# Big Data Hadoop and Spark Developer

Lab Guide
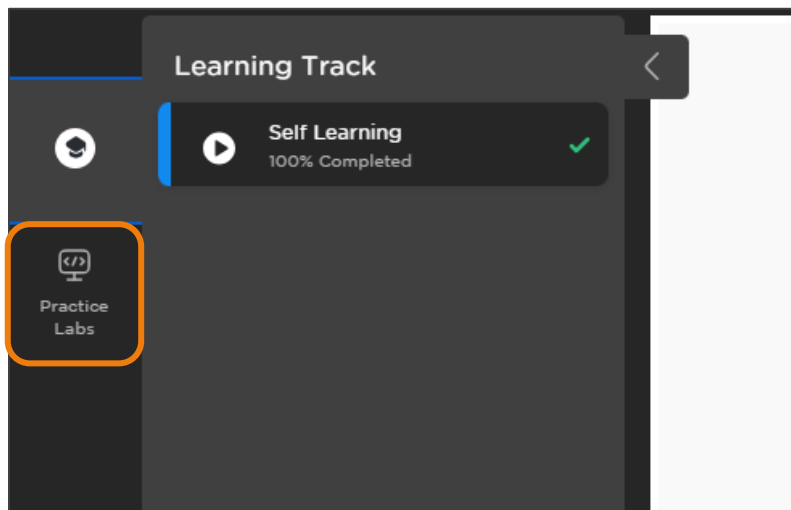
simpl¡learn

Get Certified. Get Ahead.

**Note:** The screenshots are only for your reference. Your LMS may look different depending on your course content.
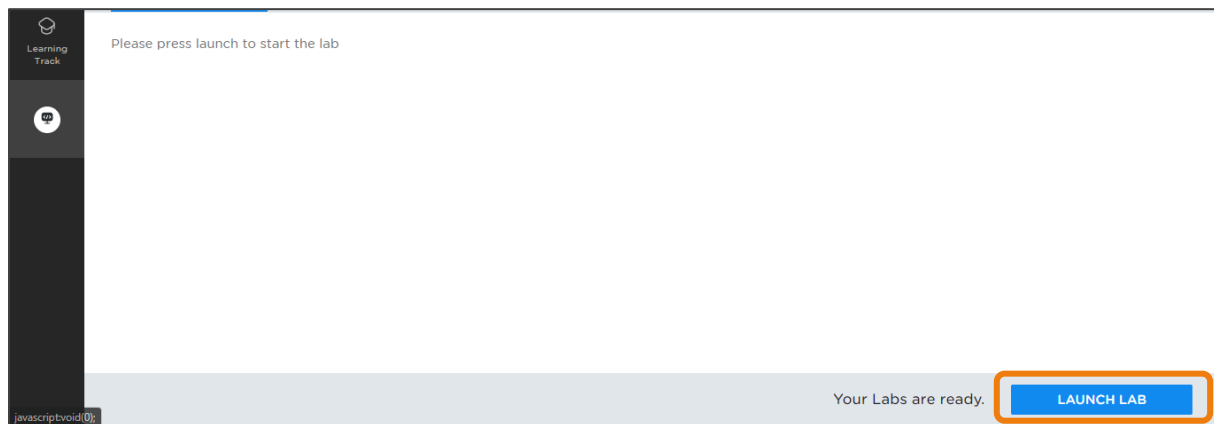
To execute the demos included in this course, follow the below steps:
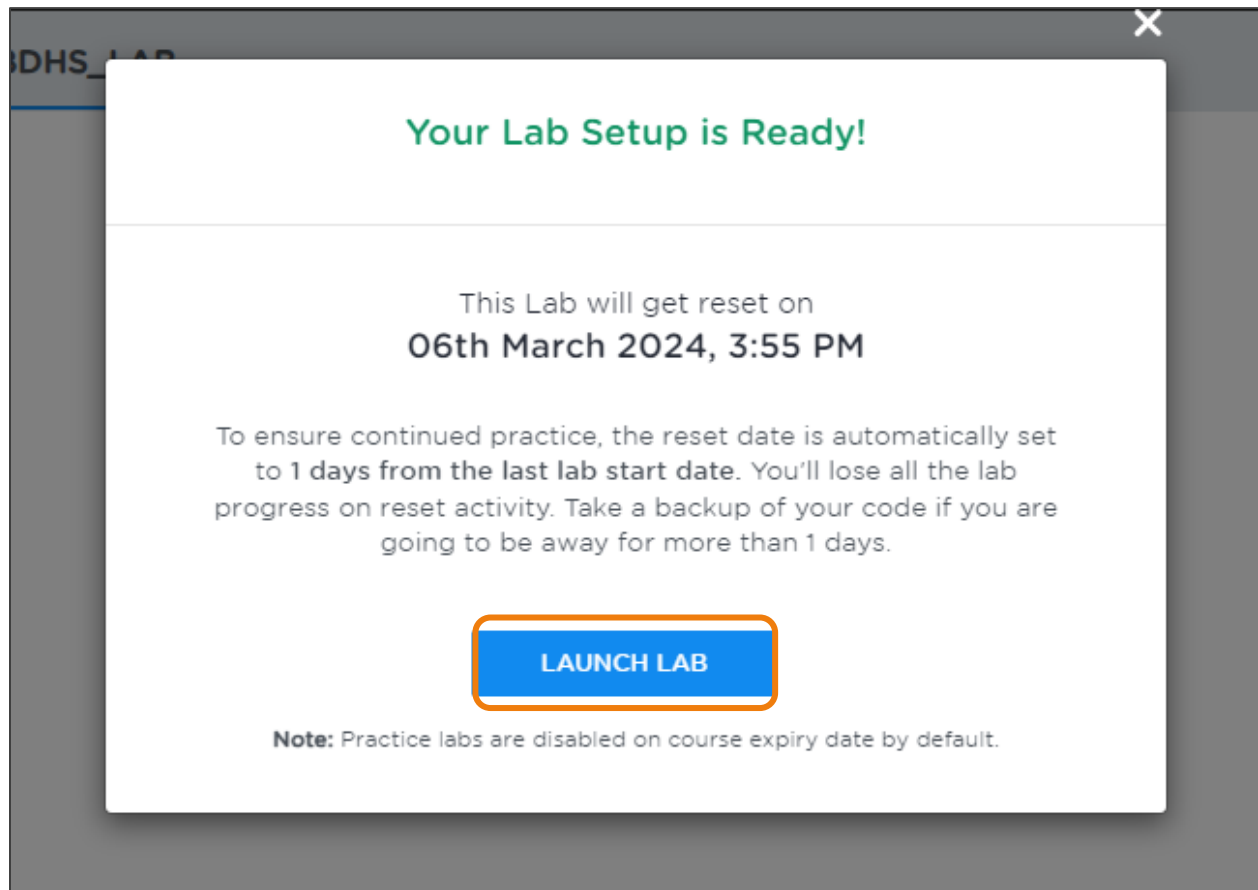
**Step 1:** Log in to the Simplilearn LMS

