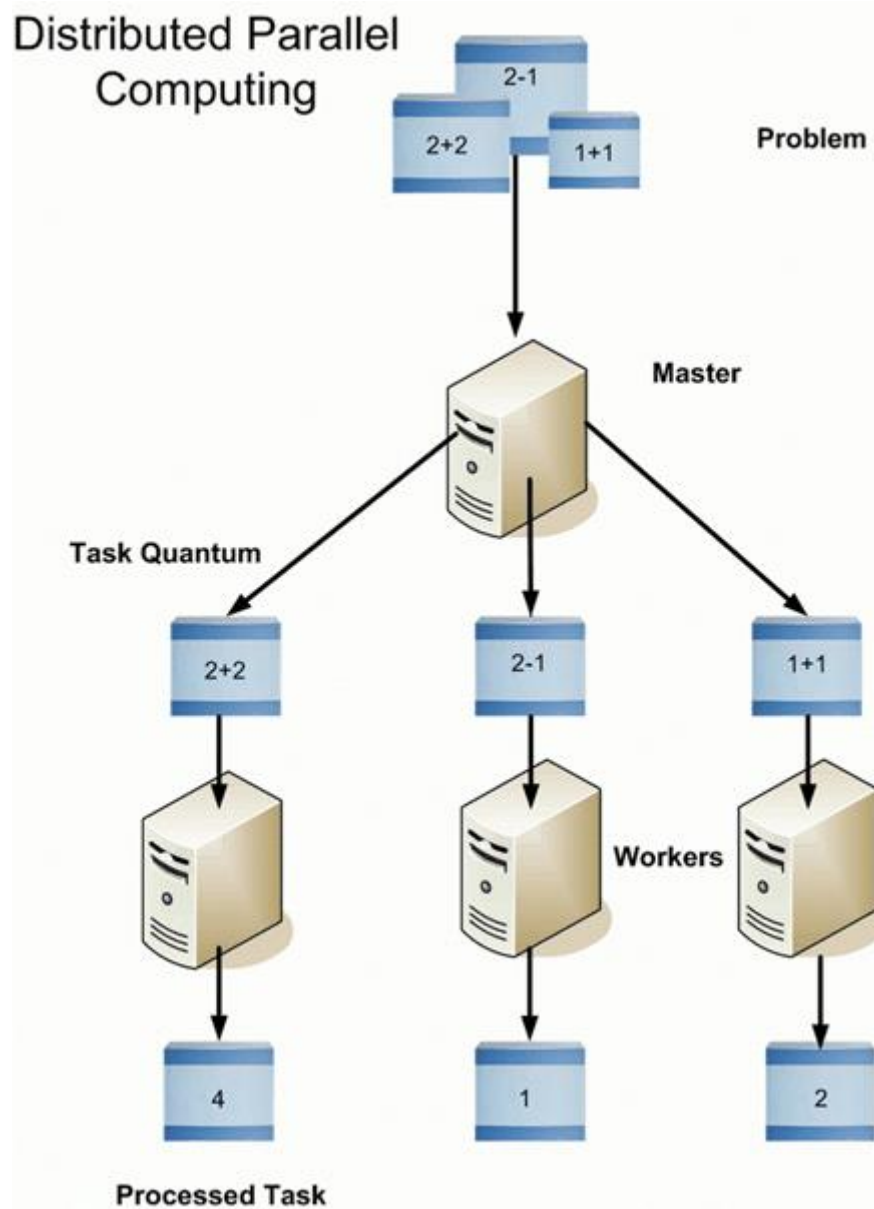


## Introduction

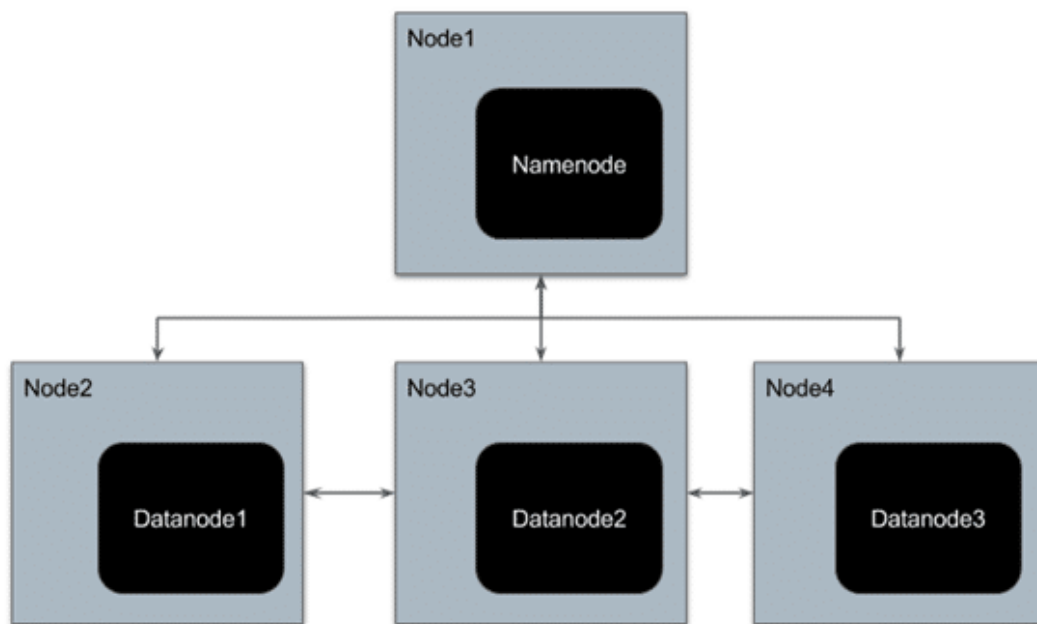
Hadoop is a software framework from Apache Software Foundation that is used to store and process Big Data. It has two main components; Hadoop Distributed File System (HDFS), its storage system and MapReduce, is its data processing framework. Hadoop has the capability to manage large datasets by distributing the dataset into smaller chunks across multiple machines and performing parallel computation on it .



