# 1

# SET UP AND CONFIGURATION HADOOP USING CLOUDERA CREATING A HDFS SYSTEM WITH MINIMUM 1 NAME NODE AND 1 DATA NODES HDFS COMMANDS
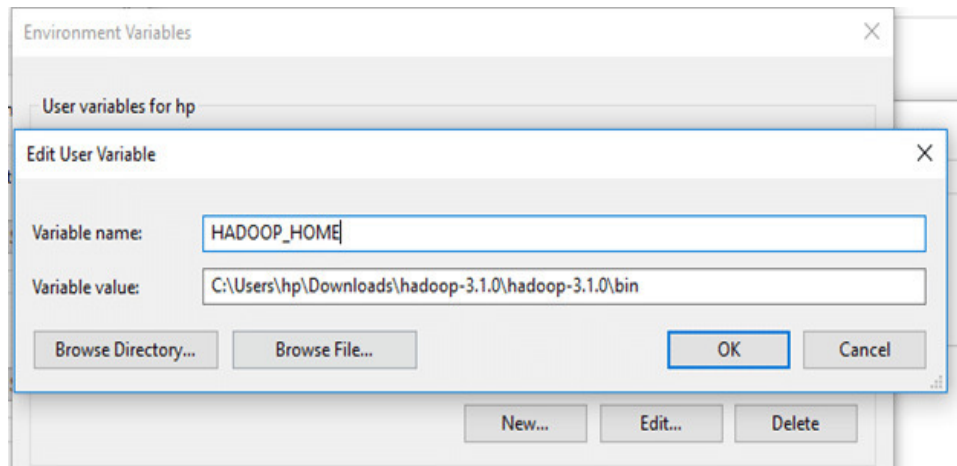
**Unit Structure :**

## 1.1 OBJECTIVES

Hadoop file system stores the data in multiple copies. Also, it's a cost-effective solution for any business to store their data efficiently. HDFS Operations acts as the key to open the vaults in which you store the data to be available from remote locations. This chapter describes how to set up and edit the deployment configuration files for HDFS

## 1.2 PREREQUISITE: TO INSTALL HADOOP, YOU SHOULD HAVE JAVA VERSION 1.8 IN YOUR SYSTEM.
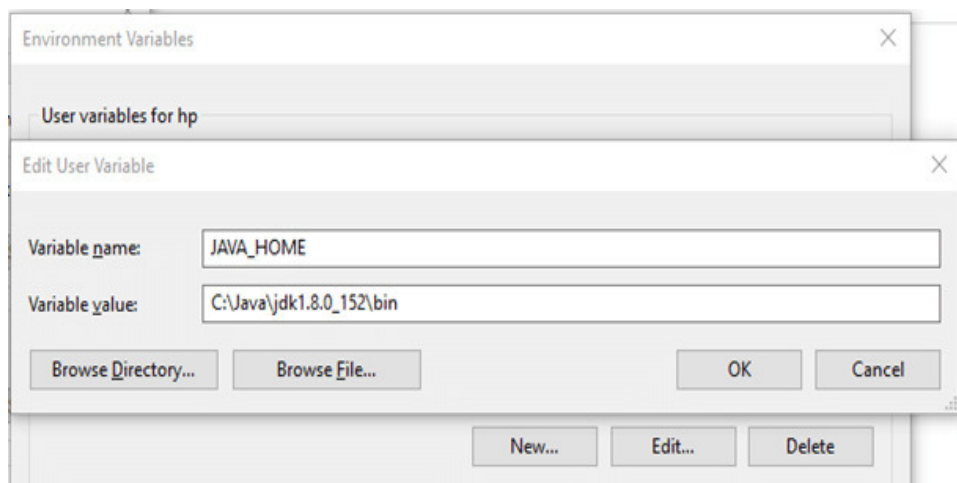
Check your java version through this command on command prompt

