# PART 2:

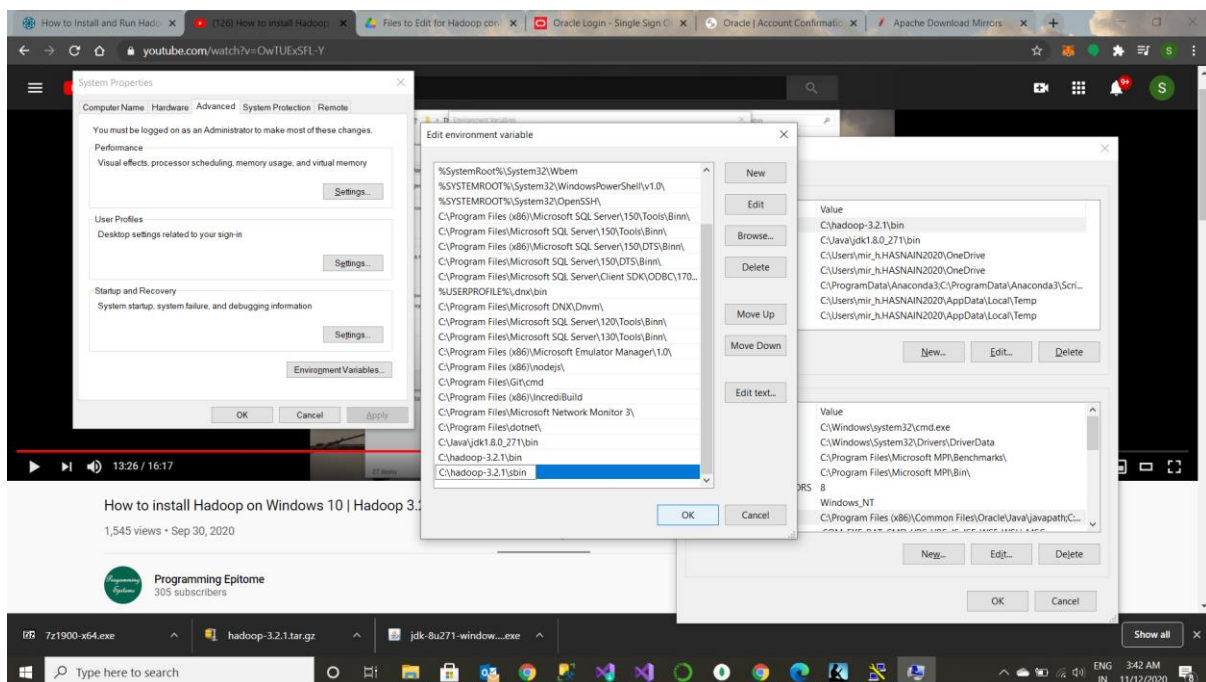## Getting Started with NoSql Systems on HDFS

**Fathima Syeda**

## Platform Setup:

Hadoop 3.2.1 is installed in this lab, but in order to do that we must have JAVA SDK-8 installed on our system , so we install that first.

Hadoop would be setup as a single node system in pseudo-distributed mode

We download Hadoop from- https://hadoop.apache.org/releases.html

Once Hadoop is installed , we must set the path and environment variables for it as:



After that 5 configuration files in the C:\hadoop-3.2.1\etc\hadoop folder, viz

core-site.xml ,mapred-site.xml, yarn -site.xml , hdfs-site.xml and hadoop-env.cmd are edited.

a) File C:/Hadoop-3.2.1/etc/hadoop/core-site.xml, paste below xml paragraph and save this file.

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

b) C:/Hadoop-3.2.1/etc/hadoop/mapred-site.xml, paste below xml paragraph and save this file.

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

```
 <value>yarn</value>
 </property>
</configuration>
```

c) Create folder "data" under "C:\Hadoop-3.2.1"

 1) Create folder "datanode" under "C:\Hadoop-3.2.1\data"

