

Monday

Big Data Practicals:

(i) Cloudera Quick Start VM is a very much CPU intensive atleast specify the 2CPU and increase the RAM size. (5GBRAM)

which gives the cloudera hadoop distribution in a single node cluster setup, which can be used for working & learning about different distribution. On a cloudera with hdfs and other hadoop ecosystem components.

In CQVM, we have one scm server and one scm agent which will start the system.

with Admin control accessible : we have to check if we have an access to the cluster. using the following commands.

In terminal

\$ hostname

\$ hdfs dfs -ls /

(There is connection establishment b/w scm server and scm agent (on cloudera) which take care of cluster.)

we can also check with following code.

\$ service cloudera-scm-server status

SEPTEMBER 2019						
Su	Mo	Tu	We	Th	Fr	Sa
	1	2	3	4	5	6
	8	9	10	11	12	13
	15	16	17	18	19	20
	22	23	24	25	26	27
	29	30				

Tuesday

The above code would tell you to use Cloudera Express.

10 Login as a root:

11 \$ su

password: Cloudera

12 \$ service cloudera-scm-server status

The above all the codes tells you hdfs working fine, now close the terminal

Now, click on the launch Cloudera Express icon, you will get the command

sudo /home/cloudera/cloudera-managed --force

Copy the above code and paste into the New Terminal window

6 \$ sudo /home/cloudera/cloudera-managed --force --express

This code would tell you shutdown all the services of Cloudera and restart the services, after which Only you will be able to access Admin console.

Now, the deployment has been configured and Wednesday and restarted the cloudera manager. It gives you Quick start admin control. using username and pwd.

Open the Browser, Change ~~http://~~ into

quickstart.cloudera:7180

click of \downarrow agree \rightarrow default port that shows Admin console.

In the Admin console window

enter username & pwd

cloudera — user name

cloudera — password.

Hue is a web interface which allows work with hdfs and depending on hue,

Remove the unused services

After, click on Actions

\downarrow

Restart

\rightarrow close (let it happen in background)

Click on Hosts

\downarrow
All hosts

\rightarrow As of now we have only one host.

SEPTEMBER 2019						
Su	Mo	Tu	We	Th	Fr	Sa
	2	3	4	5	6	7
	9	10	11	12	13	14
	16	17	18	19	20	21
	23	24	25	26	27	28
	30					

Thursday

(4)

HDFS (Hadoop Distributed File System)

- a framework for the analysis and transformation of very large data sets using mapreduce paradigm.

- Important Characteristics: Hadoop is the partitioning of data and computation across many (thousands) of hosts, and executing application computations in parallel close to their data.

Basic Commands in cloudera:

In Terminal 1

```
$ ls
$ pwd
$ mkdir Sudhakar
$ cd Sudhakar
$ pwd
$ cd ..
$ ls
$ cat name.txt
```

In Terminal 2

```
$ hadoop fs -ls /
```

```
$ pwd
```

```
$ hadoop fs -ls /usr
```

\$ hadoop fs -put employee.txt /usr
 (used to copy the file from host to hdfs)

pwd - present working directory

SEPTEMBER
2019



Week 38 | 263 • 102

20

\$ sudo jps (used for my name node's Friday
these by not)

\$ pwd (need to check where I am)

/home/cloudera

\$ cd /home

\$ ls

\$ cd ..

\$ ls /home/cloudera

\$ ls /home/cloudera -ltr

\$ ls /home/cloudera/Downloads -ltr

\$ cat /home/cloudera/Downloads/products.csv

hadoop

\$ hadoop fs -ls

Old approach
Using same thing

\$ hadoop fs -ls /

↳ list out directory list (root)

\$ hdfs dfs -ls /

\$ hdfs dfs -ls / foldername.