**Java -version**

Create a new user variable. Put the Variable_name as HADOOP_HOME and Variable_value as the path of the bin folder where you extracted hadoop.
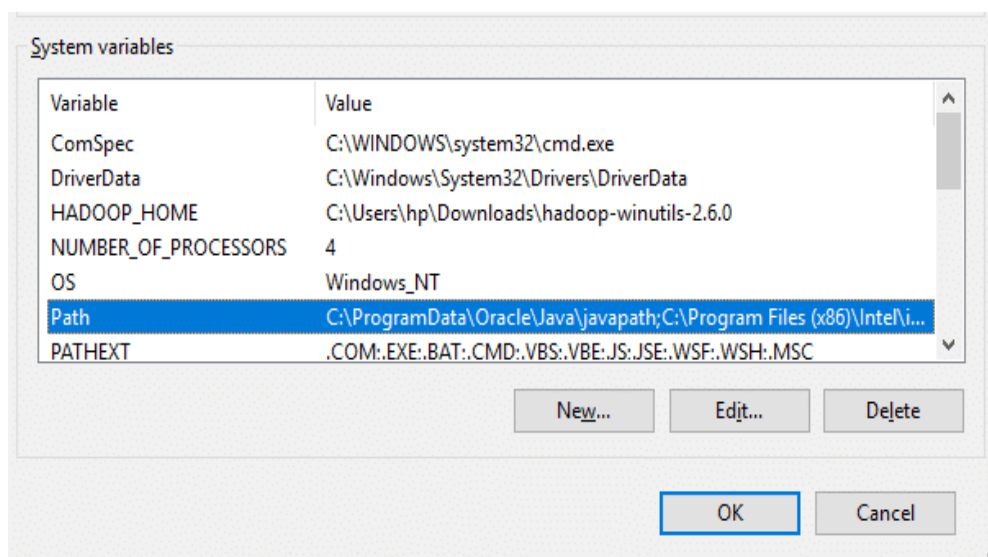
Likewise, create a new user variable with variable name as JAVA_HOME and variable value as the path of the bin folder in the Java directory.
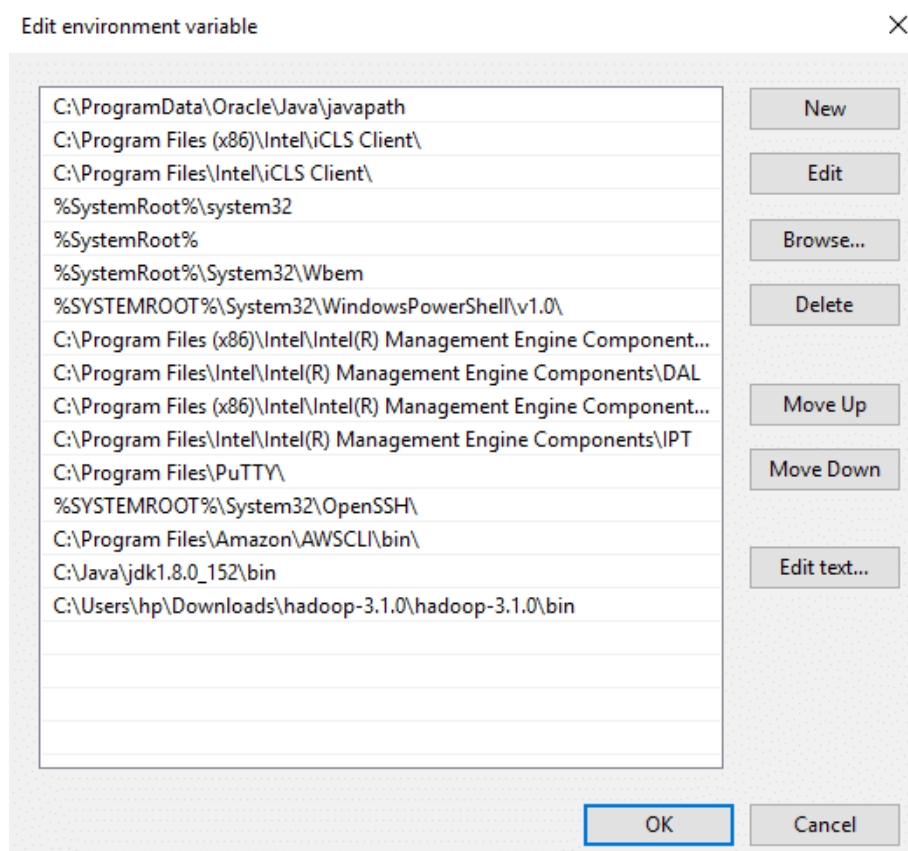


Now we need to set Hadoop bin directory and Java bin directory path in system variable path.

Edit Path in system variable

Click on New and add the bin directory path of Hadoop and Java in it.

Set up and Configuration
Hadoop using Cloudera
creating a HDFS System with
Minimum 1 Name Node
and 1 Data Nodes
HDFS Commands

## 1.3 GUI CONFIGURATIONS

Now we need to edit some files located in the hadoop directory of the etc folder where we installed hadoop. The files that need to be edited have been highlighted.

