# Assisted Practice 17: Data Exploration

**Problem Scenario:** Perform a data exploration and a descriptive analysis on the US companies' dataset.

**Objective:** In this demonstration, you will explore different commands to perform data exploration and descriptive analysis in PySpark.

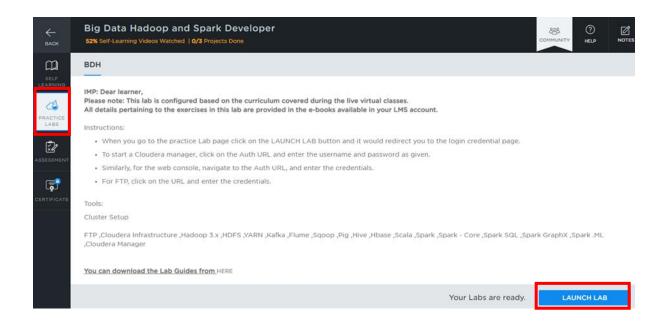**Dataset Name: "Fortune 500 Companies US.csv"**

**Steps to Perform:**

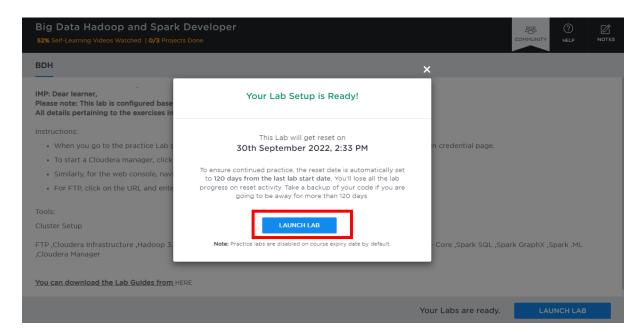**Step 1:** Download the dataset named **"Fortune 500 Companies US.csv"** from the course resources section

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

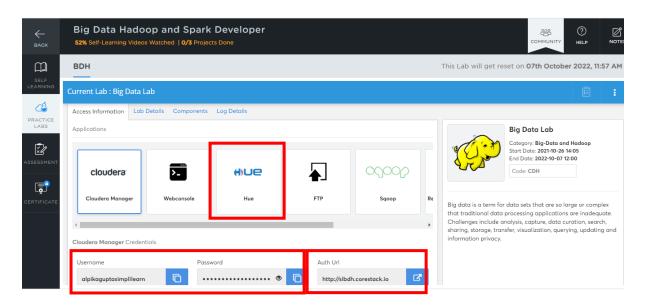**Step 3:** Open the course **"Big Data Hadoop and Spark developer"**

**Step 4:** On the left side, click on the **"PRACTICE LABS"** tab and then click on the **"LAUNCH LAB"** button
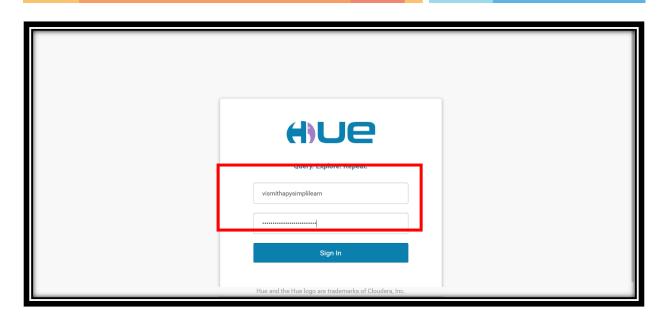
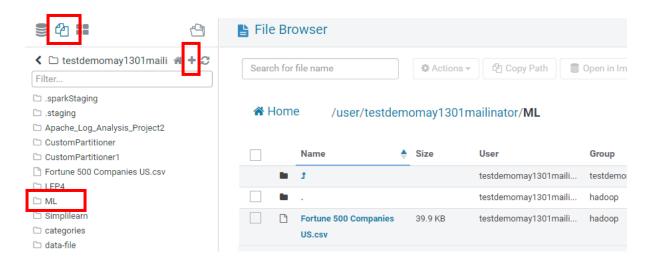**Step 5:** Again, click on the **"LAUNCH LAB"** button