\$ hdfs dfs -cat / foldername / filename, tkt

September 2019
Mo Tu We Th Fr Sa
1 2 3 4 5 6 7
8 9 10 11 12 13 14
15 16 17 18 19 20 21
22 23 24 25 26 27 28
29 30 31
To see contents of the file

21

~~It local host \$ sudo mkdir /foldername
Create directory~~
Week 38

SEPTEMBER
2019

Saturday

(b)

9 Creating the new directories in hadoop:

10 \$ hdfs dfs -mkdir /foldername

11 List out the directory.

12 \$ hdfs dfs -ls /

1 Copy files from local host to hdfs.

In local host:
list out files:

4 \$ ls /home/Cloudera/Downloads

5 \$ hadoop CopyFromLocal

6

7

22 Sunday

Understanding Hadoop command

In cloudera, we have 2 environments

- 1. Base environment (Linux)
- 2. distributed file (Hadoop)

\$ hadoop fs -ls /

It will list out all the files in hadoop

\$ hdfs dfs -ls /

This one is also to list out all the files in the hadoop.

\$ hadoop fs -mkdir /test3

↳ This command is used to create directory in hadoop environment.

\$ hadoop fs -ls /

\$ clear

Copy a file from host to distributed system

\$ ls

\$ cd test3

// with in test3 or file, create one test.txt

\$ ls

test.txt

\$ hadoop fs -put test.txt /test3

Tuesday

267 • 098 | Week 39

Check whether file is there or not.

9

\$ hadoop fs -ls /test3

10

Equivalent command

11

\$ hadoop fs -copyFromLocal test.txt /test4

12

Here test4 folder is not available, though by default it will create a directory.

2

\$ hadoop fs -cat /test4

3

This command will show you text written in test4 directory.

5

\$ cat test.txt

4

move a file from host to hadoop.

7

\$ hadoop fs -mkdir /test6

\$ hadoop fs -ls /

\$ hadoop fs -moveFromLocal test.txt /test6

\$ ls

\$ hadoop fs -ls /test6

SEPTEMBER
2019

⑨

Week 39 | 268 • 097

25

// we want to get a file from hadoop ^{Wednesday}

\$ hadoop fs -get /testb/test.txt

\$ ls

// In hadoop environment, we want to copy a file from one location to another location.

\$ hadoop fs -mkdir /test7

// move a testb/test.txt.

\$ hadoop fs -cp /testb/test.txt /test7

\$ hadoop fs -ls /test7

hadoop fs -ls

hadoop fs -ls

SEPTEMBER 2019

Su	Mo	Tu	We	Th	Fr	Sa
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30					

26

Po

SEPTEMBER

2019

Thursday

269 • 096 | Week 39

Dift. b/w HDFS and liaison

9 HDFS → is a stateless i.e. doesn't remember working directory, each time we have to write full path.

10
11 → Cannot modify the file, But we can delete and copy one
12 or write file.

1 \$ hadoop fs

In Linux

2 \$ hadoop fs -ls /

| \$ ls

3 ⇒ That '/' denotes root directory. On the HDFS.

4 If you want to go to the particular directory use this following command.

5 \$ hadoop fs -ls /usr

6 Create our own directory :

7 \$ hadoop fs -mkdir /mydata.

mydata folder will be created.

\$ hadoop fs -mkdir /mydata/testfold

Creating sub folder within mydata.

Friday

creating a text file in linux.

\$ gedit testfile.txt

\$ ls *.txt // list out all txt file

// save this file in hadoop

\$ hadoop fs -put testfile.txt /mydata/testfolder/

\$ hadoop fs -ls /mydata/testfolder/

// to get a file from hadoop to linux.

\$ hadoop fs -get /mydata/testfolder/

testfile.txt newfile.txt

\$ ls *.txt

SEPTEMBER 2019						
Su	Mo	Tu	We	Th	Fr	Sa
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28

28

(11)

271 • 094 | Week 39

Saturday

SEPTEMBER

2019

(12)

Hive setup in Cloudera

9

Step 1 : Create a folder on your
Desktop e.g: example.

10

Step 2 : Open a terminal window.