**1. Edit the file core-site.xml in the hadoop directory. Copy this xml property in the configuration in the file**

<configuration>

  <property>

    <name>fs.defaultFS</name>

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

  </property>

</configuration>

**2. Edit mapred-site.xml and copy this property in the configuration**

<configuration>

  <property>

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

    <value>yarn</value>

  </property>

</configuration>

**3. Create a folder 'data' in the hadoop directory**



**4. Create a folder with the name 'datanode' and a folder 'namenode' in this data directory**

Set up and Configuration
Hadoop using Cloudera
creating a HDFS System with
Minimum 1 Name Node
and 1 Data Nodes
HDFS Commands

**5.    Edit the file hdfs-site.xml and add below property in the configuration**

Note: The path of namenode and datanode across value would be the path of the datanode and namenode folders you just created.

\<configuration\>

  \<property\>

    \<name\>dfs.replication\</name\>

    \<value\>1\</value\>

  \</property\>

  \<property\>

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

    \<value\>C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\data\namenode\</value\>

  \</property\>

  \<property\>

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

    \<value\>                    C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\data\datanode\</value\>

  \</property\>

\</configuration\>

**6.    Edit the file yarn-site.xml and add below property in the 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\>

**7.    Edit hadoop-env.cmd and replace %JAVA_HOME% with the path of the java folder where your jdk 1.8 is installed**



**8.    Hadoop needs windows OS specific files which does not come with default download of hadoop.**

To include those files, replace the bin folder in hadoop directory with the bin folder provided in this github link.

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

Download it as zip file. Extract it and copy the bin folder in it. If you want to save the old bin folder, rename it like bin_old and paste the copied bin folder in that directory.



Check whether hadoop is successfully installed by running this command on cmd-

Set up and Configuration
Hadoop using Cloudera
creating a HDFS System with
Minimum 1 Name Node
and 1 Data Nodes
HDFS Commands

## hadoop –version

## Format the NameNode

Formatting the NameNode is done once when hadoop is installed and not for running hadoop filesystem, else it will delete all the data inside HDFS. Run this command-

## hdfs namenode –format

Now change the directory in cmd to sbin folder of hadoop directory with this command,

Start namenode and datanode with this command –

## start-dfs.cmd

Two more cmd windows will open for NameNode and DataNode

Now start yarn through this command-

## start-yarn.cmd

Note: Make sure all the 4 Apache Hadoop Distribution windows are up n running. If they are not running, you will see an error or a shutdown message. In that case, you need to debug the error.

To access information about resource manager current jobs, successful and failed jobs, go to this link in browser-

http://localhost:8088/cluster

To check the details about the hdfs (namenode and datanode),
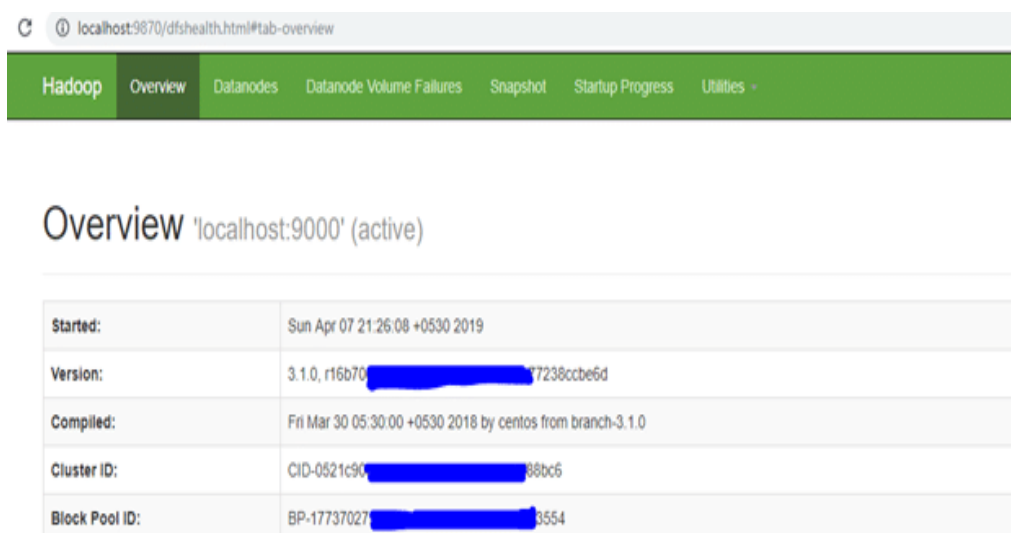
http://localhost:9870/

# Hadoop HDFS Commands

With the help of the HDFS commands, we can perform Hadoop HDFS file operations like changing the file permissions, viewing the file contents, creating files or directories, copying file/directory from the local file system to HDFS or vice-versa, etc.

