# DataBricks Exercise Manual

## Exercise 1

## Creating your DataBricks Community Workspace

1. Sign up DataBricks community using the below URL
   https://databricks.com/try-databricks?_ga=2.231253720.636689393.1585744376-1730487577.1573900075

2. Please provide your Work Email ID while signing up



3. Once you sign up, you will be taken to the page in which you have to select Free Trial or Community Edition. Please select "Get Started" option under the community Edition.

4. After you clicked "Get Started" within few minutes, you would be receiving an email in your Inbox for verification. Please verify the link provided in the Email. You will be directed to a page in which you have to provide your password for login.



5. After providing your password, you will be taken to your DataBricks Workspace

# Exercise 2

## Creating a cluster in your DataBricks Workspace

1. Click on Cluster option at the left side of your DataBricks Workspace



2. You will find the list of available cluster. Currently you will have no cluster available, Please click on "Create Cluster" option available at the top.



3. Then you need to provide the cluster name. Make sure the Runtime is 6.4 and you will be noticing the cluster size is 1 driver with 15.3GB of RAM, 2 Cores and 1 DataBricks Unit. Finally click on create cluster button.

Create Cluster

**New Cluster**   [Cancel]   [Create Cluster]   **0 Workers:** 0.0 GB Memory, 0 Cores, 0 DBU
                                                **1 Driver:** 15.3 GB Memory, 2 Cores, 1 DBU ⓘ

Cluster Name

[MyFirstCluster]

Databricks Runtime Version ⓘ

[Runtime: 6.4 (Scala 2.11, Spark 2.4.5)                    | ∨]

[New] This Runtime version supports only Python 3.

Instance

Free 15GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours.
For more configuration options, please upgrade your Databricks subscription.

**Instances**   Spark

Availability Zone ⓘ

[us-west-2c                                                | ∨]

4. Your cluster will be in pending state for some time and then it will get running status

Clusters

[ **+ Create Cluster** ]

▼ All-Purpose Clusters

| Name | State | Nodes | Driver | Worker | Runtime | Creator |
|------|-------|-------|--------|--------|---------|---------|
| ⟳ MyFirstCluster | Pending ⓘ | 0 | Community O… | Community O… | 6.4 (includes Ap… | navaneeth@d |

▼ Job Clusters

Clusters

[ **+ Create Cluster** ]

▼ All-Purpose Clusters

| Name | State | Nodes | Driver | Worker | Runtime | Creator |
|------|-------|-------|--------|--------|---------|---------|
| ● MyFirstCluster | Running | 1 (0 spot) | Community O… | Community O… | 6.4 (includes Ap… | navaneeth@d |

▼ Job Clusters

# Exercise 3

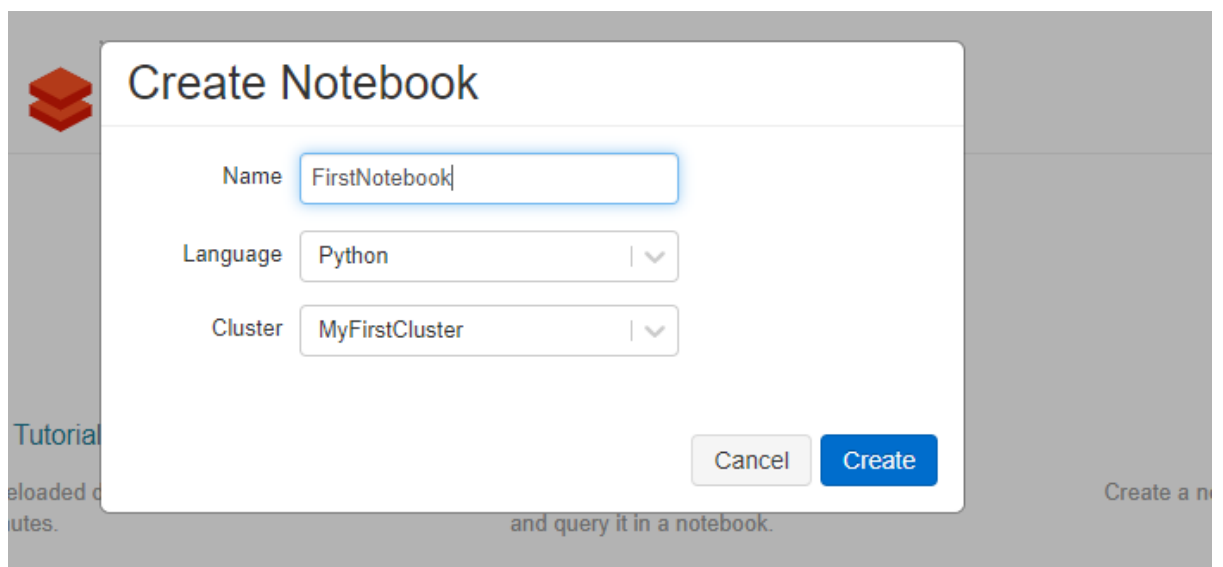## Creating a notebook and fundamentals of notebook in DataBricks Workspace

