# HINDUSTHAN COLLEGE OF ARTS & SCIENCE
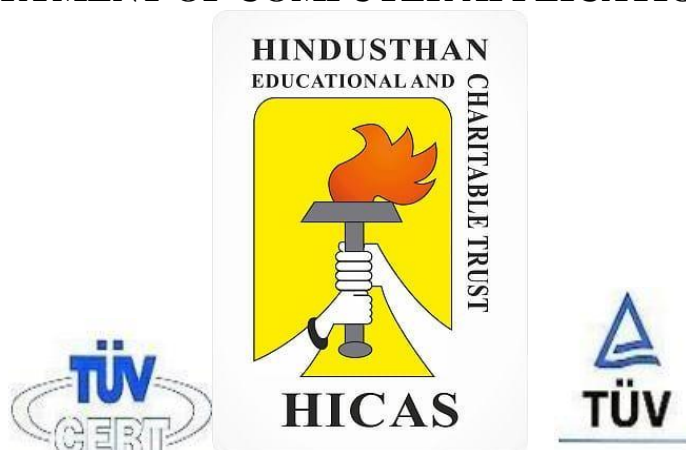
# (Autonomous)

An Autonomous Institution – Affiliated to Bharathiar University

(ISO 9001 – 2001 Certificate Instituation)

Behind Nava India, Coimbatore – 641028.

## DEPARTMENT OF COMPUTER APPLICATIONS (PG)

**MASTER OF COMPUTER APPLICATIONS**

**PRACTICAL RECORD**

**23MCP24 – PRACTICAL: BIG DATA ANALYTICS**

**NAME** : _____

**REGISTER NO** : _____

**CLASS** : _____

**SEMESTER** : _____

**YEAR** : _____

# HINDUSTHAN COLLEGE OF ARTS & SCIENCE

# (Autonomous)

An Autonomous Institution – Affiliated to Bharathiar University

(ISO 9001 – 2001 Certificate Instituation)

Behind Nava India, Coimbatore – 641028.

## DEPARTMENT OF COMPUTER APPLICATIONS (PG)

### <u>CERTIFICATE</u>

Certificate that this is a bonafide record of **Big Data Analytics (23MCP24)** done by

_____ Register No: _____ during the

academic year of 2024-2025.

**STAFF-IN CHARGE**                                                               **DIRECTOR**

Submitted for the Bharathiar University Practical Examination held on _____ at Hindusthan College of Arts & Science, Coimbatore – 641028.

**INTERNAL EXAMINER**                                             **EXTERNAL EXAMINER**

Date:

Place: Coimbatore

# CONTENTS

| PROGRAM NO: 01 DATE: | Installation of Hadoop | PAEG NO: |
|---|---|---|

**AIM:**

**SETTING UP AND INSTALLING HADOOP:**

**Prerequisites:**

**1. Install Java Development Kit (JDK):**

- Hadoop requires Java, so ensure that JDK is installed. Download it from Oracle or OpenJDK.
- Set the JAVA_HOME environment variable.

**Steps:**

- Download and install JDK 11 or later.
- Set JAVA_HOME:
  - Open Control Panel > System > Advanced system settings.
  - Click Environment Variables.
  - In the System variables section, click New and add:
    - **Variable Name:** JAVA_HOME
    - **Variable Value:** C:\Program Files\Java\jdk-11
  - Also, add Java to Path by editing the Path variable and adding %JAVA_HOME%\bin.

**2. Install WinRAR or 7-Zip:**

- To extract the Hadoop binary package you will download later.