11

Type the command :

12

password : Cloudera.

2

Step 3 : \$ cd /usr/lib/hive/conf/

3

\$ sudo gedit hive-site.xml

4

Change the file path to Desktop

5

\$ sudo !hive

6

29 Sunday

Data to Analytics:

view of

HDFS, Map reduce, Hive, Sqoop, and spark.

\$ hadoop fs // to see what are all the commands are there

\$ touch testing // make a file (empty file in local filesystem)

\$ hadoop fs -put testing /user/cloudera

\$ hadoop fs -ls /user/cloudera // list out copy file

\$ vi wordcount.txt // create a file

\$ hadoop fs -put wordcount.txt /user/cloudera

\$ sudo find / -name "*hadoop*examples*.jar"

Select the

/usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar

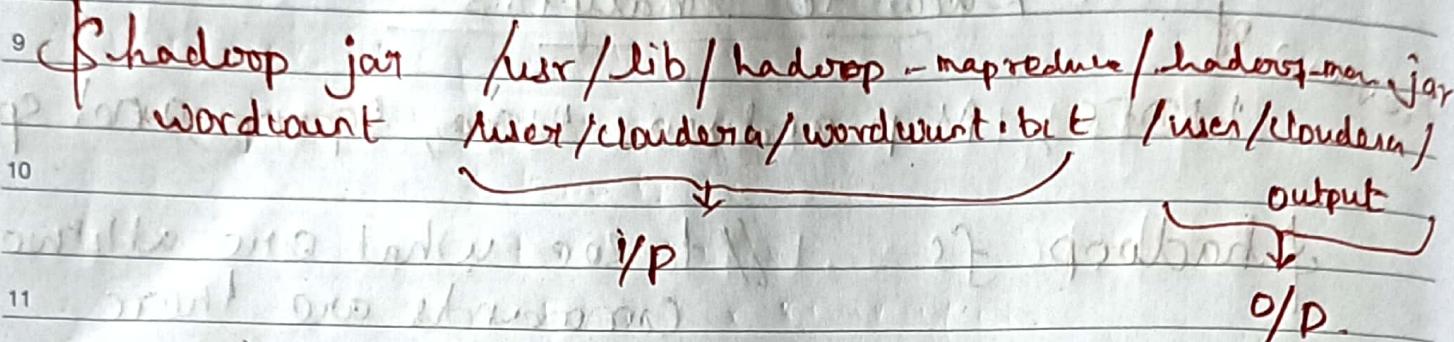
copy the link and paste it

\$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar

// This will provide all the applications that part of jar files

TEMBER 2019						
Su	Mo	We	Th	Fr	Sa	Su
1	Tu	We	Th	Fr	Sa	Su
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

Tuesday



Now, the mapreduce job will run, it will give you count of each word in wordcount.txt.

1 \$ hadoop fs -ls /user/cloudera/output

2 \$ hadoop fs -cat /user/cloudera/output/part-r-00000

3 // Count of each word

4 Share folders b/w VM & Host.

5 Devices → shared folder → shared folder setting

6 machine folder

7 click file menu (F)

with in window

folder path : C:\Projects (windows folder path)
 folder name : Projects (your path)

Check Automount

make permanent.

Ok

OCTOBER

19

touch z → is used to create empty file

dd

15

Week 40 | 275 • 090

02

Run the mount command:

Wednesday

Open a Terminal

Login as a root

\$ su

\$ cd /home/Cloudera/workspace

\$ mkdir projects

\$ ls -l

\$ mount -t vboxsf projects /home/Cloudera/workspace/projects

Shared folder name.

target

OCT

Difference b/w LFs and HDFS:

[LFs]

Linux

[HDFS]

Special commands.

Basic commands:

\$ su root

pwd: cloudera

\$ jps // Java processes

It will show you HDFS components are running

OCTOBER
2019

Mo Tu We Th Fr Sa

1 2 3 4 5

6 7 8 9 10 11

12 13 14 15 16

17 18 19 20 21

22 23 24 25 26

27 28 29 30 31

Hadoop-Doces

Thursday

9 Main components:

- Node manager

- Name node

- Data node

- Resource manager.

