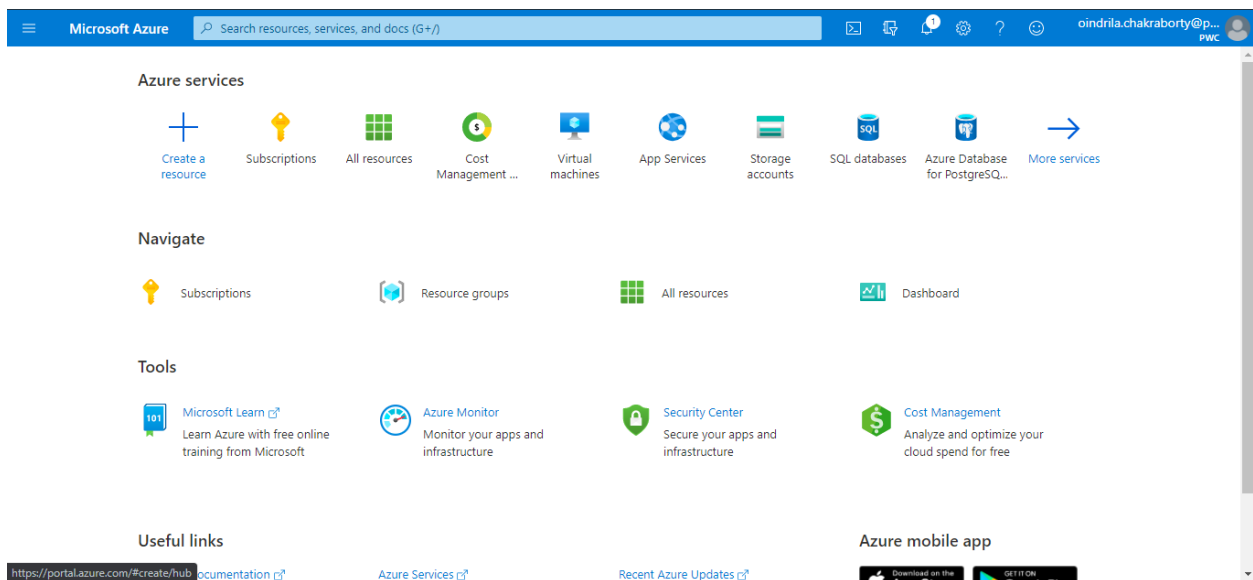# DATABRICKS NOTEBOOK WITH INPUT PARAMETER

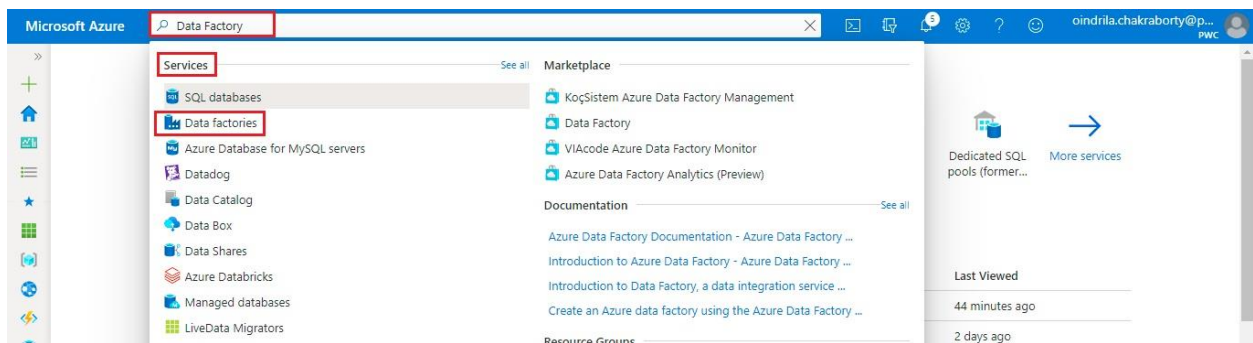## Run Databricks Notebook with Parameter from Azure Data Factory Notebook Activity

- An *Azure Data Factory Pipeline* can be *created* that *executes* a *Databricks Notebook against* the *Databricks Job Cluster*. The *Pipeline* can *also pass Azure Data Factory Parameters* to the *Databricks Notebook during execution*.
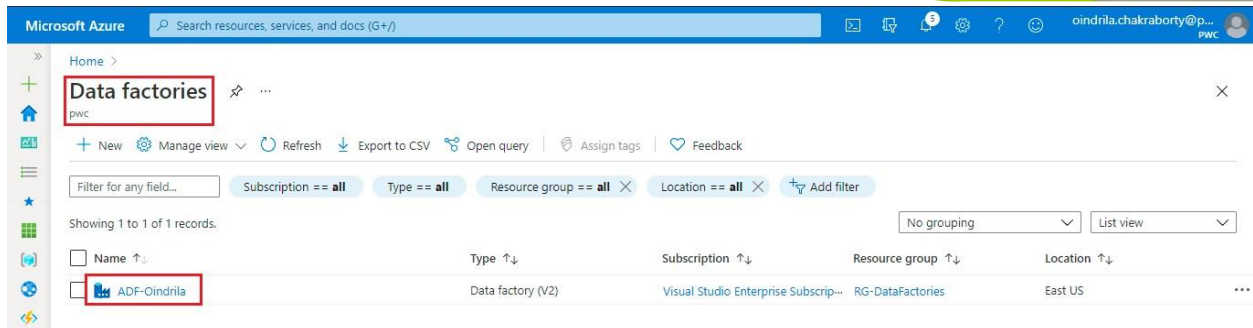- <u>Create an Azure Databricks Linked Service</u> -
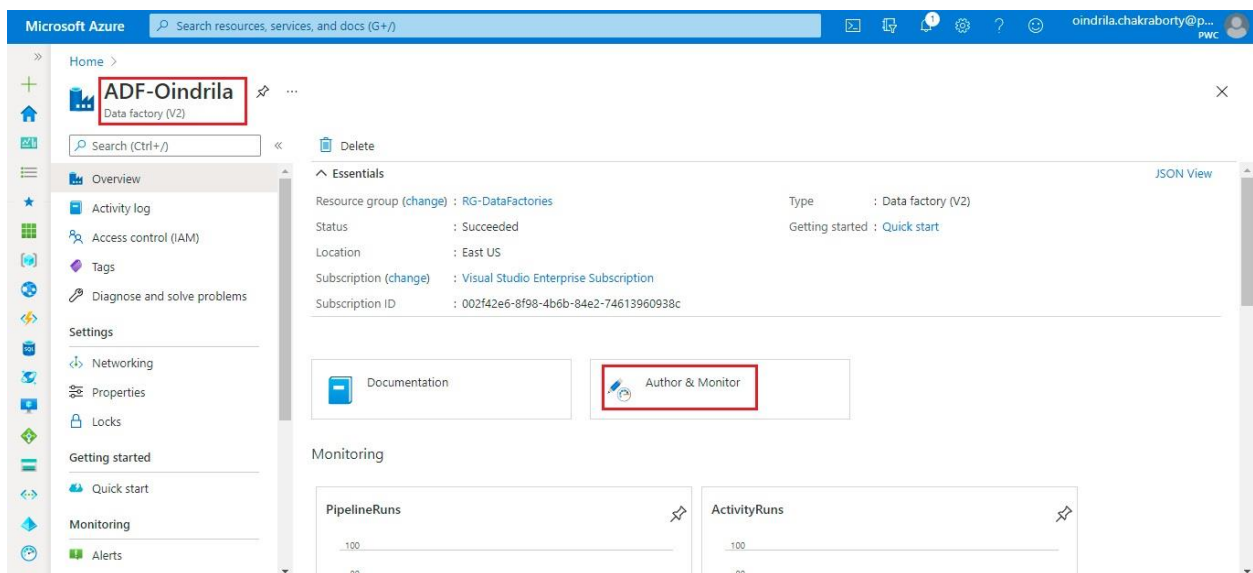    - ➤ <u>Step 1</u> - *Open* the *Azure portal* (*portal.azure.com*).



