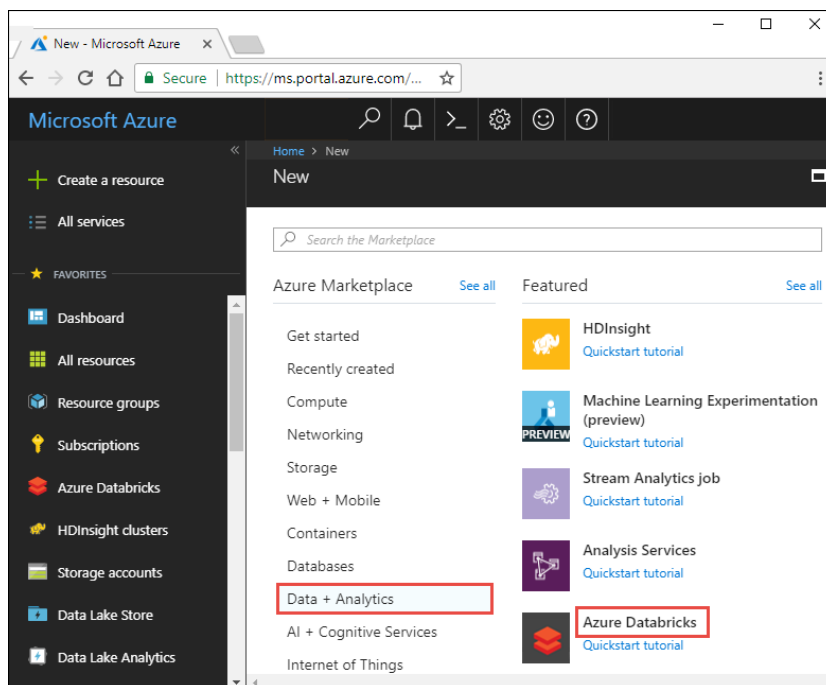# Azure Data Bricks

## LAB Overview

This lab introduces you how to create Azure Data Bricks and launch Spark SQL job using them.

---

## Task 1: Create an Azure Databricks workspace

1. Sign in the Azure portal at
   https://portal.azure.com
2. In the Azure portal, select **Create a resource > Data + Analytics > Azure Databricks**.



3. Under **Azure Databricks Service**, provide the values to create a Databricks workspace:

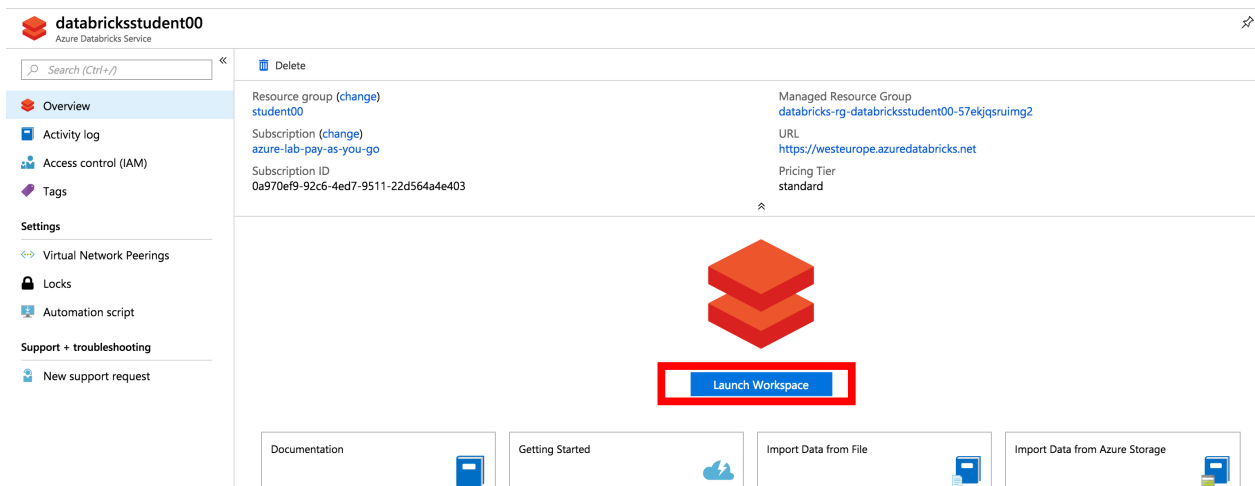   - **Workspace name:** databricksstudentXX
   - **Location:** West Europe

- **Subscription:** XXXXXX
- **Resource Group:**
  - **Create New**: studentXX
- **PricingTier:** Standard

When you finish, click on button **Create.**

---

# Task 2: Create a Spark in Databricks

In this section you will learn how to upload some blob file in Azure Storage Account using Azure Portal.