 2) Create folder "namenode" under "C:\Hadoop-3.2.1\data" data

d) Edit file C:\Hadoop-3.2.1/etc/hadoop/hdfs-site.xml, paste below xml paragraph and save this file.

```
<configuration>
 <property>
<name>dfs.replication</name>
 <value>1</value>
 </property>
 <property>
 <name>dfs.namenode.name.dir</name>
 <value>C:\hadoop-3.2.1\data\namenode</value>
 </property>
 <property>
 <name>dfs.datanode.data.dir</name>
 <value>C:\hadoop-3.2.1\data\datanode</value>
 </property>
</configuration>
```

e) Edit file C:/Hadoop-3.2.1/etc/hadoop/yarn-site.xml, paste below xml paragraph and save this file.

```
<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>
```
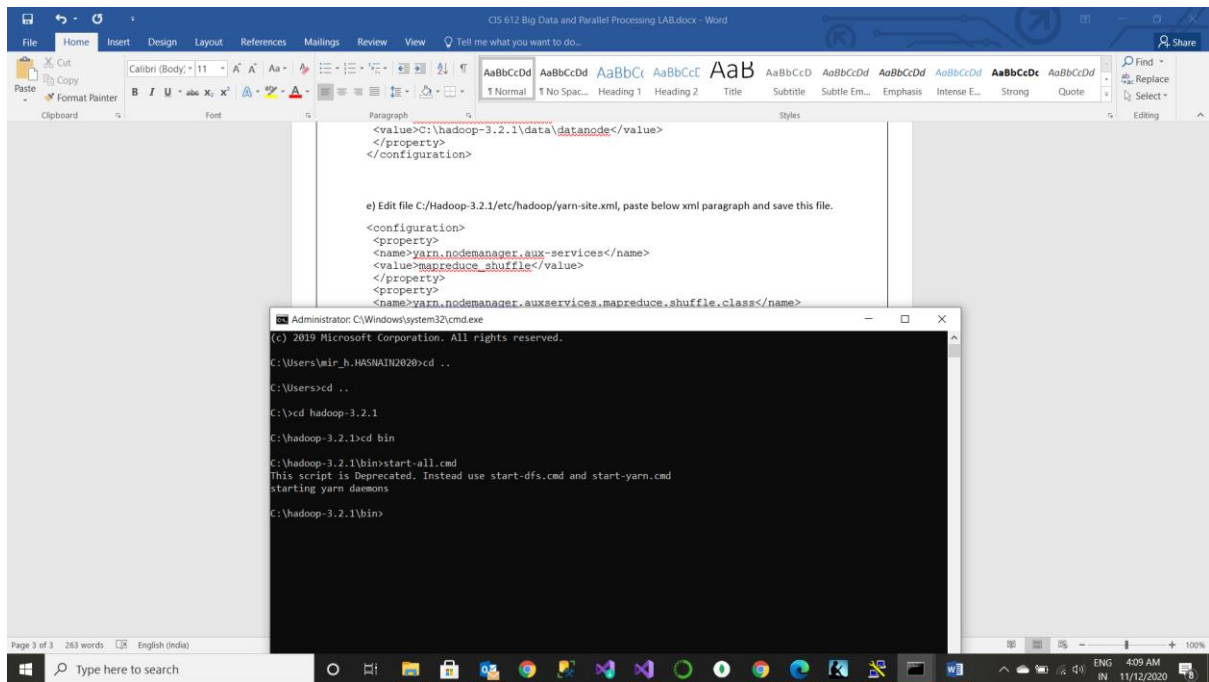
```
f) save the java path in the hadoop-env.cmd file as the path of the
java sdk's bin folder.
```