## Overview of HDFS

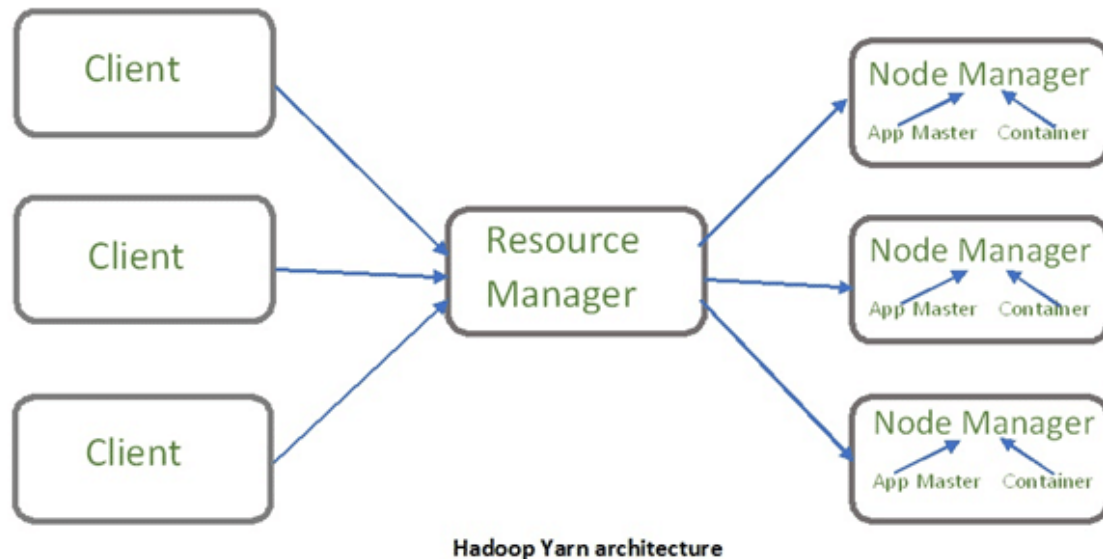
Hadoop is an essential component of the Big Data industry as it provides the most reliable storage layer, HDFS, which can scale massively. Companies like Yahoo and Facebook use HDFS to store their data.

HDFS has a master-slave architecture where the master node is called NameNode and slave node is called DataNode. The NameNode and its DataNodes form a cluster. NameNode acts like an instructor to DataNode while the DataNodes store the actual data.



source: Hasura

There is another component of Hadoop known as YARN. The idea of Yarn is to manage the resources and schedule/monitor jobs in Hadoop. Yarn has two main components, Resource Manager and Node Manager. The resource manager has the authority to allocate resources to various applications running in a cluster. The node manager is responsible for monitoring their resource usage (CPU, memory, disk) and reporting the same to the resource manager.



source: GeeksforGeeks

To understand the Hadoop architecture in detail, refer this blog –

## Advantages of Hadoop

- 1. Economical** – Hadoop is an open source Apache product, so it is free software. It has hardware cost associated with it. It is cost effective as it uses commodity hardware that are cheap machines to store its datasets and not any specialized machine.
- 2. Scalable** – Hadoop distributes large data sets across multiple machines of a cluster. New machines can be easily added to the nodes of a cluster and can scale to thousands of nodes storing thousands of terabytes of data.
- 3. Fault Tolerance** – Hadoop, by default, stores 3 replicas of data across the nodes of a cluster. So if any node goes down, data can be retrieved from other nodes.
- 4. Fast** – Since Hadoop processes distributed data parallelly, it can process large data sets much faster than the traditional systems. It is highly suitable for batch processing of data.

**5. Flexibility** – Hadoop can store structured, semi-structured as well as unstructured data. It can accept data in the form of textfile, images, CSV files, XML files, emails, etc

**6. Data Locality** – Traditionally, to process the data, the data was fetched from the location it is stored, to the location where the application is submitted; however, in Hadoop, the processing application goes to the location of data to perform computation. This reduces the delay in processing of data.

**7. Compatibility** – Most of the emerging big data tools can be easily integrated with Hadoop like Spark. They use Hadoop as a storage platform and work as its processing system.

## **Hadoop Deployment Methods**

**1. Standalone Mode** – It is the default mode of configuration of Hadoop. It doesn't use hdfs instead, it uses a local file system for both input and output. It is useful for debugging and testing.

**2. Pseudo-Distributed Mode** – It is also called a single node cluster where both NameNode and DataNode resides in the same machine. All the daemons run on the same machine in this mode. It produces a fully functioning cluster on a single machine.

**3. Fully Distributed Mode** – Hadoop runs on multiple nodes wherein there are separate nodes for master and slave daemons. The data is distributed among a cluster of machines providing a production environment.

## **Hadoop Installation on Windows 10**

As a beginner, you might feel reluctant in performing cloud computing which requires subscriptions. While you can install a virtual machine as well in your system, it requires allocation of a large amount of RAM for it to function smoothly else it would hang constantly.

You can install Hadoop in your system as well which would be a feasible way to learn Hadoop.

We will be installing single node pseudo-distributed hadoop cluster on windows 10.