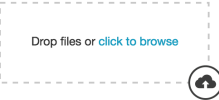1. In the Azure portal, go to the Databricks workspace that you created, and then click Launch Workspace.

2. You are redirected to the Azure Databricks portal. From the portal, click Cluster.



3. In the New cluster page, provide the values to create a cluster.
- **Cluster Name**: studentxxCluster
- **Databricks Runtime Version**: 4.0
- **Terminate after**: 120



Next click on **Create Cluster**.
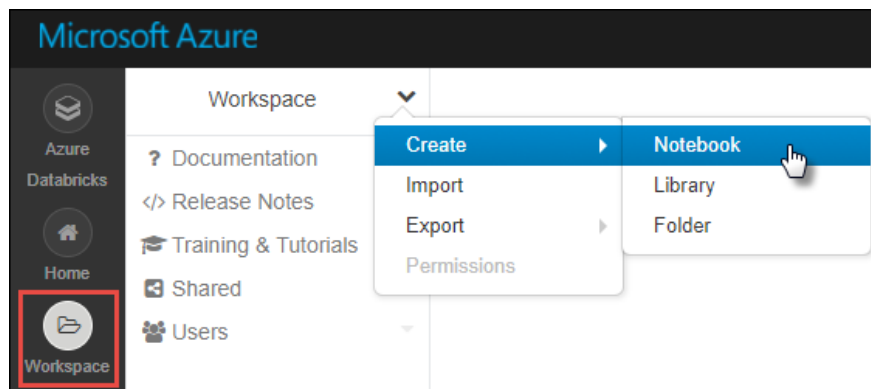
## Task 3: Download a sample data file

Download a sample JSON data file and save it into Azure blob storage.

1. Download this sample JSON data file from Github onto your local computer. Right-click and save as to save the raw file locally.
   https://raw.githubusercontent.com/Azure/usql/master/Examples/Samples/Data/json/radiowebsite/small_radio_json.json
2. Go to Azure Storage Account **studentXX**.
3. Open the storage account in the Azure portal.
4. Select **Blobs**.
5. Select **+ Container** to create a new empty container.
6. Provide a **Name** for the container **databricks**.
7. Select **Private (non anonymous access)** access level.
8. Once the container is created, select the container name.
9. Select the **Upload** button.
10. On the **Files** page, select the **Folder icon** to browse and select the sample file small_radio_json.json for upload.
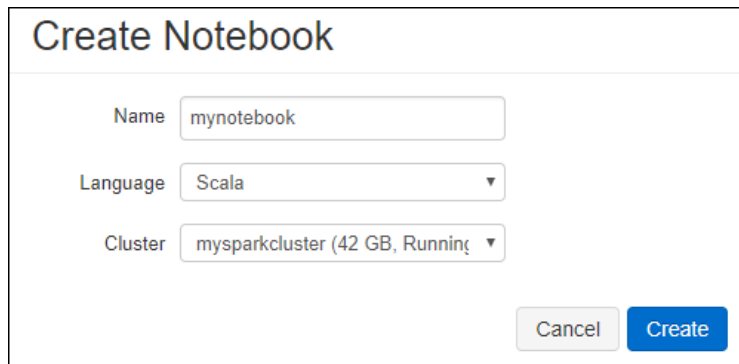11. Select **Upload** to upload the file.

## Task 4: Run a Spark SQL job

In this section you learn how to management of Azure Storage Account using Azure Storage Explorer.