- ➤ <u>Step 2</u> - *Type* "*Data Factory*" in the *global search bar* in the *home page* of the *Azure portal*. *Click on* the *second search result* "*Data factories*" under "*Services*" in the *left side*.
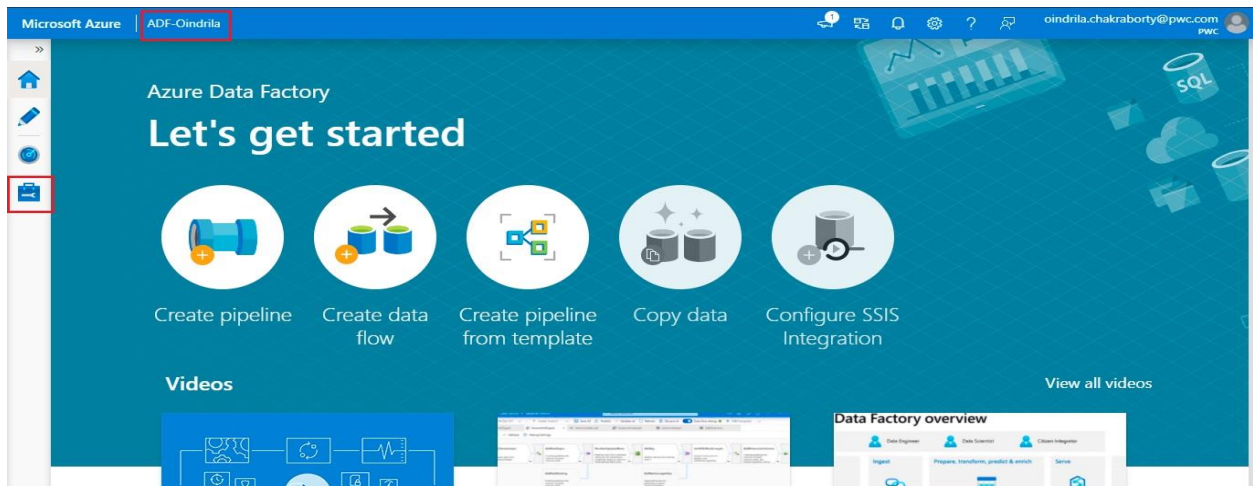


- ➤ <u>Step 3</u> - *In* the "*Data factories*" *page*, *click on* the *Azure Data Factory resource* "*ADF-Oindrila*".
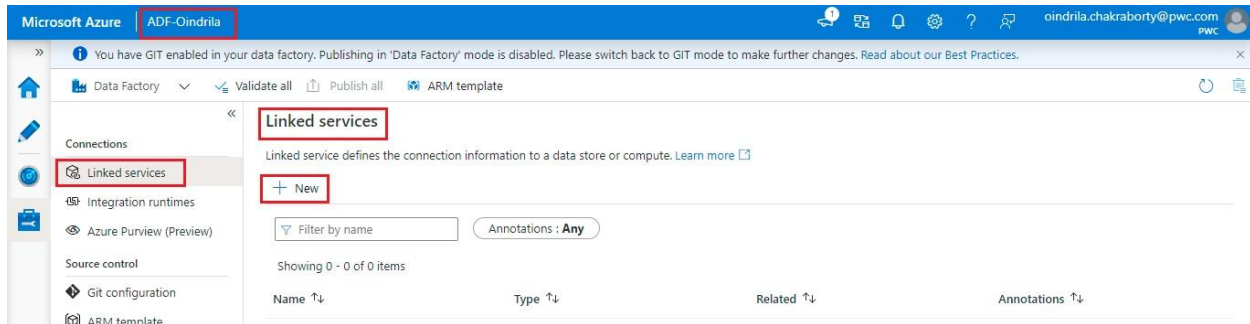
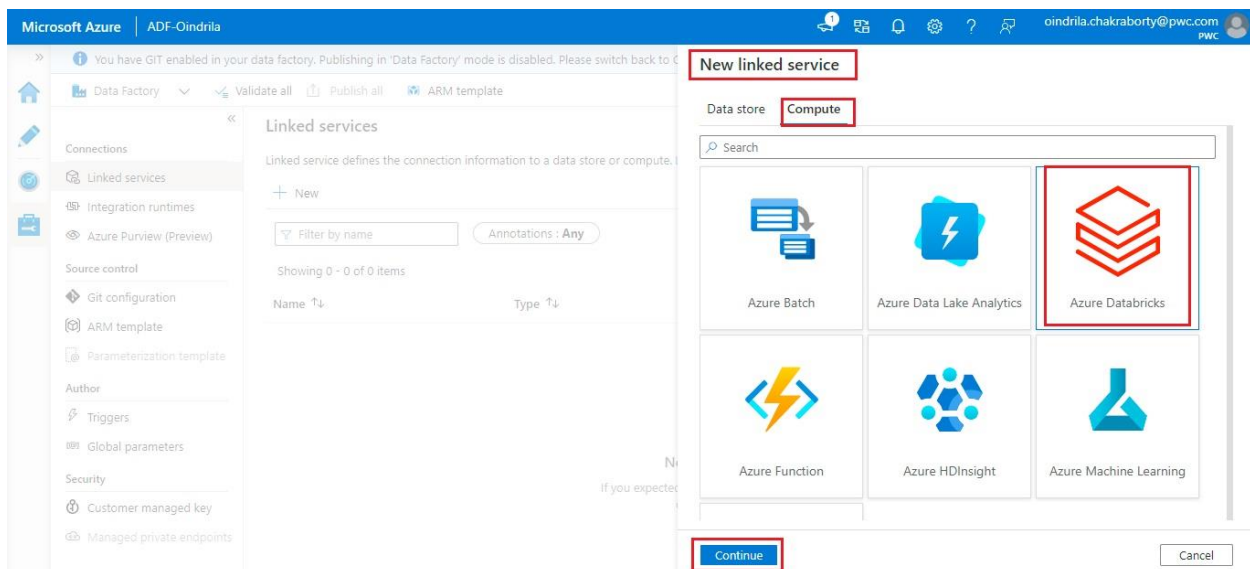> **_Step 4_** - **_Click on_** the "**_Author & Monitor_**" link.



