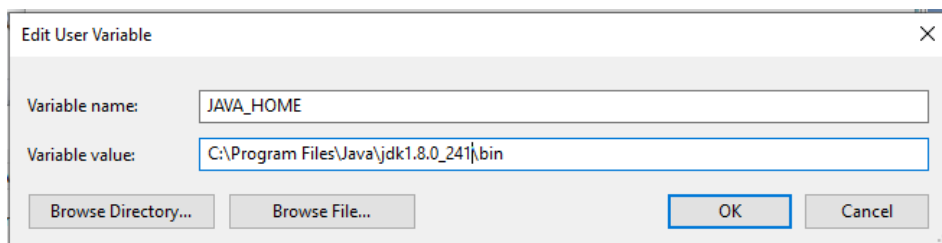
Step 3: Setting up the HADOOP_HOME variable

Use windows environment variable setting for Hadoop Path setting.

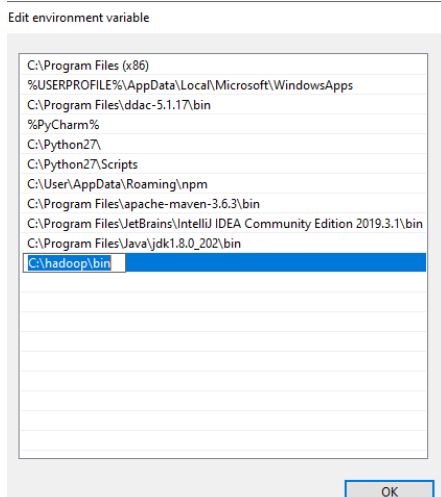


Step 4: Set JAVA_HOME variable

Use windows environment variable setting for Hadoop Path setting.



Step 5: Set Hadoop and Java bin directory path



Step 6: Hadoop Configuration :

For Hadoop Configuration we need to modify Six files that are listed below-

1. Core-site.xml
2. Mapred-site.xml
3. Hdfs-site.xml
4. Yarn-site.xml
5. Hadoop-env.cmd
6. Create two folders datanode and namenode

Step 6.1: Core-site.xml configuration

```
<configuration>
```

```
  <property>
```

```
    <name>fs.defaultFS</name>
```

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

```
  </property>
```

```
</configuration>
```

Step 6.2: Mapred-site.xml configuration

```
<configuration>
```

```
  <property>
```

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

```
    <value>yarn</value>
```

```
  </property>
```

```
</configuration>
```

Step 6.3: Hdfs-site.xml configuration

```
<configuration>
```

```
  <property>
```

```
    <name>dfs.replication</name>
```

```
    <value>1</value>
```

```
  </property>
```

```
  <property>
```

```
    <name>dfs.namenode.name.dir</name>
```

```
    <value>C:\hadoop-2.8.0\data\namenode</value>
```

```
  </property>
```

```
  <property>
```

```

    <name>dfs.datanode.data.dir</name>
    <value>C:\hadoop-2.8.0\data\datanode</value>
  </property>
</configuration>

```

Step 6.4: Yarn-site.xml configuration

```

<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>

```

Step 6.5: Hadoop-env.cmd configuration

Set "JAVA_HOME=C:\Java" (On C:\java this is path to file jdk.18.0)

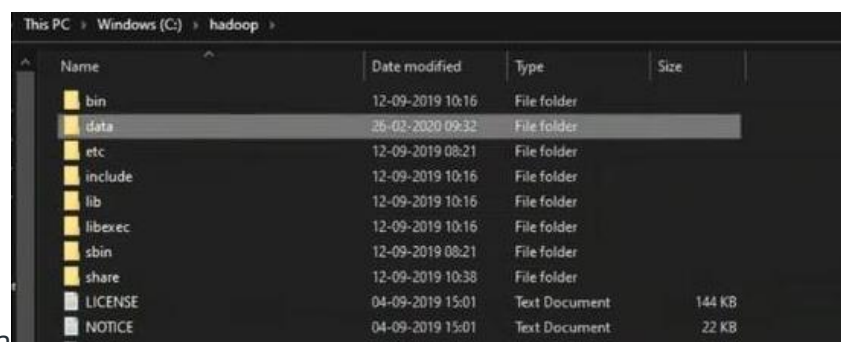
```

16  @rem Set Hadoop-specific environment variables here.
17
18
19  @rem The only required environment variable is JAVA_HOME. All others are
20  @rem optional. When running a distributed configuration it is best to
21  @rem set JAVA_HOME in this file, so that it is correctly defined on
22  @rem remote nodes.
23
24  @rem The java implementation to use. Required.
25  set JAVA_HOME=%JAVA_HOME%
26

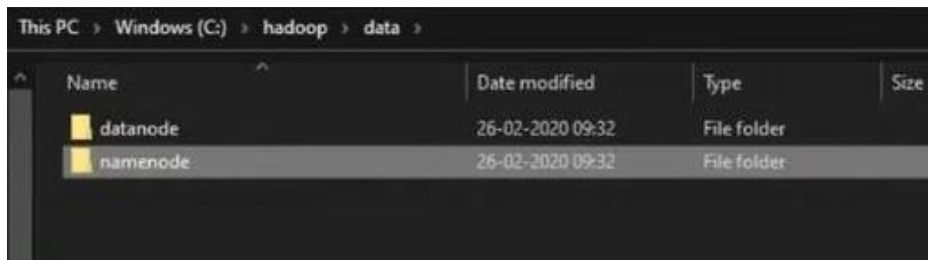
```

