## 2CS702 Big Data Analytics

## Lab-3 Task

## Submitted by: Labdhi Sheth 18BCE101

**Aim:** Setup single-node Hadoop cluster and apply HDFS commands on single-node Hadoop Cluster.

**Methodology followed:**
The installation process as shared in the drive was followed.

**Output:**
CMD run as administrator
1. Checking the java and Hadoop version

```
Microsoft Windows [Version 10.0.19043.1165]
(c) Microsoft Corporation. All rights reserved.

C:\WINDOWS\system32>cd ..

C:\Windows>cd ..

C:\>d:

D:\>cd D:\nirma\7th sem\Big Data Analytics\labwork

D:\nirma\7th sem\Big Data Analytics\labwork>java -version
java version "16.0.2" 2021-07-20
Java(TM) SE Runtime Environment (build 16.0.2+7-67)
Java HotSpot(TM) 64-Bit Server VM (build 16.0.2+7-67, mixed mode, sharing)

D:\nirma\7th sem\Big Data Analytics\labwork>Hadoop -version
Error: Could not find or load main class sem\Big

D:\nirma\7th sem\Big Data Analytics\labwork>hadoop namenode -format
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
2021-09-14 12:22:19,944 INFO namenode.NameNode: STARTUP_MSG:
/************************************************************
STARTUP_MSG: Starting NameNode
STARTUP_MSG:   host = LAPTOP-0Q3JVTL1/192.168.2.5
STARTUP_MSG:   args = [-format]
STARTUP_MSG:   version = 3.2.1
STARTUP_MSG:   classpath = D:\hadoop-3.2.1\etc\hadoop;D:\hadoop-3.2.1\share\hadoop\common;D:\hadoop-3.2.1\share\hadoop\common\lib\accessors-smart-1.2.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\animal-sniffer-an
notations-1.17.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\asm-5.0.4.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\audience-annotations-0.5.0.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\avro-1.7.7.jar;D:\hadoop-3.
2.1\share\hadoop\common\lib\checker-qual-2.5.2.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-beanutils-1.9.3.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-cli-1.2.jar;D:\hadoop-3.2.1\share\hadoop\com
mon\lib\commons-codec-1.11.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-collections-3.2.2.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-compress-1.18.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\comm
ons-configuration2-2.1.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-io-2.5.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-lang3-3.7.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-logging-1.1.3
.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-math3-3.1.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-net-3.6.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\commons-text-1.4.jar;D:\hadoop-3.2.1\share
\hadoop\common\lib\curator-client-2.13.0.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\curator-framework-2.13.0.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\curator-recipes-2.13.0.jar;D:\hadoop-3.2.1\share\hadoop\c
ommon\lib\dnsjava-2.1.7.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\error_prone_annotations-2.2.0.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\failureaccess-1.0.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\gson-2.
2.4.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\guava-27.0-jre.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\hadoop-annotations-3.2.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\hadoop-auth-3.2.1.jar;D:\hadoop-3.2
.1\share\hadoop\common\lib\htrace-core4-4.1.0-incubating.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\httpclient-4.5.6.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\httpcore-4.4.10.jar;D:\hadoop-3.2.1\share\hadoop\
common\lib\j2objc-annotations-1.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jackson-annotations-2.9.8.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jackson-core-2.9.8.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\
jackson-core-asl-1.9.13.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jackson-databind-2.9.8.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jackson-jaxrs-1.9.13.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jackson-map
per-asl-1.9.13.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jackson-xc-1.9.13.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\javax.servlet-api-3.1.0.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jaxb-api-2.2.11.jar;D:
\hadoop-3.2.1\share\hadoop\common\lib\jaxb-impl-2.2.3-1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jcip-annotations-1.0-1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jersey-core-1.19.jar;D:\hadoop-3.2.1\share\h
adoop\common\lib\jersey-json-1.19.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jersey-server-1.19.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jersey-servlet-1.19.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetti
on-1.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetty-http-9.3.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetty-io-9.3.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetty-security-9.3
.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetty-server-9.3.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetty-servlet-9.3.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jett
y-util-9.3.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetty-webapp-9.3.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jetty-xml-9.3.24.v20180605.jar;D:\hadoop-3.2.1\share\hadoop\common\li
b\jsch-0.1.54.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\json-smart-2.3.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jsp-api-2.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jsr305-3.0.0.jar;D:\hadoop-3.2.1\share
\hadoop\common\lib\jsr311-api-1.1.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\jul-to-slf4j-1.7.25.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\kerb-admin-1.0.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\kerb-c
lient-1.0.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\kerb-common-1.0.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\kerb-core-1.0.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\kerb-crypto-1.0.1.jar;D:\hadoop-3
.2.1\share\hadoop\common\lib\kerb-identity-1.0.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\kerb-server-1.0.1.jar;D:\hadoop-3.2.1\share\hadoop\common\lib\kerb-simplekdc-1.0.1.jar;D:\hadoop-3.2.1\share\hadoop\co
```