Once the configurations files have been edited and saved, the Hadoop configuration is completed successfully.
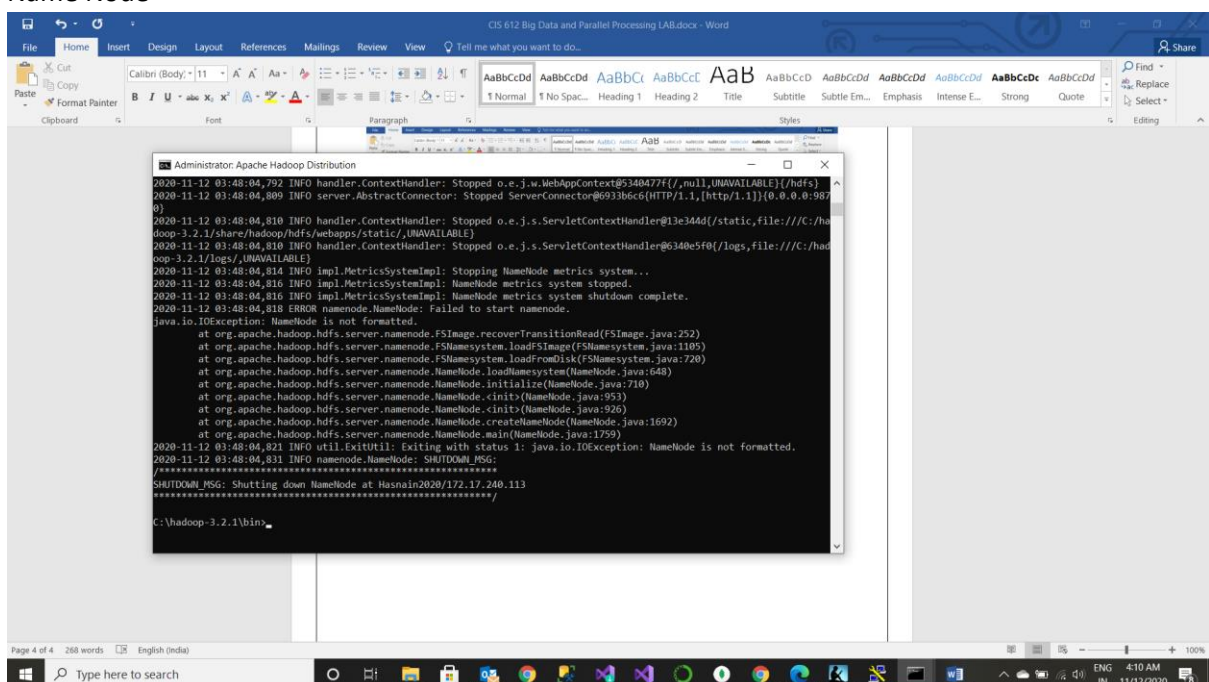
**Running the Hadoop  single node system**

Go to the cmd prompt