Before starting with the HDFS command, we have to start the Hadoop services.

In this practical, we have mentioned the Hadoop HDFS commands with their usage, examples, and description.

## 1. version

**Hadoop HDFS version Command Usage:**

hadoop version

The Hadoop fs shell command **version** prints the Hadoop version.



```
dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop version
Hadoop 3.1.2
Source code repository https://github.com/apache/hadoop.git -r 1019dde65bcf12e05ef48ac71e84550d589e5d9a
Compiled by sunilg on 2019-01-29T01:39Z
Compiled with protoc 2.5.0
From source with checksum 64b8bdd4ca6e77cce75a93eb09ab2a9
This command was run using /home/dataflair/hadoop-3.1.2/share/hadoop/common/hadoop-common-3.1.2.jar
dataflair@admin1-All-Series:~$
```

## 2. mkdir

**Hadoop HDFS mkdir Command Usage:**

hadoop fs –mkdir /path/directory_name

we create a new directory named directory_name  in HDFS using the **mkdir** command.



```
dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -mkdir /newDataFlair
dataflair@admin1-All-Series:~$ hadoop fs -ls /
Found 3 items
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 10:38 /DataFlair
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 10:39 /dataflair
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 10:41 /newDataFlair
dataflair@admin1-All-Series:~$
```

## 3. ls

**Hadoop HDFS ls Command Usage:**

hadoop fs -ls /path

**Hadoop HDFS ls Command Description:**

The Hadoop fs shell command **ls** displays a list of the contents of a directory specified in the path provided by the user. It shows the name, permissions, owner, size, and modification date for each file or directories in the specified directory.

```
                                   dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -ls -R /
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 11:30 /DataFlair
-rw-r--r--   2 dataflair supergroup         56 2020-01-29 11:30 /DataFlair/copytest
-rw-r--r--   2 dataflair supergroup          0 2020-01-29 10:44 /DataFlair/file1
-rw-r--r--   2 dataflair supergroup         39 2020-01-29 10:52 /DataFlair/sample
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 14:42 /dataflair
-rw-r--r--   1 dataflair supergroup          0 2020-01-29 12:54 /dataflair/file1
-rw-r--r--   1 dataflair supergroup       3346 2020-01-29 14:51 /dataflair/test
dataflair@admin1-All-Series:~$
```

**4. put**

**Hadoop HDFS put Command Usage:**

haoop fs -put <localsrc> <dest>

**Hadoop HDFS put Command Example:**

Here in this example, we are trying to copy localfile1 of the local file system to the Hadoop filesystem.

```
                                   dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -put ~/localfile1 /filefromlocal
dataflair@admin1-All-Series:~$
```

**Hadoop HDFS put Command Description:**

The Hadoop fs shell command **put** is similar to the **copyFromLocal**, which copies files or directory from the local filesystem to the destination in the Hadoop filesystem.

**5. copyFromLocal**

**Hadoop HDFS copyFromLocal Command Usage:**

hadoop fs -copyFromLocal <localsrc> <hdfs destination>

**Hadoop HDFS copyFromLocal Command Example:**

Here in the below example, we are trying to copy the 'test1' file present in the local file system to the newDataFlair directory of Hadoop.

```
                                   dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -copyFromLocal ~/test1 /newDataFlair/copytest
dataflair@admin1-All-Series:~$
```

```
                                   dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -copyFromLocal ~/test1 /newDataFlair/copytest
dataflair@admin1-All-Series:~$ hadoop fs -cat /newDataFlair/copytest
Hello from DataFlair

Welcome to HDFS Command Tutuorial
dataflair@admin1-All-Series:~$
```

This command copies the file from the local file system to HDFS.
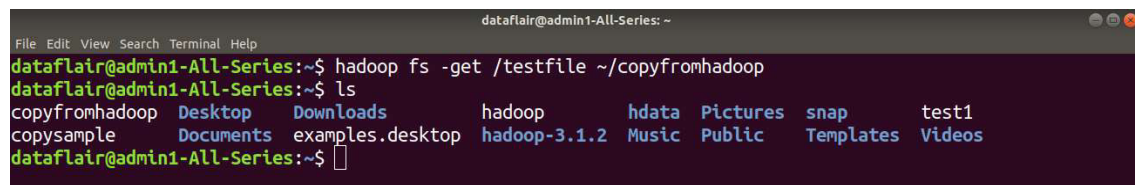
## 6. get

**Hadoop HDFS get Command Usage:**

hadoop fs -get <src> <localdest>

**Hadoop HDFS get Command Example:**

In this example, we are trying to copy the 'testfile' of the hadoop filesystem to the local file system.

