# Setting up Azure Databricks Workspace & Cluster

## Exercise 0 – Setup Azure Data Lake Gen2 Account
1. Complete Lab 1 – Working with Azure Data Lake Gen2 account.
2. Upload Files if not already uploaded.

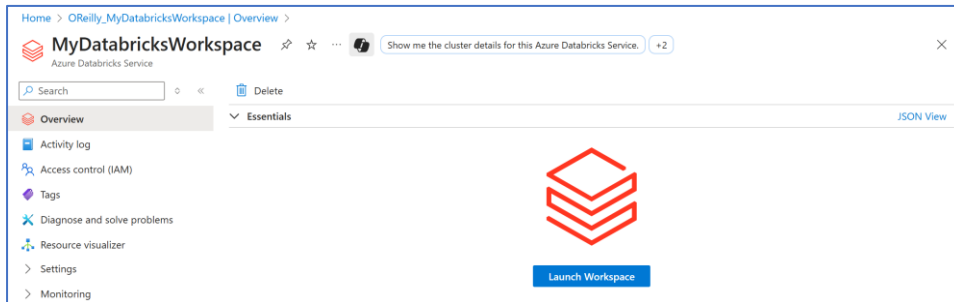## Exercise 1 – Setup Azure Databricks Workspace
1. Go to Azure portal (portal.azure.com).

2. In the search bar, search for Azure Databricks. And select it.

3. Click on **Create**.

4. Fill up the properties to create account
   a. [Basics Tab]
      i. Select subscription
      ii. Select resource group that you created in Exercise 0.
      iii. Provide a unique name
      iv. Select region of your choice (example – East US 2)
      v. Select pricing tier as Premium
      vi. Click **Review + Create**.
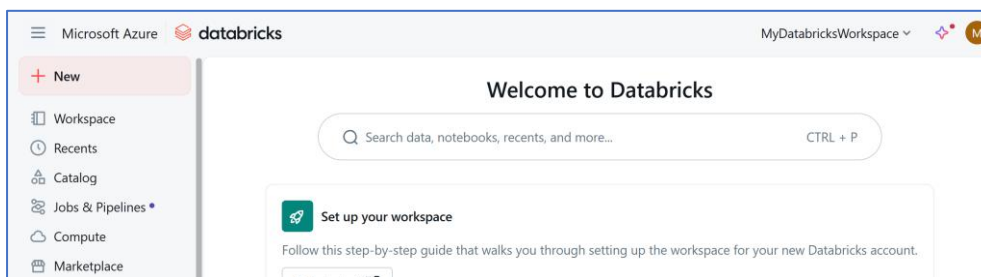


   b. Click **Create**.
      - This will take a few minutes to deploy.

# Exercise 2 – Launch Databricks Workspace & Create Cluster

1. Open Azure Databricks instance created in the previous step.

2. Click on Launch workspace, to open Databricks UI.



3. In the workspace, from left pane, go to **Compute** tab.



4. Under **All-purpose compute**, Click on **Create Compute** to create a cluster.

5. Fill up cluster properties and click on **Create**. This will take few minutes to setup a single node cluster.
   a. Compute name: Demo Cluster
   b. Databricks Runtime: Select the latest runtime with LTS (long-term support).
   c. Photon acceleration: Disable
   d. Node type: Standard_DS3_v2 (if this type is not available, select any other)
   e. Single node: Enable
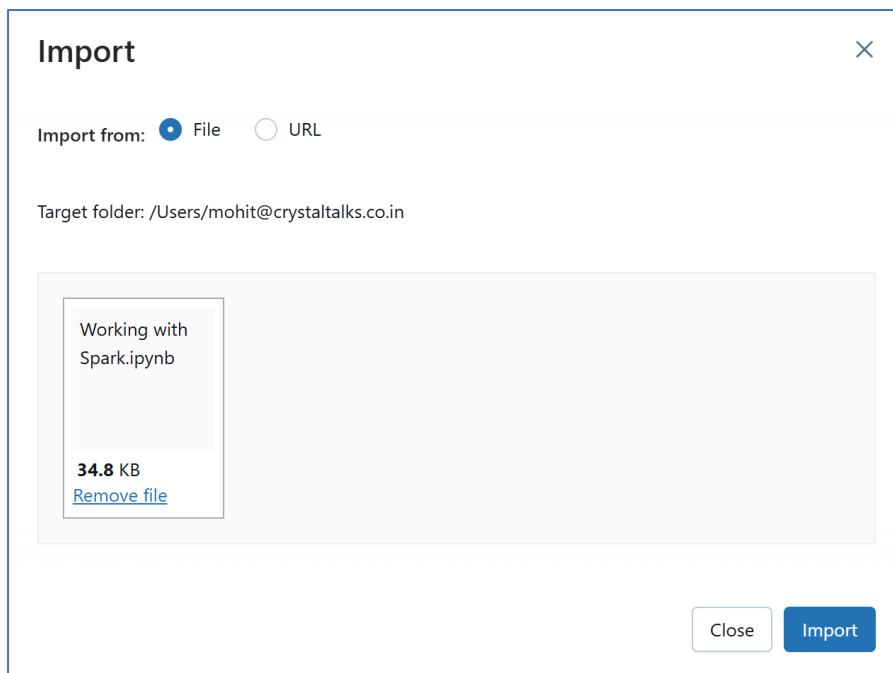   f. Terminate minutes: 30 minutes

*[Note]: If you want to setup multi-node cluster, deselect single node option from UI.*
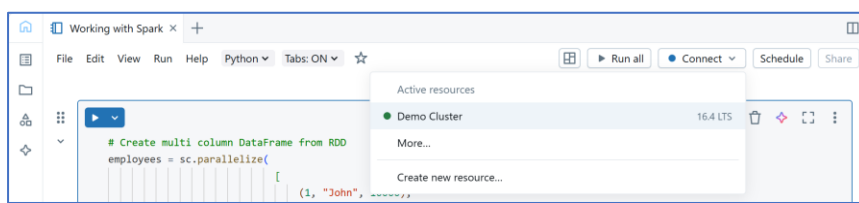
## Exercise 3 – Import Notebook & Run Commands

1. Download notebook – **Working with Spark.py** from GitHub repository.

2. Once the cluster is ready, from left pane, go to **Workspace** tab.

3. Click on Workspace folder → Users folder → Your user account.

4. Right-click on the account folder and click **Import**.



5. Upload the notebook - **Working with Spark.py** and click **Import**.

6. Open the notebook and connect to cluster (Demo Cluster) that you previously created.



7. Run the commands 1 to 12.