# Azure Project Steps/Demo

**Obj:** Use Data Flow to move data to Azure Data Lake Storage or Azure Blob Storage. Use Databricks to read and write data from and to these data stores, and also perform advanced analytics and machine learning tasks.



**Setting up the Linked Services**

**Git configuration/integration to Git Repo
from organisation settings
after adding new Repo in Git**



**1- setting the data type for each column**



**2- Data preview**



**3- Dataset Description/Statistics**

**Revenue** 1.2f



Not Null    Null

| | |
|---|---|
| Not Null | 64 |
| Percentile 25 | 2.2367074E7 |
| Standard Deviation | 11641498.36 |
| Null | 32 |
| Percentile 50 | 3.1650092E7 |
| Average | 32360452.23 |
| Percentile 75 | 3.9831564E7 |
| Variance | 135524484023262.06 |
| Maximum | 58756472 |
| Minimum | 14021480 |

## Data cleaning

- Check the nan values from each column using the statistics above
- Using Filter data flow using the expression: !isNull(Revenue)



1- using Filter dataflow to remove the nan values using builder expression



**Revenue** 1.2f

| | |
|---|---|
| Not Null | 64 |
| Percentile 25 | 2.2367074E7 |
| Standard Deviation | 11641498.36 |
| Null | 0 |
| Percentile 50 | 3.1650092E7 |
| Average | 32360452.23 |
| Percentile 75 | 3.9831564E7 |
| Variance | 135524484023262.06 |
| Maximum | 58756472 |
| Minimum | 14021480 |

**Sales_quantity** 123

| | |
|---|---|
| Not Null | 64 |
| Percentile 25 | 14959 |
| Standard Deviation | 6591.29 |
| Null | 0 |
| Percentile 50 | 18339 |
| Average | 19197.38 |
| Percentile 75 | 22786 |
| Variance | 43445067.70 |
| Maximum | 38069 |
| Minimum | 8314 |

**2- Columns after filtering has zero Null values**

**Getting the Avg of all columns per year using the aggregate transformation**

three sinks are added:
afterSortingandFiltering
AddingColumnYear
yearlyAvgRevenue