**Hadoop HDFS get Command Description:**
The Hadoop fs shell command get copies the file or directory from the Hadoop file system to the local file system.

```
                                    dataflair@admin1-All-Series: ~                              ⊖ ⊟ ⊗
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -get /testfile ~/copyfromhadoop
dataflair@admin1-All-Series:~$ ls
copyfromhadoop  Desktop    Downloads      hadoop      hdata Pictures  snap       test1
copysample      Documents  examples.desktop hadoop-3.1.2  Music Public    Templates Videos
dataflair@admin1-All-Series:~$ 
```

## 7. copyToLocal
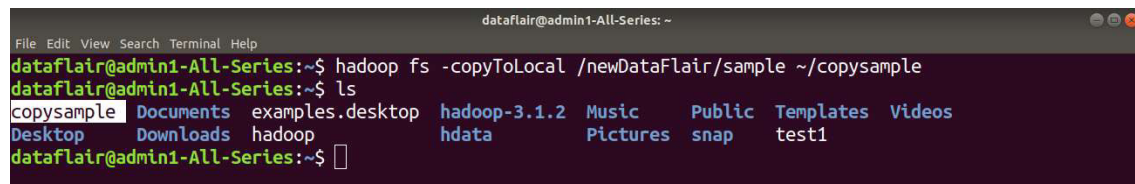
**Hadoop HDFS copyToLocal Command Usage:**

hadoop fs -copyToLocal <hdfs source> <localdst>

**Hadoop HDFS copyToLocal Command Example:**

Here in this example, we are trying to copy the 'sample' file present in the newDataFlair directory of HDFS to the local file system.

**adoop HDFS copyToLocal Description:**

**copyToLocal** command copies the file from HDFS to the local file system.

```
                                    dataflair@admin1-All-Series: ~                              ⊖ ⊟ ⊗
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -copyToLocal /newDataFlair/sample ~/copysample
dataflair@admin1-All-Series:~$ ls
copysample  Documents  examples.desktop  hadoop-3.1.2  Music    Public  Templates  Videos
Desktop     Downloads  hadoop            hdata         Pictures snap    test1
dataflair@admin1-All-Series:~$ 
```

## 8. cat

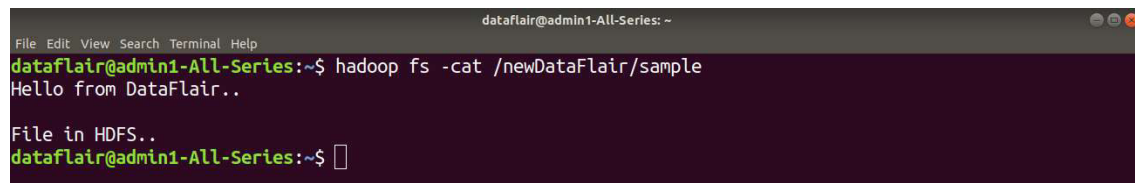**Hadoop HDFS cat Command Usage:**

hadoop fs –cat /path_to_file_in_hdfs

**Hadoop HDFS cat Command Example:**

Here in this example, we are using the cat command to display the content of the 'sample' file present in newDataFlair directory of HDFS.

**Hadoop HDFS cat Command Description:**

The **cat** command reads the file in HDFS and displays the content of the file on console or stdout.

```
                              dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -cat /newDataFlair/sample
Hello from DataFlair..

File in HDFS..
dataflair@admin1-All-Series:~$ 
```