> **_Step 5_** - **_Click on_** the "**_Manage_**" button.

➢ **Step 6** - *Under* the *"Connections" left menu, click on* the *menu option "Linked services"*. Then, *click on* the *"+ New" button* to *add* a *new Linked Service*.



➢ **Step 7** - *In* the *"New Linked Service" window, click on* the *"Compute" tab*. Then, *select* the *"Azure Databricks"* and *click on* the *"Continue" button*.



➢ **Step 8** - *In* the *"New Linked Service (Azure Databricks)" window, provide* *"LS_databricks_workspace_001"* in the *"Name" textbox*.
*Select* the *option "From Azure Subscription" from* the *dropdown "Account selection method"*.
*Select* the *proper Azure Subscription from* the *dropdown "Azure subscription"*.
*Select* the *appropriate Databricks Workspace* that will *run* the *Notebooks, from* the *dropdown "Databricks workspace"*, e.g., *"databricks-workspace-001"*.
*Select* the *radio button option "New job cluster" for* the *property "Select cluster"*.

The *Databricks Workspace URL* is *auto populated*.

Select *"Access Token" from* the *dropdown "Authentication type"*.

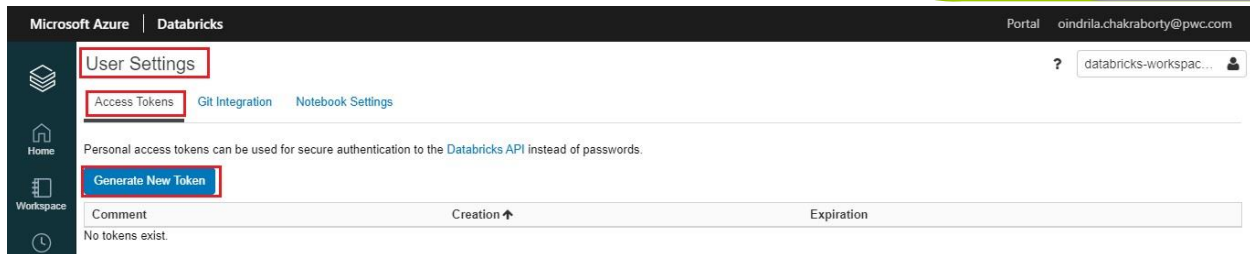To *authenticate* to and *access Databricks REST APIs*, the *Databricks Personal Access Token*, or, *Databricks Password* can be *used*. It is *recommended* to *use Databricks Personal Access Token*. *Token-Based Authentication* is *enabled by default for all Databricks Accounts*, *launched after January 2018*. The *number* of *Databricks Personal Access Tokens per User* is *limited* to *600 per Workspace*. *Following* are the *steps* to *generate* the *Databricks Personal Access Token* for the *selected Databricks Workspace* -
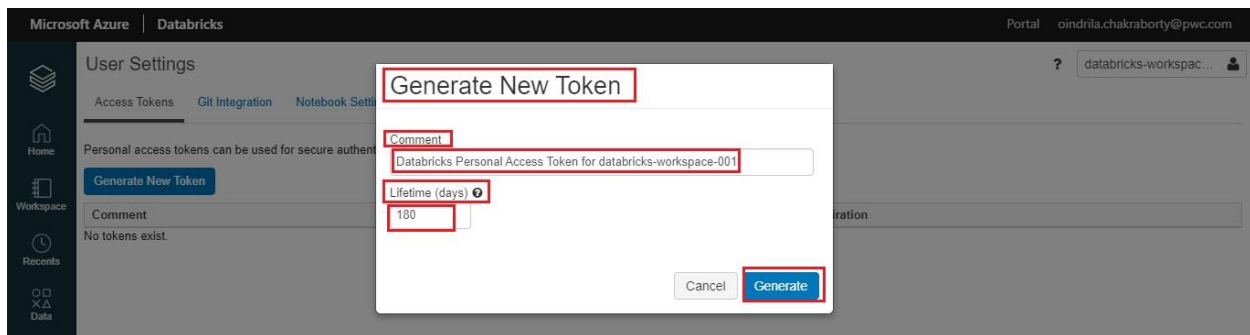
    ✓  <u>Step 8.1</u> - *Click on* the *User Profile Icon* in the *upper right corner* of the *Databricks Workspace*. *Click on* the *menu option "User Settings"*.
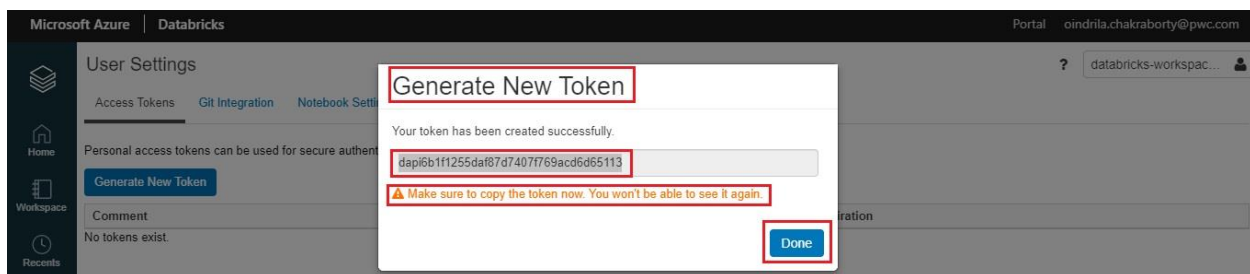


