# Big Data Hadoop and Spark Developer

Lesson-End Project Solution
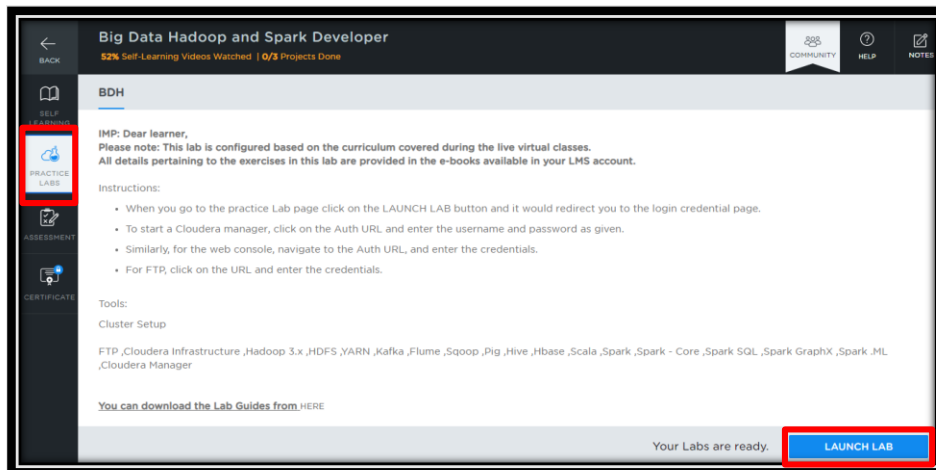
**simpli·learn**

Get Certified. Get Ahead.

