# Assisted Practice 16.2: UDF with DataFrame

**Problem Scenario:** Create a DataFrame and a Python function to convert it into UDF

**Objective:** In this demonstration, you will create a built-in function using UDF to convert the first letter of every word into uppercase.

**Tasks to Perform:**

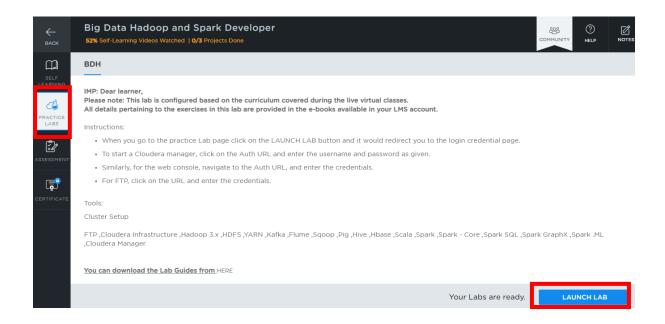1. Import required packages and create a DataFrame with two columns (S_No, Name)
2. Create a function convertCase() which will convert the first letter of every word into a capital letter
3. Convert a Python function to a PySpark UDF
4. Apply the convertUDF function to a DataFrame column as a built-in function

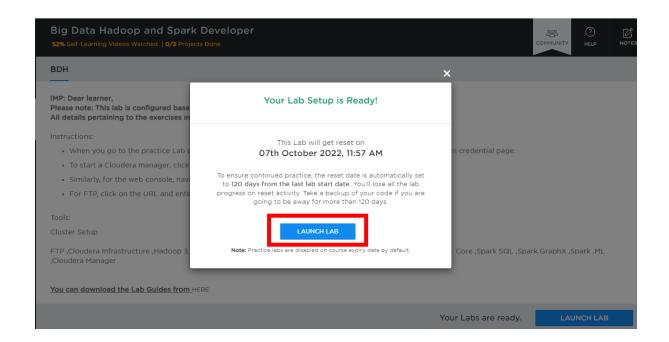**Steps to Perform:**

**Step 1:** Log in to your LMS account

**Step 2:** Open the course "**Big Data Hadoop and Spark developer**"

**Step 3:** On the left side, click on the "**PRACTICE LABS**" tab and click on the "**LAUNCH LAB**" button

**Step 4:** Again, click on the "**LAUNCH LAB**" button



**Step 5**: Click on "**Webconsole**" and click on the "**Auth Url**"

**Step 6:** Copy the **Username** and the **Password** provided to log in to the Webconsole

**Step 7:** 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 8:** Login into the PySpark shell using the below command:

**Command:** pyspark3

```
bdh-cluster2-edgenode10 login: testdemomay1301mailinator
Password:
Last login: Wed Jun  8 17:10:24 on pts/1
     | |            | |
     | |            | |
    _| |     __  _ _| |___  ___  __   __
   /   | |  /  \ | | |  /   _\ | ' .  __ '_|
  |  (   ) |  |_| | |  (   | | | | | (_| |
   \__|_|\__/ \_,_|\_,_|\_\___| \_,_|
==================================================
     *                  :

Password for testdemomay1301mailinator@BDH-ENV.GNE4-RUTX.CLOUDERA.SITE:
[testdemomay1301mailinator@bdh-cluster2-edgenode10 ~]$ pyspark3
```

**Step 9**: Import the required libraries and create a DataFrame with two columns (S_No, Name)

```
>>> import pyspark
>>> from pyspark.sql import SparkSession
>>> from pyspark.sql.functions import col, udf
>>> from pyspark.sql.types import StringType
```

```
>>> spark = SparkSession.builder.appName('SparkByExamples.com').getOrCreate()
>>> columns = ["S_No","Name"]
>>> data = [("1", "stephan petit"),
...        ("2", "audrey smith"),
...        ("3", "ray sanders")]
>>>
>>> df = spark.createDataFrame(data=data,schema=columns)
>>> df.show(truncate=False)
+----+-------------+
|S_No|Name         |
+----+-------------+
|1   |stephan petit|
|2   |audrey smith |
|3   |ray sanders  |
+----+-------------+
```

**Step 10:** Creates a function convertCase() which takes a string parameter and converts the first letter of every word to capital letter

```
>> def convertCase(str):
.       resStr=""
.       arr = str.split(" ")
.       for x in arr:
.           resStr= resStr + x[0:1].upper() + x[1:len(x)] + " "
.       return resStr
.
```

**Step 11:** Convert the function convertCase() to convertUDF

```
...
>>> convertUDF = udf(lambda z: convertCase(z))
>>>
```

**Step 12:** Now, you can use convertUDF on a DataFrame column as a regular built-in function

```
>>> df.select(col("S_No"), \
...      convertUDF(col("Name")).alias("Name") ) \
...    .show(truncate=False)
+----+--------------+
|S_No|Name          |
+----+--------------+
|1   |Stephan Petit |
|2   |Audrey Smith  |
|3   |Ray Sanders   |
+----+--------------+
```