- Go to the course
- Click on **Practice Labs** on the left



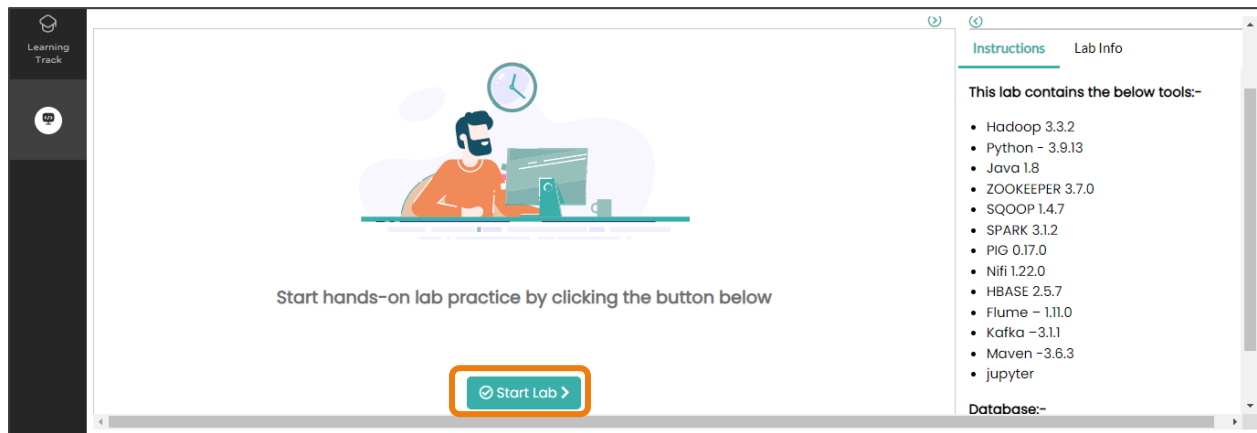**Step 2:** Click on **LAUNCH LAB**

**Step 3:** A small screen will pop up in the middle of your screen with important information about the lab. Again, click on the **LAUNCH LAB** button.
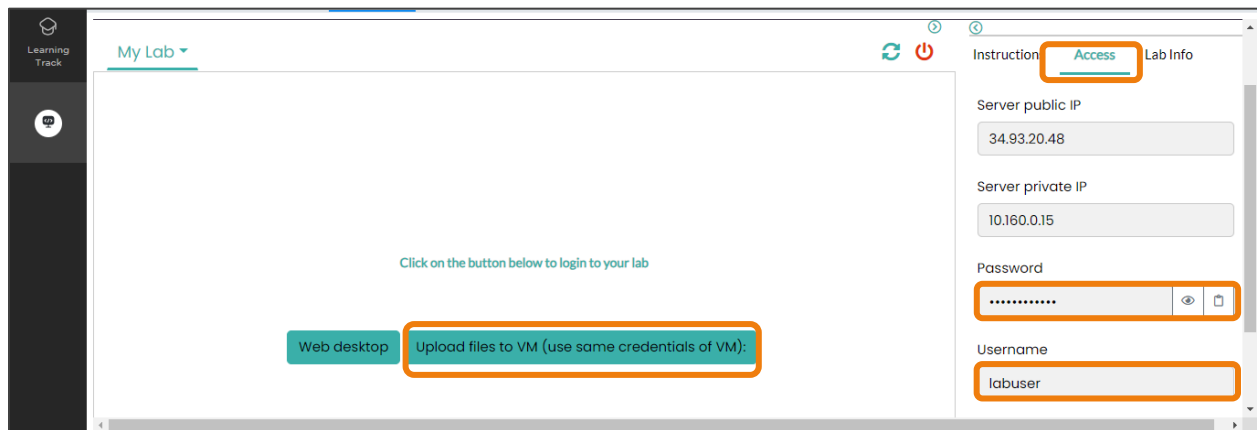


**Note:** It will take about three to five minutes for the lab environment to load.
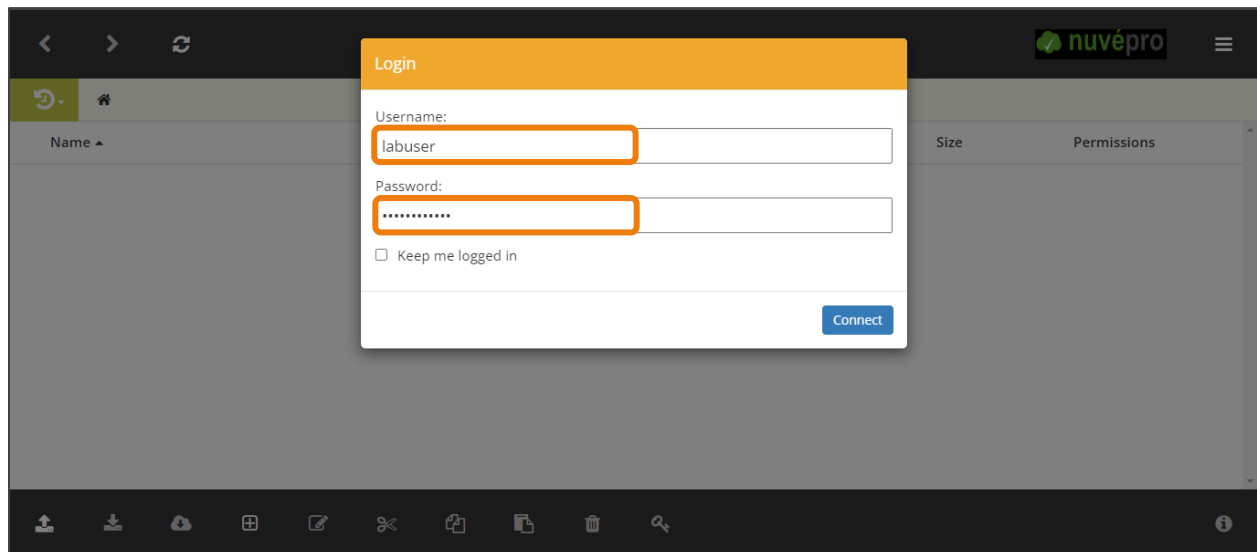
**Step 4:** Click on the **Start Lab** button



**Step 5:** Once the environment has loaded, log in to the **FTP (Upload files to VM)** using the specified **Username** and **Password** from the **Access** tab
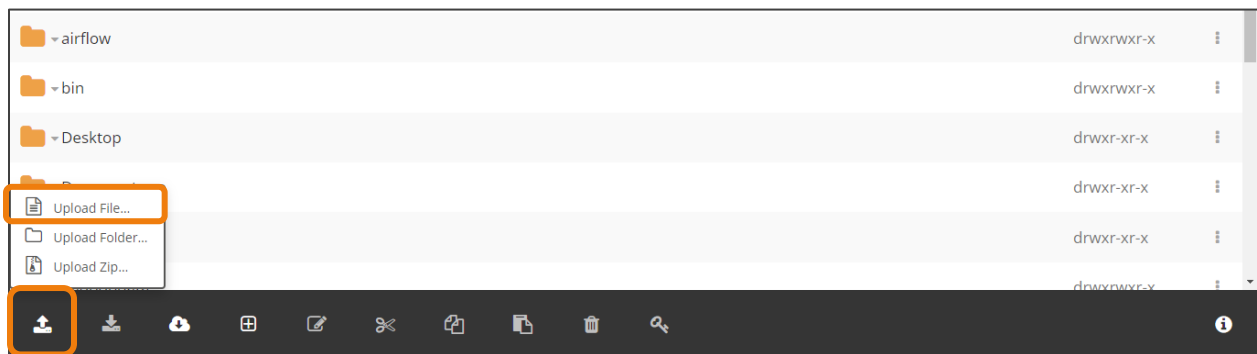
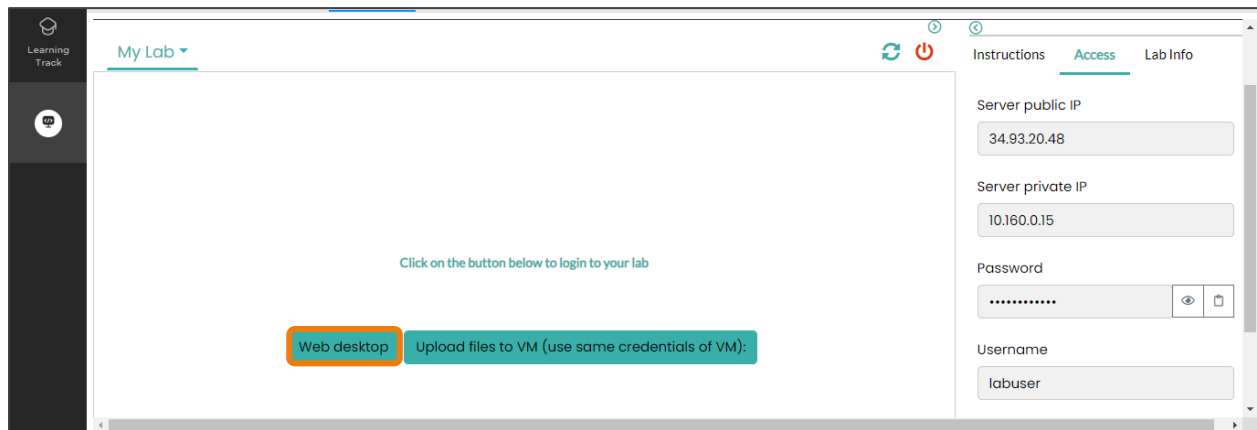**Step 6:** You will now be directed to the login screen where you can enter your **Username** and **Password**