and go to the Hadoop-3.2.1/bin folder and type start-all.cmd to start all the nodes in the Hadoop system.
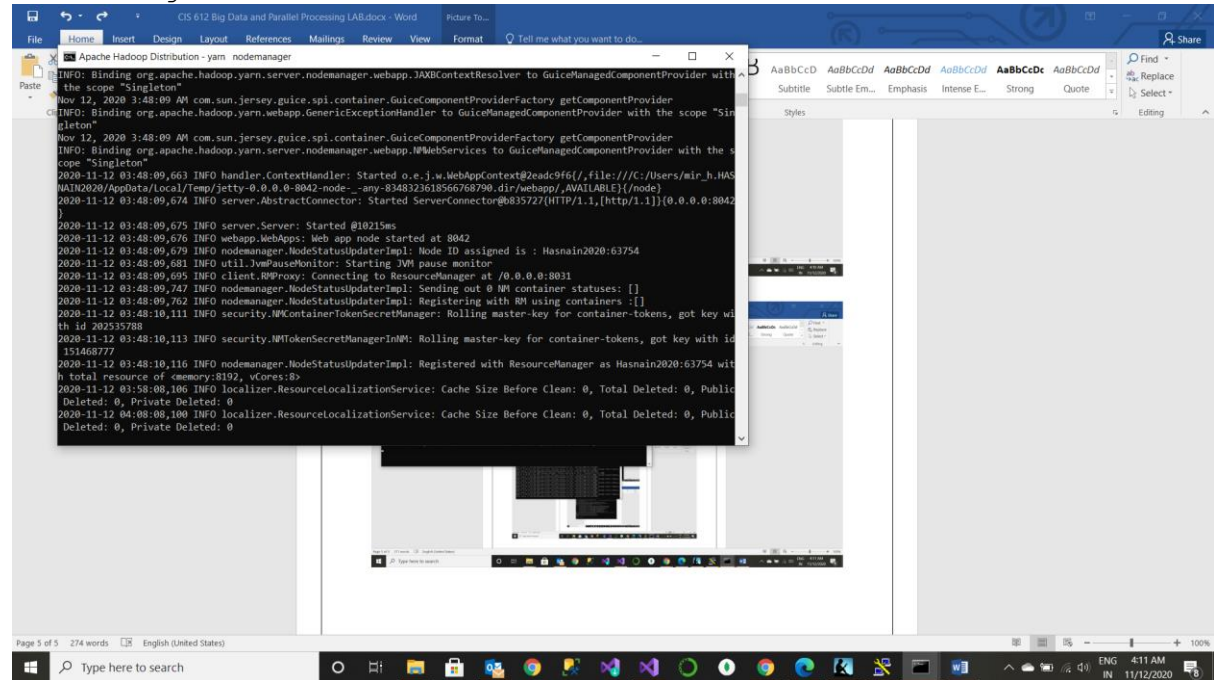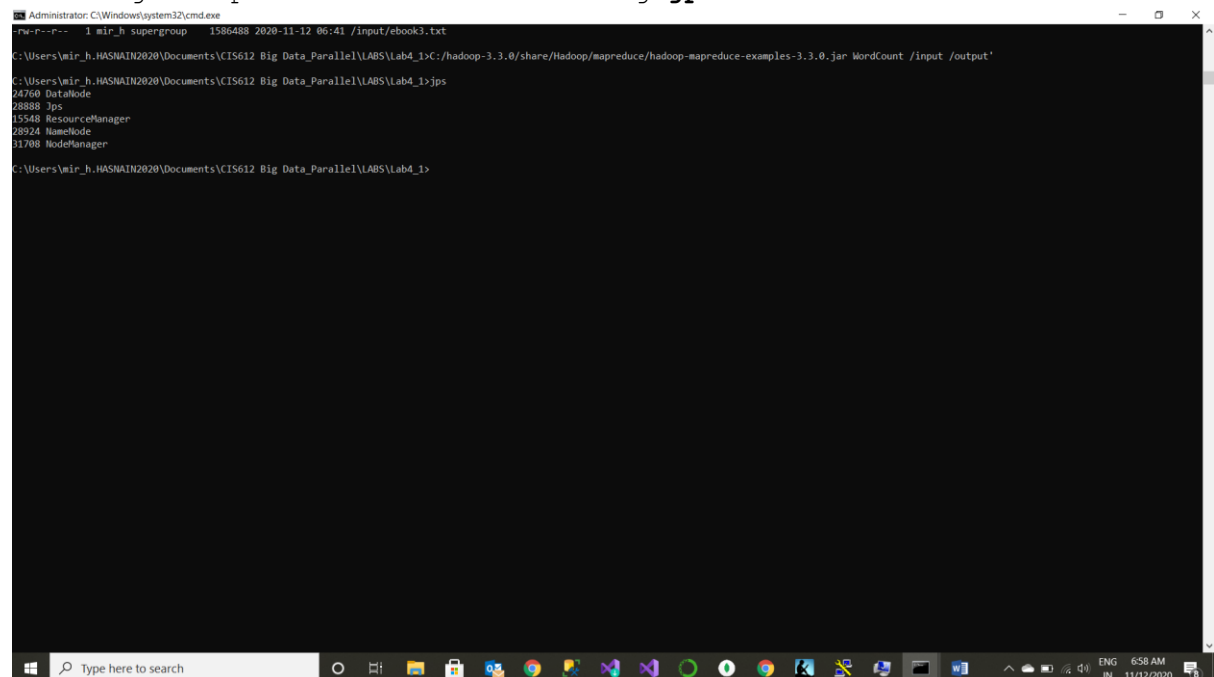
This would start the following