**3. Install SSH (Optional for pseudo/fully distributed mode):**

- You will need an SSH client (such as PuTTY) for fully distributed setups if multiple machines are involved. However, it's optional for single-machine setups.
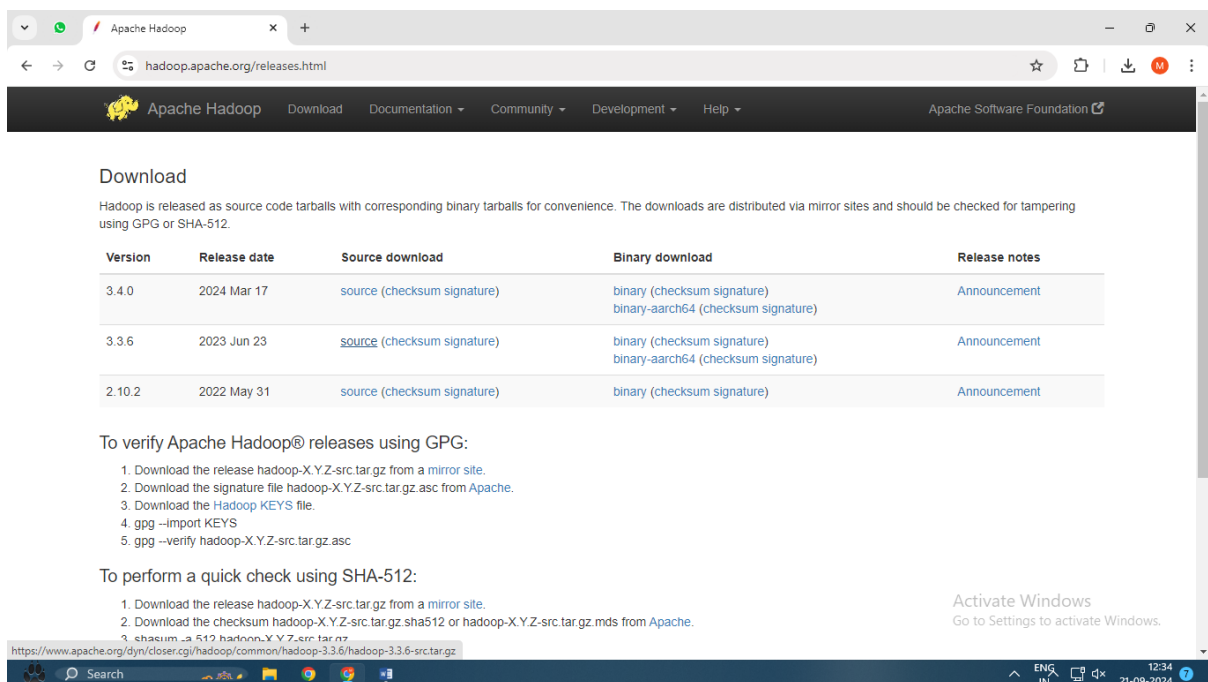
# 1. Standalone Mode

Standalone mode requires no Hadoop daemon services (such as NameNode or DataNode). It runs directly on the local filesystem.

**Steps:**

**1. Download Hadoop:**

- Download the binary release for Hadoop from Apache Hadoop Releases.
- Extract the .tar.gz file to a directory like C:\hadoop.

**2. Configure Environment Variables:** Add Hadoop to the system's PATH.

- o Open Control Panel > System > Advanced system settings > Environment Variables.
- o Add the following variables:
  - **HADOOP_HOME:** C:\hadoop
  - Edit the Path variable and add %HADOOP_HOME%\bin.

**3. Test Installation:** Open a new command prompt and run:

hadoop version

**4. Running a MapReduce Job:** You can run a sample MapReduce job in standalone mode:

hadoop jar %HADOOP_HOME%

\share\hadoop\mapreduce\hadoop-mapreduce-examples-*.jar wordcount input output

## 2. Pseudo-Distributed Mode

Pseudo-distributed mode runs all of Hadoop's daemons (NameNode, DataNode, ResourceManager, NodeManager) on a single machine, but simulates a distributed cluster.

**Steps:**

**1. Configure Hadoop:** You'll need to modify several XML configuration files located in the C:\hadoop\etc\hadoop directory.

**2. Configure core-site.xml:** Modify the file C:\hadoop\etc\hadoop\core-site.xml:

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

**3. Configure hdfs-site.xml:** Modify the file C:\hadoop\etc\hadoop\hdfs-site.xml:

```
<configuration>
```

```xml
  <property>

    <name>dfs.replication</name>

    <value>1</value> <!-- Since it's running on a single machine -->

  </property>

  <property>

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

    <value>file:/C:/hadoop_data/hdfs/namenode</value>

  </property>

  <property>

    <name>dfs.datanode.data.dir</name>

    <value>file:/C:/hadoop_data/hdfs/datanode</value>

  </property>

</configuration>
```

**4. Configure mapred-site.xml:** Modify the file C:\hadoop\etc\hadoop\mapred-site.xml (first copy it from the template):

cp C:\hadoop\etc\hadoop\mapred-site.xml.template C:\hadoop\etc\hadoop\mapred-site.xml