1. In the left pane, click Workspace. From the Workspace drop-down, click Create, and then click Notebook.
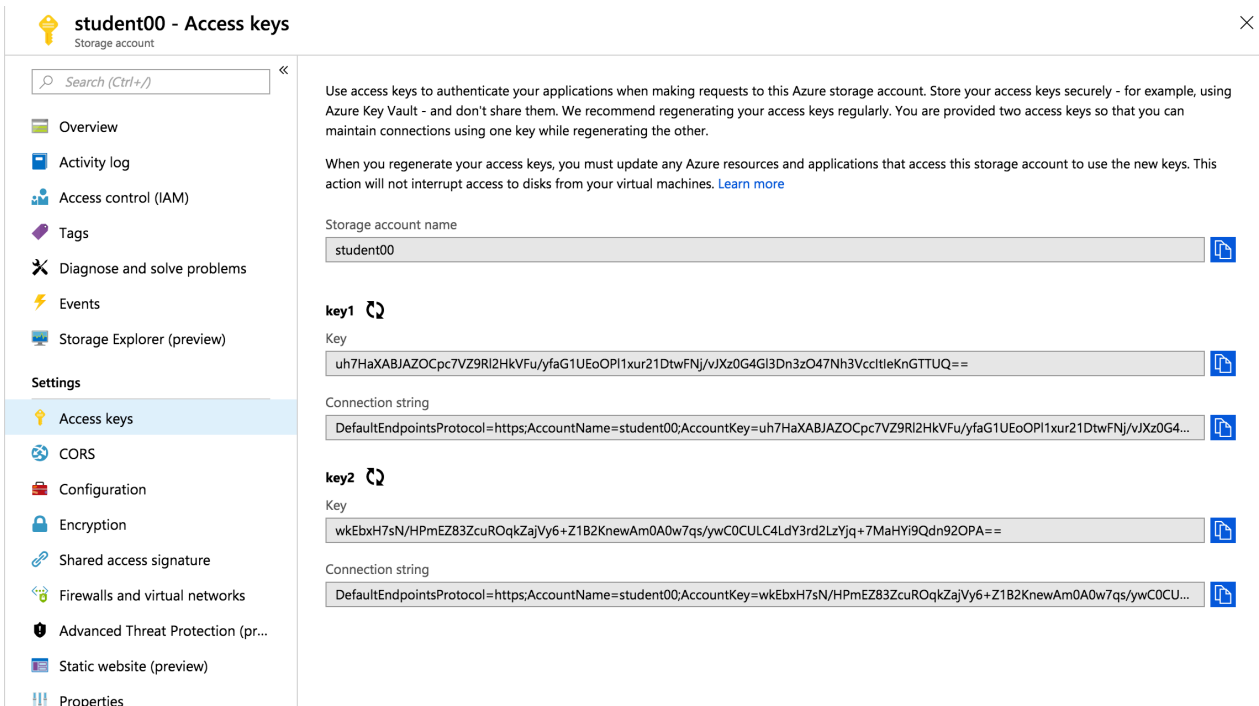


2. In the **Create Notebook** dialog box, enter a name, select **Scala** as the language, and select the Spark cluster that you created earlier.



3. Click **Create**.

4. Mount the storage account with DBFS. The Azure Storage account path is mounted to /mnt/mypath.
5. Paste the code to notebook:
   **dbutils.fs.mount(source = "wasbs://databricks@{YOUR STORAGE ACCOUNT NAME}.blob.core.windows.net/", mountPoint = "/mnt/mypath",  extraConfigs = Map("fs.azure.account.key.{YOUR STORAGE ACCOUNT NAME}.blob.core.windows.net" -> "{YOUR STORAGE ACCOUNT ACCESS KEY}"))**
6. Open Azure Storage account **studentXX** and replace values from step 5 {YOUR STORAGE ACCOUNT NAME} with storage account name and {YOUR STORAGE ACCOUNT ACCESS KEY} using key1 value.