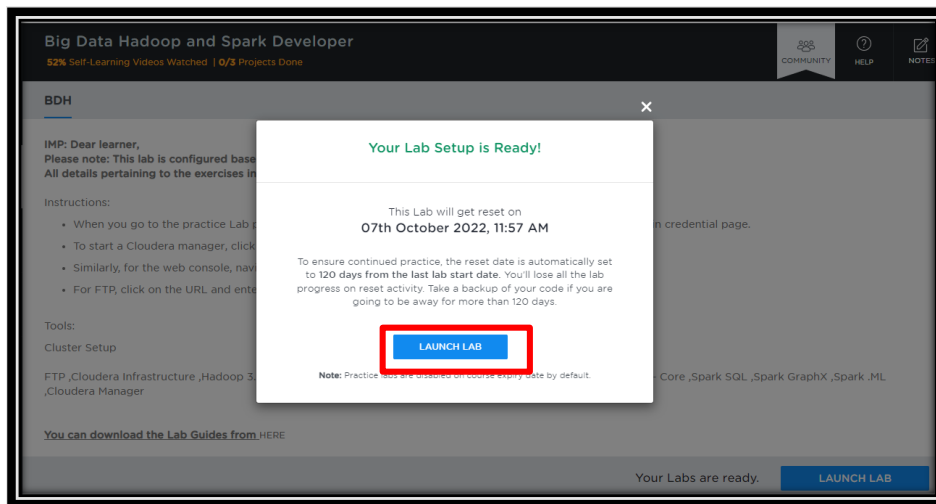# Post Office Data Analysis Using Hive

**Steps to Perform:**

**Step 1:** Download the lesson 6 dataset from the course resources and upload it to **"HDFS"**

1.1 Open the course "**Big Data Hadoop and Spark Developer**"

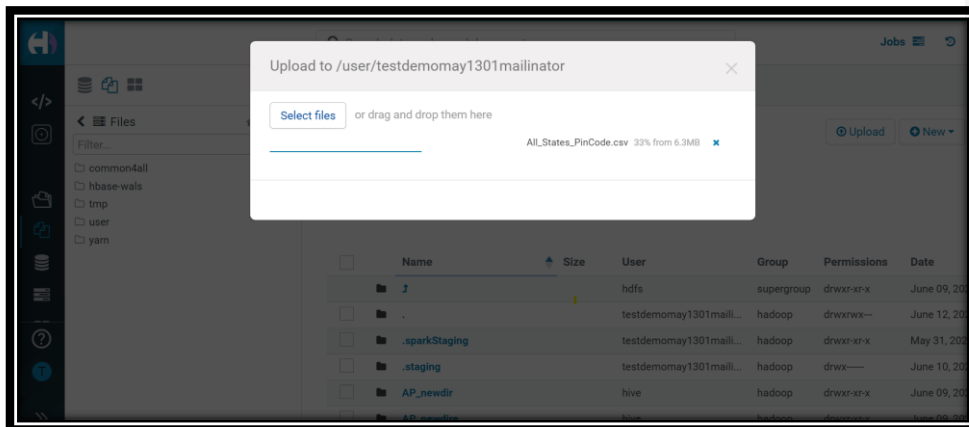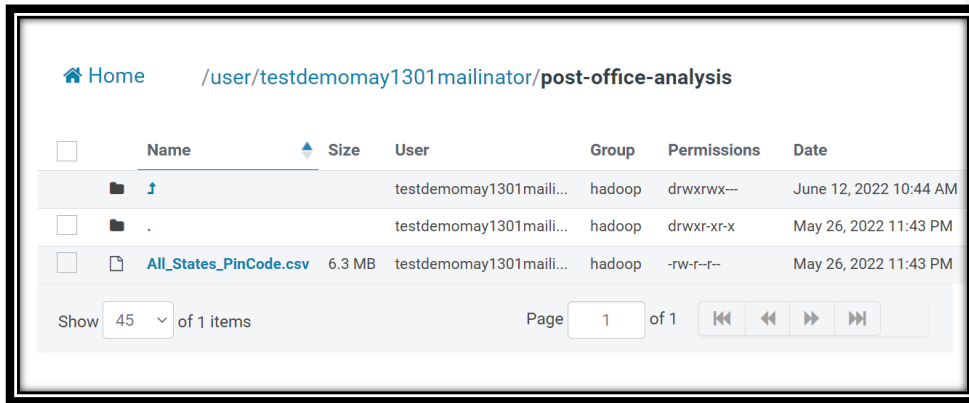1.2 Click on the "**PRACTICE LABS**" tab on the left side and select "**LAUNCH LAB**" on



**1.3** Click on the "**LAUNCH LAB**" button

1.4 Log in to the **"HUE"** lab

1.5 Click on **"HDFS"** and upload the downloaded dataset

| | Name | | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|---|---|
| ☐ | 📁 | ⬆ | | testdemomay1301maili... | hadoop | drwxrwx--- | June 12, 2022 10:44 AM |
| ☐ | 📁 | . | | testdemomay1301maili... | hadoop | drwxr-xr-x | May 26, 2022 11:43 PM |
| ☐ | 📄 | All_States_PinCode.csv | 6.3 MB | testdemomay1301maili... | hadoop | -rw-r--r-- | May 26, 2022 11:43 PM |

🏠 Home  /user/testdemomay1301mailinator/**post-office-analysis**

Show 45 ▾ of 1 items   Page 1 of 1 ⏮ ◀◀ ▶▶ ⏭

**Step 2**: Create a table on the Hive editor

2.1 Create a database

**Command:**
create database lesson6_lep;



```
1 create database lesson6_lep;
2
```

```
INFO  : Starting task [Stage-0:DDL] in serial mode
INFO  : Completed executing command(queryId=hive_202206121709
54_3eabdd46-10f6-4275-ac7c-3c9d0a05fac7); Time taken: 0.062 s
econds
INFO  : OK
```

✔ Success.

2.2 Use the created database

**Command:**
use lesson6_lep;



```
1  create database lesson6_lep;
2
3  use lesson6_lep;
4
```

```
INFO  : Starting task [Stage-0:DDL] in serial mode
INFO  : Completed executing command(queryId=hive_202206121711
01_b3e04345-17da-470b-a514-371826f74d54); Time taken: 0.007 s
econds
INFO  : OK
```
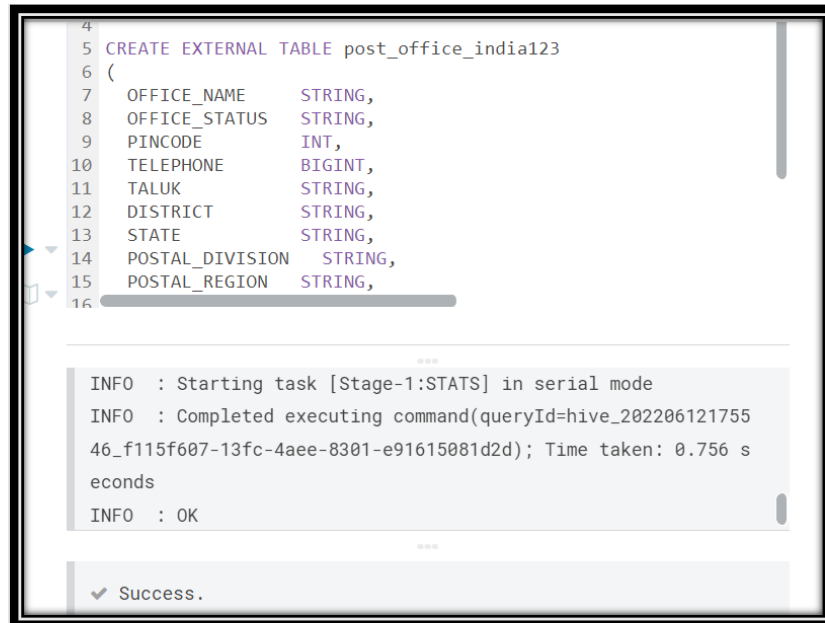
✔ Success.