Writing the order of the three sinks

**Creating new pipeline with our dataflow setting the sinks to run in parallel**



**Debugging the pipeline**



**Dataflow Debugging Dashboard after enabling the debug**

**the 3 Sinks have been added to the azure storage after debugging the pipeline**

## New trigger

Name *

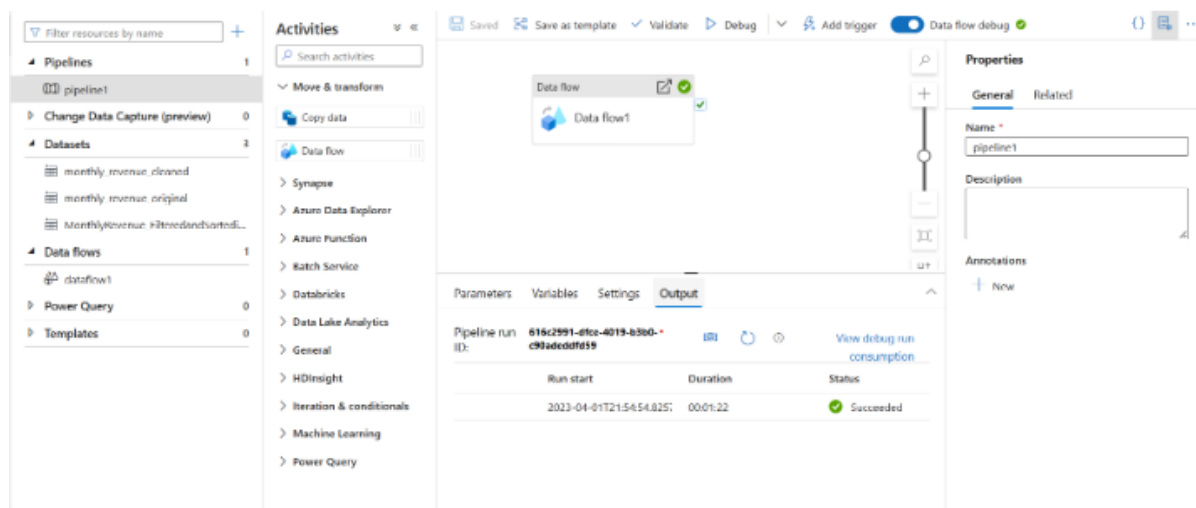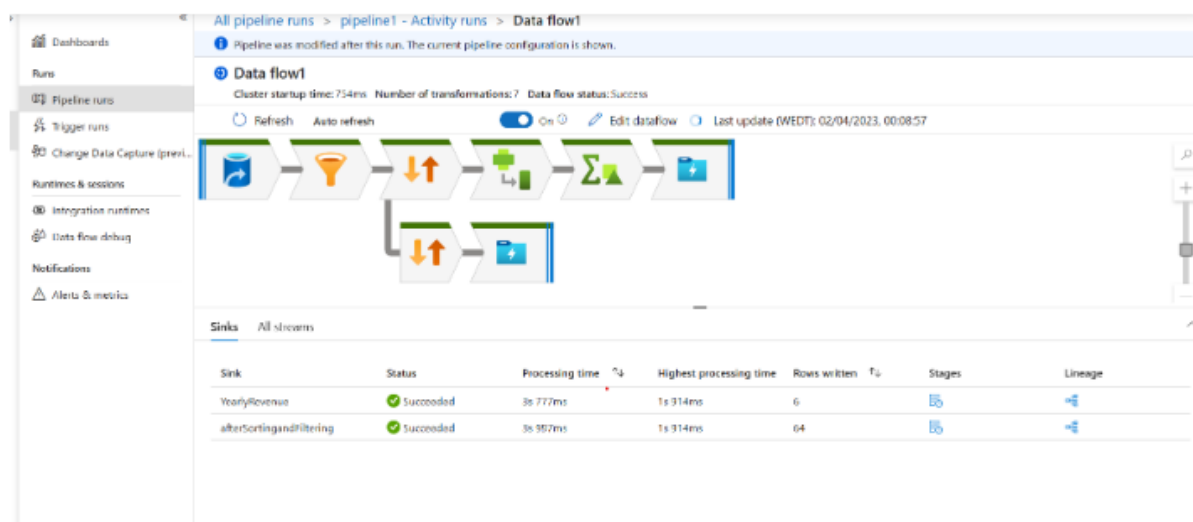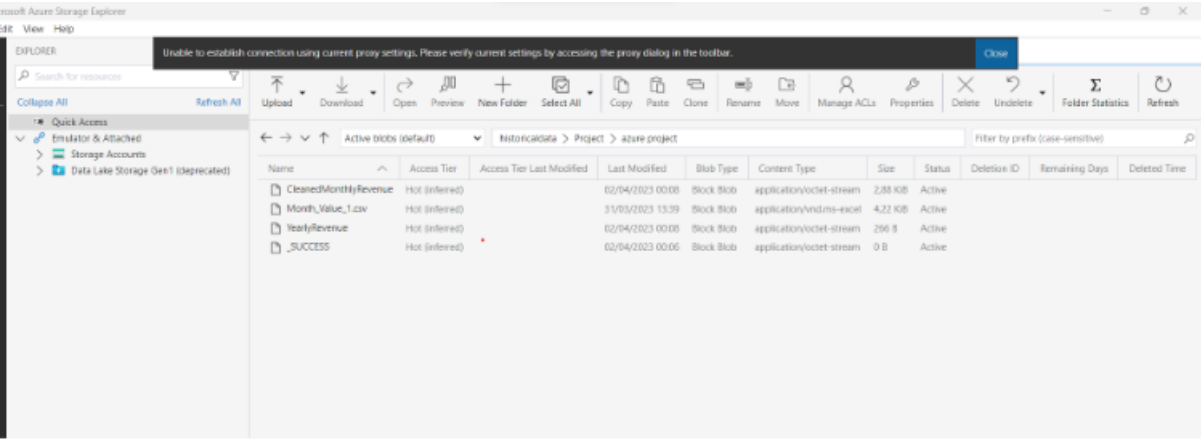UpdatingMonthRevenue