12 Others are supporting components.

1 \$ su cloudera // come back to normal user.

2 \$ clear.

3 How we can check the health i.e. How nodes are running.

4 fsck → check the ~~not~~ health of hdfs.

5 \$ hdfs fsck /

6 ✓ \$ hdfs dfs -ls /

7 ✓ \$ hdfs dfs -mkdir /bigdatatesting.

✓ \$ hdfs dfs -ls /

touchz → create a new file on HDFS with size 0 bytes.

\$ hdfs dfs -touchz /bigdatatesting/test.dat

OCTOBER
19

vi editor for local file system

17

Week 40 | 277 • 088

04

du \Rightarrow checks the size of a file. Friday

\$ hdfs dfs -du -s /bigdata/testing/test.dat

How we can content to the existing file:

appendToFile \rightarrow appends the contents to the file which is present on HDFS.

\$ hdfs dfs -appendToFile - /bigdata/testing/test.dat

. testing per

\$ hdfs dfs -du -s /bigdata/testing/test.dat

To display the contents of the file.

using cat command present in HDFS.

\$ hdfs dfs -cat /bigdata/testing/test.dat

CopyFromLocal (put)

copy To Local (get)

\$ cat >> test1.dat

testing hadoop.

test1 \$ ls test1.dat

Mo	Tu	We	Th	Fr	Sa
1	2	3	4	5	
7	8	9	10	11	12
14	15	16	17	18	19
21	22	23	24	25	26
28	29	30	31		

05

278 • 087 | Week 40

OCTOBER
2019

Saturday

18

9 \$ hdfs dfs - copyFromLocal test1.dat /bigdata/testing

10 \$ hdfs dfs -ls /bigdata/testing.

11 \$ cat >> test2.dat

12 ---

1 \$ hdfs dfs -put test2.dat /bigdata/testing.

2 \$ hdfs dfs -ls /bigdata/testing.

3
4 Name node : is the master daemon that
5 maintains and manages data nodes
6 and it records only the metadata

7 no. of frags
filename,
file permission,
which rack the file
is there.

06 Sunday

Data node : are slave daemons (or) processes
which runs on each slave
machine & it stores actual data

CTOBER
19

ab

19

Week 41 | 280 • 085

07

Monday

\$ hadoop fs -test -d destination

directory name.

// test command used for conditional checking, we can check some condition

-d means whatever path I am specifying after -d, is it is directory (or) ~~not~~ file

\$ echo \$? // print results.

if it returning zero(0) means it is directory otherwise file.

\$ hadoop fs -test -e destination

The given directory is existing (or) not.

\$ hadoop fs -test -f destination

file (or) not.

\$ hadoop fs -test -z destination/sy z.txt

// this would test the file size is zero or not.

\$ echo \$?

TOBER

2019

	Mo	Tu	We	Th	Fr	Sa
1	2	3	4	5		
7	8	9	10	11	12	
14	15	16	17	18	19	
21	22	23	24	25	26	
28	29	30	31			

// move from local to hdfs

\$ hadoop fs -moveFromLocal a2.txt destination

\$ cat a2.txt

08

281 • 084 | Week 41

Tuesday

29

OCTOBER
2019

9

//merge two file contents

10

Before that create two text files in any directory using appendToFile

\$ hadoop fs -cat destination / xyz.txt
directory name

1000-1000

hadoop fs -cat destination/sig71.txt

2

\$ hadoop fs -getmerge -nl destination/partition
\$ hadoop fs -put destination/partition newFile

// Verify whether megabots not

3

\$ cat Desktop/merge.txt

4

- Desktop / mega bit
Where we have to store new file name

5 Checksum Command:

6

checks integrity of one file

whether file is modified or not.

We can check their hash value.