**Note:** Once you are successfully logged in, you will be redirected to the below page as shown in Step 7.

**Step 7:** Click on the **Upload File** button to upload the dataset
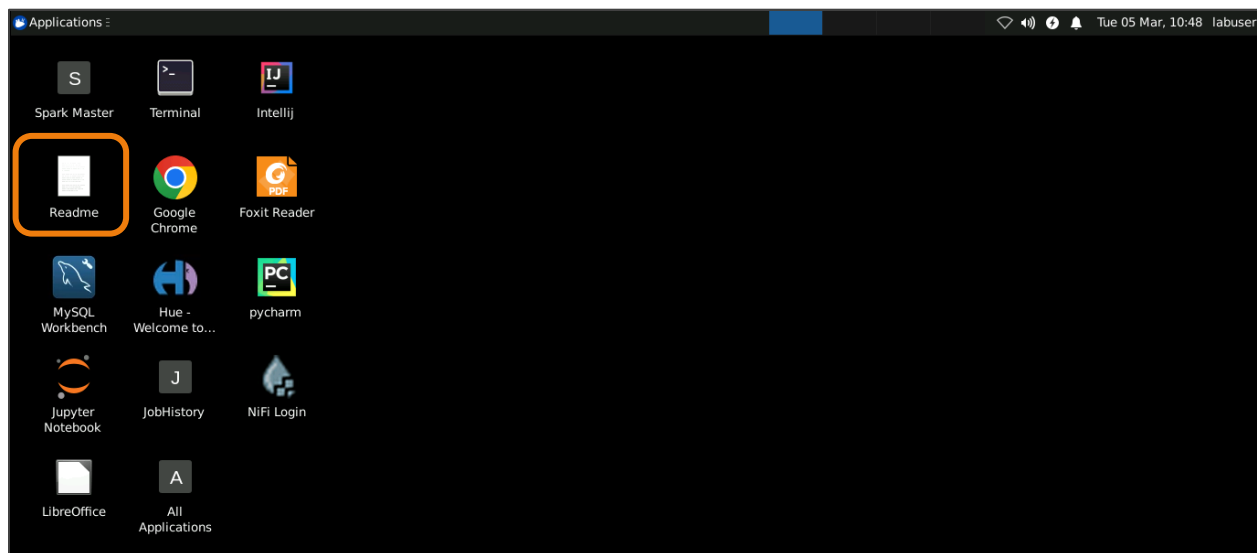


**Note: Datasets** are accessible in the lab for online use. Alternatively, for offline work, you can download and extract them from the **Reference Materials** in your LMS.

**Step 8:** To access the terminal, HUE, or any other provided service, click on **Web desktop**



**Step 9:** Click on the **Readme** file to discover the installed tools, their URL/path, and the Username and Password for different services

**Step 10:** Copy the **Username** and **Password** provided in Readme file to log in to **HUE**

```
=============
hadoop,spark,hbase :-

sudo systemctl start allservice.service
sudo systemctl stop allservice.service

hue:-

sudo systemctl start hue.service
sudo systemctl stop hue.service


To check hadoop daemons:-
---------------------
command :- jps

To start job history server :-
--------------------
go to cd /opt/hadoop/sbin/
              ./mr-jobhistory-daemon.sh start historyserver

click on jobhisory desktop icon

1) Hue :-

Url :- http://localhost:8132/
username :- admin
password :- admin
```

**Step 11:** Click on **HUE**

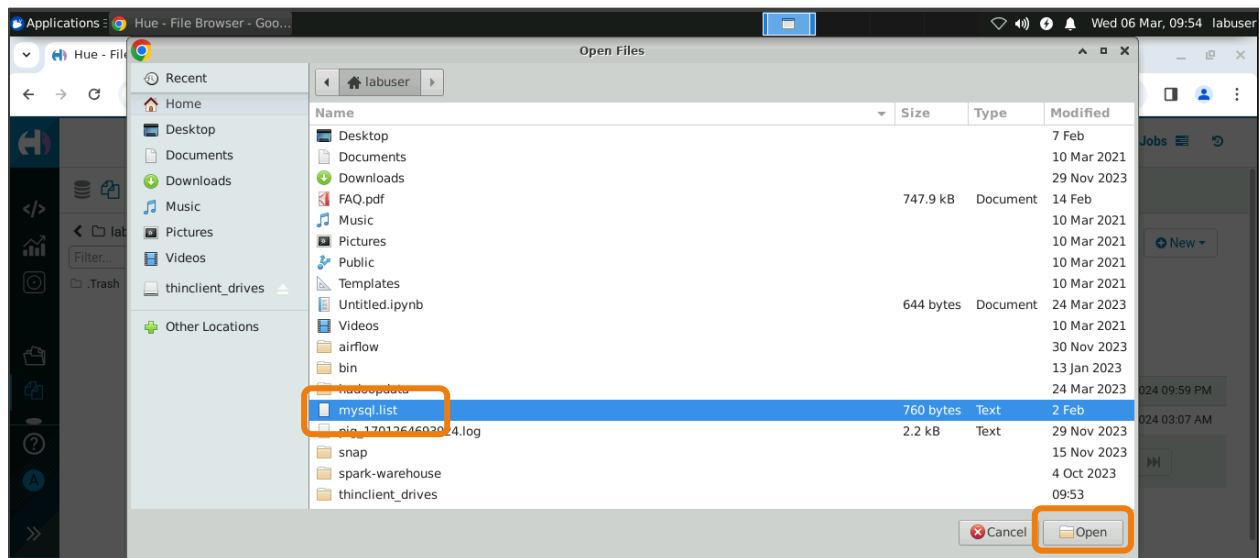**Step 12:** Enter the **Username** and **Password** you copied to log in to **HUE**



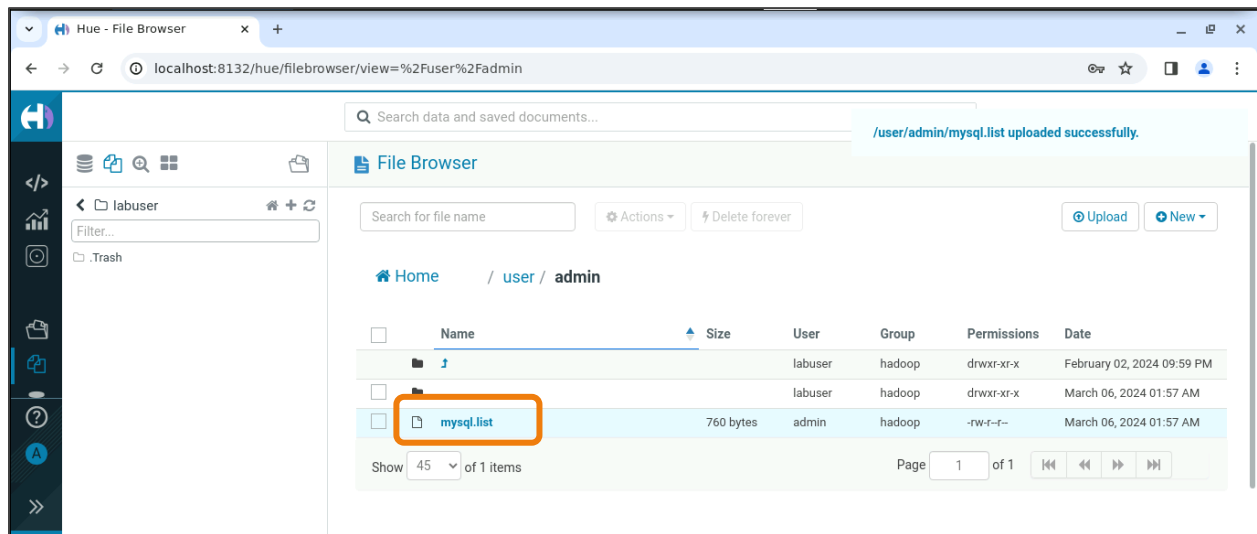**Note:** You will be navigated to the dashboard as shown below:

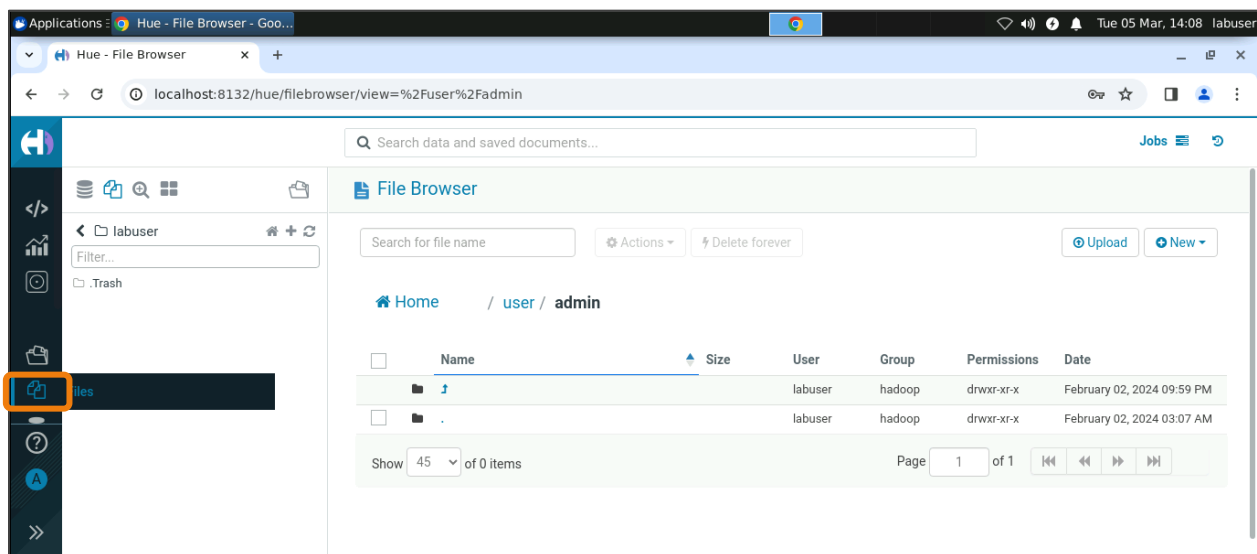**Step 13:** Click on the shown **+** icon to upload the dataset in HUE



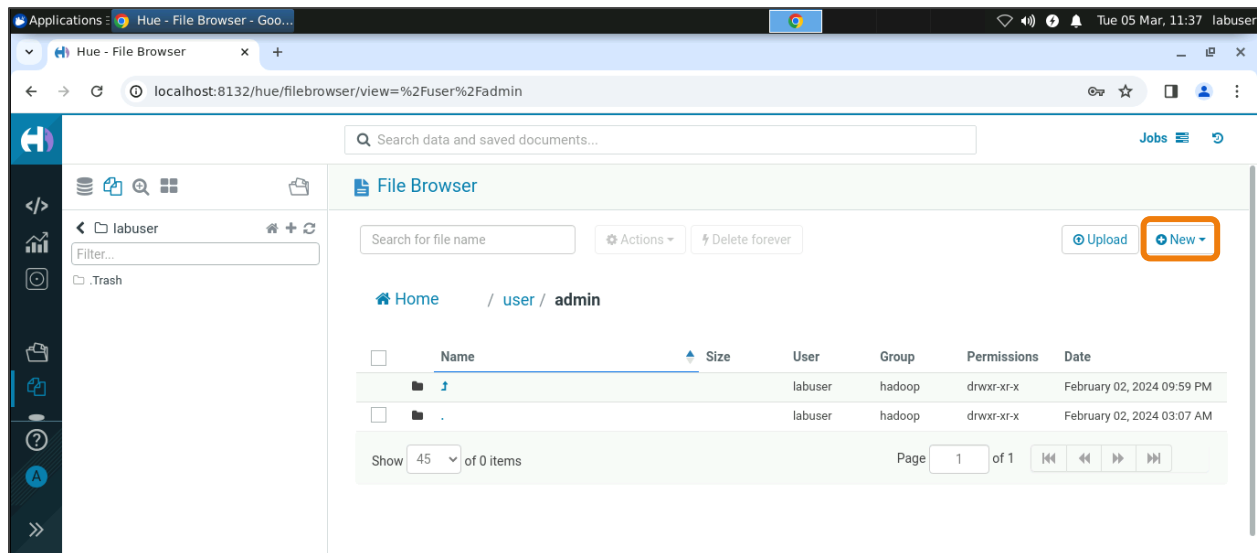**Step 14:** Choose the dataset you want to upload, and then click on the **Open** button

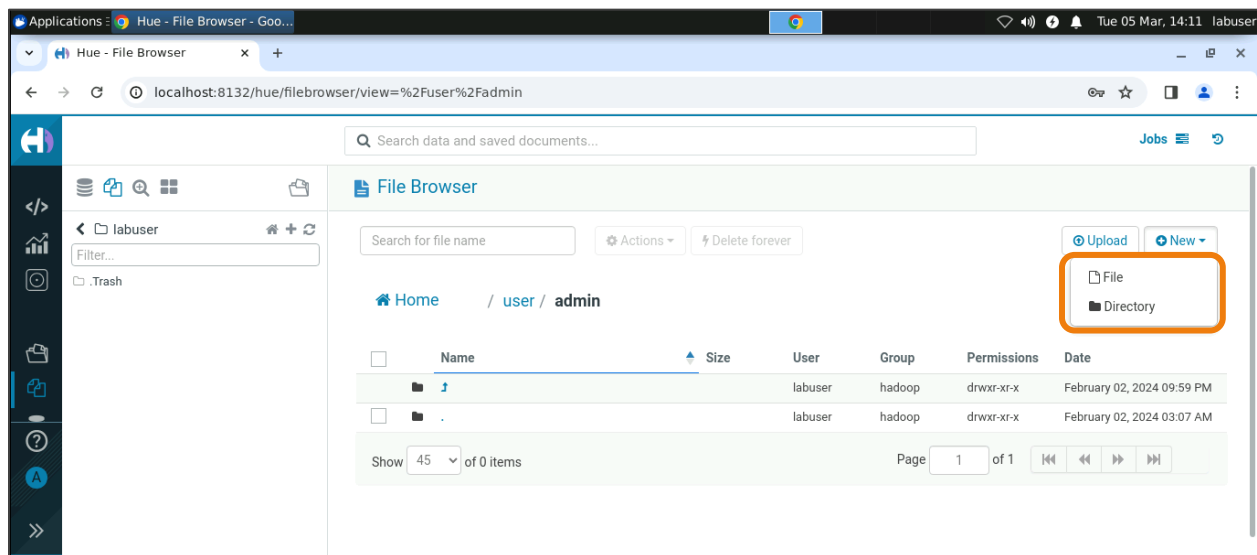**Note:** After clicking the **Open** button, you will be able to see the uploaded dataset as shown below:
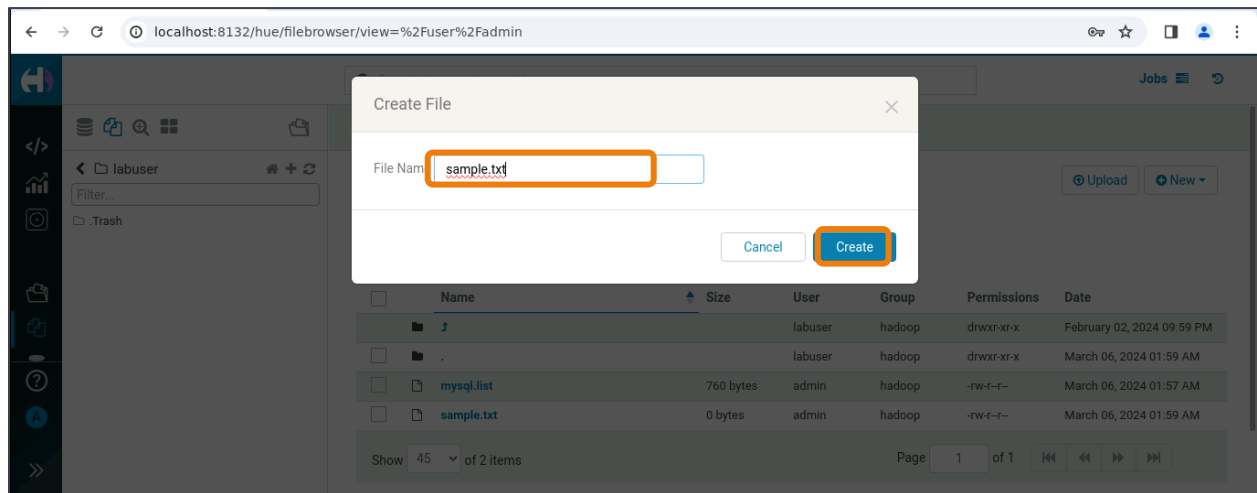