1.  On the DataBricks home page, you can use either, New Notebook options or Create a Blank Notebook option to create a notebook
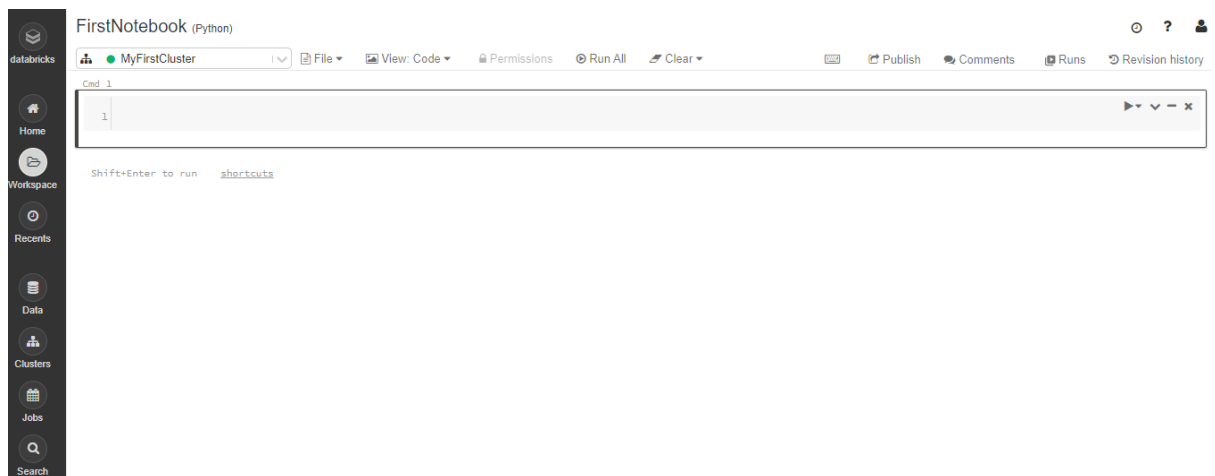


2.  Provide a name for the notebook, Select Python, Scala, SQL or R based upon the execution environment and also select eh cluster.

    For current practice select "Python" and the cluster you have created recently.

    Then click on "Create" Option.

