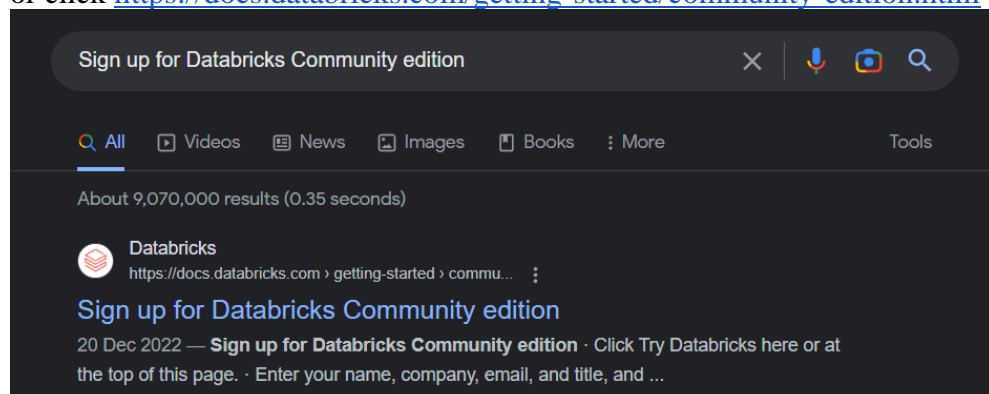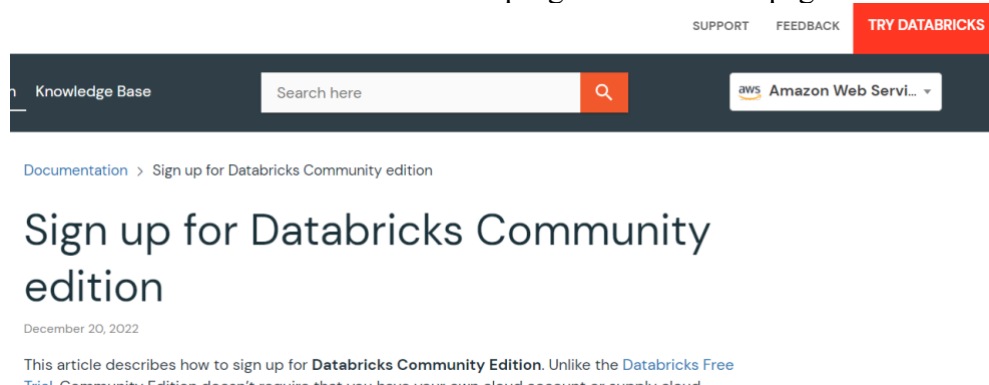# 1 Databricks Guide

## 1.1  Databricks Community Edition Registration

To get started, we would be using the free version, databricks community edition and you would have to register for it.

1) Google "Sign up for Databricks Community edition" and click on the first link or click https://docs.databricks.com/getting-started/community-edition.html



2) click "TRY DATABRICKS" on the top right corner of the page



3) Fill in the necessary details and click continue to create an account
4) You will reach the page below and since we are using the Community Edition, please click on the portion highlighted below

5) Validate your email address and login to databricks

## 1.2 Create a Cluster

You need to create a cluster to run your notebooks.