2.3 Create the table

**Command:**
CREATE EXTERNAL TABLE post_office_india123
(
     OFFICE_NAME     STRING,
     OFFICE_STATUS    STRING,
     PINCODE     INT,
     TELEPHONE    BIGINT,
     TALUK    STRING,
     DISTRICT    STRING,
     STATE    STRING,
     POSTAL_DIVISION   STRING,
     POSTAL_REGION   STRING,
     POSTAL_CIRCLE   STRING
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','

STORED AS TEXTFILE;

```
4
5  CREATE EXTERNAL TABLE post_office_india123
6  (
7    OFFICE_NAME      STRING,
8    OFFICE_STATUS    STRING,
9    PINCODE          INT,
10   TELEPHONE        BIGINT,
11   TALUK            STRING,
12   DISTRICT         STRING,
13   STATE            STRING,
14   POSTAL_DIVISION   STRING,
15   POSTAL_REGION    STRING,
16
```

```
INFO  : Starting task [Stage-1:STATS] in serial mode
INFO  : Completed executing command(queryId=hive_202206121755
46_f115f607-13fc-4aee-8301-e91615081d2d); Time taken: 0.756 s
econds
INFO  : OK
```

✔ Success.

2.4 Load the CSV file into the table

**Command:**
LOAD DATA INPATH '/user/testdemomay1301mailinator/All_States_PinCode.csv'
INTO TABLE post_office_india123;

```
12    DISTRICT         STRING,
13    STATE            STRING,
14    POSTAL_DIVISION   STRING,
15    POSTAL_REGION    STRING,
16    POSTAL_CIRCLE    STRING
17  )
18  ROW FORMAT DELIMITED
19  FIELDS TERMINATED BY ','
20  STORED AS TEXTFILE;
21
22  LOAD DATA INPATH '/user/testdemomay1301mailinator/hivedemo1/Al
23
```

```
INFO  : Starting task [Stage-1:STATS] in serial mode
INFO  : Completed executing command(queryId=hive_202206121755
46_f115f607-13fc-4aee-8301-e91615081d2d); Time taken: 0.756 s
econds
INFO  : OK
```
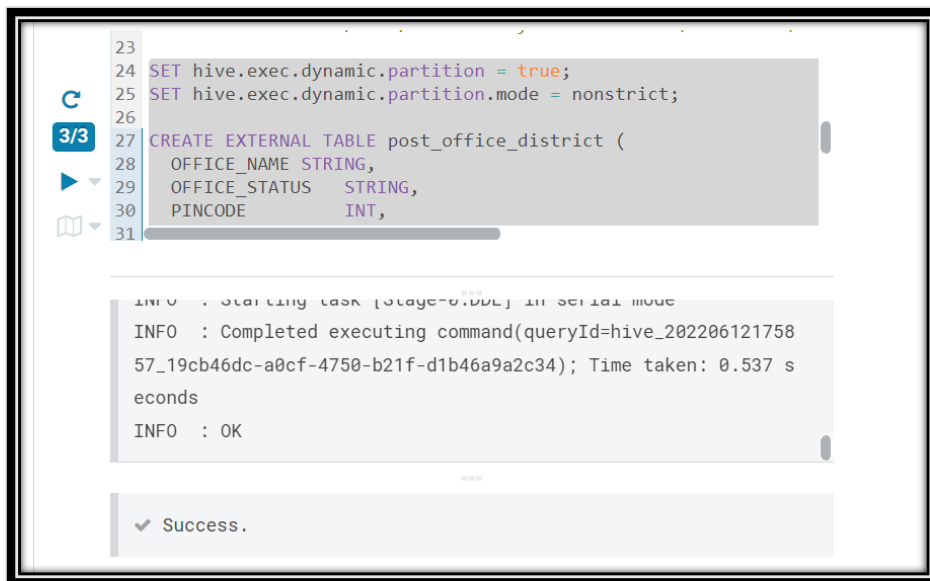
✔ Success.

2.5 Create a partitioned table to fetch data easily

Configure the Hive to support dynamic partition creation and enter the following set commands.

