# Assisted Practice 15.4: Create an RDD in Spark

**Problem Scenario:** Create an RDD with a real-world retail business dataset of different categories.

**Objective:** In this demonstration, you will read the data from HDFS and print the distinct categories.

**Dataset Name:** **"part-m-00000"**

**Tasks to Perform:**

1.  Download the "part-m-00000" dataset from the categories folder from the course resource section and upload it into the HDFS using "Hue"
2.  Login into the "webconsole" and open the PySpark shell
3.  Create an RDD using textFile and update the path of the dataset
4.  Create a lambda function that will split the line
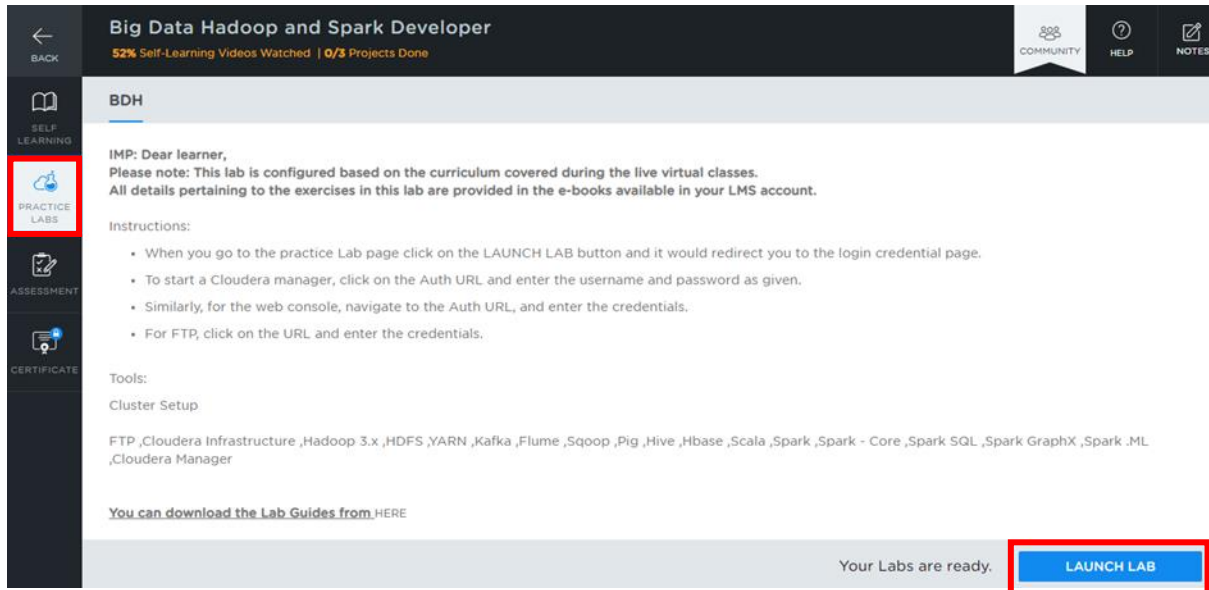5.  Print each element using collect() method

**Steps to Perform:**

**Step 1:** Download the dataset with the name **"part-m-00000"** from the course resources section.

**Step 2:** Log in to your LMS account

**Step 3:** Open the course **"Big Data Hadoop and Spark developer"**

**Step 4:** On the left side, click on the **"PRACTICE LABS"** tab and click on the **"LAUNCH LAB"** button

**Big Data Hadoop and Spark Developer**
52% Self-Learning Videos Watched | 0/3 Projects Done

**BDH**

IMP: Dear learner,
Please note: This lab is configured based on the curriculum covered during the live virtual classes.
All details pertaining to the exercises in this lab are provided in the e-books available in your LMS account.

Instructions:

- When you go to the practice Lab page click on the LAUNCH LAB button and it would redirect you to the login credential page.
- To start a Cloudera manager, click on the Auth URL and enter the username and password as given.
- Similarly, for the web console, navigate to the Auth URL, and enter the credentials.
- For FTP, click on the URL and enter the credentials.

