

Getting Started with Azure Databricks

Azure Databricks is a fast, easy and collaborative Spark based analytics service. It is used to accelerate big data analytics, artificial intelligence, performant data lakes, interactive data science, machine learning and collaboration. You will discover the Azure Databricks environment and the main topics around it: workspace, cluster, notebook.

To begin, you need to have access to an Azure Databricks workspace with an interactive cluster. If you do not have a workspace and/or the required cluster, follow the instructions below. Otherwise, you can skip to the **Upload data** section below.

Create Azure Databricks resources

To use Azure Databricks, you first need to deploy an Azure Databricks workspace in your Azure subscription and create a cluster on which you will run notebooks and code. You can then upload the data and notebooks to experiment with in your workspace.

Deploy an Azure Databricks workspace

1. In the Azure portal 'https://portal.azure.com', create a new **Azure Databricks** resource, specifying the following settings:
 - **Subscription:** Choose the Azure Subscription in which to deploy the workspace.
 - **Resource Group:** Create a new resource group.
 - **Workspace Name:** Provide a name for your workspace.
 - **Region:** Select a location near you for deployment. For the list of regions supported by Azure Databricks, see [Azure services available by region](#).
 - **Pricing Tier:** Standard
2. Wait for the workspace to be created. Workspace creation takes a few minutes. During workspace creation, the portal displays the Submitting deployment for Azure Databricks tile on the right side. You may need to scroll right on your dashboard to see the tile. There is also a progress bar displayed near the top of the screen. You can watch either area for progress.

Deploy an Azure
Databricks
workspace

[Create a cluster](#)

[Upload data](#)

[Import
Databricks
notebooks](#)

Create a cluster


1. When your Azure Databricks workspace resource has been created, go to it in the portal, and select **Launch Workspace** to open your Databricks workspace in a new tab, signing in if prompted.
2. In the left-hand menu of your Databricks workspace, select **Compute**, and then select + **Create Cluster** to add a new cluster with the following configuration:
 - **Name:** Enter a unique name.
 - **Cluster Mode:** Single Node
 - **Databricks Runtime Version:** Select the **ML** edition of the latest available version of the runtime. Ensure that the version selected:
 - Does **not** use a GPU
 - Includes Scala > **2.11**
 - Includes Spark > **3.0**
 - **Terminate after:** 120 minutes of inactivity
 - **Node Type:** Standard_DS3_v2
3. Wait for your cluster to be created, which may take several minutes. The cluster will start automatically, and eventually the spinning *Pending* indicator next to the cluster name will change to a solid green circle to indicate a status of *Running*.

Upload data

1. Download `https://raw.githubusercontent.com/MicrosoftLearning/dp-090-databricks-ml/master/data/nyc-taxi.csv` to your computer, saving it as **nyc-taxi.csv** in any folder.

2. On the **Data** page in the Databricks Workspace, select the option to **Create Table**.
3. In the **Files** area, select **browse** and then browse to the **nyc-taxi.csv** file you downloaded.
4. After the file is uploaded to the workspace, select **Create Table with UI**. Then select your cluster and select **Preview Table**.
5. Specify the following table attributes:
 - **Table Name:** nyc_taxi
 - **Create in Database:** default
 - **File Type:** CSV
 - **Column Delimiter:** , (comma)
 - **First row is header:** checked
 - **Infer schema:** checked
 - **Multi-line:** unchecked
6. Set the appropriate data type for each column: Locate the **passengerCount** column. In the drop-down menu, set the column to **INT**.
 - passengerCount: **INT**
 - tripDistance: **DOUBLE**
 - hour_of_day: **INT**
 - day_of_week: **INT**
 - month_num: **INT**
 - normalizeHolidayName: **STRING**
 - isPaidTimeOff: **BOOLEAN**
 - snowDepth: **DOUBLE**
 - precipTime: **DOUBLE**
 - precipDepth: **DOUBLE**
 - temperature: **DOUBLE**
 - totalAmount: **DOUBLE**
7. Click **Create Table**.
8. After the table has been created, view it in the workspace.

Import Databricks notebooks

1. In the Azure Databricks Workspace, using the command bar on the left, select **Workspace**. Then select **Users**, and  **your_user_name**.
2. In the blade that appears, select the downwards pointing chevron (▼) next to your name, and select **Import**.
3. On the **Import Notebooks** dialog, import the notebook archive from the following URL, noting that a folder with the archive name is created, containing one or more notebooks:

```
https://github.com/MicrosoftLearning/dp-090-databricks-ml/raw/master/01%20-%20Introduction%20to%20Azure%20Databricks.dbc
```

4. Repeat the previous step to import the following notebook archives, noting that a folder is created for each archive as it is imported.

```
https://github.com/MicrosoftLearning/dp-090-databricks-ml/raw/master/02%20-%20Training%20and%20Evaluating%20Machine%20Learning%20Models.dbc
```

```
https://github.com/MicrosoftLearning/dp-090-databricks-ml/raw/master/03%20-%20Managing%20Experiments%20and%20Models.dbc
```

```
https://github.com/MicrosoftLearning/dp-090-databricks-ml/raw/master/04%20-%20Integrating%20Azure%20Databricks%20and%20Azure%20Machine%20Learning.dbc
```

Explore Azure Databricks

In this exercise, you will discover the Azure Databricks environment.

1. In the **01 - Introduction to Azure Databricks** folder in your workspace, open the **Getting Started with Azure Databricks** notebook.
2. In the top left dropdown menu, choose your cluster to attach your notebook to that cluster. *(Alternatively, you will be prompted to attach a cluster when running the first cell in an unattached notebook).*
3. Read the notes in the notebook, running each code cell in turn.

Clean-up

If you're finished working with Azure Databricks for now, in Azure Databricks workspace, on the **Compute** page, select your cluster and select **■ Terminate** to shut it down. Otherwise, leave it running for the next exercise.