    ✓  <u>Step 8.2</u> - *Go to* the *"Access Tokens" tab*, and, *click on* the *"Generate New Token" button*.

✓ *Step 8.3* - *In* the "*Generate New Token*" *pop up*, *provide* a *comment*, e.g., "*Databricks Personal Access Token for databricks-workspace-001*" in the *textbox "Comment"*, and, *provide* "*180*" in the *textbox "Lifetime (days)"*. *Finally*, *click on* the "*Generate*" *button*.
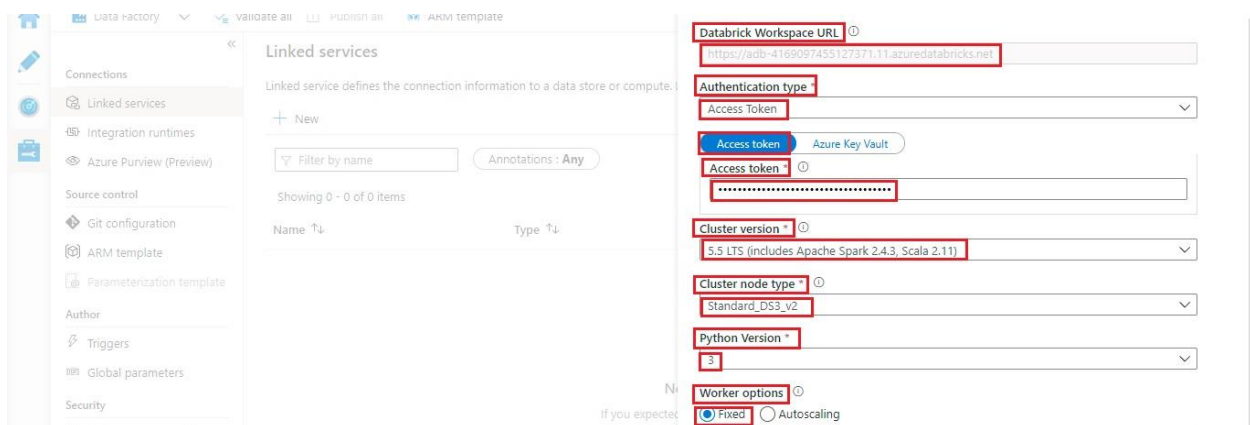


✓ *Step 8.4* - *Copy* the *generated Token*, i.e., "*dapi6b1f1255daf87d7407f769acd6d65113*", and *save* it in a *secured location*, *as*, the *Token* will *not* be *displayed again*. *Click on* the "*Done*" *button*.
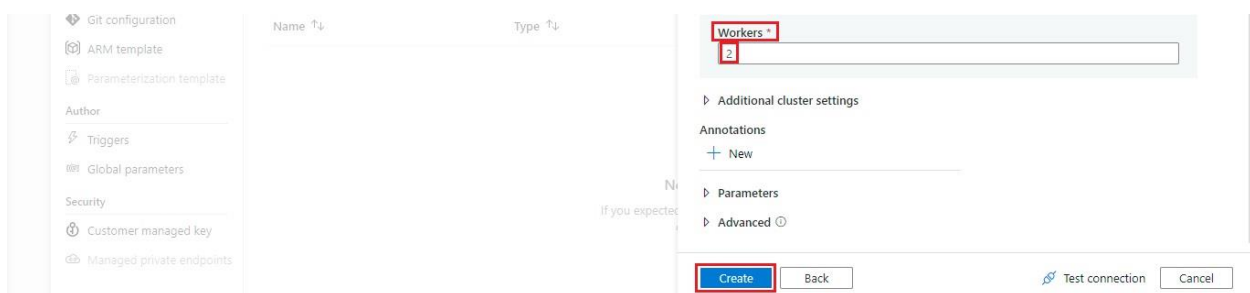


✓ *Step 8.5* - The *Databricks Personal Access Token for* the *Databricks Workspace "databricks-workspace-001*" is *created*.
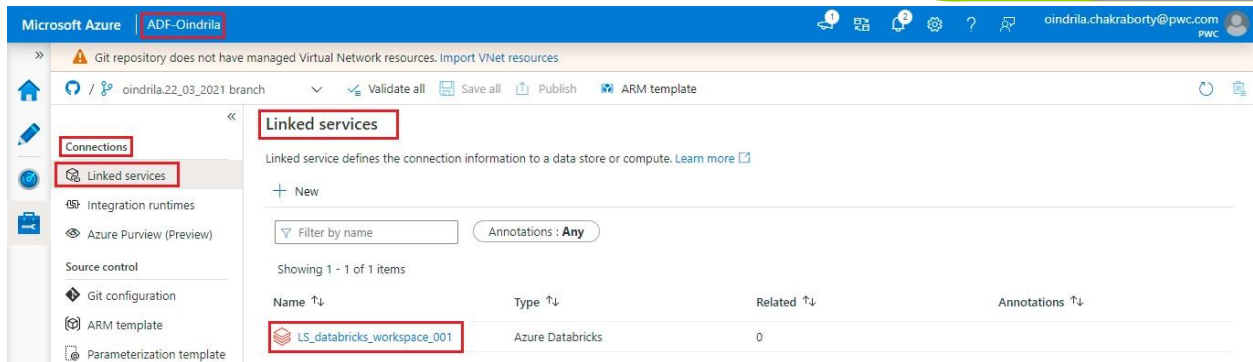
Provide the created Databricks Personal Token *"dapi6b1f1255daf87d7407f769acd6d65113"* in the *textbox "Access token"*. *Select* the *option "5.5 LTS (includes Apache Spark 2.4.3, Scala 2.11)" from* the *dropdown "Cluster version"*.