Tools:

Cluster Setup

FTP ,Cloudera Infrastructure ,Hadoop 3.x ,HDFS ,YARN ,Kafka ,Flume ,Sqoop ,Pig ,Hive ,Hbase ,Scala ,Spark ,Spark - Core ,Spark SQL ,Spark GraphX ,Spark .ML ,Cloudera Manager

You can download the Lab Guides from HERE

Your Labs are ready.    **LAUNCH LAB**

**Step 5:** Again, click on the **"LAUNCH LAB"** button
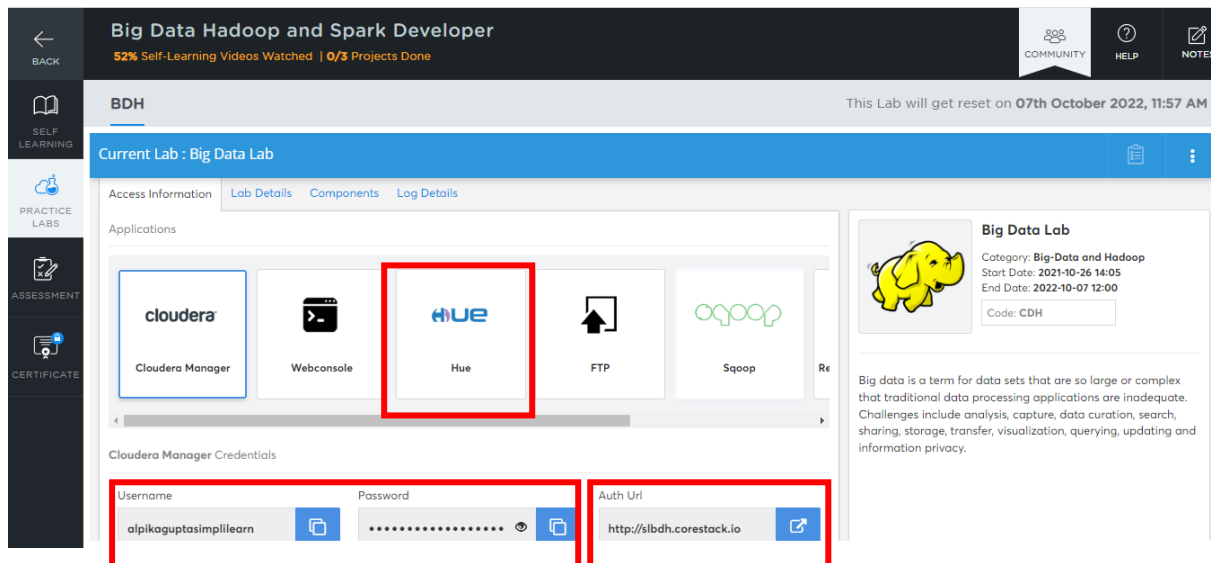


**Your Lab Setup is Ready!**

This Lab will get reset on