**Prerequisite:** To install Hadoop, you should have Java version 1.8 in your system.

Check your java version through this command on command prompt

1 java -version

```
Command Prompt
Microsoft Windows [Version 10.0.17134.648]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\hp>java -version
java version "1.8.0_152"
Java(TM) SE Runtime Environment (build 1.8.0_152-b16)
Java HotSpot(TM) 64-Bit Server VM (build 25.152-b16, mixed mode)

C:\Users\hp>
```

If java is not installed in your system, then –


Go this link –

<https://www.oracle.com/technetwork/java/javase/downloads/jdk8-downl...>

Accept the license,

Java SE Development Kit 8u201 Demos and Samples Downloads		
You must accept the <a href="#">Oracle BSD License</a> . to download this software.		
<input type="radio"/> Accept License Agreement <input checked="" type="radio"/> Decline License Agreement		
Product / File Description	File Size	Download
Linux ARM 32 Hard Float ABI	9.05 MB	<a href="#">jdk-8u201-linux-arm32-vfp-hflt-demos.tar.gz</a>
Linux ARM 64 Hard Float ABI	9.06 MB	<a href="#">jdk-8u201-linux-arm64-vfp-hflt-demos.tar.gz</a>
Linux x86	56.13 MB	<a href="#">jdk-8u201-linux-i586-demos.rpm</a>
Linux x86	55.98 MB	<a href="#">jdk-8u201-linux-i586-demos.tar.gz</a>
Linux x64	56.23 MB	<a href="#">jdk-8u201-linux-x64-demos.rpm</a>
Linux x64	56.08 MB	<a href="#">jdk-8u201-linux-x64-demos.tar.gz</a>
Mac OS X	56.25 MB	<a href="#">jdk-8u201-macosx-x86_64-demos.zip</a>
Solaris SPARC 64-bit	12.2 MB	<a href="#">jdk-8u201-solaris-sparcv9-demos.tar.Z</a>
Solaris SPARC 64-bit	8.46 MB	<a href="#">jdk-8u201-solaris-sparcv9-demos.tar.gz</a>
Solaris x64	12.19 MB	<a href="#">jdk-8u201-solaris-x64-demos.tar.Z</a>
Solaris x64	8.42 MB	<a href="#">jdk-8u201-solaris-x64-demos.tar.gz</a>
Windows x86	56.96 MB	<a href="#">jdk-8u201-windows-i586-demos.zip</a>
Windows x64	56.98 MB	<a href="#">jdk-8u201-windows-x64-demos.zip</a>











Download the file according to your operating system. Keep the java folder directly under the local disk directory (C:\Java\jdk1.8.0\_152) rather than in Program Files (C:\Program Files\Java\jdk1.8.0\_152) as it can create errors afterwards.

This PC > Local Disk (C:) > Java		
Name	Date modified	Type
 jdk1.8.0_152	4/7/2019 2:54 PM	File folder

After downloading java version 1.8, download hadoop version 3.1 from this link –

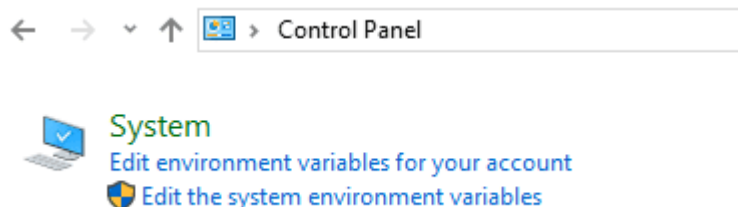
<https://archive.apache.org/dist/hadoop/common/hadoop-3.1.0/hadoop-3...>

Extract it to a folder.

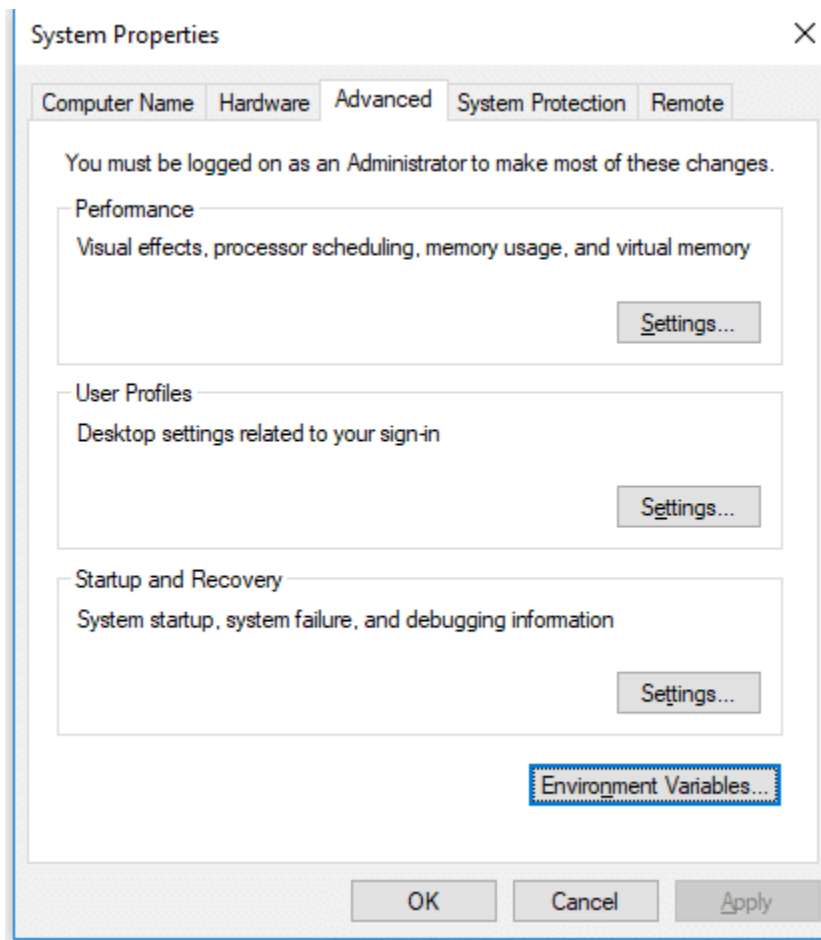
PC > Downloads > hadoop-3.1.0 > hadoop-3.1.0				
Name	Date modified	Type	Size	
 bin	4/7/2019 8:24 PM	File folder		
 etc	4/7/2019 8:24 PM	File folder		
 include	4/7/2019 8:24 PM	File folder		
 lib	4/7/2019 8:24 PM	File folder		
 libexec	4/7/2019 8:24 PM	File folder		
 sbin	4/7/2019 8:24 PM	File folder		
 share	4/7/2019 8:16 PM	File folder		
 LICENSE	3/21/2018 11:27 PM	Text Document	144 KB	
 NOTICE	3/21/2018 11:27 PM	Text Document	22 KB	
 README	3/21/2018 11:27 PM	Text Document	2 KB	