*Select* the *option "Standard_DS3_V2" under* the *"General purpose" category from* the *dropdown "Cluster node type"*.

*Select* the *option "3" from* the *dropdown "Python Version"*.

*Select* the *radio button option "Fixed" for* the *property "Worker options"*.



Provide *"2"* in the *textbox "Workers"*, and, *finally*, *click on* the *"Create" button*.
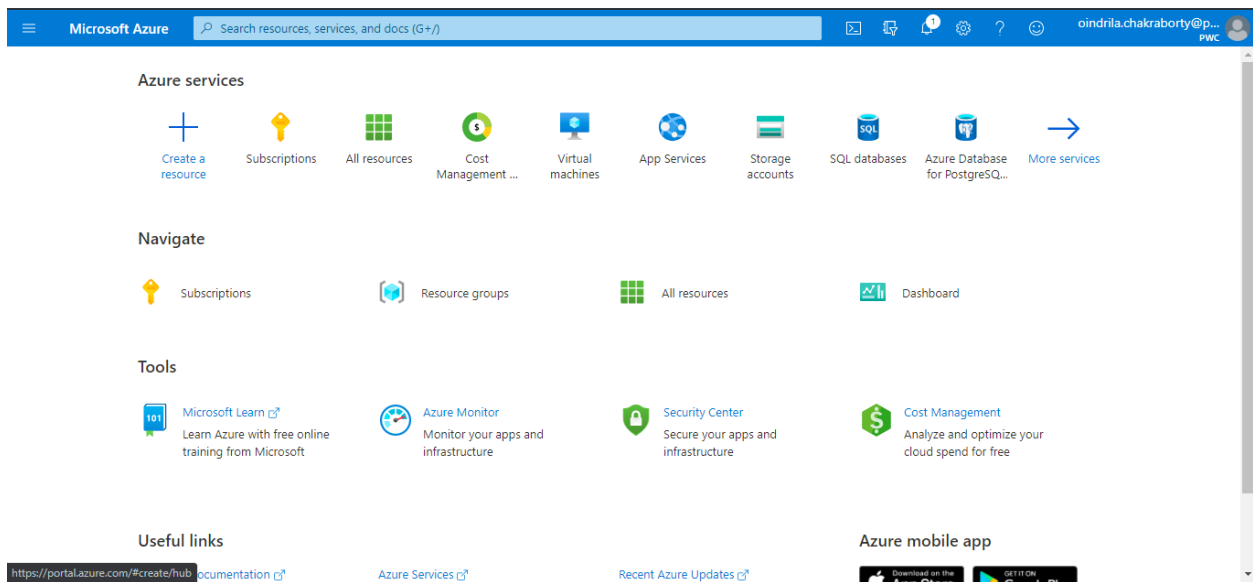


> ➤ **Step 9** - The *Linked Service "LS_databricks_workspace_001"* is *created successfully*.
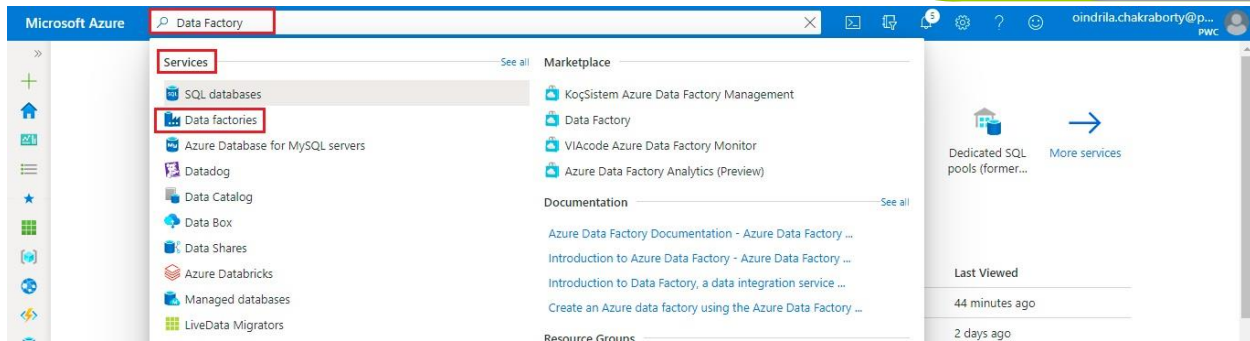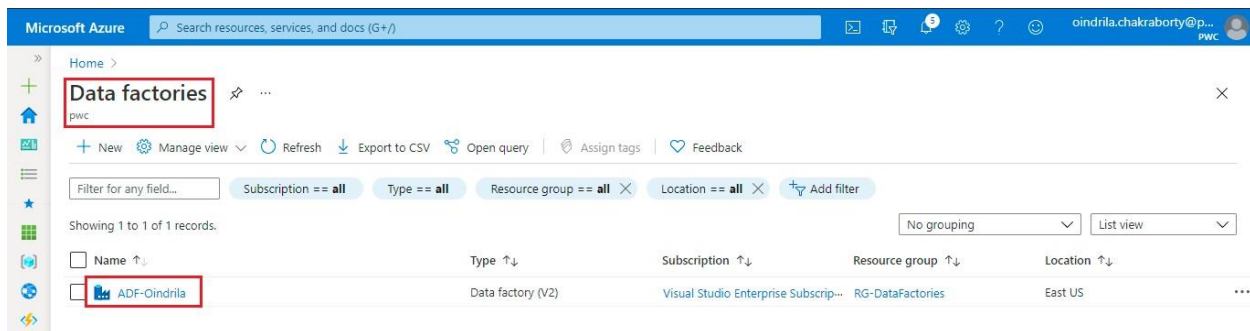
**Create a Pipeline -**