SET hive.exec.dynamic.partition = true;
SET hive.exec.dynamic.partition.mode = nonstrict;

CREATE EXTERNAL TABLE post_office_district (
        OFFICE_NAME STRING,
        OFFICE_STATUS    STRING,
        PINCODE            INT,
        TELEPHONE   BIGINT,
        TALUK          STRING,
        DISTRICT STRING,
        POSTAL_DIVISION   STRING,
        POSTAL_REGION    STRING,
        POSTAL_CIRCLE    STRING
)

PARTITIONED BY (STATE STRING)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;

```
23
24  SET hive.exec.dynamic.partition = true;
25  SET hive.exec.dynamic.partition.mode = nonstrict;
26
27  CREATE EXTERNAL TABLE post_office_district (
28    OFFICE_NAME STRING,
29    OFFICE_STATUS   STRING,
30    PINCODE         INT,
31
```

```
INFO  : Starting task [Stage-0:DDL] in serial mode
INFO  : Completed executing command(queryId=hive_202206121758
57_19cb46dc-a0cf-4750-b21f-d1b46a9a2c34); Time taken: 0.537 s
econds
INFO  : OK
```

✔ Success.

insert overwrite table post_office_district partition (STATE) select * from post_office_india123;

```
39 ROW FORMAT DELIMITED
40 FIELDS TERMINATED BY ','
41 STORED AS TEXTFILE;
42
43
44
45
46 insert overwrite table post_office_district partition (STATE)
47
```

```
INFO : Tez session was closed. Reopening...
INFO : Session re-established          application_1654575316638_0410
INFO : Session re-established.
INFO : Status: Running (Executing on YARN cluster with App i
d application_1654575316638_0410)
```

Query History          Saved Queries

a few seconds ago    ⇥    insert overwrite table post_office_di:
                          select * from post_office_india123

a minute ago              CREATE EXTERNAL TABLE post_office_di:

---

```
39 ROW FORMAT DELIMITED
40 FIELDS TERMINATED BY ','
41 STORED AS TEXTFILE;
42
43
44
45
46 insert overwrite table post_office_district partition (STATE)
47
```

```
INFO : Starting task [Stage-3:STATS] in serial mode
INFO : Completed executing command(queryId...        application_1654575316638_0410
50_f20a923b-b0a4-4083-80c6-7cd0933e986d); Time taken: 88.917
seconds
INFO : OK
```
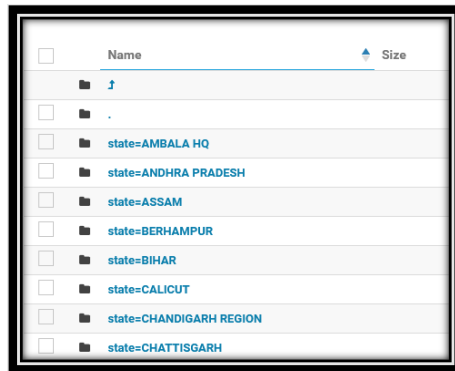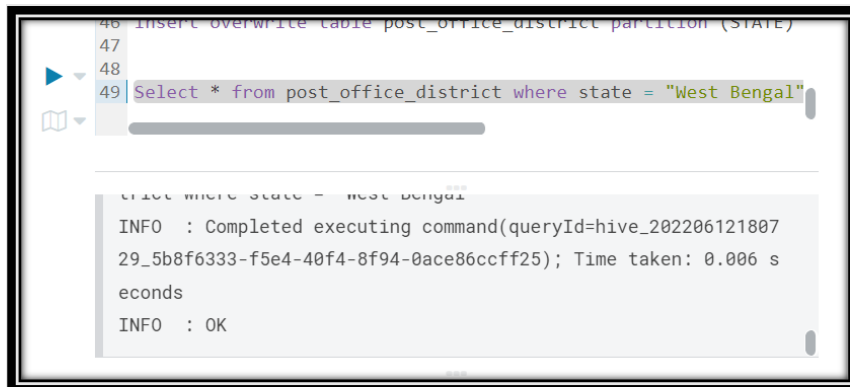
✔ Success.

**Note**: The above command changes the setting only for a single session. This will create multiple partitions in the HDFS directory with each state as a subfolder.

| | Name | Size |
|---|---|---|
| ☐ | ⬆ | |
| ☐ | . | |
| ☐ | state=AMBALA HQ | |
| ☐ | state=ANDHRA PRADESH | |
| ☐ | state=ASSAM | |
| ☐ | state=BERHAMPUR | |
| ☐ | state=BIHAR | |
| ☐ | state=CALICUT | |
| ☐ | state=CHANDIGARH REGION | |
| ☐ | state=CHATTISGARH | |

2.7 Run the query where the state is West Bengal

Select * from post_office_district where state = 'West Bengal';

```
46  insert overwrite table post_office_district partition (STATE)
47
48
49 |Select * from post_office_district where state = "West Bengal"
```

```
trict where state = "West Bengal"
INFO  : Completed executing command(queryId=hive_202206121807
29_5b8f6333-f5e4-40f4-8f94-0ace86ccff25); Time taken: 0.006 s
econds
INFO  : OK
```

In comparison to the previous runs, which took 30 seconds to complete, this will return the results in seconds.