## Setup System Environment Variables

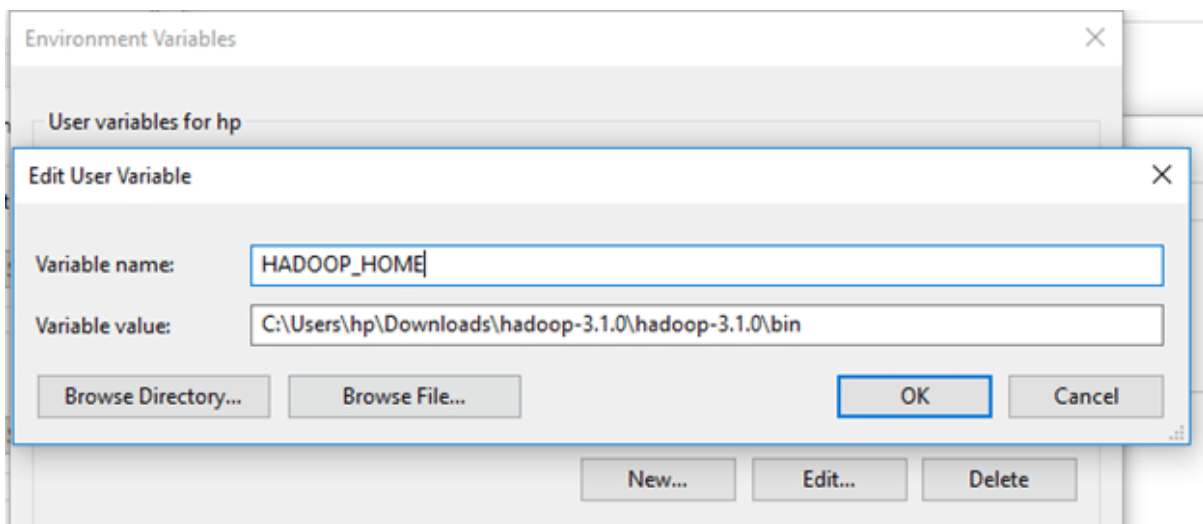
Open control panel to edit the system environment variable



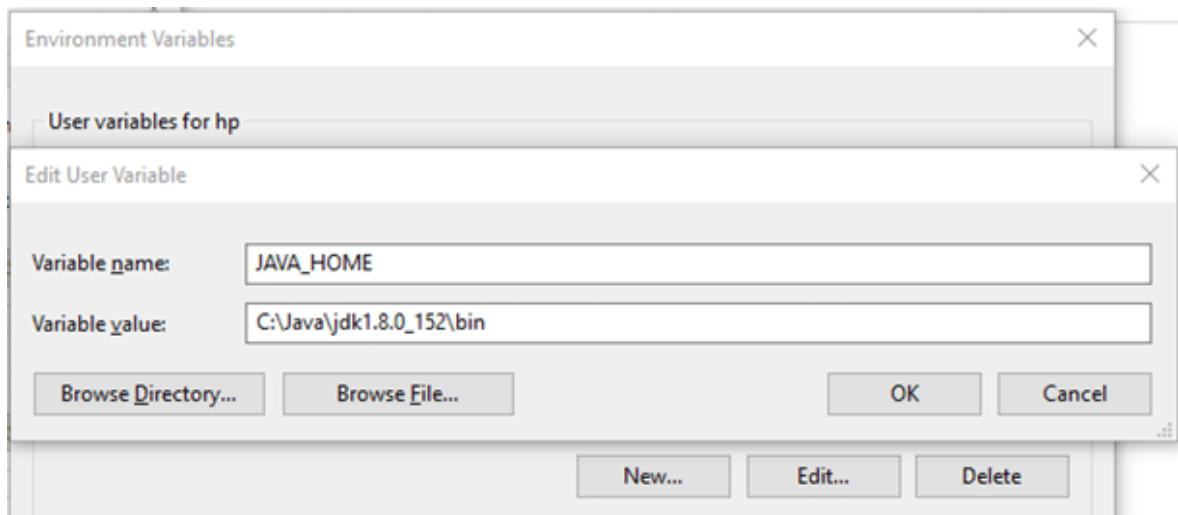
Go to environment variable in system properties



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.