➢ *Step 1 - **Open** the **Azure portal** (portal.azure.com).*
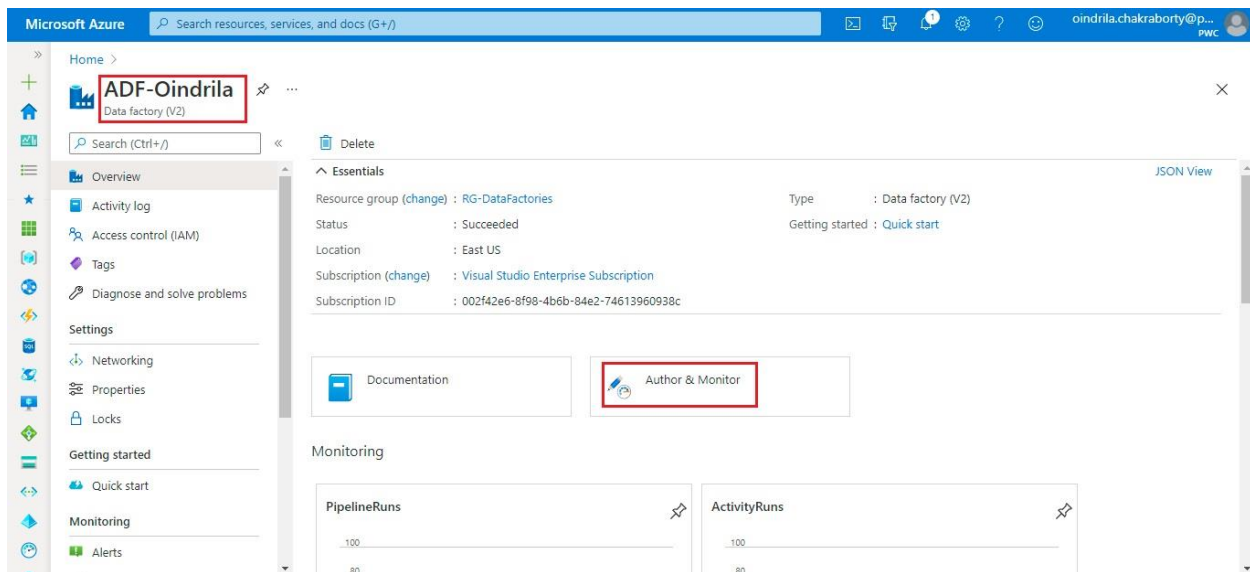


➢ *Step 2 - **Type** "**Data Factory**" in the **global search bar** in the **home page** of the **Azure portal**. **Click on** the **second search result** "**Data factories**" under "**Services**" in the **left side**.*
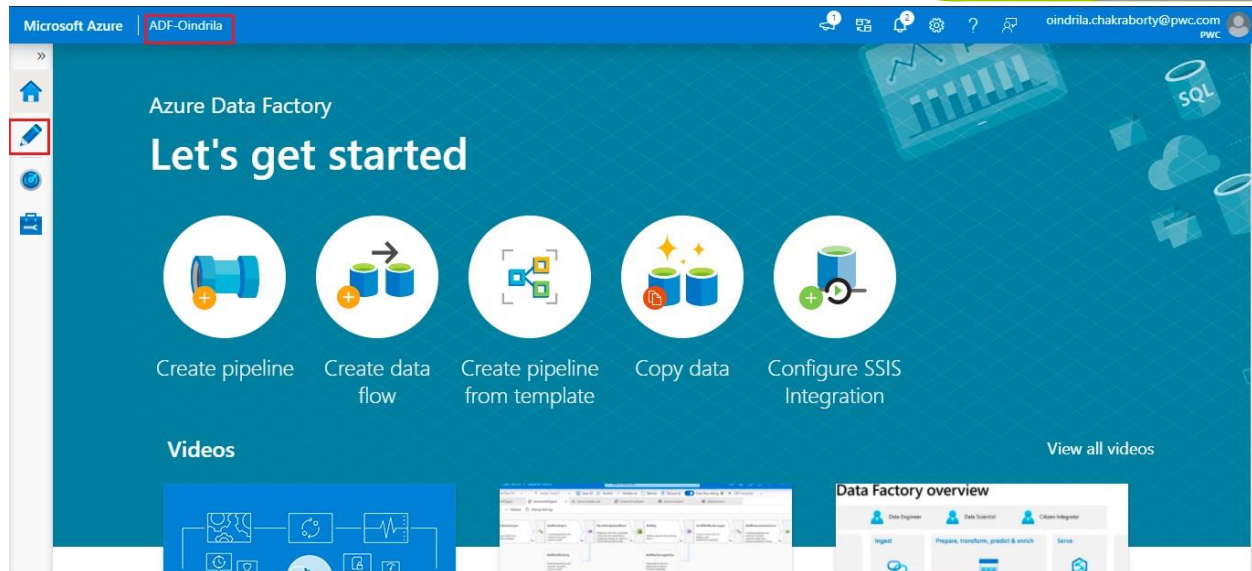
> ➤ **Step 3** - *In* the *"Data factories" page, click on* the *Azure Data Factory resource "ADF-Oindrila".*
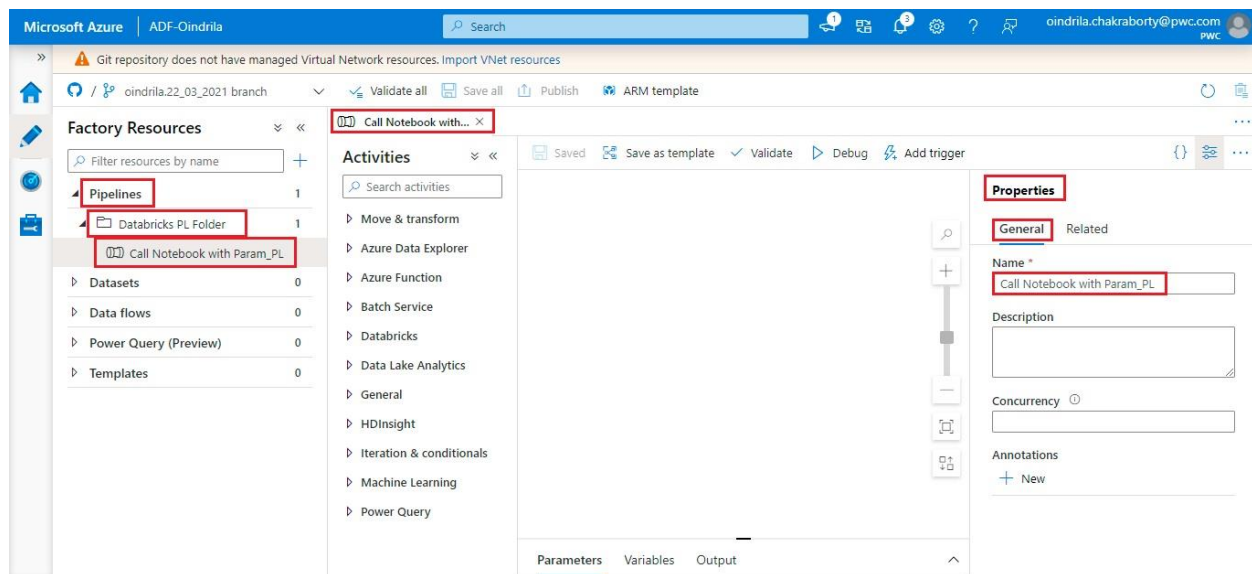