## 2. Starting dfs and yarn



```
2021-09-14 12:22:24,272 INFO namenode.FSDirectory: XATTR serial map: bits=24 maxEntries=16777215
2021-09-14 12:22:24,288 INFO util.GSet: Computing capacity for map INodeMap
2021-09-14 12:22:24,288 INFO util.GSet: VM type       = 64-bit
2021-09-14 12:22:24,289 INFO util.GSet: 1.0% max memory 889 MB = 8.9 MB
2021-09-14 12:22:24,290 INFO util.GSet: capacity      = 2^20 = 1048576 entries
2021-09-14 12:22:24,290 INFO namenode.FSDirectory: ACLs enabled? false
2021-09-14 12:22:24,291 INFO namenode.FSDirectory: POSIX ACL inheritance enabled? true
2021-09-14 12:22:24,291 INFO namenode.FSDirectory: XAttrs enabled? true
2021-09-14 12:22:24,292 INFO namenode.NameNode: Caching file names occurring more than 10 times
2021-09-14 12:22:24,300 INFO snapshot.SnapshotManager: Loaded config captureOpenFiles: false, skipCaptureAccessTimeOnlyChange: false, snapshotDiffAllowSnapRootDescendant: true, maxSnapshotLimit: 65536
2021-09-14 12:22:24,303 INFO snapshot.SnapshotManager: SkipList is disabled
2021-09-14 12:22:24,309 INFO util.GSet: Computing capacity for map cachedBlocks
2021-09-14 12:22:24,309 INFO util.GSet: VM type       = 64-bit
2021-09-14 12:22:24,310 INFO util.GSet: 0.25% max memory 889 MB = 2.2 MB
2021-09-14 12:22:24,310 INFO util.GSet: capacity      = 2^18 = 262144 entries
2021-09-14 12:22:24,322 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.window.num.buckets = 10
2021-09-14 12:22:24,322 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.num.users = 10
2021-09-14 12:22:24,323 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.windows.minutes = 1,5,25
2021-09-14 12:22:24,327 INFO namenode.FSNamesystem: Retry cache on namenode is enabled
2021-09-14 12:22:24,328 INFO namenode.FSNamesystem: Retry cache will use 0.03 of total heap and retry cache entry expiry time is 600000 millis
2021-09-14 12:22:24,338 INFO util.GSet: Computing capacity for map NameNodeRetryCache
2021-09-14 12:22:24,338 INFO util.GSet: VM type       = 64-bit
2021-09-14 12:22:24,339 INFO util.GSet: 0.029999999329447746% max memory 889 MB = 273.1 KB
2021-09-14 12:22:24,340 INFO util.GSet: capacity      = 2^15 = 32768 entries
2021-09-14 12:22:24,517 INFO namenode.FSImage: Allocated new BlockPoolId: BP-2023101307-192.168.2.5-1631602344469
2021-09-14 12:22:24,655 INFO common.Storage: Storage directory \tmp\hadoop-labdh\dfs\name has been successfully formatted.
2021-09-14 12:22:24,738 INFO namenode.FSImageFormatProtobuf: Saving image file \tmp\hadoop-labdh\dfs\name\current\fsimage.ckpt_0000000000000000000 using no compression
2021-09-14 12:22:24,997 INFO namenode.FSImageFormatProtobuf: Image file \tmp\hadoop-labdh\dfs\name\current\fsimage.ckpt_0000000000000000000 of size 400 bytes saved in 0 seconds .
2021-09-14 12:22:25,044 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
2021-09-14 12:22:25,054 INFO namenode.FSImage: FSImageSaver clean checkpoint: txid=0 when meet shutdown.
2021-09-14 12:22:25,055 INFO namenode.NameNode: SHUTDOWN_MSG:
/************************************************************
SHUTDOWN_MSG: Shutting down NameNode at LAPTOP-0Q3JVTL1/192.168.2.5
************************************************************/

D:\nirma\7th sem\Big Data Analytics\labwork>start-dfs.cmd

D:\nirma\7th sem\Big Data Analytics\labwork>start-yarn.cmd
starting yarn daemons

D:\nirma\7th sem\Big Data Analytics\labwork>jps
20068 NameNode
20644 Jps
22628 ResourceManager
2872 NodeManager
21532 DataNode

D:\nirma\7th sem\Big Data Analytics\labwork>
```

# 3. Checking the clusters and working on localhost

| | |
|---|---|
| **Configured Capacity:** | 465.76 GB |
| **Configured Remote Capacity:** | 0 B |
| **DFS Used:** | 148 B (0%) |
| **Non DFS Used:** | 96.76 GB |
| **DFS Remaining:** | 369 GB (79.23%) |
| **Block Pool Used:** | 148 B (0%) |
| **DataNodes usages% (Min/Median/Max/stdDev):** | 0.00% / 0.00% / 0.00% / 0.00% |
| **Live Nodes** | 1 (Decommissioned: 0, In Maintenance: 0) |
| **Dead Nodes** | 0 (Decommissioned: 0, In Maintenance: 0) |
| **Decommissioning Nodes** | 0 |
| **Entering Maintenance Nodes** | 0 |
| **Total Datanode Volume Failures** | 0 (0 B) |
| **Number of Under-Replicated Blocks** | 0 |
| **Number of Blocks Pending Deletion (including replicas)** | 0 |
| **Block Deletion Start Time** | Tue Sep 14 12:22:37 +0530 2021 |
| **Last Checkpoint Time** | Tue Sep 14 12:22:24 +0530 2021 |
| **Enabled Erasure Coding Policies** | RS-6-3-1024k |