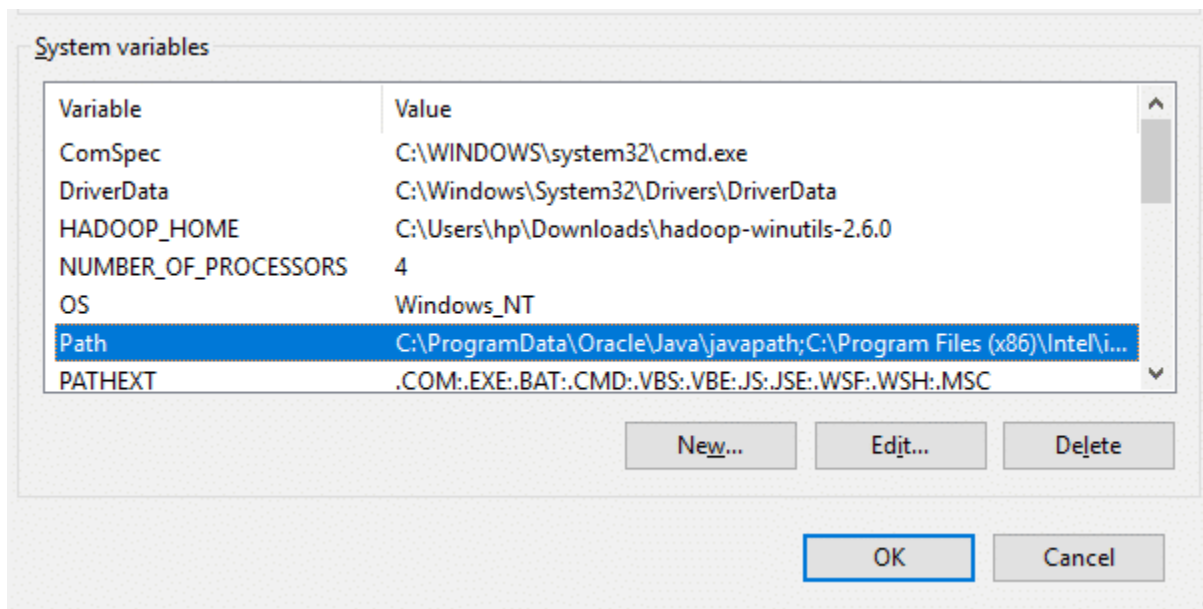


Likewise, create a new user variable with variable name as JAVA\_HOME and variable value as the path of the bin folder in the Java directory.



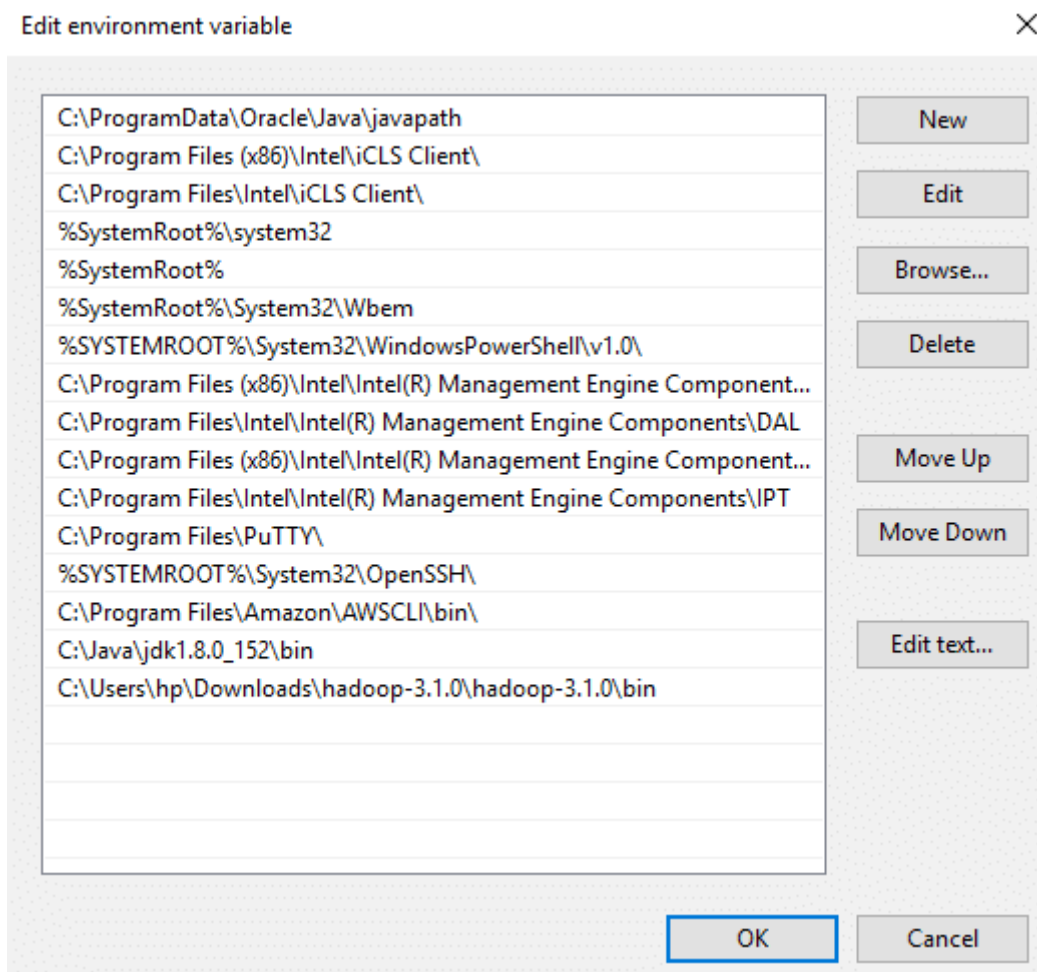
Now we need to set Hadoop bin directory and Java bin directory path in system variable path.

Edit Path in system variable



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


























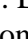






## Configurations

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.