7

↳ shadow of - checksum destination / syn. bit.

CTOBER 09 mysql -u root -h localhost -p -P 3306
09 default port no.
09 21 282 • 083 09

mapreduce

How to create and reduce mapreduce program over cloudera environment.

Mysql is running (or) not on cloudera environment

\$ sudo service mysqld status

\$ mysql -u root -h localhost -p

pwd: cloudera ↗ db name

mysql > show databases;

// list out all the db

mysql > use retail_db;

✓ db name

// get onto the db

mysql > show tables;

// list the tables.

mysql > select count(1) from orders;

mysql > describe orders;

// list out all the field & type.

mysql > exit

OCTOBER 2019						
Su	Mo	We	Th	Fr	Sa	Su
1	2	3	4	5		
8	9	10	11	12		
15	16	17	18	19		
22	23	24	25	26		
29	30	31				

10

283 • 082 | Week 41

g2

OCTOBER
2019

Thursday

9 Sqoop: is a language that transfer the relational data (OR) SQL Databases to HDFS

10 \$ mysql -u root -P,

11 mysql> \pwd; Cloudera // enter into mysql;

12 mysql> show databases;

1 mysql> create database myfirstdb;

2 // creating database

mysql> use myfirstdb;

3 // going to use the db

4 mysql> create table mytable (ID int, name
5 varchar(20), address varchar(20));

6 // ~~Table~~ creation Table creation.

7 mysql> describe mytable;

mysql> insert into mytable values (1, "Adi", "USA"),
(2, "Ram", "UK"), (3, "Rey", "India");

mysql> select * from mytable;

drop table mytable;

OCTOBER
19

IP address : IP address
terminal 100
hostname - I
Week 41 284 • 081
23 Friday
10.0.2.15
local host
IP address of VM depends on machine

In next Terminal:

Login as a root:

To run a Sqoop

\$ sqoop import --connect jdbc:mysql://10.0.0.41:3306

/myfirstdata --username root --password cloudera

table db name

--table mytable --target-dir=/user/cloudera/

myfirstdata -m 1
No. of maps

\$ hadoop fs -cat /user/cloudera/myfirstdata

\$ hadoop fs -ls /user/cloudera/myfirstdata.

\$ hadoop fs -cat /user/cloudera/myfirstdata/part-m-00000

mysql to Sqoop:

Exporting Data from HDFS to MySQL Sqoop

OCTOBER 2019						
Mo	Tu	We	Th	Fr	Sa	Su
1	2	3	4	5		
7	8	9	10	11	12	
14	15	16	17	18	19	
21	22	23	24	25	26	
28	29	30	31			

Login as root

Sudo -i

Sqoop export --connect jdbc:mysql://

Saturday

9 Use Sqoop to copy a relational db table to HDFS

10 \$ mysql --user root --password cloudera movieLens

11 mysql > show databases;

12 mysql > show use movieLens;

13 mysql > show tables;

14 mysql > describe movie;

15 mysql > select * from movie limit 5;

16 mysql > select count(*) from movie;

17 Copy me all records from mysql to HDFS

18 \$ Sqoop import |

> --connect jdbc:mysql://localhost/movieLens |

> --username root -p |

19 Sunday

> --fields-terminated-by '\t' \

> --warehouse-dir /mydata \

> --table movie

Monday

\$ hadoop fs -ls /mydata

\$ hadoop fs -ls /mydata/movie

// view the contents.