Description

Type *

Schedule

Start date * ⓘ

4/2/2023, 12:13:51 PM

Time zone * ⓘ

Coordinated Universal Time (UTC)

Recurrence * ⓘ

Every 1                              Month(s)

∨ Advanced recurrence options

◉ Month days   ◯ Week days

Select day(s) of the month to execute

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|

**adding trigger to be triggered monthly**

## New trigger

◉ Month days   ◯ Week days

Select day(s) of the month to execute

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|----|----|----|------|----|----|----|
| 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| 15 | 16 | 17 | 18 | 19 | 20 | 21 |
| 22 | 23 | 24 | 25 | 26 | 27 | 28 |
| 29 | 30 | 31 | Last | | | |

Execute at these times  ⓘ

Hours

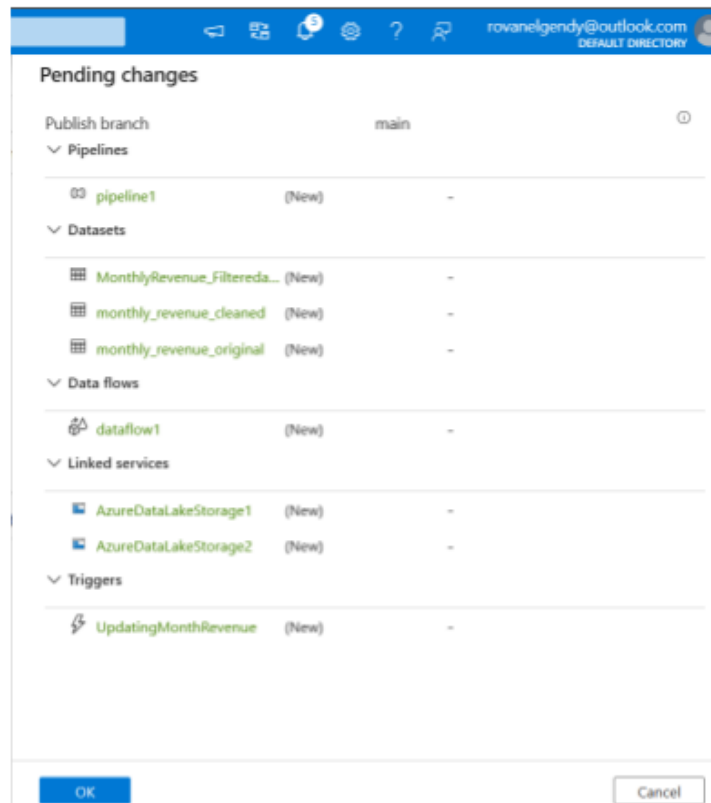Minutes

Schedule execution times
12:13

☐ Specify an end date

Annotations

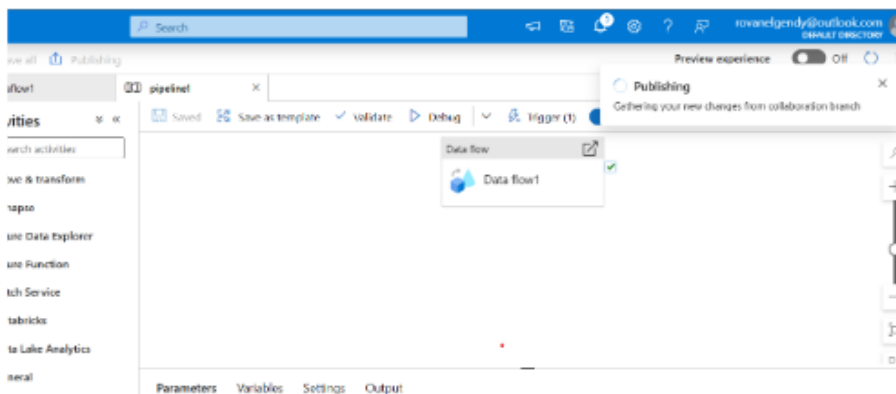+ New

Start trigger ⓘ

☑ Start trigger on creation

[ OK ]   [ Cancel ]

**Setting the trigger**
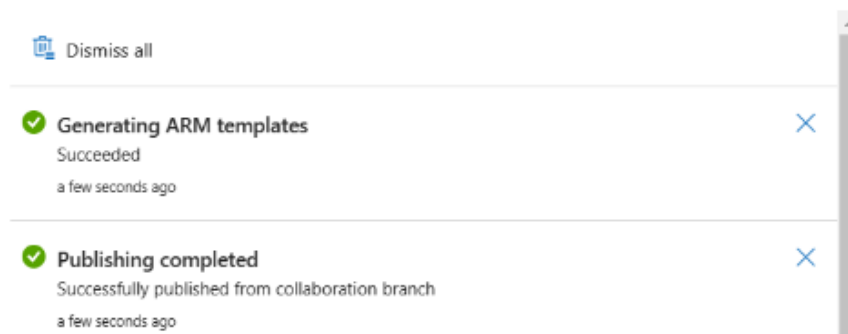
## Setting the trigger



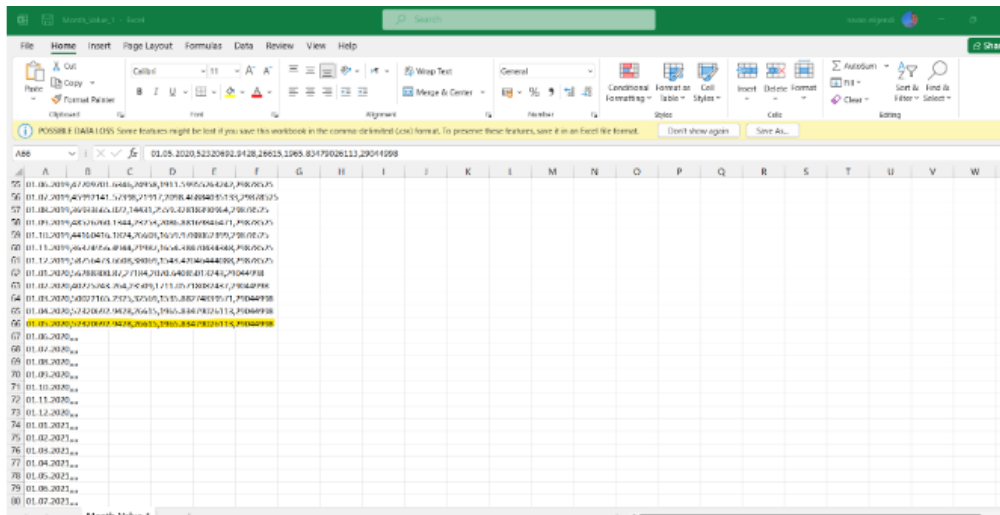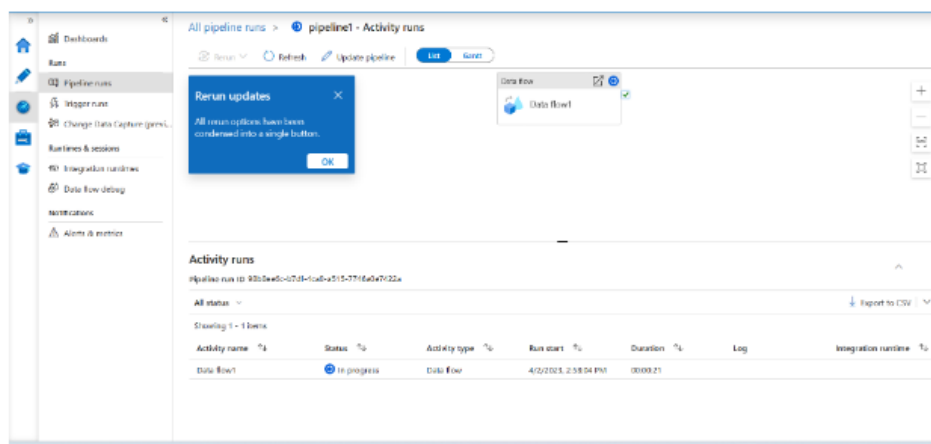## Saved and to be published