3. An empty notebook will be created



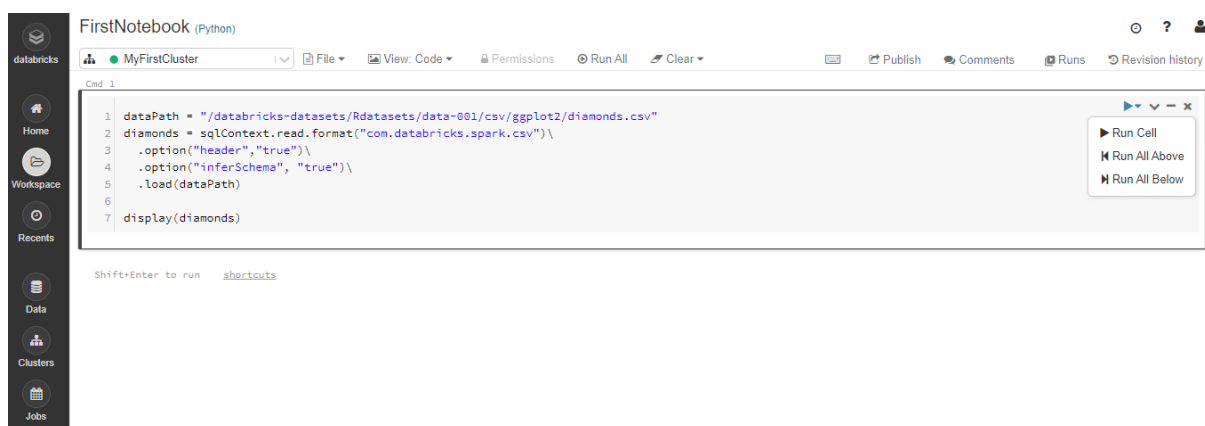4. Add the below lines of code in the cell provided in the databrick

```
dataPath = "/databricks-datasets/Rdatasets/data-001/csv/ggplot2/diamonds.csv"
diamonds = sqlContext.read.format("com.databricks.spark.csv")\
  .option("header","true")\
  .option("inferSchema", "true")\
  .load(dataPath)

display(diamonds)
```

We are using the sample data provided by the databricks for practice
In the above code, we are parsing a CSV, identifying the schema and displaying the data.
We will be learning RDD, DataFrame, SparkSQL in next topics, this is just for demo purpose of notebook, we will be learning more clear explanations in upcoming topic.



5. Click on Run Cell option to execute the code. You will see the spark cluster running the job. Once the job execution is done, you will get the result like below

6. At the bottom of the notebook, we have options for different data visualization options, you can select any one of that now, you may witness meaning less graph. We can see graph with meaningful data while we practice spark.