> ➤ **Step 4** - *Click on* the *"Author & Monitor" link.*
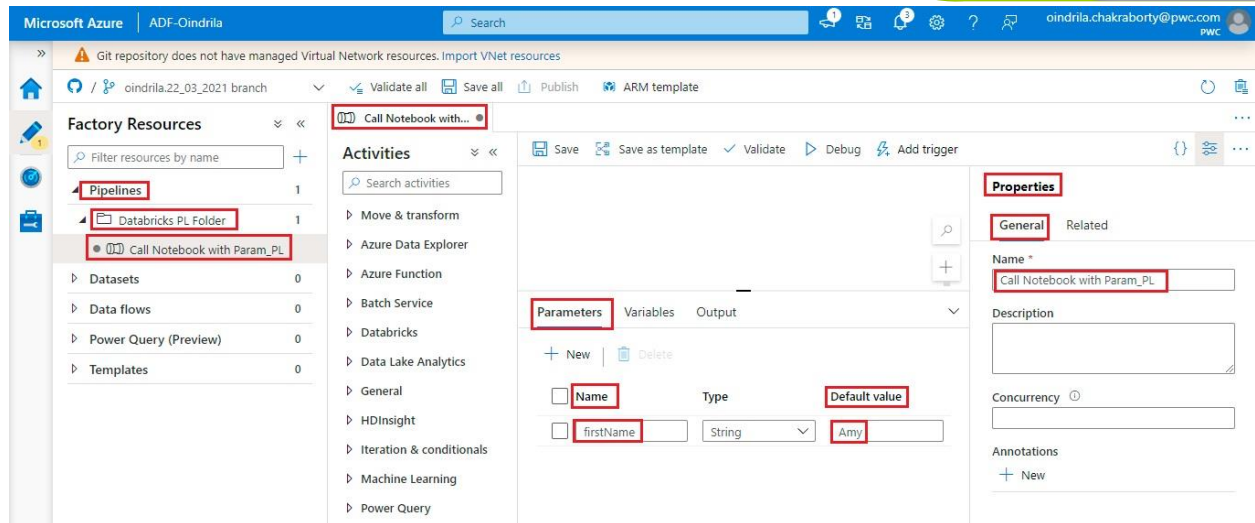


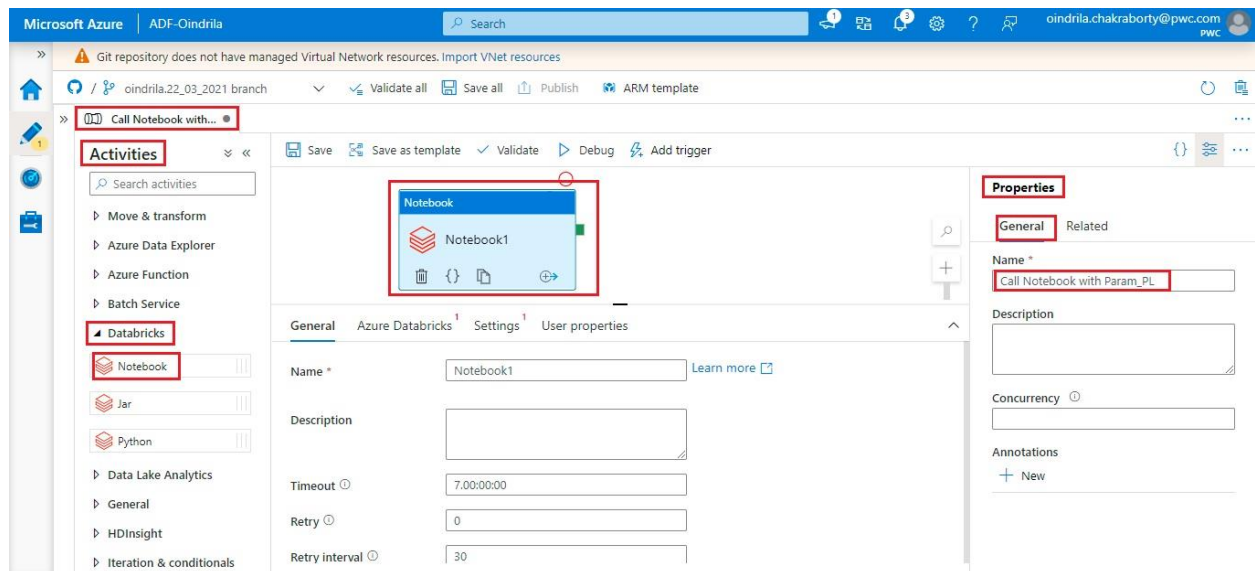> ➤ **Step 5** - *Click on* the *"Author" button.*

➢ **Step 6** - *Create* a *Pipeline* "*Call Notebook with Param_PL*".