**30th September 2022, 2:33 PM**

To ensure continued practice, the reset date is automatically set to **120 days from the last lab start date**. You'll lose all the lab progress on reset activity. Take a backup of your code if you are going to be away for more than 120 days.
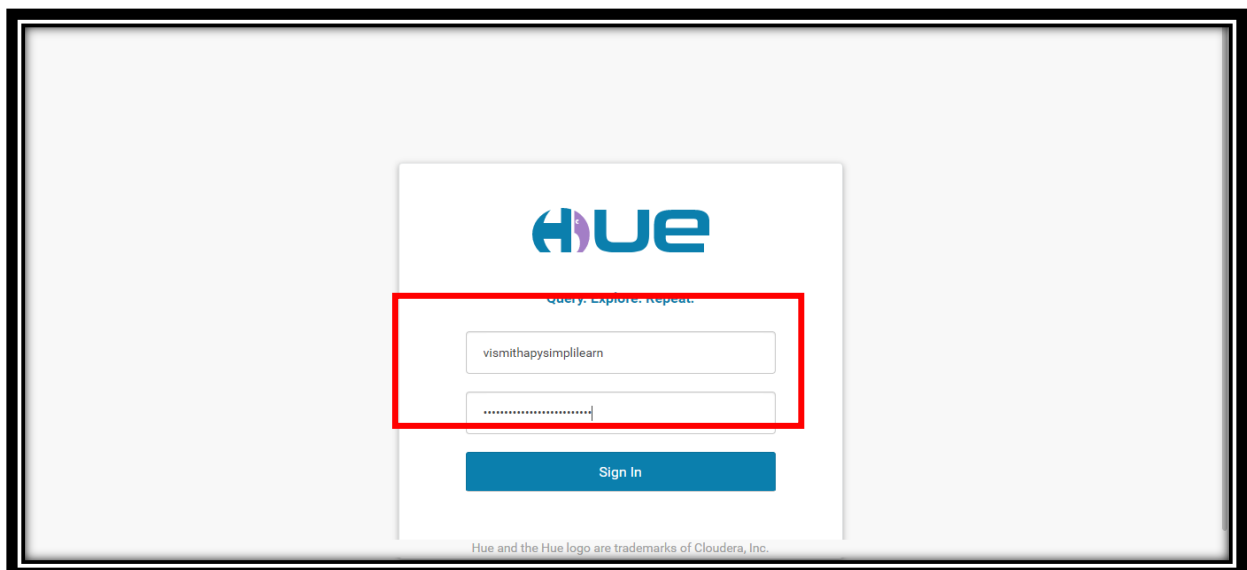
**LAUNCH LAB**

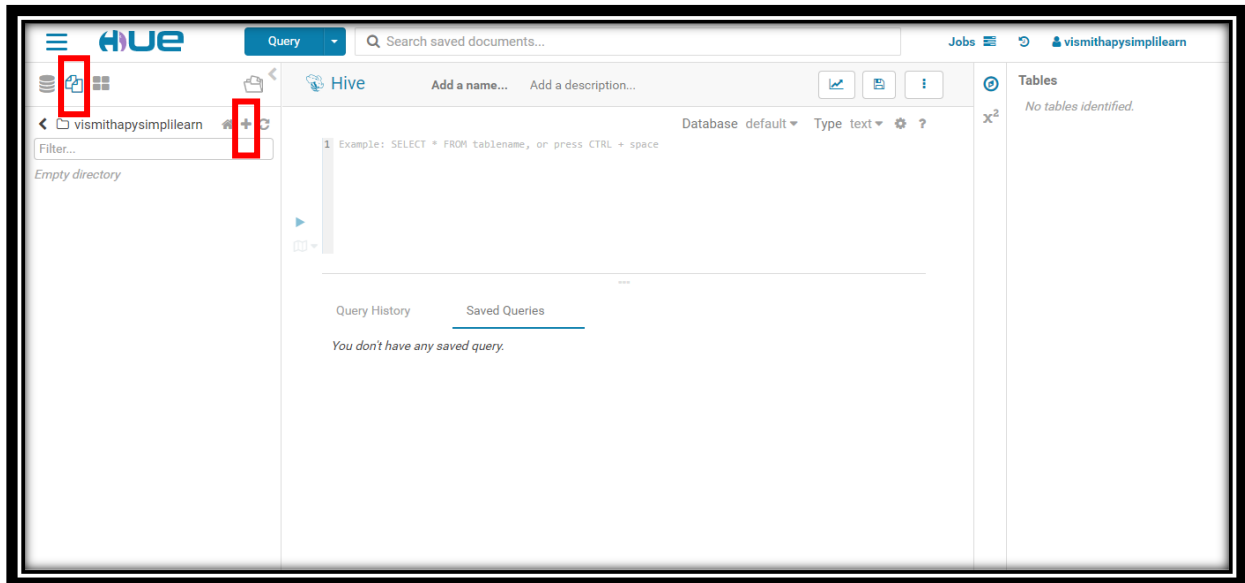**Note:** Practice labs are disabled on course expiry date by default.