**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"**
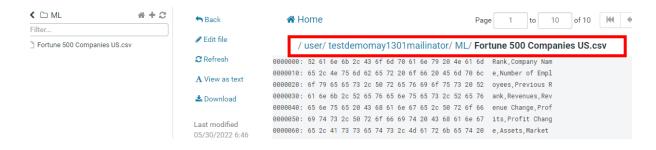


**Step 7:** Paste the **"Username"** and the **"Password"** on the log in window and click on **"Sign In"**

**Step 8**: Create a directory named **"ML"** and click on the **"HDFS"** icon and then on the **"+"** symbol to upload the dataset
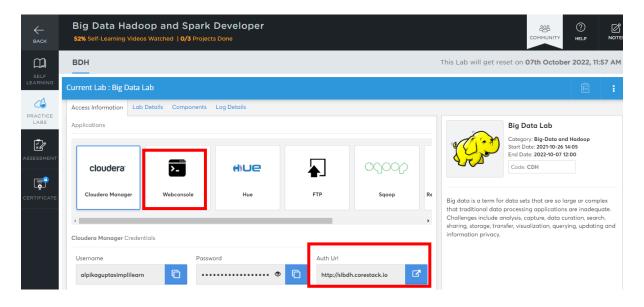


**Step 9:** Copy the path of the dataset that has been uploaded

**Step 10:** Go back to the lab window and click on the "**Webconsole**" and 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 "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**: Import the necessary modules

**Command:**

from pyspark import SparkConf, SparkContext

from pyspark.sql import SQLContext

**Step 15:** Create a Spark Session, and then create a DataFrame from a CSV file to load data

**Note:** The path should be provided to the ML folder.

**Command:**

sc = SparkContext = SparkSession \

    .builder \

    .appName("Simlilearn Examples") \

.getOrCreate() \

.sparkContext

companydata = spark.read.option("header", "true") \

.option("inferSchema","true") \

.csv("/user/testdemomay1301mailinator/ML")

```
>>> from pyspark import SparkContext
>>> from pyspark.sql import SparkSession
>>>
>>> # Create Spark Session.
... sc = SparkContext = SparkSession \
...     .builder \
...     .appName("Simplilearn Examples") \
...     .getOrCreate() \
...     .sparkContext
>>>
>>> companydata = spark.read.option("header", "true") \
...     .option("inferSchema", "true") \
...     .csv("/user/testdemomay1301mailinator/ML/")
```

**Step 16:** View the loaded data using the below command:

**Command:**

companydata.take(2)

```
>>> companydata.take(2)
[Row(Rank=1, Company Name='Walmart', Number of Employees='23,00,000', Previous Rank=1, Revenues='$4,85,873', Revenue Change='0.8%', Prof
nge='-7.2%', Assets='$1,98,825', Market Value='$2,18,619'), Row(Rank=2, Company Name='Berkshire Hathaway', Number of Employees='3,67,700
s='$2,23,604', Revenue Change='6.1%', Profits='$24,074.0', Profit Change='0.0%', Assets='$6,20,854', Market Value='$4,11,035')]
```

**Step 17:** To check the data type of every column of a DataFrame and to print the schema of the DataFrame in a tree format, you can use the below commands:

**Command:**

companydata.cache()

companydata.printSchema()

```
>>> companydata.cache()
DataFrame[Rank: int, Company Name: string, Number of Employees: string, Previous Rank: int, Revenues: string, Revenue Change: string
e: string, Assets: string, Market Value: string]
```

```
>>> companydata.printSchema()
root
 |-- Rank: integer (nullable = true)
 |-- Company Name: string (nullable = true)
 |-- Number of Employees: string (nullable = true)
 |-- Previous Rank: integer (nullable = true)
 |-- Revenues: string (nullable = true)
 |-- Revenue Change: string (nullable = true)
 |-- Profits: string (nullable = true)
 |-- Profit Change: string (nullable = true)
 |-- Assets: string (nullable = true)
 |-- Market Value: string (nullable = true)
```

**Step 18:** To perform a descriptive analysis of the company data you will use the below command:

**Command:**

companydata.describe()

```
>>> companydata.describe()
DataFrame[summary: string, Rank: string, Company Name: string, Number of Employees: string,
s: string, Profit Change: string, Assets: string, Market Value: string]
```

companydata.describe().toPandas().transpose()

```
>>> companydata.describe().toPandas().transpose()
                          0                  1                    2              3                  4
summary               count               mean               stddev            min                max
Rank                    500              250.5     144.4818327679989              1                500
Company Name            500               None                 None             3M      salesforce.com
Number of Employees     500               None                 None        1,00,300            98,800
Previous Rank           492  257.1117886178862    154.04809767869145              1                761
Revenues                500               None                 None     $1,00,288           $94,595
Revenue Change          500               None                 None              -             94.5%
Profits                 500               None                 None      $1,006.0            -$97.0
Profit Change           500               None                 None              -             99.7%
Assets                  500               None                 None     $1,00,245           $95,377
Market Value            500               None                 None     $1,00,595                 -
```