**Step 15:** Click on the **Files** button to create a new file or directory within HUE
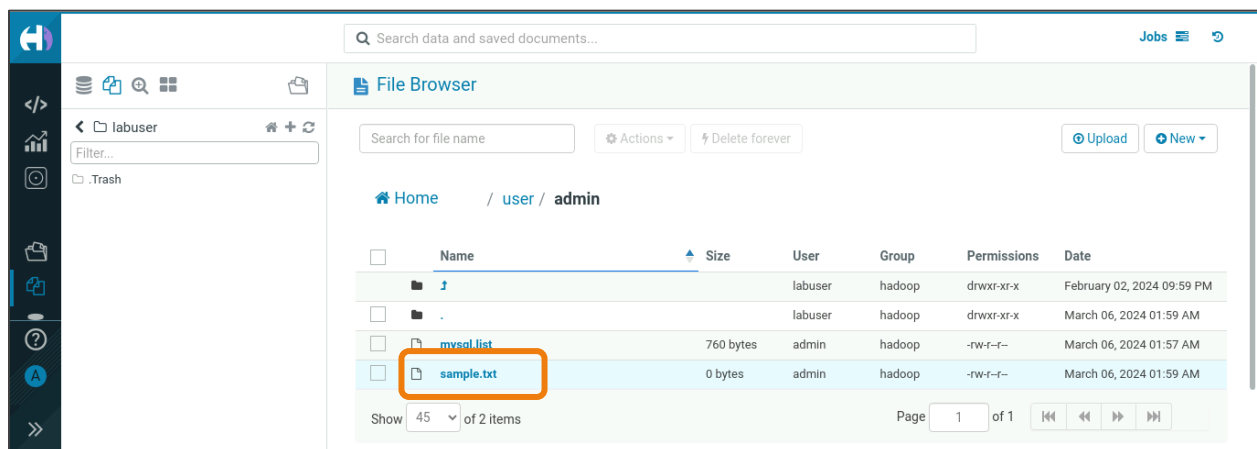
**Step 16:** Click on the **New** button
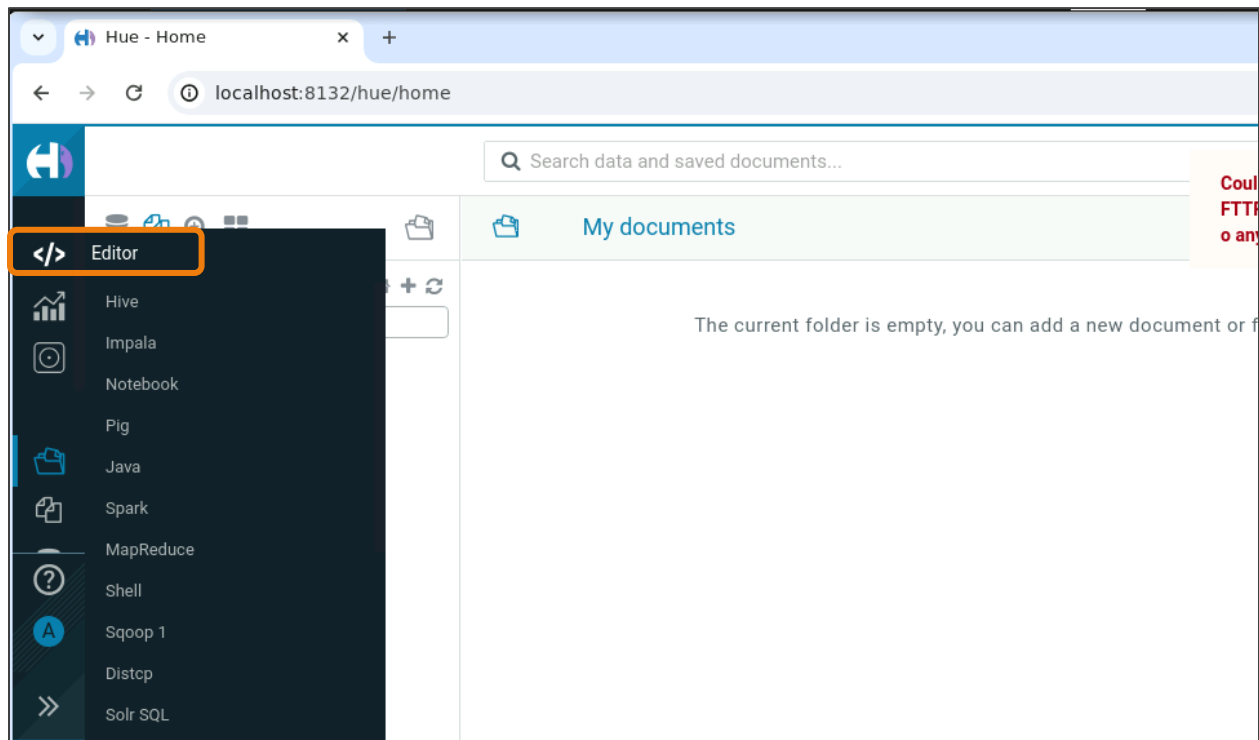