7. At the top, you have publish option. You can publish the notebook for future reference.
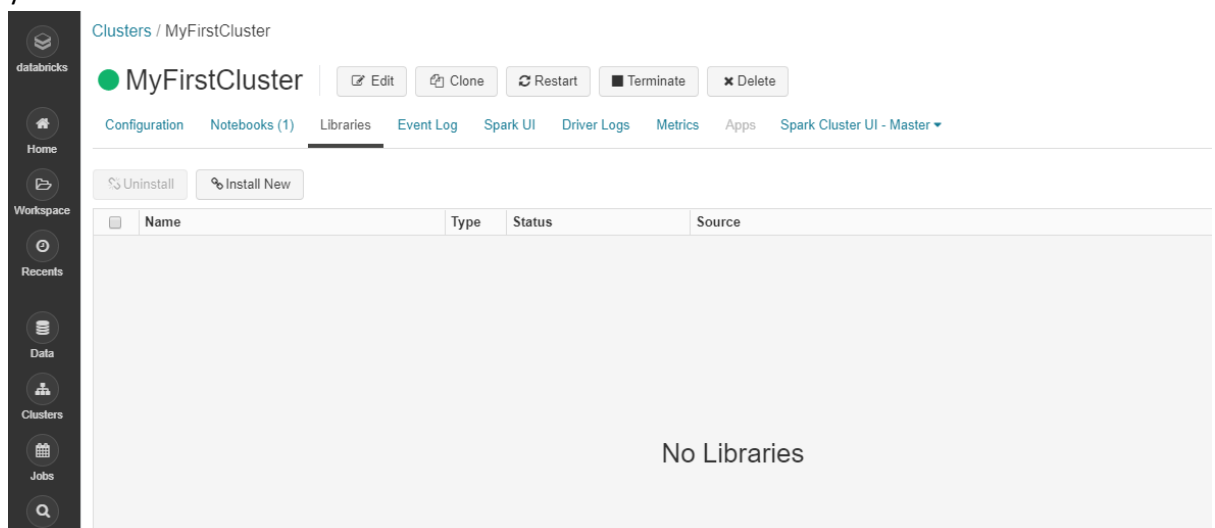
8. You can also see the revision history in notebook about the previous cell executions
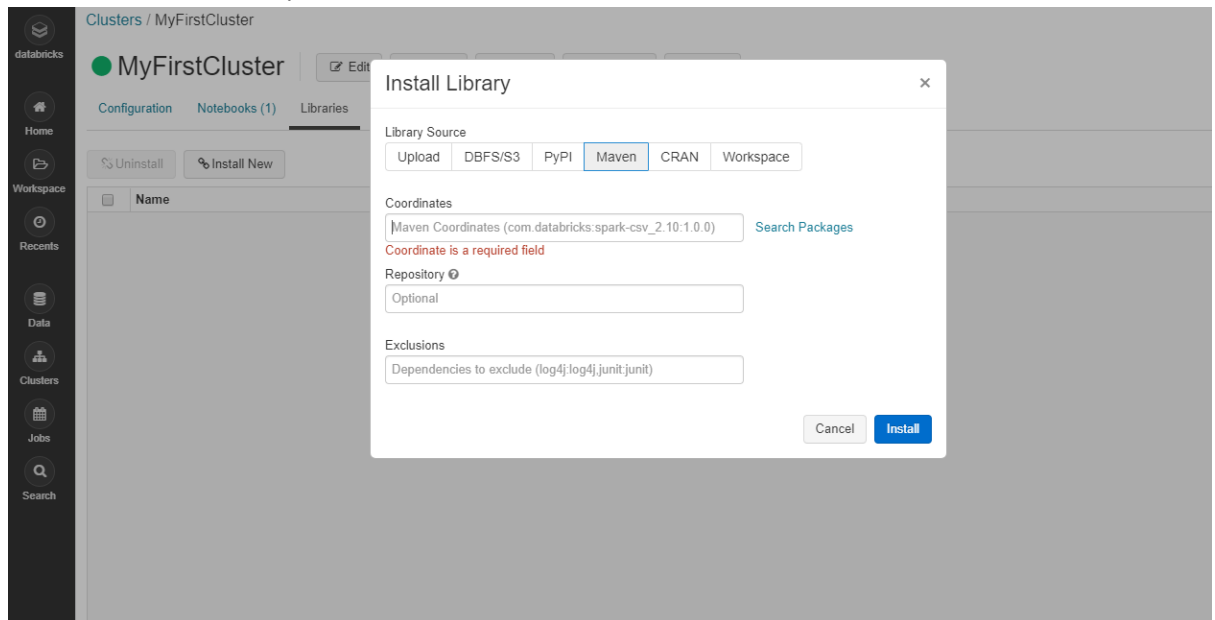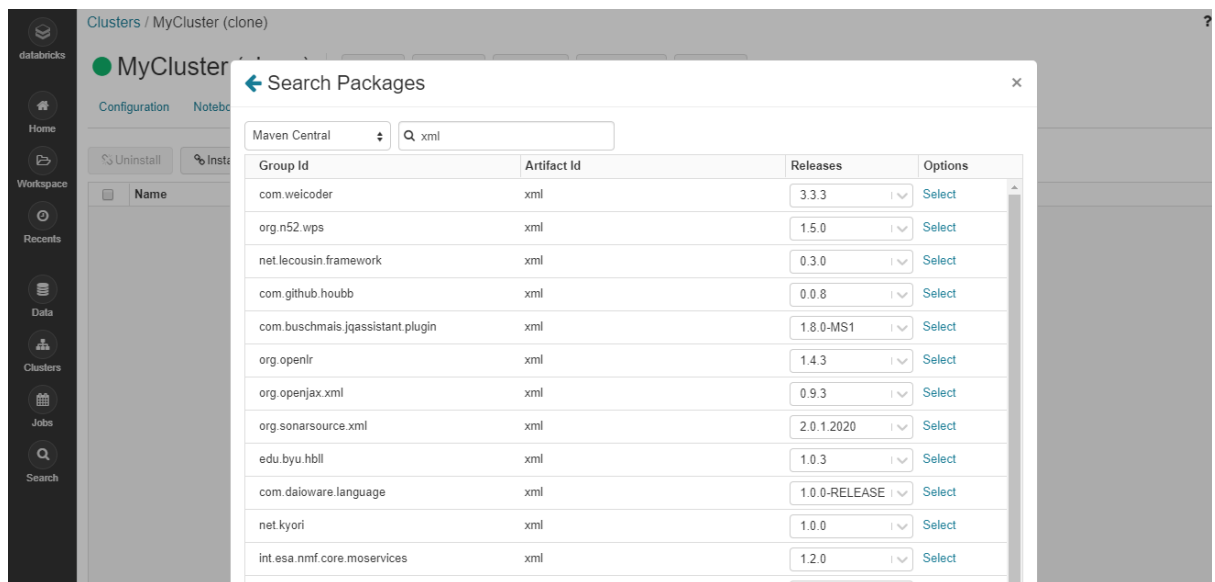