**Step 6**: Click on **"Hue"** and click on the **"Auth Url"** to upload the dataset and copy the "**Username**" and the "**Password**" provided to log in to the **"Hue"**
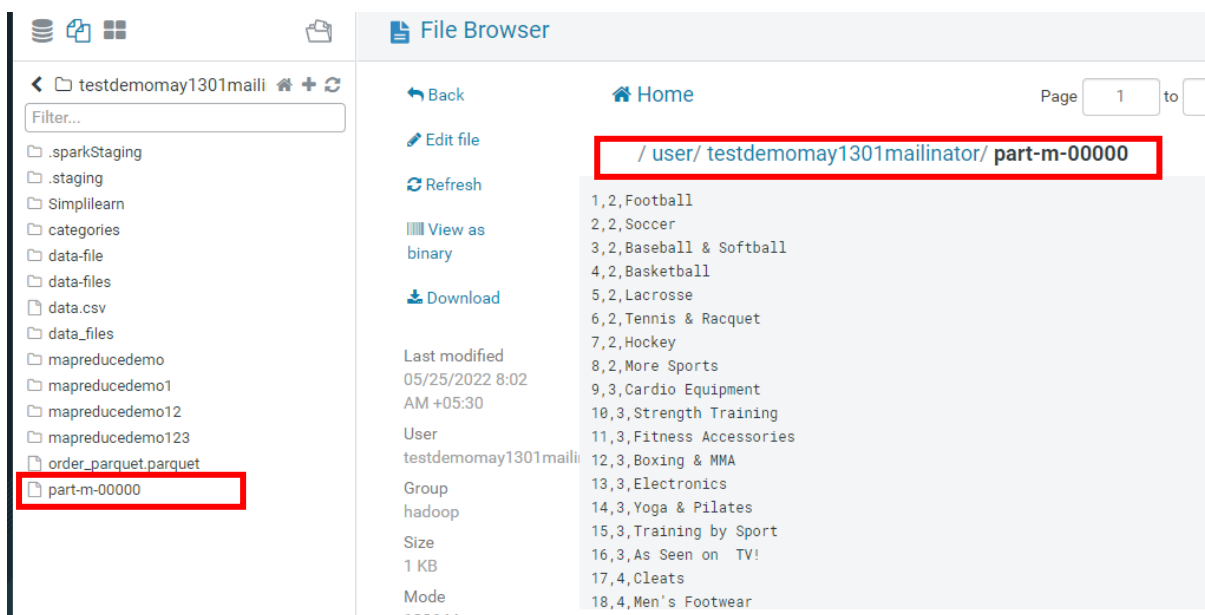


**Step 7:** Paste the **"Username"** and the **"Password"** on the login window and click on sign in



**Step 8**: Click on **"HDFS"** icon and click on the **"+"** symbol to upload the dataset

**Step 9**: Select the downloaded dataset file and upload it to **"HDFS."** In addition, by right-clicking, copy the path from the dataset that has been uploaded.



**Step 10:** Go back to the lab window and click on "**Webconsole**" and click on the **"Auth Url"**

**Step 11:** Copy the "**Username**" and the "**Password**" provided to log in to the

**"Webconsole"**



**Step 12:** Paste the "**Username**" and the "**Password**" on the console and click on

enter.

Note: The password will not be visible when pasted on the console

**Step 13:** Enter the **"PySpark"** console by running the below command.

      **Command:**

pyspark3



**Step 14**: Create an RDD using textFile and update the path of the dataset.

**Command:**

catRDD = spark.sparkContext.textFile("/ user/testdemomay1301mailinator/ part-m-00000 ")

**Step 15:** Create a lambda function that will split the line.

**Command:**

catRDD = catRDD.map(lambda line: line.split(",")[2]) \.distinct()

**Step 16:** Print each element using the collect() method. You will see the output as shown below:

**Command:**

for ele in catRDD.collect():

    print(ele)

```
>>>
>>> catRDD = spark.sparkContext.textFile("/user/testdemomay1301mailinator/part-m-00000")
>>> catRDD = catRDD.map(lambda line: line.split(",")[2]) \
...     .distinct()
22/05/25 02:42:25 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your
22/05/25 02:42:25 WARN shortcircuit.DomainSocketFactory: The short-circuit local reads featur
>>>
>>> for ele in catRDD.collect():
...     print(ele)
...
Football
Soccer
Baseball & Softball
Lacrosse
Hockey
Cardio Equipment
Electronics
```