**Step 17:** Click on either **File** or **Directory** based on your needs

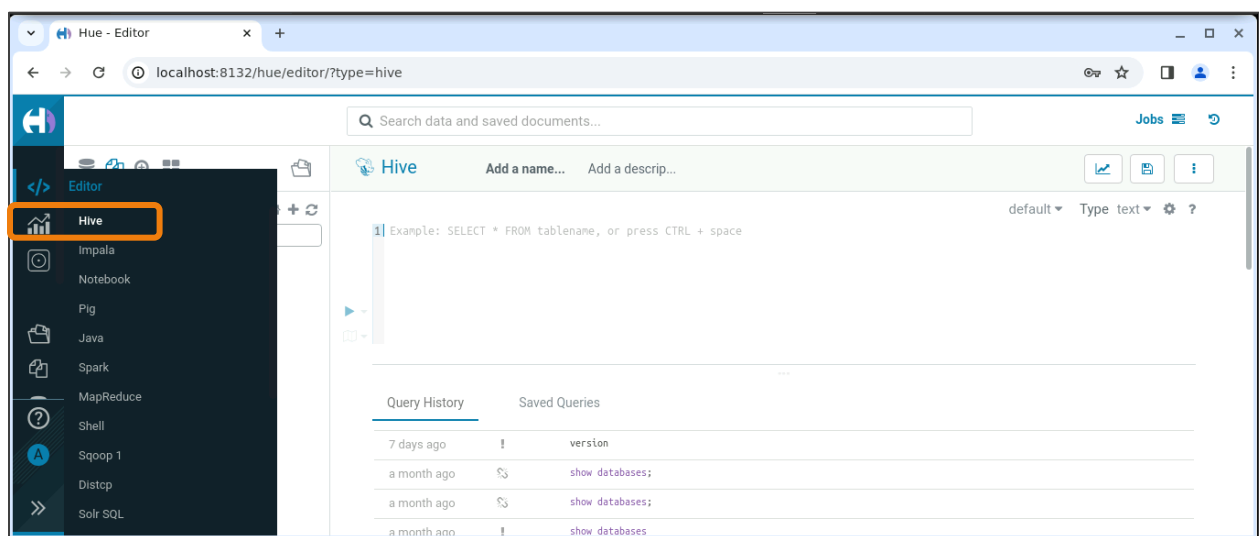**Step 18:** Enter the file name and click on the **Create** button



**Note:** After clicking the **Create** button, you will be able to see the new file as shown below:

**Step 19:** Select the **Editor** button situated on the left side of the HUE dashboard to compose the **Hive query**



**Step 20:** Click on **Hive**

**Step 21:** Type the command into the editor, and then click on the **triangle** button located on the left side of the editor to execute the command. The output will be displayed in the **Results** column.
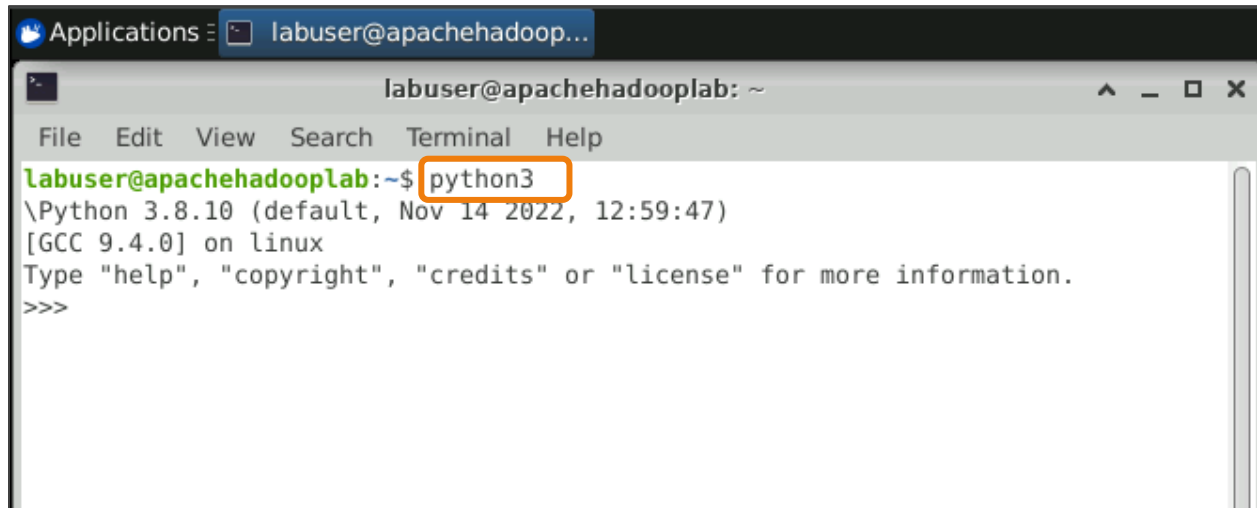
**Step 22:** Click on **Terminal** to open the terminal window



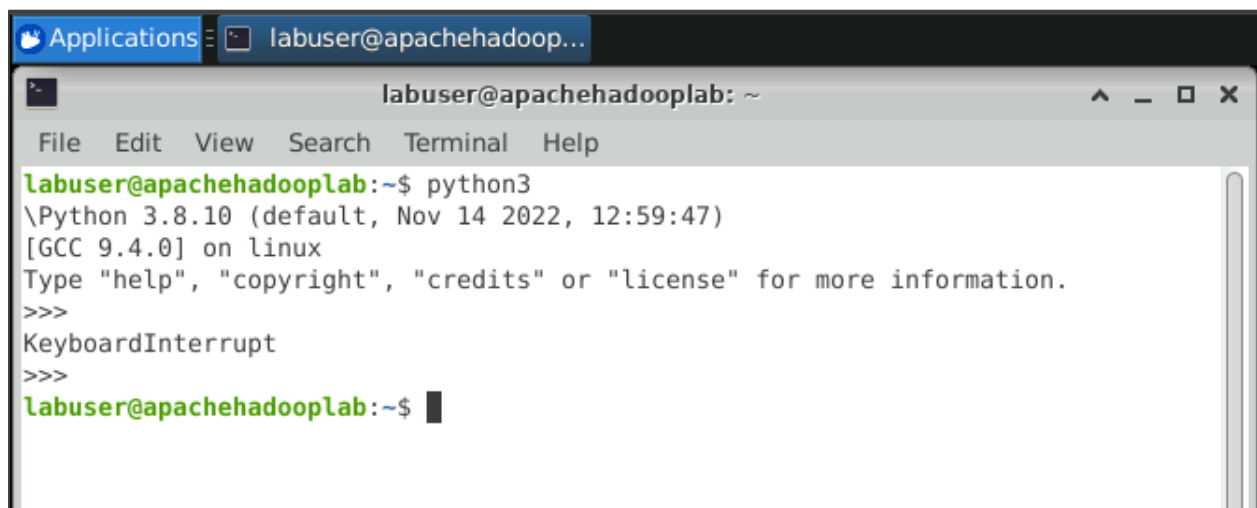**Note:** You will be able to see the new terminal file as shown below:

**Step 23:** To access the **PySpark console**, execute the following command:
**pyspark3**

**Step 24:** To access the **Python shell**, execute the following command:
**python3**



**Step 25:** Press **Ctrl+d** to exit the Python shell. This action will not be visibly displayed.

**Step 26:** To enter the **vi editor** and to write any Python file or txt file, use the command below:

**vi sample.py**

or

**vi sample.txt**
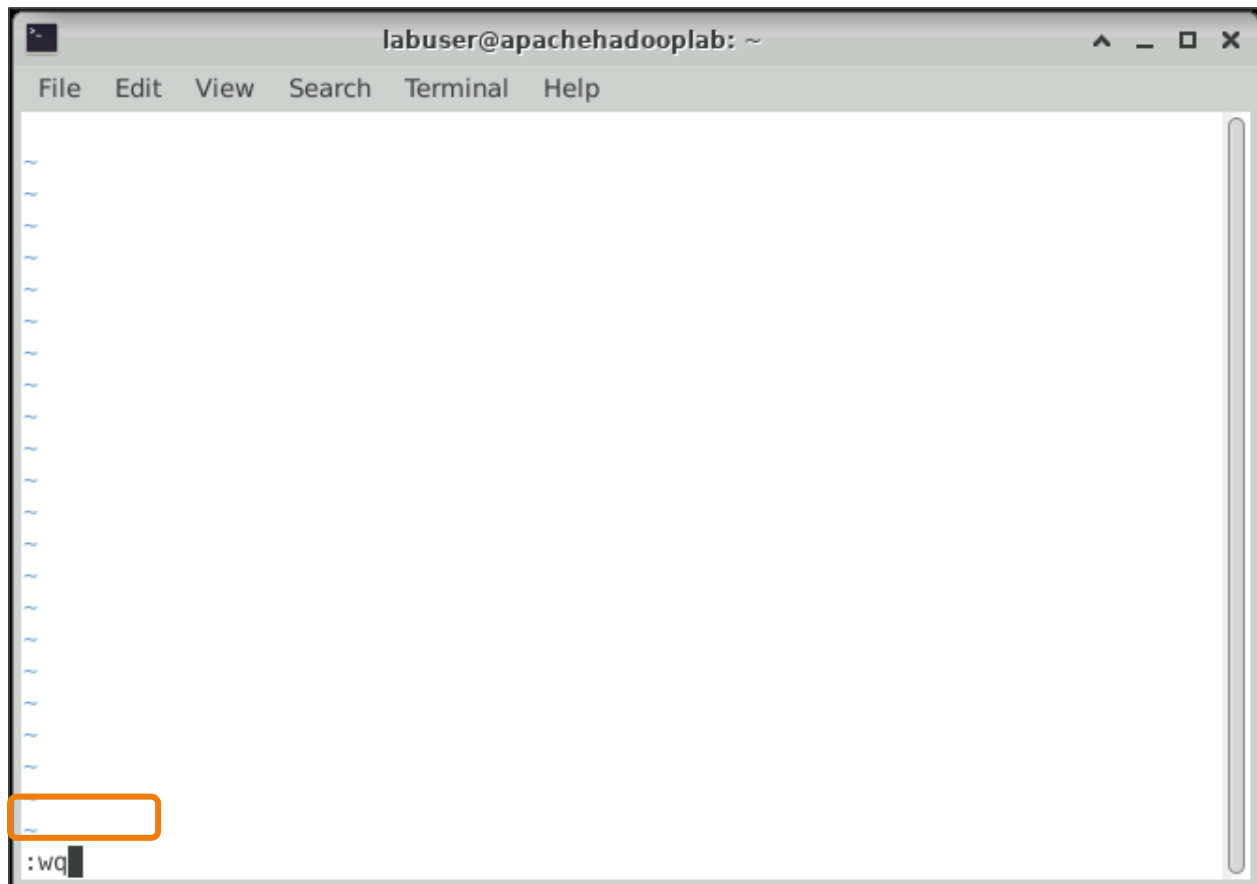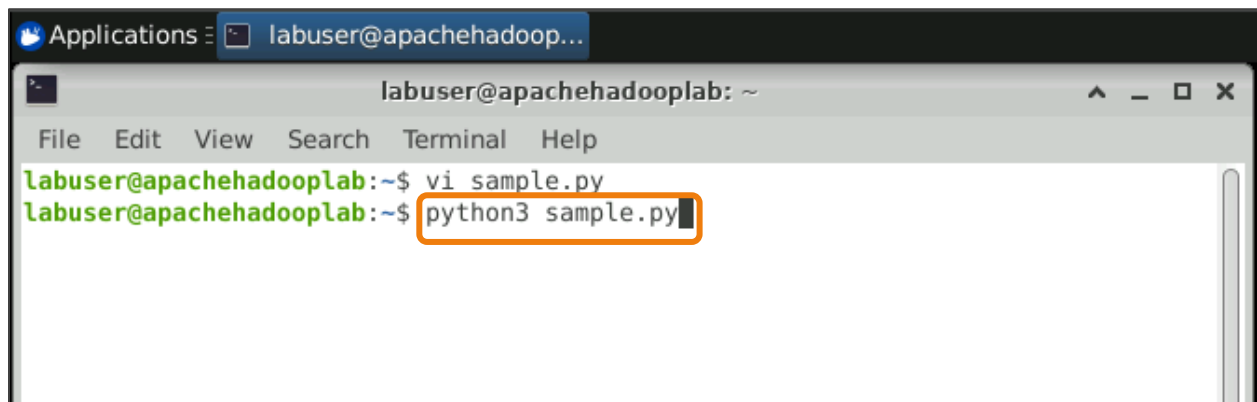
**Step 27:** Click on **i** on your keyboard to enter the **insert mode**