**9. mv**

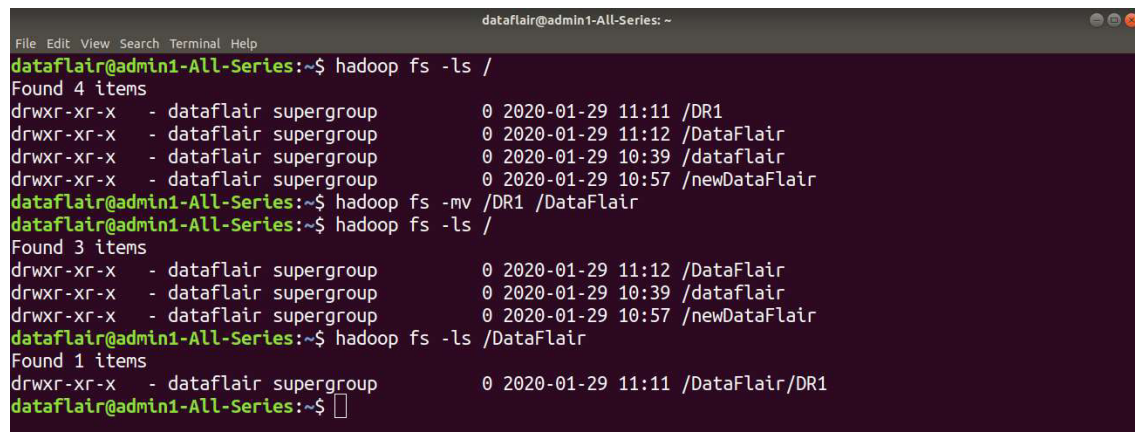**Hadoop HDFS mv Command Usage:**

hadoop fs -mv <src> <dest>

**Hadoop HDFS mv Command Example:**

In this example, we have a directory 'DR1' in HDFS. We are using **mv** command to move the DR1 directory to the DataFlair directory in HDFS.

**Hadoop HDFS mv Command Description:**

The HDFS mv command moves the files or directories from the source to a destination within **HDFS**.

```
                              dataflair@admin1-All-Series: ~
File Edit View Search Terminal Help
dataflair@admin1-All-Series:~$ hadoop fs -ls /
Found 4 items
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 11:11 /DR1
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 11:12 /DataFlair
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 10:39 /dataflair
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 10:57 /newDataFlair
dataflair@admin1-All-Series:~$ hadoop fs -mv /DR1 /DataFlair
dataflair@admin1-All-Series:~$ hadoop fs -ls /
Found 3 items
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 11:12 /DataFlair
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 10:39 /dataflair
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 10:57 /newDataFlair
dataflair@admin1-All-Series:~$ hadoop fs -ls /DataFlair
Found 1 items
drwxr-xr-x   - dataflair supergroup          0 2020-01-29 11:11 /DataFlair/DR1
dataflair@admin1-All-Series:~$ 
```

**10. cp**
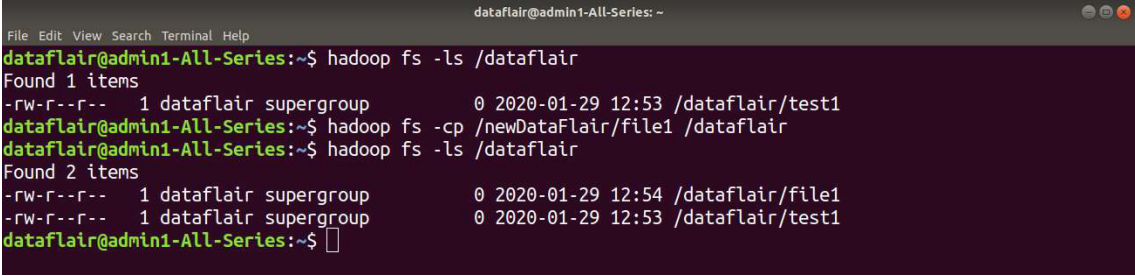
**Hadoop HDFS cp Command Usage:**

hadoop fs -cp <src> <dest>

**Hadoop HDFS cp Command Example:**

In the below example we are copying the 'file1' present in newDataFlair directory in HDFS to the dataflair directory of HDFS.

**Hadoop HDFS cp Command Description:**

The **cp** command copies a file from one directory to another directory within the HDFS.

```
                              dataflair@admin1-All-Series: ~
File  Edit  View  Search  Terminal  Help
dataflair@admin1-All-Series:~$ hadoop fs -ls /dataflair
Found 1 items
-rw-r--r--   1 dataflair supergroup          0 2020-01-29 12:53 /dataflair/test1
dataflair@admin1-All-Series:~$ hadoop fs -cp /newDataFlair/file1 /dataflair
dataflair@admin1-All-Series:~$ hadoop fs -ls /dataflair
Found 2 items
-rw-r--r--   1 dataflair supergroup          0 2020-01-29 12:54 /dataflair/file1
-rw-r--r--   1 dataflair supergroup          0 2020-01-29 12:53 /dataflair/test1
dataflair@admin1-All-Series:~$ 
```