Name Node

Data Node



```
2020-11-12 04:09:52,430 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 2 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:09:55,431 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 3 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:09:58,434 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 4 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:01,438 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 5 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:04,440 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 6 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:07,445 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 7 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:10,448 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 8 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:13,452 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 9 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:15,457 WARN datanode.DataNode: Problem connecting to server: localhost/127.0.0.1:9000
2020-11-12 04:10:23,460 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 0 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:26,464 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 1 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:29,469 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 2 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:32,474 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 3 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:35,476 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 4 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
2020-11-12 04:10:38,479 INFO ipc.Client: Retrying connect to server: localhost/127.0.0.1:9000. Already tried 5 time(s);
retry policy is RetryUpToMaximumCountWithFixedSleep(maxRetries=10, sleepTime=1000 MILLISECONDS)
```

Resource Manager



```
2020-11-12 03:48:07,960 INFO pb.RpcServerFactoryPBImpl: Adding protocol org.apache.hadoop.yarn.server.api.ResourceTracke
rPB to the server
2020-11-12 03:48:07,970 INFO ipc.Server: IPC Server Responder: starting
2020-11-12 03:48:07,970 INFO ipc.Server: IPC Server listener on 8031: starting
2020-11-12 03:48:07,985 INFO util.JvmPauseMonitor: Starting JVM pause monitor
2020-11-12 03:48:08,004 INFO ipc.CallQueueManager: Using callQueue: class java.util.concurrent.LinkedBlockingQueue, queu
eCapacity: 5000, scheduler: class org.apache.hadoop.ipc.DefaultRpcScheduler, ipcBackoff: false.
2020-11-12 03:48:08,020 INFO ipc.Server: Starting Socket Reader #1 for port 8030
2020-11-12 03:48:08,039 INFO pb.RpcServerFactoryPBImpl: Adding protocol org.apache.hadoop.yarn.api.ApplicationMasterProt
ocolPB to the server
2020-11-12 03:48:08,045 INFO ipc.Server: IPC Server Responder: starting
2020-11-12 03:48:08,045 INFO ipc.Server: IPC Server listener on 8030: starting
2020-11-12 03:48:08,165 INFO ipc.CallQueueManager: Using callQueue: class java.util.concurrent.LinkedBlockingQueue, queu
eCapacity: 5000, scheduler: class org.apache.hadoop.ipc.DefaultRpcScheduler, ipcBackoff: false.
2020-11-12 03:48:08,168 INFO ipc.Server: Starting Socket Reader #1 for port 8032
2020-11-12 03:48:08,172 INFO pb.RpcServerFactoryPBImpl: Adding protocol org.apache.hadoop.yarn.api.ApplicationClientProt
ocolPB to the server
2020-11-12 03:48:08,173 INFO ipc.Server: IPC Server Responder: starting
2020-11-12 03:48:08,173 INFO ipc.Server: IPC Server listener on 8032: starting
2020-11-12 03:48:08,181 INFO resourcemanager.ResourceManager: Transitioned to active state
2020-11-12 03:48:10,084 INFO resourcemanager.ResourceTrackerService: NodeManager from node Hasnain2020(cmPort: 63754 htt
pPort: 8042) registered with capability: <memory:8192, vCores:8>, assigned nodeId Hasnain2020:63754
2020-11-12 03:48:10,090 INFO rmnode.RMNodeImpl: Hasnain2020:63754 Node Transitioned from NEW to RUNNING
2020-11-12 03:48:10,118 INFO capacity.CapacityScheduler: Added node Hasnain2020:63754 clusterResource: <memory:8192, vCo
res:8>
2020-11-12 03:58:07,868 INFO scheduler.AbstractYarnScheduler: Release request cache is cleaned up
```

Node Manager



Checking the port of the daemons using **jps** command



Open the browser and type localhost:9870 anf localhost:8086/clusters to see the Hadoop connection and the nodes running on it .

## Overview 'localhost:9000' (✓active)

| Started: | Thu Nov 12 06:24:46 +0530 2020 |
|---|---|
| Version: | 3.3.0, raa96f1871bfd858f9bac59cf2a81ec470da649af |
| Compiled: | Tue Jul 07 00:14:00 +0530 2020 by brahma from branch-3.3.0 |
| Cluster ID: | CID-0ac3c6ac-bc64-4296-8f3b-362fab35b811 |
| Block Pool ID: | BP-1846466883-172.17.240.113-1605142393221 |

## 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).

Heap Memory used 78.58 MB of 306.5 MB Heap Memory. Max Heap Memory is 889 MB.

Non Heap Memory used 48.26 MB of 49.86 MB Commited Non Heap Memory. Max Non Heap Memory is <unbounded>.

| Configured Capacity: | 458.62 GB |
|---|---|

## All Applications

### Cluster Metrics

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

### Cluster Nodes Metrics

| Active Nodes | Decommissioning Nodes | Decommissioned Nodes | Lost Nodes |
|---|---|---|---|
| 1 | 0 | 0 | 0 |

### Scheduler Metrics

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

Show 20 ∨ entries

| ID | User | Name | Application Type | Application Tags | Queue | Application Priority | StartTime | LaunchTime | FinishTime | State | FinalStatus | Running Containers | Alloc CP VCo |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | No data available in table | | | | |

Showing 0 to 0 of 0 entries

# Part-1 Hive

Our input data is the videogame sales data



## 1. Creation of Table

**CREATE EXTERNAL TABLE** videogames

(       **Rank INT,**

      **Name STRING,**

      **Platform STRING,**

      **Year INT,**

      **Genre STRING,**

      **Publisher STRING,**

      **NA_Sales DOUBLE,**

      **EU_Sales DOUBLE,**

      **JP_Sales DOUBLE,**

      **Other_Sales DOUBLE,**

      **Global_Sales DOUBLE,**

**)**

## 2. Loading data into created table all states

Load data local inpath '/home/hduser/Desktop/videogames.csv' into table `videogames`;