**Step 28:** Click on the **ESC** key and then type **:wq** to save and exit



**Step 29:** To execute the **Python script**, run the command below:
**python3 sample.py**

**Step 30:** To enter the **Scala environment,** execute the following command:
**spark3-shell**

**Step 31:** To activate the **JobHistory server**, the initial step involves launching the **jobhistory daemon**. To achieve this, simply copy the command specified in the Readme file and execute it within your terminal.

```
Tools Installed::-
================

Systemd service:-
=============
hadoop,spark,hbase :-

sudo systemctl start allservice.service
sudo systemctl stop allservice.service

hue:-

sudo systemctl start hue.service
sudo systemctl stop hue.service



To check hadoop daemons:-
----------------------
command :- jps

To start job history server :-
--------------------
go to cd /opt/hadoop/sbin/
          ./mr-jobhistory-daemon.sh start historyserver

click on jobhisory desktop icon

1) Hue :-

Url :- http://localhost:8132/
```
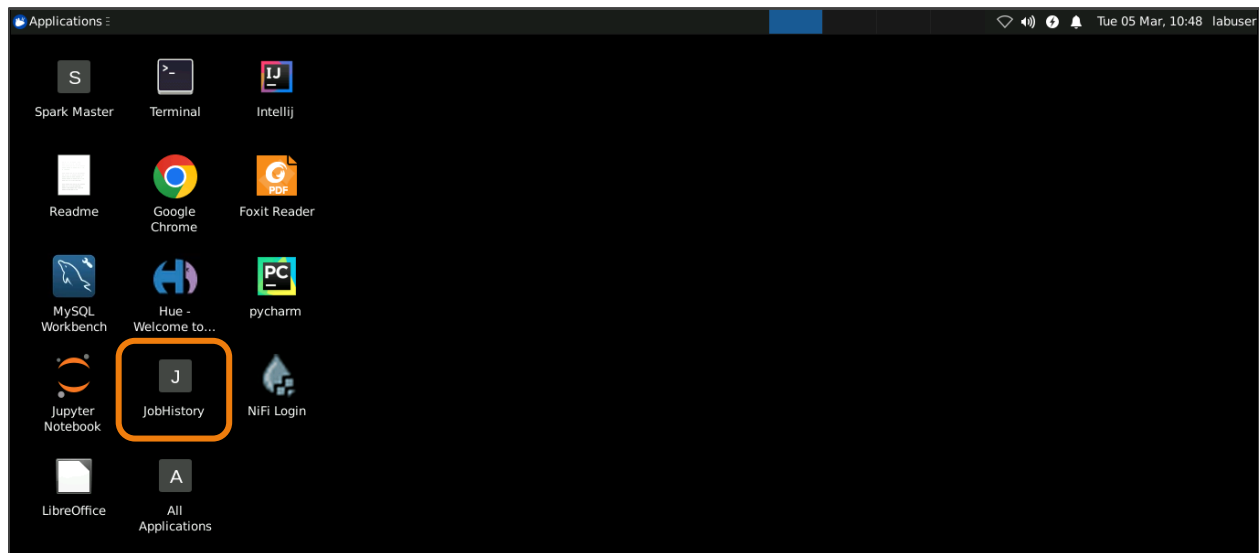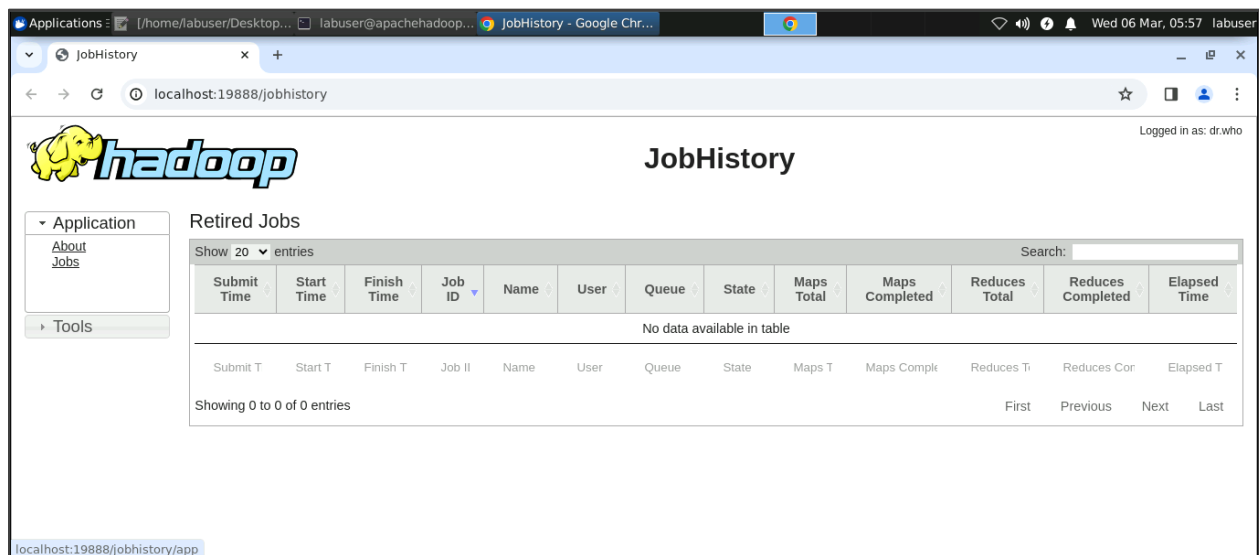
```
                    labuser@apachehadooplab: /opt/hadoop/sbin            ^ _ ▢ ✕
File   Edit   View   Search   Terminal   Help
labuser@apachehadooplab:~$ cd /opt/hadoop/sbin/
labuser@apachehadooplab:/opt/hadoop/sbin$ ./mr-jobhistory-daemon.sh start historyserver
WARNING: Use of this script to start the MR JobHistory daemon is deprecated.
WARNING: Attempting to execute replacement "mapred --daemon start" instead.
labuser@apachehadooplab:/opt/hadoop/sbin$
```
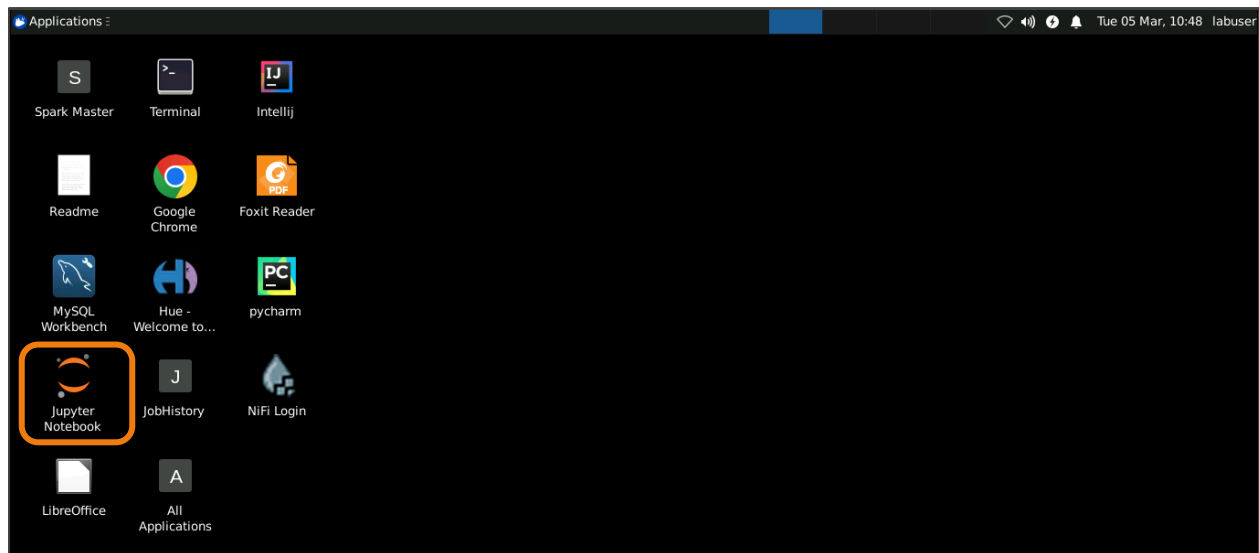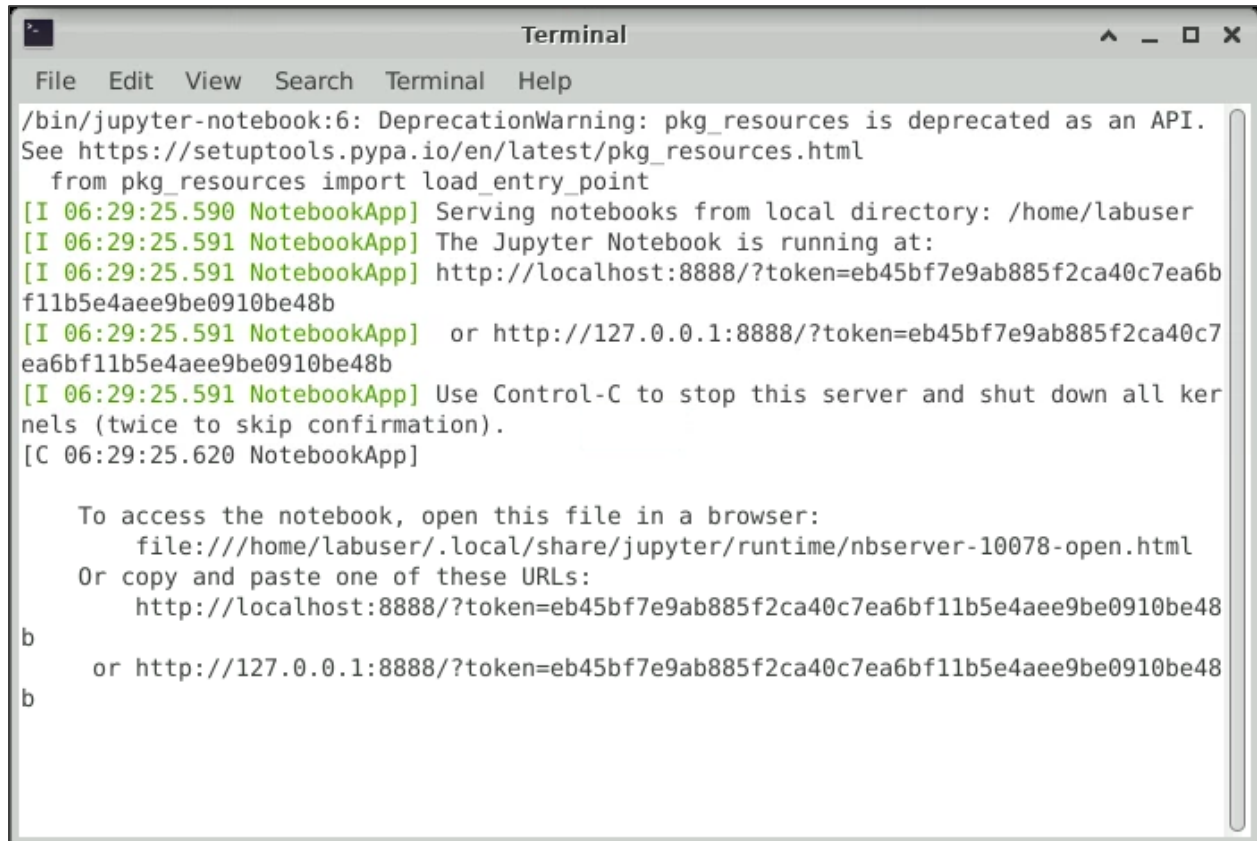
**Step 32:** Now, click on **JobHistory**



**Note:** After completing the activation of the jobhistory daemon, you can access the JobHistory dashboard to review your job history.

**Step 33:** Click on **Jupyter Notebook** to access the **Jupyter** environment

**Note:** You will be prompted with a terminal window, followed by the display of your Jupyter dashboard, where you can access your notebooks.