---

Hadoop    Overview    Utilities ▾

## DataNode on 192.168.2.5:9866

| | |
|---|---|
| **Cluster ID:** | CID-82e3b47f-b10e-4055-b5c0-e06f5248bb97 |
| **Version:** | 3.2.1, rb3cbbb467e22ea829b3808f4b7b01d07e0bf3842 |

## Block Pools

| Namenode Address | Block Pool ID | Actor State | Last Heartbeat | Last Block Report | Last Block Report Size (Max Size) |
|---|---|---|---|---|---|
| localhost:9000 | BP-2023101307-192.168.2.5-1631602344469 | RUNNING | 1s | a few seconds | 0 B (64 MB) |

## Volume Information

| Directory | StorageType | Capacity Used | Capacity Left | Capacity Reserved | Reserved Space for Replicas | Blocks |
|---|---|---|---|---|---|---|
| D:\tmp\hadoop-labdh\dfs\data | DISK | 148 B | 369 GB | 0 B | 0 B | 0 |

# All Applications

**Cluster**
- About
- Nodes
- Node Labels
- Applications
  - NEW
  - NEW_SAVING
  - SUBMITTED
  - ACCEPTED
  - RUNNING
  - FINISHED
  - FAILED
  - KILLED
- Scheduler

**Tools**

### Cluster Metrics

| Apps Submitted | Apps Pending | Apps Running | Apps Completed | Containers Running | Memory Used | Memory Total | Memory Reserved | VCores Used | VCores Total |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 B | 8 GB | 0 B | 0 | 8 |

### 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 Cluster Appl |
|---|---|---|---|---|
| Capacity Scheduler | [memory-mb (unit=Mi), vcores] | <memory:1024, vCores:1> | <memory:8192, vCores:4> | 0 |

Show 20 entries

Search:

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

Showing 0 to 0 of 0 entries

First

# 4. Mkdir command

```
D:\nirma\7th sem\Big Data Analytics\labwork>start-dfs.cmd

D:\nirma\7th sem\Big Data Analytics\labwork>start-yarn.cmd
starting yarn daemons

D:\nirma\7th sem\Big Data Analytics\labwork>jps
20068 NameNode
20644 Jps
22628 ResourceManager
2872 NodeManager
21532 DataNode

D:\nirma\7th sem\Big Data Analytics\labwork>hdfs dfs -ls /

D:\nirma\7th sem\Big Data Analytics\labwork>hdfs dfs -mkdir /input

D:\nirma\7th sem\Big Data Analytics\labwork>hdfs dfs -put d:\hadoop-3.2.1\etc\hadoop\*.xml\input
put: `.': No such file or directory: `hdfs://localhost:9000/user/labdh'

D:\nirma\7th sem\Big Data Analytics\labwork>hdfs dfs -put d:\hadoop-3.2.1\etc\hadoop\*.xml /input
2021-09-14 12:42:19,903 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:20,586 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:20,674 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:20,745 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:20,813 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:20,914 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:21,027 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:21,141 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:42:21,220 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false

D:\nirma\7th sem\Big Data Analytics\labwork>hdfs dfs -ls /
Found 1 items
drwxr-xr-x   - labdh supergroup          0 2021-09-14 12:42 /input

D:\nirma\7th sem\Big Data Analytics\labwork>hadoop jar d:\hadoop-3.2.1\share\hadoop\mapreduce\hadoop-mapreduce-examples-3.2.1.jar grep /input /output 'dfs[a-z.]+'
2021-09-14 12:44:31,079 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
2021-09-14 12:44:32,166 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/labdh/.staging/job_1631602376065_0001
2021-09-14 12:44:32,348 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:44:32,541 INFO input.FileInputFormat: Total input files to process : 9
2021-09-14 12:44:32,783 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:44:32,860 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:44:32,882 INFO mapreduce.JobSubmitter: number of splits:9
2021-09-14 12:44:33,083 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2021-09-14 12:44:33,117 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1631602376065_0001
2021-09-14 12:44:33,117 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-09-14 12:44:33,330 INFO conf.Configuration: resource-types.xml not found
2021-09-14 12:44:33,330 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2021-09-14 12:44:33,786 INFO impl.YarnClientImpl: Submitted application application_1631602376065_0001
2021-09-14 12:44:33,862 INFO mapreduce.Job: The url to track the job: http://LAPTOP-0Q3JVTL1:8088/proxy/application_1631602376065_0001/
2021-09-14 12:44:33,863 INFO mapreduce.Job: Running job: job_1631602376065_0001
```

# 5. Grep command