**Then, edit the following:**

```xml
<configuration>

  <property>

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

    <value>yarn</value>

  </property>

</configuration>
```

**5. Configure yarn-site.xml:** Modify the file C:\hadoop\etc\hadoop\yarn-site.xml:

```
<configuration>

  <property>

    <name>yarn.nodemanager.aux-services</name>

    <value>mapreduce_shuffle</value>

  </property>

</configuration>
```

**6. Format the NameNode:** Open a command prompt and run:

```
hdfs namenode -format
```

**7. Start Hadoop Daemons:** Run the following commands to start the Hadoop daemons:

```
start-dfs.cmd
```

```
start-yarn.cmd
```

**8. Test HDFS:** Verify that HDFS is running correctly:

```
hdfs dfs -mkdir /user
```

```
hdfs dfs -ls /
```

### 3. Fully Distributed Mode

Fully distributed mode is where Hadoop runs on multiple machines (master and worker nodes). For this, you'll need to configure Hadoop on each machine and ensure proper communication between them.

**Steps:**

**1. Master-Slave Setup:**

- Set up Master Node (NameNode) on one machine and Slave Nodes (DataNodes) on other machines.
- SSH setup for passwordless login between master and slave nodes may be required (for cross-machine communication).

**2. Configure core-site.xml on all machines:** On the master node and all slave nodes, configure **C:\hadoop\etc\hadoop\core-site.xml:**

```
<configuration>

  <property>

    <name>fs.defaultFS</name>

    <value>hdfs://master-node-ip:9000</value>

  </property>

</configuration>
```

**3. Configure hdfs-site.xml:** On all machines, configure **C:\hadoop\etc\hadoop\hdfs-site.xml:**

```
<configuration>

  <property>

    <name>dfs.replication</name>

    <value>3</value>

  </property>

</configuration>
```

**4. Configure workers file on Master Node:** In the file **C:\hadoop\etc\hadoop\workers,** list all the slave nodes:

```
slave-node1

slave-node2
```

**5. Set Up SSH and Communication:** Set up SSH for passwordless communication between the master and the slave nodes.

**6. Format the NameNode:** On the master node, run:

hdfs namenode -format

**7. Start Hadoop Daemons on Master and Slave Nodes:** On the master node, start the services:

start-dfs.cmd

start-yarn.cmd

**8. Verify Hadoop Cluster:** On the master node, check the cluster status:

hdfs dfsadmin -report

**RESULT:**

| PROGRAM NO: 02 DATE: | File Management in Hadoop | PAEG NO: |
|---|---|---|

**AIM:**

**ALGORITHM:**

**PROGRAM:**

**Start Hadoop Daemons:**

C:\Windows\system32>start-all

This script is Deprecated. Instead use start-dfs.cmd and start-yarn.cmd

starting yarn daemons



C:\Windows\system32>jps

2832 ResourceManager

10328 DataNode

12936 NodeManager

6504 NameNode

10012 Jps

**ENABLE HADOOP INTERFACE**

Open Browser and enable "**localhost:9870**"

## ENABLE HADOOP ALL APPLICATION

Open Browser and enable "**localhost:8088**"



## Adding files and directories:

### 1. Directory

C:\Windows\system32>hdfs dfs -mkdir /HADOOP

### 2. Files

C:\Windows\system32>hdfs dfs -put D:\example.txt /HADOOP

## Retrieving files:

C:\Windows\system32>hdfs dfs -get /HADOOP /example.txt D:\

get: `/example.txt': No such file or directory

## View Data from the file:

C:\Windows\system32>hdfs dfs -cat /HADOOP/example.txt

hi

hello everyone

how are you?

## Deleting files and directories:

### 1. File

C:\Windows\system32>hdfs dfs -rm -r /HADOOP/example.txt

Deleted /HADOOP/example.txt

### 2. Directory

C:\Windows\system32>hdfs dfs -rm -r /HADOOP

Deleted /HADOOP

**RESULT:**

| PROGRAM NO: 03<br>DATE: | **Word Count using MapReduce Paradigm** | **PAEG NO:** |
|---|---|---|

**AIM:**

**ALGORITHM:**

**PROGRAM:**

**example.txt**

hi

hello everyone

how are you?