## 3. Creation of partition table

**create table Platform _part(Rank INT,**

       **Name STRING,**

       **Year INT,**

       **Genre STRING,**

       **Publisher STRING,**

       **NA_Sales DOUBLE,**

       **EU_Sales DOUBLE,**

       **JP_Sales DOUBLE,**

       **Other_Sales DOUBLE,**

       **Global_Sales DOUBLE,**

**) PARTITIONED BY(Platform STRING);**

## 4. Set the following property on the partition
**set hive.exec.dynamic.partition.mode=nonstrict**

## 5. Loading data into partition table

**INSERT OVERWRITE TABLE Platform _part PARTITION(Platform)**

**SELECT Name, Year, Genre, Publisher ,NA_Sales , EU_Sales , JP_Sales , Other_Sales , Global_Sales from  videogames;**

6. ## Creating Bucket

In Hive, we have to enable buckets by using the **set.hive.enforce.bucketing=true;**

**Create table sample_bucket { Rank INT,**

   **Name STRING,**

   **Year INT,**

   **Platform STRING,**

   **Publisher STRING,**

   **NA_Sales DOUBLE,**

   **EU_Sales DOUBLE,**

   **JP_Sales DOUBLE,**

   **Other_Sales DOUBLE,**

**Global_Sales DOUBLE} clustered by (genre ) into 4 buckets**

**Into row format delimited**

**Fields terminated by ',';**

7. Loading data into sample_bucket
   **From videogames**
   **Insert overwrite table sample_bucket**
   **SELECT Name, Year, Platform, Publisher ,NA_Sales , EU_Sales , JP_Sales , Other_Sales ,**
   **Global_Sales**
8.

# `Part-2 MongoDb on HDFS`

Download sample data

Using the business100.json and review100.json files

- Install MongoDB on the VM

- Start MongoDB – a default configuration file is installed by yum so you can just run this to start on localhost and the default port 27017

**mongod -f /etc/mongod.conf**

- Install MongoDB Hadoop Connector
- Inserting into the mongodb collection using

**mongoimport --jsonArray --db test --collection docs --file C:\Users\mir_h.HASNAIN2020\Documents\CIS612 Big Data_Parallel\LABS\LAb4_2\business.json**
**mongoimport --jsonArray --db test --collection docs --file C:\Users\mir_h.HASNAIN2020\Documents\CIS612 Big Data_Parallel\LABS\LAb4_2\review.json**

Displaying the data in the collection using mongodb via Hadoop input format

# set up parameters for reading from MongoDB via Hadoop input format

config = {"mongo.input.uri": "mongodb://localhost:27017/YelpBusiness.business"}

inputFormatClassName = "com.mongodb.hadoop.MongoInputFormat"

keyClassName = "org.apache.hadoop.io.Text"

valueClassName = "org.apache.hadoop.io.MapWritable"


RawRDD = sc.newAPIHadoopRDD(inputFormatClassName, keyClassName, valueClassName, None, None, config)


# configuration for output to MongoDB

config["mongo.output.uri"] = "mongodb://localhost:27017/ YelpBusiness. business "

outputFormatClassName = "com.mongodb.hadoop.MongoOutputFormat"

RDD = RawRDD.values()

```
uncaught exception: SyntaxError: illegal character :
@(shell):1:0
> config = {"mongo.input.uri": "mongodb://localhost:27017/YelpBusiness.business"}
{ "mongo.input.uri" : "mongodb://localhost:27017/YelpBusiness.business" }
> inputFormatClassName = "com.mongodb.hadoop.MongoInputFormat"
com.mongodb.hadoop.MongoInputFormat
> keyClassName = "org.apache.hadoop.io.Text"
org.apache.hadoop.io.Text
> valueClassName = "org.apache.hadoop.io.MapWritable"
org.apache.hadoop.io.MapWritable
>
> RawRDD = sc.newAPIHadoopRDD(inputFormatClassName, keyClassName, valueClassName, None, None, config)
uncaught exception: ReferenceError: sc is not defined :
@(shell):1:1
>
> # configuration for output to MongoDB
uncaught exception: SyntaxError: illegal character :
@(shell):1:0
> config["mongo.output.uri"] = "mongodb://localhost:27017/ YelpBusiness. business "
mongodb://localhost:27017/ YelpBusiness. business
> outputFormatClassName = "com.mongodb.hadoop.MongoOutputFormat"
com.mongodb.hadoop.MongoOutputFormat
> RDD = RawRDD.values()
uncaught exception: ReferenceError: RawRDD is not defined :
@(shell):1:1
>
```

config["mongo.output.uri"] = "mongodb://localhost:27017/ YelpBusiness. business "

outputFormatClassName = "com.mongodb.hadoop.MongoOutputFormat"

RDD = RawRDD.values()