Step 6.6: Create datanode and namenode folders

1. Create folder "data" under "C:\Hadoop-2.8.0"
2. Create folder "datanode" under "C:\Hadoop-2.8.0\data"



3. Create folder "namenode" under "C:\Hadoop-2.8.0\data"



Step 7: Format the namenode folder

Open command window (cmd) and typing command “hdfs namenode –format”

```
C:\Users\Ravikiran>hdfs namenode -format
2020-02-26 09:42:38,498 INFO namenode.NameNode: STARTUP_MSG:
*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = LAP10P-AV5R03TS/192.168.207.1
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 3.1.3
STARTUP_MSG: classpath = C:\hadoop\etc\hadoop;C:\hadoop\share\hadoop\common;C:\hadoop\share\hadoop\common\lib\access-smart-1.2.jar;C:\hadoop\share\hadoop\common\lib\animal-sniffer-annotations-1.17.jar;C:\hadoop\share\hadoop\common\lib\asm-5.0.4.jar;C:\hadoop\share\hadoop\common\lib\audience-annotations-0.5.0.jar;C:\hadoop\share\hadoop\common\lib\avro-1.7.jar;C:\hadoop\share\hadoop\common\lib\checker-qual-2.5.2.jar;C:\hadoop\share\hadoop\common\lib\commons-beanutils-1.3.jar;C:\hadoop\share\hadoop\common\lib\commons-cli-1.2.jar;C:\hadoop\share\hadoop\common\lib\commons-codec-1.11.jar;C:\hadoop\share\hadoop\common\lib\commons-collections-3.2.2.jar;C:\hadoop\share\hadoop\common\lib\commons-compress-1.18.jar;C:\hadoop\share\hadoop\common\lib\commons-configuration2-2.1.1.jar;C:\hadoop\share\hadoop\common\lib\commons-io-2.5.jar;C:\hadoop\share\hadoop\common\lib\commons-lang-2.6.jar;C:\hadoop\share\hadoop\common\lib\commons-lang3-3.4.jar;C:\hadoop\share\hadoop\common\lib\commons-logging-1.1.3.jar;C:\hadoop\share\hadoop\common\lib\commons-math3-3.1.1.jar;C:\hadoop\share\hadoop\common\lib\commons-net-3.6.jar;C:\hadoop\share\hadoop\common\lib\curator-client-2.13.0.jar;C:\hadoop\share\hadoop\common\lib\curator-framework-2.13.0.jar;C:\hadoop\share\hadoop\common\lib\curator-recipes-2.13.0.jar;C:\hadoop\share\hadoop\common\lib\error-prone-annotations-2.2.0.jar;C:\hadoop\share\hadoop\common\lib\failureaccess-1.0.jar;C:\hadoop\share\hadoop\common\lib\gson-2.2.4.jar;C:\hadoop\share\hadoop\common\lib\guava-27.0-jre.jar;C:\hadoop\share\hadoop\common\lib\hadoop-annotations-3.1.3.jar;C:\hadoop\share\hadoop\common\lib\hadoop-auth-3.1.3.jar;C:\hadoop\share\hadoop\common\lib\htrace-core4-4.1.0-incubating.jar;C:\hadoop\share\hadoop\common\lib\httpclient-4.5.2.jar;C:\hadoop\share\hadoop\common\lib\httpcore-4.4.4.jar;C:\hadoop\share\hadoop\common\lib\j2objc-annotations-1.1.jar;C:\hadoop\share\hadoop\common\lib\jackson-annotations-2.7.8.jar;C:\hadoop\share\hadoop\common\lib\jackson-core-2.7.8.jar;C:\hadoop\share\hadoop\common\lib\jackson-core-asl-1.9.13.jar;C:\hadoop\share\hadoop\common\lib\jackson-databind-2.7.8.jar;C:\hadoop\share\hadoop\common\lib\jackson-jaxrs-1.9.13.jar;C:\hadoop\share\hadoop\common\lib\jackson-mapper-asl-1.9.13.jar;C:\hadoop\share\hadoop\common\lib\jackson-xc-1.9.13.jar;C:\hadoop\share\hadoop\common\lib\javax-servlet-api-3.1.0.jar;C:\hadoop\share\hadoop\common\lib\...
```

Step 8: Testing the setup

Open command window (cmd) and typing command “start-all.cmd”