**Command:**

C:\Windows\system32>hdfs dfs -mkdir /ex3

C:\Windows\system32>hdfs dfs -put D:\example.txt /ex3

C:\Windows\system32>hadoop jar C:\hadoop\share\hadoop\mapreduce\hadoop-mapreduce-examples-3.2.4.jar wordcount /ex3/example.txt /output

**OUTPUT:**



**RESULT:**

| PROGRAM NO: 04 DATE: | Weather Report using MapReduce | PAEG NO: |
| --- | --- | --- |

**AIM:**

**ALGORITHM:**

**PROGRAM:**

**weather_data.txt**

2024-08-01,Rainy

2024-08-02,Sunny

2024-08-03,Rainy

2024-08-04,Cloudy


weather_mapper.py

```python
import sys
for line in sys.stdin:
    line=line.strip()
    date, weather=line.split(',')
    print(f"{weather}\t1")
```


weather_reducere.py

```python
import sys
current_weather=None
current_count=0

for line in sys.stdin:
    line=line.strip()
    weather, count=line.split('\t',1)
    count=int(count)

    if current_weather == weather:
        count_count += count

    else:
        if current_weather:
            print(f"{current_weather}\t{current_count}")
        current_weather = weather
```

```
        current_count = count
```

```
if current_weather == weather:
    print(f"{current_weather}\t{current_count}")
```

C:\Windows\system32>hdfs dfs -mkdir /wc

C:\Windows\system32>hdfs dfs -put D:\weather_data.txt /wc

C:\Windows\system32>hdfs dfs -put D:\weather_mapper.py /wc

C:\Windows\system32>hdfs dfs -put D:\weather_reducer.py /wc

C:\Windows\system32>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar ^

More?   -input /wc/weather_data.txt ^

More?   -output /weather_output ^

More?   -mapper /weather_mapper.py ^

More?   -reducer /weather_reducer.py

**OUTPUT:**



**RESULT:**

| PROGRAM NO: 05 DATE: | Matrix Multiplication using MapReduce | PAEG NO: |
| --- | --- | --- |

**AIM:**

**ALGORITHM:**

**PROGRAM:**

**MatrixA.txt**

1,2,3

4,5,6

**MatrixB.txt**

7,8

9,10

11,12

**matrix_mapper.py**

```python
import sys
for line in sys.stdin:
    line = line.strip()
    elements = line.split()
    if elements[0] == "A":
        print(f"{elements[1]}\t{elements[2]}\t{elements[3]}\tA")
    else:
        print(f"{elements[2]}\t{elements[1]}\t{elements[3]}\tB")
```

**matrix_reducer.py**

```python
import sys
from collections import defaultdict


MatrixA = defaultdict(list)
MatrixB = defaultdict(list)s


for line in sys.stdin:
    line = line.strip()
    i, j, value, Matrix = line.split()
```

```python
    if Matrix == "A":

        MatrixA[int(i)].append((int(j), int(value)))

    else:

        MatrixB[int(j)].append((int(i), int(value)))


for i in MatrixA:

    for a_col, a_value in MatrixA[i]:

        for b_row, b_value in MatrixB[a_col]:

            print(f"{i},{b_row}\t{a_value * b_value}")
```

**Command:**

C:\Windows\system32>hdfs dfs -mkdir /ex5

C:\Windows\system32>hdfs dfs -put D:\MatrixA.txt /ex5

C:\Windows\system32>hdfs dfs -put D:\MatrixB.txt /ex5

C:\Windows\system32>hdfs dfs -put D:\matrix_mapper.py /ex5

C:\Windows\system32>hdfs dfs -put D:\matrix_reducer.py /ex5

C:\Windows\system32>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar ^

More?   -input /ex5/MatrixA.txt ^

More?   -input /ex5/MatrixB.txt ^

More?   -output /Matrix_output ^

More?   -mapper /matrix_mapper.py ^

More?   -reducer /matrix_reducer.py

**OUTPUT:**



**RESULT:**

| PROGRAM NO: 06<br>DATE: | Sales Data Report using MapReduce | PAEG NO: |
|---|---|---|

**AIM:**

**ALGORITHM:**

**PROGRAM:**

**sales_data.txt**

USA,100

India,200

USA,150

UK,50

**sales_mapper.py**

```python
import sys
for line in sys.stdin:
    line = line.strip()
    country, sales = line.split(',')
    print(f"{country}\t{sales}")
```