**7.** Run a SQL statement in notebook to create a temporary table using data from the sample JSON data file, **small_radio_json.json**.
**%sql**

**DROP TABLE IF EXISTS radio_sample_data;**

**CREATE TABLE radio_sample_data**

**USING json**

**OPTIONS (**

 **path "/mnt/mypath/small_radio_json.json"**

**)**

8. Select data from temporary table using command:

**%sql**

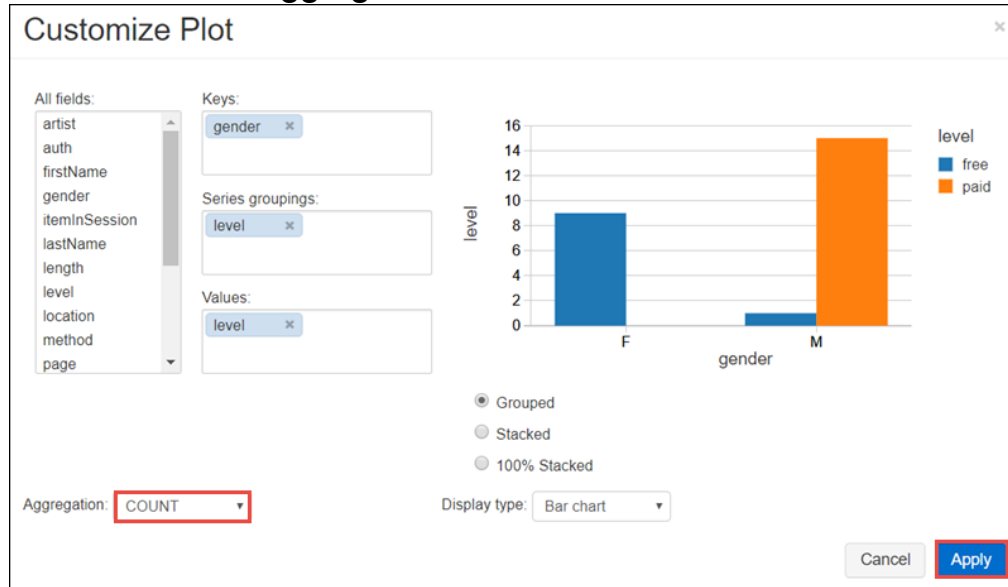**SELECT * from radio_sample_data**

9. You will see a tabular output.

| artist | auth | firstName | gender | itemInSession | lastName | length | level |
|--------|------|-----------|--------|---------------|----------|--------|-------|
| El Arrebato | Logged In | Annalyse | F | 2 | Montgomery | 234.57914 | free |
| Creedence Clearwater Revival | Logged In | Dylann | M | 9 | Thomas | 340.87138 | paid |
| Gorillaz | Logged In | Liam | M | 11 | Watts | 246.17751 | paid |
| null | Logged In | Tess | F | 0 | Townsend | null | free |
| Otis Redding | Logged In | Margaux | F | 2 | Smith | 135.57506 | free |

10. You now create a visual representation of this data to show for each gender, how many users have free accounts and how many are paid subscribers. From the bottom of the tabular output, click the Bar chart icon, and then click Plot Options.
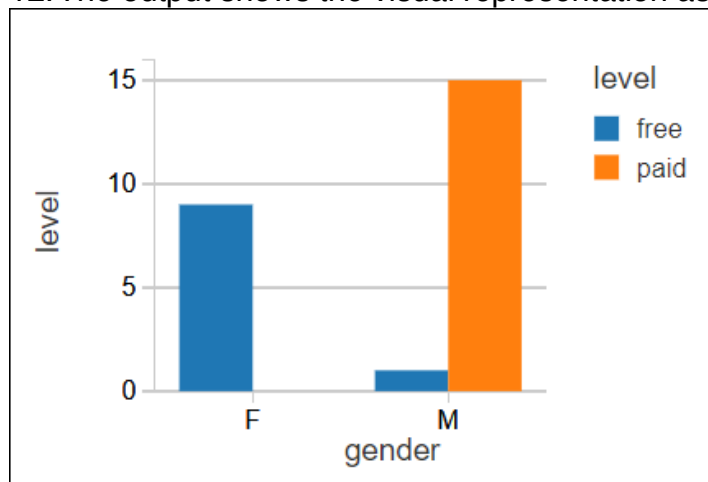
11. In Customize Plot, drag-and-drop values as shown in the screenshot.
    a. Set **Keys** to **gender**.
    b. Set **Series groupings** to **level**.
    c. Set **Values** to **level**.
    d. Set **Aggregation** to **COUNT**.



Click **Apply**.

12. The output shows the visual representation as depicted in the following screenshot.



13. Stop the cluster.