PC > Downloads > hadoop-3.1.0 > hadoop-3.1.0 > etc > hadoop

Name	Date modified	Type	Size
 <b>core-site</b>	3/30/2018 5:31 AM	XML Document	1 KB
 <b>hadoop-env</b>	3/30/2018 5:31 AM	Windows Comma...	4 KB
 hadoop-env.sh	3/30/2018 5:52 AM	SH File	16 KB
 hadoop-metrics2.properties	3/30/2018 5:31 AM	PROPERTIES File	4 KB
 <b>hadoop-policy</b>	3/30/2018 5:31 AM	XML Document	11 KB
 hadoop-user-functions.sh.example	3/30/2018 5:31 AM	EXAMPLE File	4 KB
 <b>hdfs-site</b>	3/30/2018 5:33 AM	XML Document	1 KB
 https-env.sh	3/30/2018 5:33 AM	SH File	2 KB
 https-log4j.properties	3/30/2018 5:33 AM	PROPERTIES File	2 KB
 https-signature.secret	3/30/2018 5:33 AM	SECRET File	1 KB
 <b>https-site</b>	3/30/2018 5:33 AM	XML Document	1 KB
 <b>kms-acls</b>	3/30/2018 5:31 AM	XML Document	4 KB
 kms-env.sh	3/30/2018 5:31 AM	SH File	2 KB
 kms-log4j.properties	3/30/2018 5:31 AM	PROPERTIES File	2 KB
 <b>kms-site</b>	3/30/2018 5:31 AM	XML Document	1 KB
 log4j.properties	3/30/2018 5:31 AM	PROPERTIES File	14 KB
 <b>mapred-env</b>	3/30/2018 5:44 AM	Windows Comma...	1 KB
 mapred-env.sh	3/30/2018 5:44 AM	SH File	2 KB
 mapred-queues.xml.template	3/30/2018 5:44 AM	TEMPLATE File	5 KB
 <b>mapred-site</b>	3/30/2018 5:44 AM	XML Document	1 KB
 ssl-client.xml.example	3/30/2018 5:31 AM	EXAMPLE File	3 KB
 ssl-server.xml.example	3/30/2018 5:31 AM	EXAMPLE File	3 KB
 user_ec_policies.xml.template	3/30/2018 5:33 AM	TEMPLATE File	3 KB
 workers	3/30/2018 5:31 AM	File	1 KB
 <b>yarn-env</b>	3/30/2018 5:43 AM	Windows Comma...	3 KB
 yarn-env.sh	3/30/2018 5:43 AM	SH File	6 KB
 yarnservice-log4j.properties	3/30/2018 5:43 AM	PROPERTIES File	3 KB
 <b>yarn-site</b>	3/30/2018 5:43 AM	XML Document	1 KB

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

```
1 /span>configuration>
2 /span>property>
3 /span>name>fs.defaultFS/span>/name>
4 /span>value>hdfs://localhost:9000</value>
5 /span>/property>
6 /span>/configuration>
```

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

```
1 /span>configuration>
2 /span>property>
3 /span>name>mapreduce.framework.name/span>/name>
```

```
4 /span>value>yarn/span>/value>
5 /span>/property>
6 /span>/configuration>
```

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

PC > Downloads > hadoop-3.1.0 > hadoop-3.1.0

Name	Date modified	Type	Size
bin	4/7/2019 8:24 PM	File folder	
data	4/7/2019 8:34 PM	File folder	
etc	4/7/2019 8:24 PM	File folder	
include	4/7/2019 8:24 PM	File folder	
lib	4/7/2019 8:24 PM	File folder	
libexec	4/7/2019 8:24 PM	File folder	
sbin	4/7/2019 8:24 PM	File folder	
share	4/7/2019 8:16 PM	File folder	
LICENSE	3/21/2018 11:27 PM	Text Document	144 KB
NOTICE	3/21/2018 11:27 PM	Text Document	22 KB
README	3/21/2018 11:27 PM	Text Document	2 KB

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

PC > Downloads > hadoop-3.1.0 > hadoop-3.1.0 > data

Name	Date modified	Type
datanode	4/7/2019 8:35 PM	File folder
namenode	4/7/2019 8:35 PM	File folder

4. 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.

```
1 /span>configuration>
2 /span>property>
3 /span>name>dfs.replication/span>/name>
4 /span>value>1/span>/value>
```

```

5  /span>/property>
6  /span>property>
7    /span>name>dfs.namenode.name.dir/span>/name>
8    /span>value>C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-
9  3.1.0\data\namenode/span>/value>
10 /span>/property>
11 /span>property>
12   /span>name>dfs.datanode.data.dir/span>/name>
13   /span>value> C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-
14 3.1.0\data\datanode/span>/value>
    /span>/property>
    /span>/configuration>

```

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

```

1  /span>configuration>
2  /span>property>
3    /span>name>yarn.nodemanager.aux-services/span>/name>
4    /span>value>mapreduce_shuffle/span>/value>
5  /span>/property>
6  /span>property>
7    /span>name>yarn.nodemanager.auxservices.mapreduce.shuffle.class/span>/name>
8  /span>value>org.apache.hadoop.mapred.ShuffleHandler/span>/value>
9  /span>/property>
10 /span>/configuration>

```

6. Edit hadoop-env.cmd and replace %JAVA\_HOME% with the path of the java folder where your jdk 1.8 is installed

```
hadoop-env - Notepad
File Edit Format View Help

@rem Set Hadoop-specific environment variables here.

@rem The only required environment variable is JAVA_HOME. All others are
@rem optional. When running a distributed configuration it is best to
@rem set JAVA_HOME in this file, so that it is correctly defined on
@rem remote nodes.

@rem The java implementation to use. Required.
set JAVA_HOME=C:\Java\jdk1.8.0_152

@rem The jsvc implementation to use. Jsvc is required to run secure datanodes.
@rem set JSVC_HOME=%JSVC_HOME%

@rem set HADOOP_CONF_DIR=

@rem Extra Java CLASSPATH elements. Automatically insert capacity-scheduler.
if exist %HADOOP_HOME%\contrib\capacity-scheduler (
  if not defined HADOOP_CLASSPATH (
    set HADOOP_CLASSPATH=%HADOOP_HOME%\contrib\capacity-scheduler\*.jar
  ) else (
    set HADOOP_CLASSPATH=%HADOOP_CLASSPATH%;%HADOOP_HOME%\contrib\capacity-scheduler\*.jar
  )
)
```

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.