Step 8.1: Testing the setup:

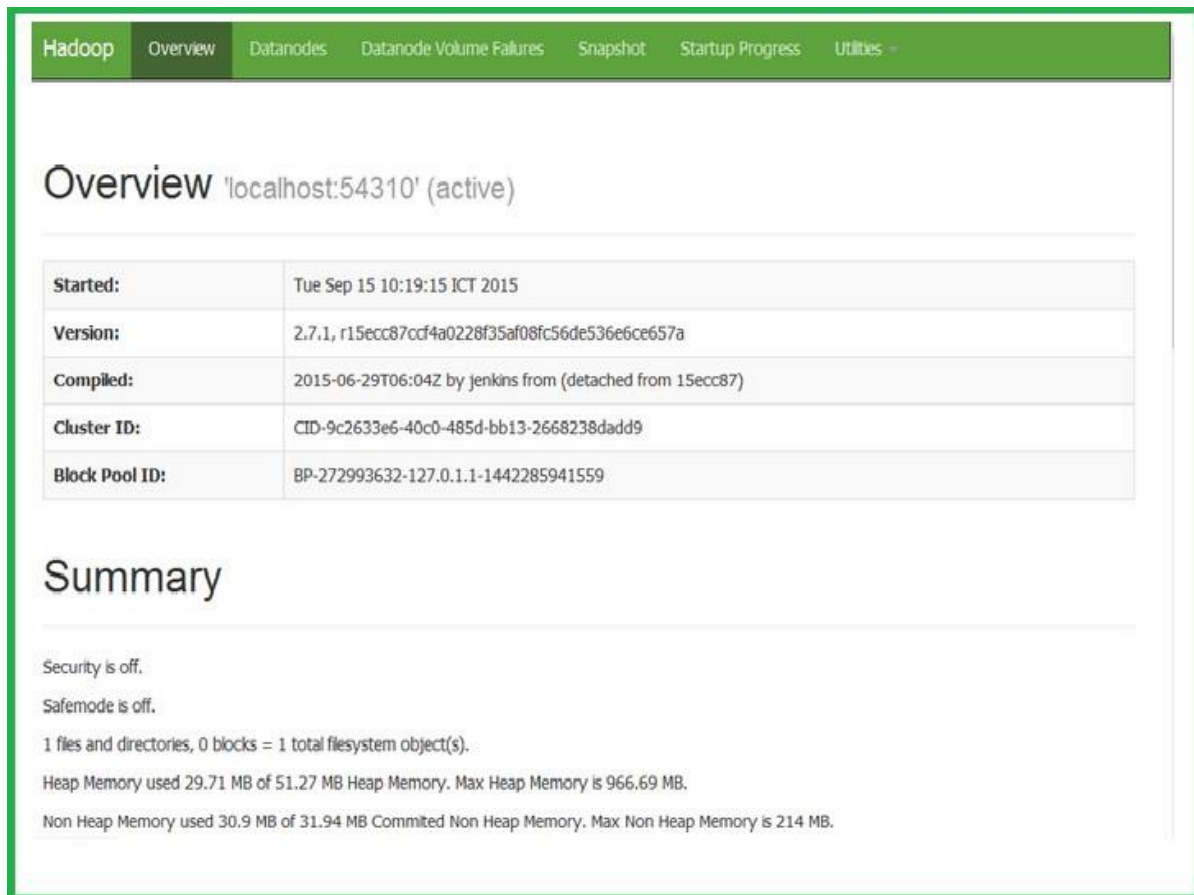
Ensure that namenode, datanode, and Resource manager are running

Step 9: Open: <http://localhost:8088>



Step 10:

Open: <http://localhost:50070>



The screenshot displays the Hadoop Overview page for a cluster named 'localhost:54310' (active). The page has a green header with navigation tabs: Hadoop, Overview, Datanodes, Datanode Volume Failures, Snapshot, Startup Progress, and Utilities. The Overview tab is selected. Below the header, the title 'Overview 'localhost:54310' (active)' is shown. A table provides key cluster information:

Started:	Tue Sep 15 10:19:15 ICT 2015
Version:	2.7.1, r15ecc87ccf4a0228f35af08fc56de536e6ce657a
Compiled:	2015-06-29T06:04Z by jenkins from (detached from 15ecc87)
Cluster ID:	CID-9c2633e6-40c0-485d-bb13-2668238dadd9
Block Pool ID:	BP-272993632-127.0.1.1-1442285941559

Below the table, a 'Summary' section provides additional details:

- Security is off.
- Safemode is off.
- 1 files and directories, 0 blocks = 1 total filesystem object(s).
- Heap Memory used 29.71 MB of 51.27 MB Heap Memory. Max Heap Memory is 966.69 MB.
- Non Heap Memory used 30.9 MB of 31.94 MB Committed Non Heap Memory. Max Non Heap Memory is 214 MB.

ANALYSIS:

Here, we have successfully installed the Hadoop.