\$ hadoop fs -cat /mydata/movie/part-m-00000

// list top 5 records only

Create and

Load data in Hive table

first create sample file, the field names are id, name, Dept, salary, Domain

~ Sachin ~ pune ~ Product Engg ~ 10000 ~

\$ cat employee

\$ hadoop fs -put employee /user/cloudera

\$ hadoop fs -ls /user/cloudera

\$ hadoop fs -cat /user/cloudera

// display all records from local

to HDFS

TOBER 2019						
Mo	Tu	We	Th	Fr	Sa	Su
1	2	3	4	5		
7	8	9	10	11	12	
14	15	16	17	18	19	
21	22	23	24	25	26	
28	29	30	31			

15

288 • 077 | Week 42

Tuesday

In another terminal

OCTOBER
2019

26

9

\$ hive // enter into hive shell

11

@ hive> show databases;

12

hive> create database organization;

1

hive> show databases;

2

hive> use organization;

3

hive> go Show tables; // show info

4

hive> ! clear;

5

hive> create table employee (

id int,

name string,

city string,

department string,

salary int,

domain string);

row format delimited

fields terminated by '\n';

hive> Show tables;

// list out the tables

live) Select * from employee Wednesday

// it will not return the data because we have not loaded

hive) load data inpath '/user/cloudera/employee'
Overwrite into table employee;

live) Select * from employee;

hive) desc employee;

// describe table

hive) Select name, city from employee;

✓ Hive is one of the most popular Datawarehouse

✓ developed by facebook; maintained by apache hive

✓ mainly netflix and amazon.

✓ why Hive was developed?

✓ How and when it can be used?

✓ when Hive cannot be used?

✓ features of Hive

✓ Hive Vs RDBMS

OBER							2019
Su	Mo	Tu	We	Th	Fr	Sa	
1	2	3	4	5			
7	8	9	10	11	12		
14	15	16	17	18	19		
21	22	23	24	25	26		
28	29	30	31				

Thursday

✓ Tree can be used in OLAP;

✓ Scalable & flexible

✓ used in multidimensional data

Cannot be used:

— It is not a relational database

— It cannot be used for OLTP

12: developing global DB structures

1: developing model of travel (cont)

2: developing model of travel (cont)

3: global structures

4: developing model of travel (cont)

5

5: developing model of travel (cont)

6

6: developing model of travel (cont)

7

7: normalizing without primary

8: normalizing with primary

9: normalizing with primary

10: normalizing with primary

11: normalizing with primary

12: normalizing with primary

OCTOBER
2019

28

① 29

Week 43 | 296 • 069

23

Wednesday

Big Data Practicals

Cloudera!

Word count Prg:

It will count each and every words as many times occurred on the file.

\$ jps // return all the running nodes

\$ pwd // present working directory // filename

\$ hdfs dfs -put/Downloads/input /

// Copy the input file from local file system to hadoop system.

\$ hdfs dfs -ls /

// return all the files in hadoop

Otherwise we can check with the user interface

on URL → localhost:9870

↑
utilities

↓
Browse Directory

↑
we can see the
input file

OCTOBER		2019				
Su	Mo	Tu	We	Th	Fr	Sa
	1	2	3	4	5	
	6	7	8	9	10	11
	13	14	15	16	17	18
	20	21	22	23	24	25
	27	28	29	30	31	

24

297 • 068 | Week 43

129

30

OCTOBER
2019

Thursday

9 `cd /bin/hadoop jar share/hadoop/mapreduce/
10 hadoop-mapreduce-examples-3.0.0.jar wordcount`

11 /input /input-output
12 filename directory name

We can check the output either User Interface or Terminal

User interface: use the local host address,

localhost:9870 / Explorer.html # /

Utilities

Browse Directory

input_output // file name

~~palt-τ-00000~~

Select & Head the file

You will see the effect

OCTOBER
2019

(31)

Week 43 | 298 • 067

25

Friday

In Terminal:

\$ hdfs dfs -cat /input_output/part-r-00000

// This would return the output in terminal.

Hive

Create and load data in Hive table.

\$ cat employee

\$ hadoop fs -put employee /user/cloudera

\$ hadoop fs -~~cp~~ -ls /user/cloudera

\$ hadoop fs -cat /user/cloudera/employee

Another Terminal:

\$ hive

hive> show databases;

· > create database organization;

OCTOBER 2019						
Su	Mo	Tu	We	Th	Fr	Sa
	1	2	3	4	5	
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30	31		