# Exercise 4

# Adding a library in the cluster

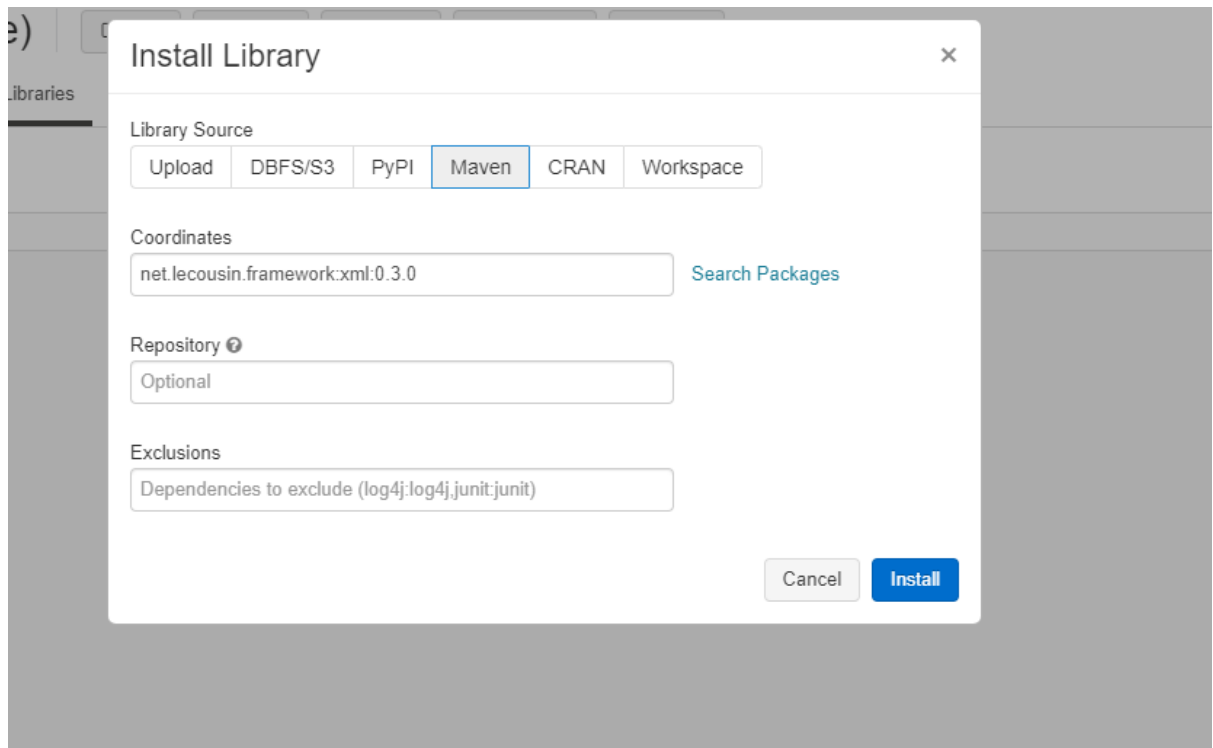1. Click on the cluster option from the left side bar. Click on the libraries option available when you scroll over.

2. Click on "Install New" option available and choose maven



3. Click on "Search Packages" and wait for some time. The libraries will be retrieved for display. Then you can select the library.



4. Once the library is selected, click on install option

5. The library will be available in the cluster in sometime