PC > Downloads > hadoop-3.1.0 > hadoop-3.1.0				
Name	Date modified	Type	Size	
bin	4/7/2019 8:40 PM	File folder		
bin_old	4/7/2019 8:24 PM	File folder		
data	4/7/2019 8:35 PM	File folder		
etc	4/7/2019 8:24 PM	File folder		
include	4/7/2019 8:24 PM	File folder		
lib	4/7/2019 8:24 PM	File folder		
libexec	4/7/2019 8:24 PM	File folder		
sbin	4/7/2019 8:24 PM	File folder		
share	4/7/2019 8:16 PM	File folder		
LICENSE	3/21/2018 11:27 PM	Text Document	144 KB	
NOTICE	3/21/2018 11:27 PM	Text Document	22 KB	
README	3/21/2018 11:27 PM	Text Document	2 KB	

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

1 hadoop version

```

Microsoft Windows [Version 10.0.17134.648]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\hp>hadoop version
Hadoop 3.1.0
Source code repository https://github.com/apache/hadoop -r 16b70619a24cdcf5d3b0fcf4b58ca77238ccbe6d
Compiled by centos on 2018-03-30T00:00Z
Compiled with protoc 2.5.0
From source with checksum 14182d20c972b3e2105580a1ad6990
This command was run using /C:/Users/hp/Downloads/hadoop-3.1.0/hadoop-3.1.0/share/hadoop/common/hadoop-common-3.1.0.jar
C:\Users\hp>

```

Since it doesn't throw error and successfully shows the hadoop version, that means hadoop is successfully installed in the system.

## 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-

1 hdfs namenode –format

It would appear something like this –

```
C:\Users\hp>hdfs namenode -format
2019-04-07 21:12:06,653 INFO namenode.NameNode: STARTUP_MSG:
/*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG:  host = DESKTOP-S9AL7B0/[REDACTED]
STARTUP_MSG:  args = [-format]
STARTUP_MSG:  version = 3.1.0
*****/
2019-04-07 21:12:08,941 INFO snapshot.SnapshotManager: SkipList is disabled
2019-04-07 21:12:08,941 INFO util.GSet: Computing capacity for map cachedBlocks
2019-04-07 21:12:08,941 INFO util.GSet: VM type = 64-bit
2019-04-07 21:12:08,941 INFO util.GSet: 0.25% max memory 889 MB = 2.2 MB
2019-04-07 21:12:08,941 INFO util.GSet: capacity = 2^18 = 262144 entries
2019-04-07 21:12:08,957 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.window.num.buckets = 10
2019-04-07 21:12:08,957 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.num.users = 10
2019-04-07 21:12:08,973 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.windows.minutes = 1,5,25
2019-04-07 21:12:08,973 INFO namenode.FSNamesystem: Retry cache on namenode is enabled
2019-04-07 21:12:08,973 INFO namenode.FSNamesystem: Retry cache will use 0.03 of total heap and retry cache entry exp
2019-04-07 21:12:08,973 INFO util.GSet: Computing capacity for map NameNodeRetryCache
2019-04-07 21:12:08,973 INFO util.GSet: VM type = 64-bit
2019-04-07 21:12:08,988 INFO util.GSet: 0.029999999329447746% max memory 889 MB = 273.1 KB
2019-04-07 21:12:08,988 INFO util.GSet: capacity = 2^15 = 32768 entries
2019-04-07 21:12:13,586 INFO namenode.FSImage: Allocated new BlockPoolId: BP-1773702794-192.168.56.1-1554651733554
2019-04-07 21:12:13,696 INFO common.Storage: Storage directory C:\hadoop-3.1.0\data\namenode has been successfully fo
2019-04-07 21:12:13,727 INFO namenode.FSImageFormatProtobuf: Saving image file C:\hadoop-3.1.0\data\namenode\current\
on
2019-04-07 21:12:13,887 INFO namenode.FSImageFormatProtobuf: Image file C:\hadoop-3.1.0\data\namenode\current\fsimage
n 0 seconds .
2019-04-07 21:12:14,046 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
2019-04-07 21:12:14,062 INFO namenode.NameNode: SHUTDOWN_MSG:
/*****
SHUTDOWN_MSG: Shutting down NameNode at DESKTOP-S9AL7B0/[REDACTED]
*****/
C:\Users\hp>
```

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

(Note: Make sure you are writing the path as per your system)

1 cd C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin

Start namenode and datanode with this command –

1 start-dfs.cmd

```
C:\Users\hp>cd C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>start-dfs.cmd
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>
```