➢ **Step 7** - *Create* a *Pipeline Parameter* by the *name* "*firstName*", *with* the *Default Value* "*Amy*". This *Pipeline Parameter* will be *passed* to the *Databricks Notebook*.

> ➢ _**Step 8**_ - _In_ the "_Activities_" _toolbox, expand_ "_Databricks_". _Drag_ the "_Notebook_" _Activity from_ the "_Activities_" _toolbox_ to the _Pipeline Designer surface_.



> ➢ _**Step 9**_ - _In_ the "_Notebook_" _Activity window, switch to_ the "_Azure Databricks_" _tab_. _Select_ the _created Linked Service_ "_LS_databricks_workspace_001_" in the _dropdown_ "_Databricks linked service_".

> ➢ **Step 10** - *Create* a *new Databricks Notebook* by the *name* "*Call ADF Notebook*", *inside* the *folder* "*ADF_Notebook_Folder*", in the *Databricks Workspace* "*databricks-workspace-001*".



> ➢ **Step 11** - *In* the "*Settings*" *tab*, *browse* to *select* the *path* of the *Databricks Notebook*. The *path* of the *Databricks Notebook* is "*/Users/oindrila.chakraborty@pwc.com/ADF_Notebook_Folder/Call ADF Notebook*".

> ➢ **Step 12** - **Add** the *Pipeline Parameter* "*firstName*" to the "*Notebook*" *Activity*. **Expand** the *property* "*Base Parameters*". **Click on** the "*+ New*" *link*. **Provide** "*firstName*" in the *textbox* "*Name*".



**Click on** the *textbox* "*Value*". This *enables* the *link* "*Add dynamic content [Alt+P]*" *to be displayed just below* the *textbox*. **Click** on that *link*. **Select** the *Pipeline Parameter* "*firstName*", and, *click on* the "*Finish*" *button*.

The *Value* of the *Parameter* of the *"Notebook"* *Activity* *"firstName"* is the *expression* *"@pipeline().parameters.firstName"*.



> ➢ **Step 13** - *Save* the *Pipeline* and *Validate*. Then *run* the *Pipeline, by clicking on* the *"Debug"* *link*.

➕ **Verify the Output** - *Open* the *Azure Databricks Workspace "databricks-workspace-001"*. *Click on* the *"Clusters"* *left menu link*.

Switch to the "*Job Clusters*" *tab*. The **Status** of the **Job** can be "*Pending*", "*Execution*", "*Running*", or, "*Terminated*". **Since**, the **Job** is **finished executing**, its **Status** is **displayed** as "*Terminated*".



Click on the **Job Name**, then, *click on* the **button** "*Go To Job Run*".



On successful run, the **parameter passed**, and, the **output** of the **Notebook** can be **validated**.

ADF_ADF-Oindrila_Call Notebook with Param_PL_Notebook1_99f3748e-65ab-4426-88ef-d8a2ee00c826 (Run id: 10)   ?   databricks-workspac... 👤

View: Code    ☁ Export to HTML

### ADF_ADF-Oindrila_Call Notebook with Param_PL_Notebook1_99f3748e-65ab-4426-88ef-d8a2ee00c826 (Run id: 10)   ✖ Delete

**Started:** 2021-03-22 04:16:09 IST
**Duration:** 1m 25s
**Status:** Succeeded
**Run ID:** 10
**Task:** Notebook at /Users/oindrila.chakraborty@pwc.com/ADF_Notebook_Folder/Call ADF Notebook
▶ **Parameters:**
    {"firstName":"Amy"}
**Cluster:** Driver: Standard_DS3_v2, Workers: Standard_DS3_v2, 2 workers, 5.5 LTS (includes Apache Spark 2.4.3, Scala 2.11) - View Spark UI / Logs / Metrics

## Output

```
customerDf = spark.read.options\
                    (
                        header = "true",\
                        delimiter = "|"\
                    )\
                .csv("dbfs:/mnt/datalakegen2oindrila/CSV-Data/customer.dat")

display(customerDf)
```

ADF_ADF-Oindrila_Call Notebook with Param_PL_Notebook1_99f3748e-65ab-4426-88ef-d8a2ee00c826 (Run id: 10)   ?   databricks-workspac... 👤

⬇

Command took 21.58 seconds

```
from pyspark.sql.functions import col

firstName = dbutils.widgets.get("firstName")
customerWithFirstNameDf = customerDf.where(col("c_first_name").contains(firstName))
display(customerWithFirstNameDf)
```

▶ 🔢 customerWithFirstNameDf: pyspark.sql.dataframe.DataFrame = [c_customer_sk: string, c_customer_id: string ... 16 more fields]

| | c_customer_sk ▲ | c_customer_id ▲ | c_current_cdemo_sk ▲ | c_current_hdemo_sk ▲ | c_current_addr_sk ▲ | c_first_shipto_date_sk ▲ | c_first_sales_date_sk ▲ | c_ |
|---|---|---|---|---|---|---|---|---|
| 1 | 2 | AAAAAAAACAAAAAAA | 819667 | 1461 | 31655 | 2452318 | 2452288 | Dr. |
| 2 | 588 | AAAAAAAAMECAAAAA | 650147 | 4116 | 25285 | 2449080 | 2449050 | Mi |
| 3 | 821 | AAAAAAAAFDDAAAAA | 1273654 | 5124 | 33794 | 2450778 | 2450748 | Dr. |
| 4 | 962 | AAAAAAAACMDAAAAA | 1345722 | 5686 | 44097 | 2451508 | 2451478 | Mr |
| 5 | 1482 | AAAAAAAAKMFAAAAA | 1072096 | 1326 | 2909 | 2449121 | 2449091 | Mi |
| 6 | 2431 | AAAAAAAAPHJAAAAA | 1487079 | 2979 | 22126 | 2452371 | 2452341 | Dr. |
| 7 | 2485 | AAAAAAAAFLJAAAAA | 731621 | 1735 | 23002 | 2449135 | 2449105 | Ms |
| 8 | 3396 | AAAAAAAAEENAAAAA | 809859 | 5854 | 31849 | 2449687 | 2449657 | Mi |

Showing all 244 rows.