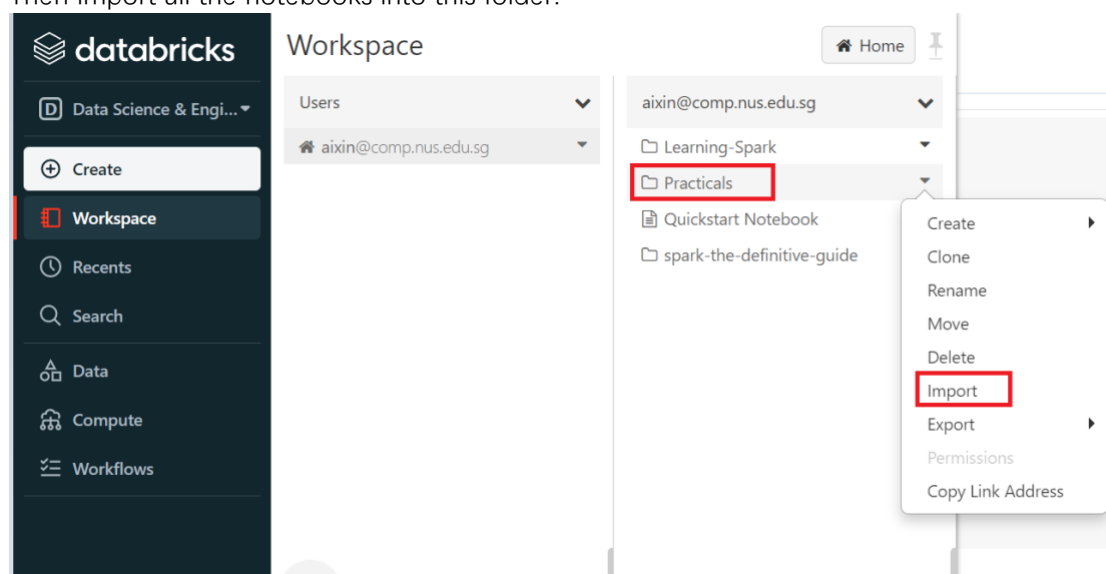Note: the created cluster will automatically terminate after an idle period of two hours for the Community Edition. Hence you would need to re-create a new cluster to run your code if your cluster has terminated.
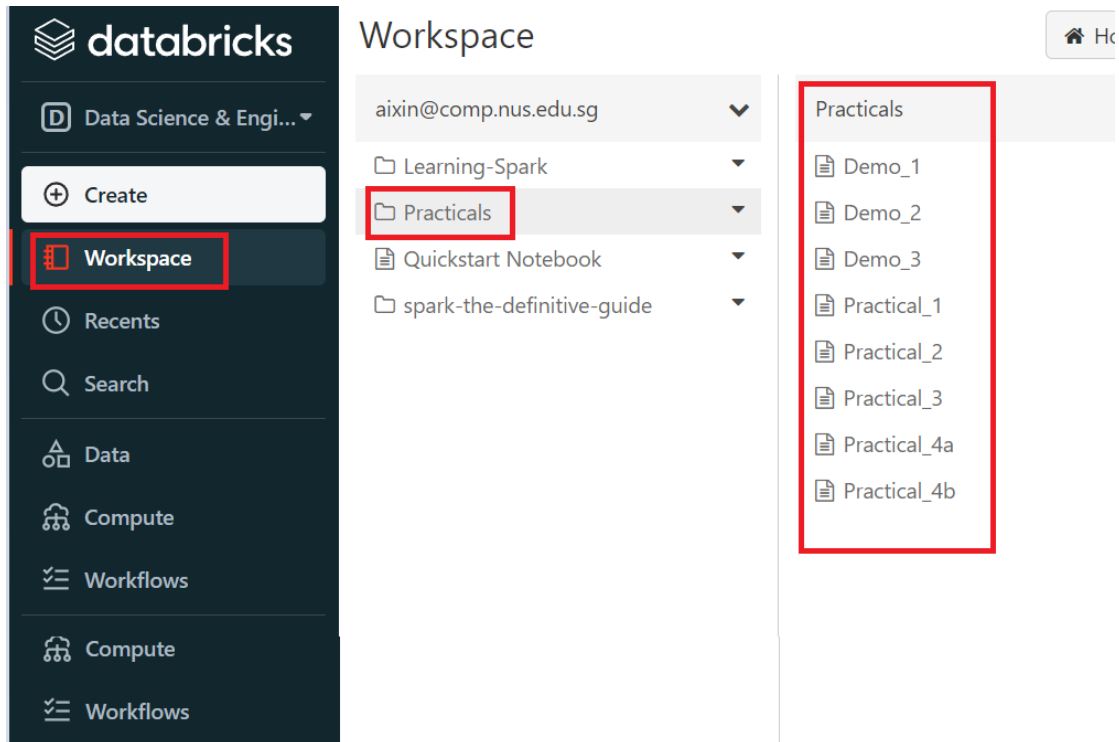
## 1.3   Upload Notebooks

First, create a folder, e.g. Practicals (or any other name you like).



Then import all the notebooks into this folder.

Open Demo_1 and attach the notebook to your active cluster, e.g. Test.



Now you can start run the notebook cell by cell. You can also access Spark UI for each Job by click "View".