Two more cmd windows will open for NameNode and DataNode

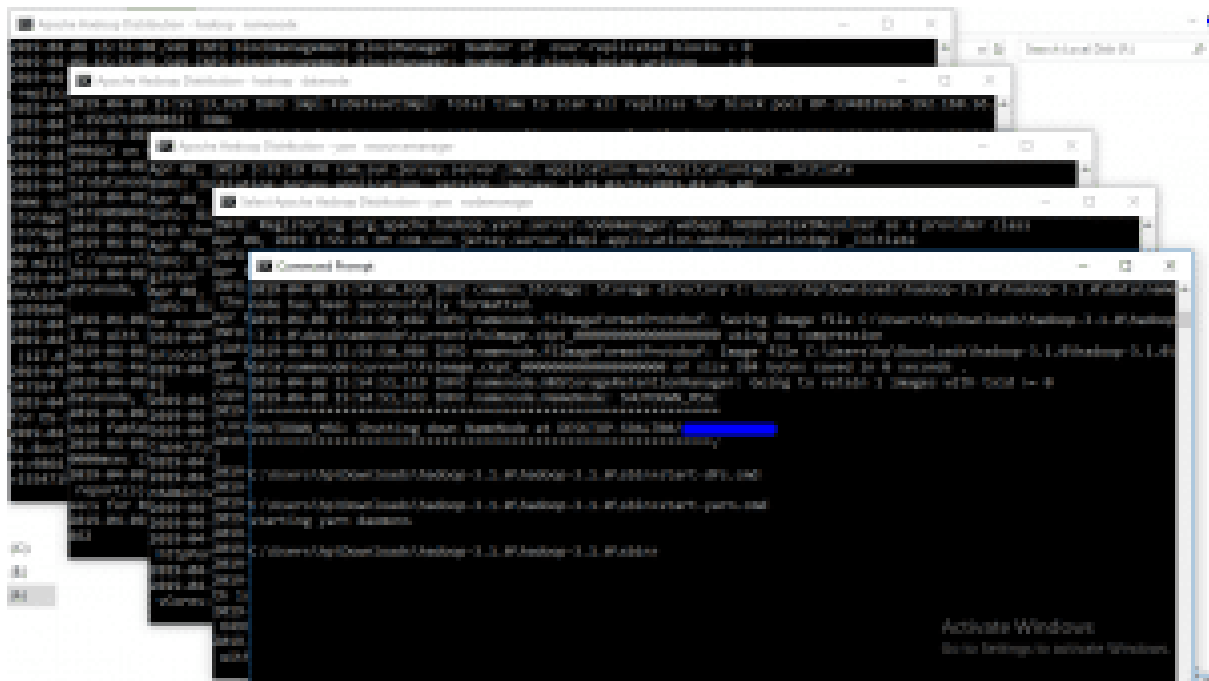
Now start yarn through this command-

1 start-yarn.cmd

```
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>start-yarn.cmd
starting yarn daemons

C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>
```

Two more windows will open, one for yarn resource manager and one for yarn node manager.



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>



localhost:8088/cluster

## All Applications

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	Vcores Used
0	0	0	0	0	0 B	8 GB	0 B	0

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy Nodes	Rebooted Nodes
1	0	0	0	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation	Maximum Capacity
Capacity Scheduler	[memory-mb (unit=Mi), vcores]	<memory:1024, vCores:1>	<memory:8192, vCores:4>	0

Show 20 entries

ID	User	Name	Application Type	Queue	Application Priority	StartTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU Vcores	Allocated Memory MB	Reserved CPU Vcores	Reserved Memory MB	% of Queue	% of Cluster
No data available in table																

Showing 0 to 0 of 0 entries

To check the details about the hdfs (namenode and datanode),

Open this link on browser-

<http://localhost:9870/>

Note: If you are using Hadoop version prior to 3.0.0 – Alpha 1, then use port <http://localhost:50070/>

localhost:9870/dfshealth.html#tab-overview

## Hadoop Overview

Overview 'localhost:9000' (active)

Started:	Sun Apr 07 21:26:08 +0530 2019
Version:	3.1.0, r16b70 [redacted] 7238ccbe6d
Compiled:	Fri Mar 30 05:30:00 +0530 2018 by centos from branch-3.1.0
Cluster ID:	CID-0521c90 [redacted] 88bc6
Block Pool ID:	BP-17737027 [redacted] 3554

## Summary

Security is off.  
Safemode is off.

1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).

## Working with HDFS

I will be using a small text file in my local file system. To put it in hdfs using hdfs command line tool.

I will create a directory named 'sample' in my hadoop directory using the following command-

```
1 hdfs dfs -mkdir /sample
```

```
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>hdfs dfs -mkdir /sample
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>
```

To verify if the directory is created in hdfs, we will use 'ls' command which will list the files present in hdfs –

```
1 hdfs dfs -ls /
```

```
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>hdfs dfs -ls /
Found 1 items
drwxr-xr-x   - hp supergroup          0 2019-04-07 23:39 /sample
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>
```

Then I will copy a text file named 'potatoes' from my local file system to this folder that I just created in hdfs using copyFromLocal command-

```
1 hdfs dfs -copyFromLocal C:\Users\hp\Downloads\potatoes.txt /sample
```

```
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>hdfs dfs -copyFromLocal C:\Users\hp\Downloads\potatoes.txt /sample
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>
```

To verify if the file is copied to the folder, I will use 'ls' command by specifying the folder name which will read the list of files in that folder –

```
1 hdfs dfs -ls /sample
```

```
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>hdfs dfs -ls /sample
Found 1 items
-rw-r--r--  1 hp supergroup          3736 2019-04-07 23:39 /sample/potatoes.txt
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>
```

To view the contents of the file we copied, I will use cat command-

```
1 hdfs dfs -cat /sample/potatoes.txt
```

```
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>hdfs dfs -cat /sample/potatoes.txt
area    temp    size    storage method    texture    flavor    moistness
1       1       1       1      1      2.9      3.2      3.0
1       1       1       1      2      2.3      2.5      2.6
1       1       1       1      3      2.5      2.8      2.8
1       1       1       1      4      2.1      2.9      2.4
1       1       1       1      5      1.9      2.8      2.2
1       1       1       2      1      1.8      3.0      1.7
1       1       1       2      2      2.6      3.1      2.4
1       1       1       2      3      3.0      3.0      2.9
1       1       1       2      4      2.2      3.2      2.5
1       1       1       2      5      2.0      2.8      1.9
1       1       1       3      1      1.8      2.6      1.5
1       1       1       3      2      2.0      2.8      1.9
1       1       1       3      3      2.6      2.6      2.6
1       1       1       3      4      2.1      3.2      2.1
1       1       1       3      5      2.5      3.0      2.1
1       1       1       4      1      2.6      3.1      2.4
```

To Copy file from hdfs to local directory, I will use get command –

```
1 hdfs dfs -get /sample/potatoes.txt C:\Users\hp\Desktop\priyanka
```

```
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>hdfs dfs -get /sample/potatoes.txt C:\Users\hp\Desktop\priyanka
C:\Users\hp\Downloads\hadoop-3.1.0\hadoop-3.1.0\sbin>
```