**sales_reducer.py**

```python
import sys
current_country = None
current_sales = 0

for line in sys.stdin:
    line = line.strip()
    country, sales = line.split('\t')
    sales = int(sales)

    if current_country == country:
        current_sales += sales
    else:
        if current_country:
            print(f"{current_country}\t{current_sales}")
        current_country = country
```

```
        current_sales = sales


if current_country == country:

    print(f"{current_country}\t{current_sales}")
```

C:\Windows\system32>hdfs dfs -mkdir /ex6

C:\Windows\system32>hdfs dfs -put D:\sales_data.txt /ex6

C:\Windows\system32>hdfs dfs -put D:\sales_mapper.py /ex6

C:\Windows\system32>hdfs dfs -put D:\sales_reducer.py /ex6

C:\Windows\system32>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar ^

More?   -input /ex6/sales_data.txt ^

More?   -output /sales_output ^

More?   -mapper /sales_mapper.py ^

More?   -reducer /sales_reducer.py

**OUTPUT:**



**RESULT:**

| PROGRAM NO: 07<br>DATE: | Electrical Consumption Report using<br>MapReduce | PAEG NO: |
|---|---|---|

**AIM:**

**ALGORITHM:**

**PROGRAM:**

**electricity.txt**

2020,5000

2020,6000

2021,5500

2021,7000

**electricity_mapper.py**

```
import sys
for line in sys.stdin:
    line = line.strip()
    year, consumption = line.split(',')
    print(f"{year}\t{consumption}")
```

**electricity_reducer.py**

```
import sys
current_year = None
max_consumption = 0

for line in sys.stdin:
    line = line.strip()
    year, consumption = line.split('\t')
    consumption = int(consumption)

    if current_year == year:
        if consumption > max_consumption:
            max_consumption = consumption
    else:
        if current_year:
            print(f"{current_year}\t{max_consumption}")
```

```python
        current_year = year
        max_consumption = consumption


if current_year == year:
    print(f"{current_year}\t{max_consumption}")
```

**Command:**

C:\Windows\system32>hdfs dfs -mkdir /ex7

C:\Windows\system32>hdfs dfs -put D:\electricity.txt /ex7

C:\Windows\system32>hdfs dfs -put D:\electricity_mapper.py /ex7

C:\Windows\system32>hdfs dfs -put D:\electricity_reducer.py /ex7

C:\Windows\system32>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar ^

More?   -input /ex7/electricity.txt ^

More?   -output /electricity_output ^

More?   -mapper /electricity_mapper.py ^

More?   -reducer /electricity_reducer.py

**OUTPUT:**



**RESULT:**

| PROGRAM NO: 08 DATE: | Real-time data Analysis using MapReduce | PAEG NO: |
|---|---|---|

**AIM:**

**ALGORITHM:**

**PROGRAM:**

**twitter_data.txt**

2024-08-01,#hadoop

2024-08-01,#bigdata

2024-08-02,#hadoop

2024-08-03,#ai

**twitter_mapper.py**

```python
import sys
for line in sys.stdin:
    line = line.strip()
    date, hashtag = line.split(',')
    print(f"{hashtag}\t1")
```

**twitter_reducer.py**

```python
import sys
current_hashtag = None
current_count = 0

for line in sys.stdin:
    line = line.strip()
    hashtag, count = line.split('\t')
    count = int(count)

    if current_hashtag == hashtag:
        current_count += count
    else:
        if current_hashtag:
            print(f"{current_hashtag}\t{current_count}")
        current_hashtag = hashtag
```

```
    current_count = count


if current_hashtag == hashtag:

    print(f"{current_hashtag}\t{current_count}")
```

**Command:**

C:\Windows\system32>hdfs dfs -mkdir /ex8

C:\Windows\system32>hdfs dfs -put D:\twitter_data.txt /ex8

C:\Windows\system32>hdfs dfs -put D:\twitter_mapper.py /ex8

C:\Windows\system32>hdfs dfs -put D:\twitter_reducer.py /ex8

C:\Windows\system32>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar ^

More?   -input /ex8/twitter_data.txt ^

More?   -output /twitter_output ^

More?   -mapper /twitter_mapper.py ^

More?   -reducer /twitter_reducer.py

**OUTPUT:**



**RESULT:**