## published



## Notifications

Updating the original dataset by one row



pushing the trigger we made earlier manually
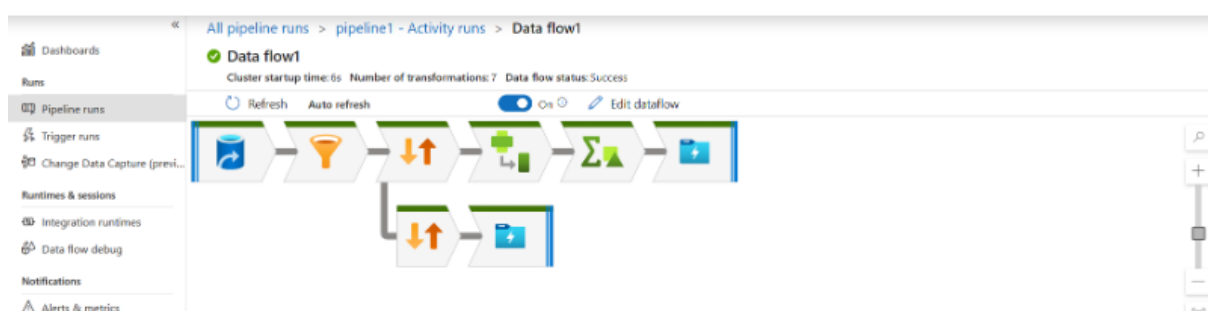


Trigger success notification

| Sink | Status | Processing time ↑↓ | Highest processing time | Rows written ↑↓ | Stages | Lineage |
|------|--------|-------------------|------------------------|-----------------|--------|---------|
| YearlyRevenue | ✅ Succeeded | 4s 189ms | 2s 212ms | 6 | | |
| afterSortingandFiltering | ✅ Succeeded | 4s 439ms | 2s 212ms | 65 | | |

**here we can see after trigger the rows read is updated to 65 rows instead of 64 rows**



**the updated row in the data preview**

You can save the pipeline as template to import later



importing the tempelate



the tempelate gallery



last step to click on use this tempelate

**You can save the pipeline as template to import later**



**importing the tempelate**



**the tempelate gallery**



**last step to click on use this tempelate**

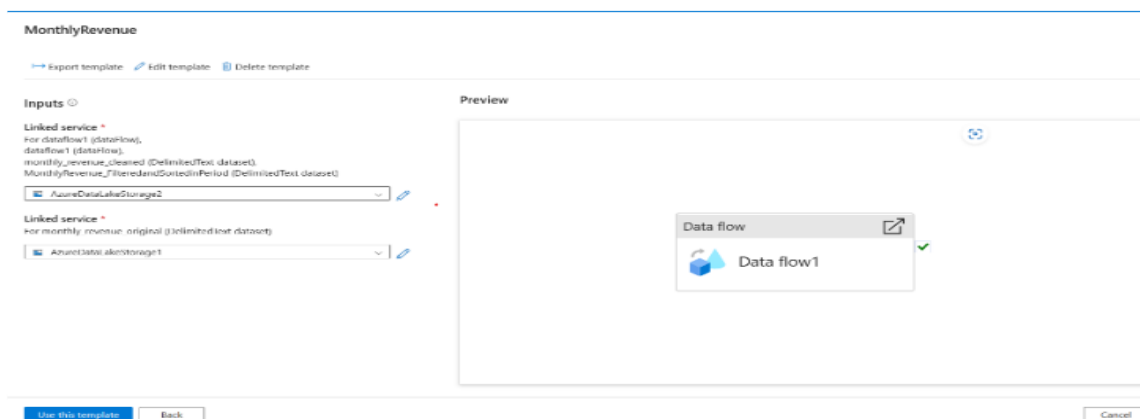**Adding Databricks Service using existing cluster and the access token**



