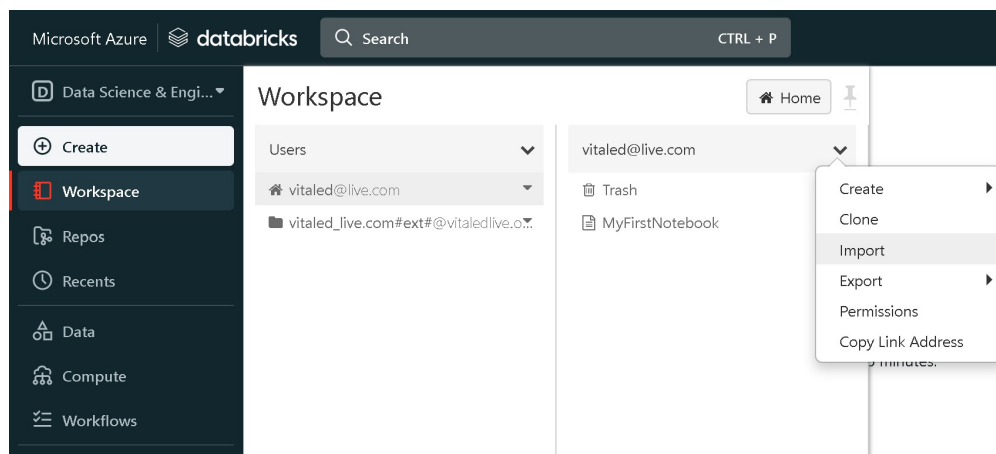# Lab 04: Create and schedule a Job

## Goal

During this lab you will learn how to create and schedule a Job in a Azure Databricks
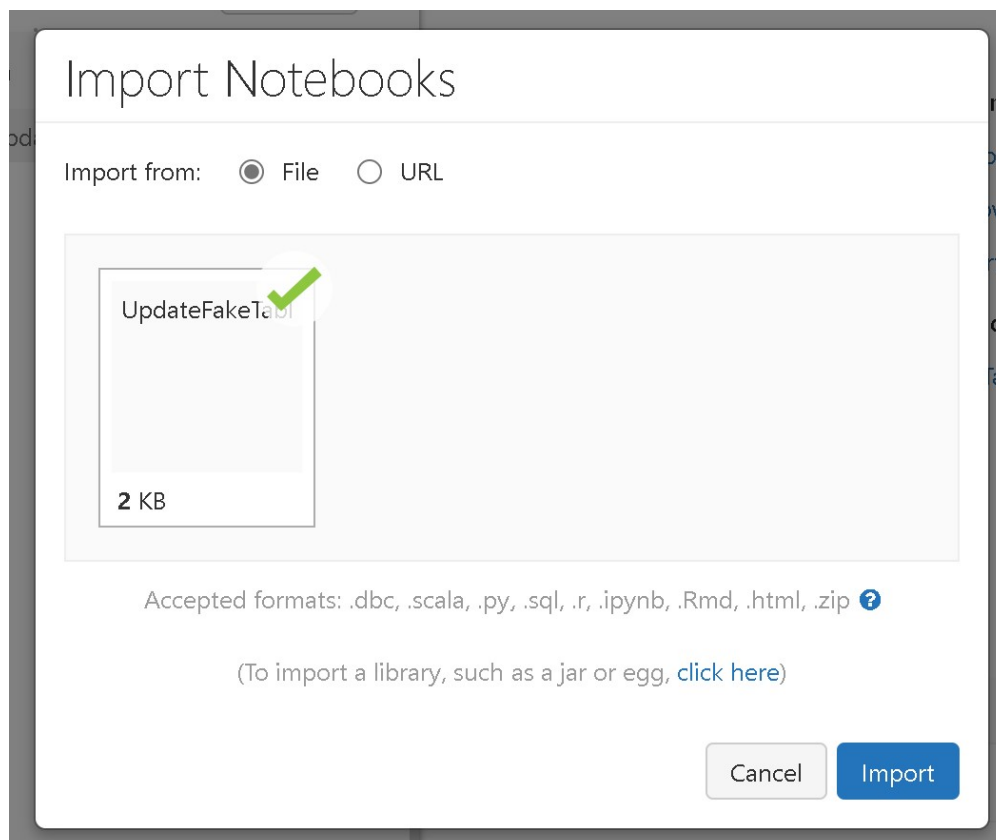
## Tasks

### Task 1: Import notebook to schedule

Click on **Workspace** navigate to your home directory and click on **Import**



Click on browse, select the file *UpdateFakeTable.ipynb* and click on import



In the notebook page click on **Run** and then **Run all**

This first execution will create the table *default.fake_table* and insert 1000 entries you can check this by looking at the result of the count statement



## Task 2: Create and Schedule a Job

Into the Azure Databrick portal from the lef-sided menu select **Create** + **Job**

In the form provide the following details:

| property | value |
| --- | --- |
| Job Name | *FirstJob* |
| Task name | *UpdateFakeTable* |
| Type | *Notebook* |
| Source | *Workspace* |
| Path | *path to the UpdateFakeTable notebook* |
| Cluster | *select cluster created in the Lab 02 |

then click on **Create**

In the Job details click on **Edit Schedule**

⋮   Run now  ⌄



Select **Scheduled** as Trigger type and **Every 1** and **minute** as other settings



## Task 3: Check job results

On the left-sided menu click on **Workflows** and then **Job Runs** after one minute you should see the succeded execution

Click on the date under the *Start time* column

You will see the result of the notebook execution the count now should be **2000**



## Task 3: Delete the Job

On the left-sided menu click on **Workflows** and then **Jobs** click on the delete icon at then end of the job entry row



This Lab has been completed!