**Adding databricks notebook -which is already made in different file- to our dataflow**

Triggering the pipeline after updating the original dataset with one more row:
- The dataflow is updated
- The notebook is updated as well.

creating workflow job in databricks with Git Repo integration to the databricks notebook located in the Repo. then trigger it.

It shows the updated workflow of the notebook and gets updated in Git Repo.

**MonthlyRevenue run**

Delete job run

Output

Dashboard: MonthlyRevenueDashboard | Export as HTML

# MonthlyRevenueDashboard

## Monthly Revenue Dataset

| | Period | Revenue | Sales_quantity | Average_cost | The_average_annual_payroll_of_the_region | Year |
|---|---|---|---|---|---|---|
| 60 | 2019-12-01 | 58756472 | 38069 | 1543 | 29878525 | 2019 |
| 61 | 2020-01-01 | 56288300 | 27184 | 2070 | 29044998 | 2020 |
| 62 | 2020-02-01 | 40225244 | 23509 | 1711 | 29044998 | 2020 |
| 63 | 2020-03-01 | 50022164 | 32569 | 1535 | 29044998 | 2020 |
| 64 | 2020-04-01 | 52320692 | 26615 | 1965 | 29044998 | 2020 |
| 65 | 2020-05-01 | 52320692 | 26615 | 1965 | 29044998 | 2020 |
| 66 | 2020-06-01 | 58756472 | 38069 | 1543 | 29878525 | 2020 |
| 67 | 2020-07-01 | 44160416 | 26603 | 1659 | 29878525 | 2020 |

## From Statistics,it is found that:

• Revenue increases through the years with increasing the

## Yearly Revenue

50M
40M

| | |
|---|---|
| Duration | 24s |
| Status | ✓ Succeeded |

**Git** ⓘ

| | |
|---|---|
| Git URL | https://github.com/RovanElgendy/Rovan-Elgendy-Wiley_s-Training |
| Branch | main |
| Commit | 806a7756 |

**Notebook**

project/Project_Code_Demo/Azure Databricks Project notebook ( 806a7756)

**Compute**

✓ Rovan Elgendy's Cluster
Driver: Standard_DS3_v2 · Workers: Standard_DS3_v2 · 0 workers · 11.3 LTS (includes Apache Spark 3.3.0, Scala 2.12)

View details | Spark UI | Logs | Metrics

---



**relgendy | Containers**
Storage account

+ Container | Change access level | Restore containers | Refresh | Delete | Give feedback

Search containers by prefix | Show deleted containers

| Name | Last modified | Public access level | Lease state | |
|---|---|---|---|---|
| $logs | 3/20/2023, 4:02:19 PM | Private | Available | ** |
| $web | 4/5/2023, 1:29:35 AM | Blob | Available | ** |
| historicaldata | 3/20/2023, 4:17:14 PM | Container | Available | ** |

Search
Overview
Activity log
Tags
Diagnose and solve problems
Access Control (IAM)
Data migration
Events
Storage browser

Data storage
Containers
File shares
Queues
Tables

**creating $web container in blob storage to store web files .html to present the workflow of the project.**