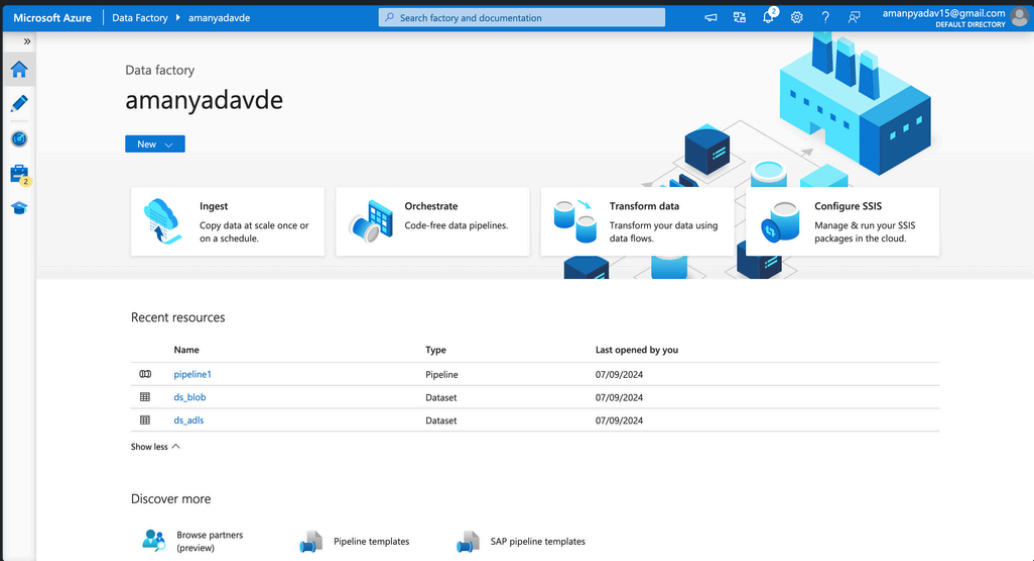# Creation of a data Pipeline

## Setting Up the Data Factory & Creating Storage Accounts

1. Create and Setup Azure Data Factory.
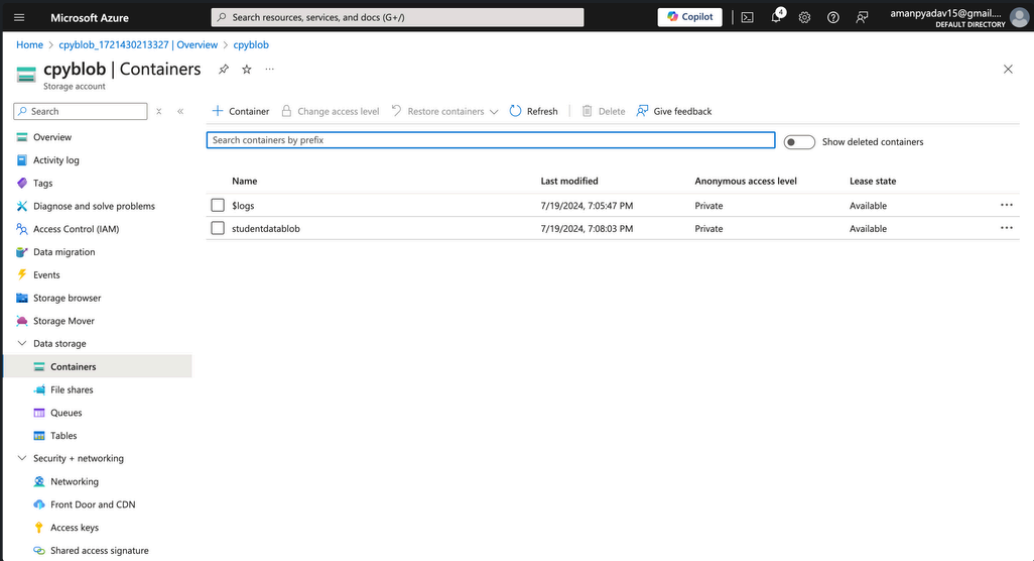


Azure Data Factory
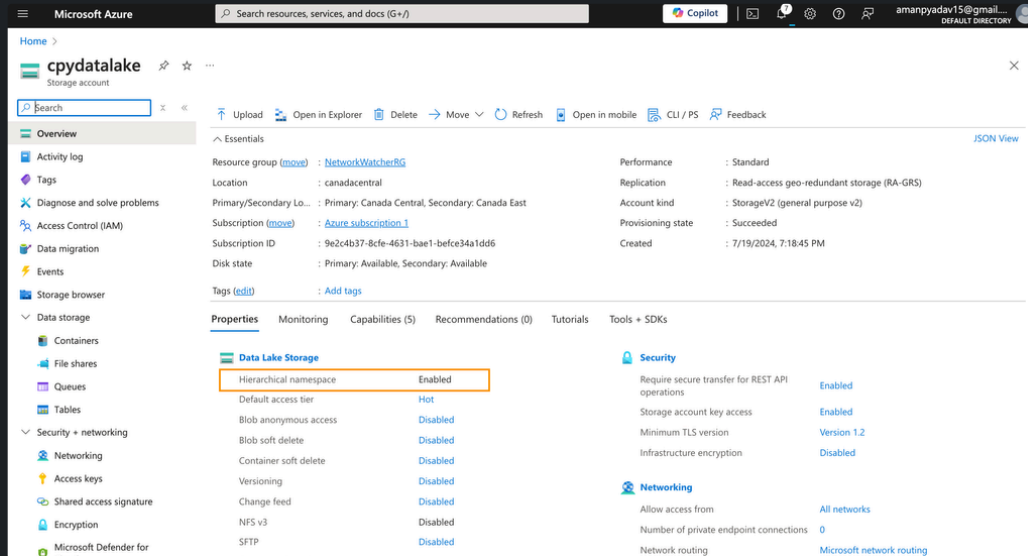
2. Create Blob Storage Account.

## Blob Storage

- We'll create first blob storage account in the name **cpyblob**. Inside blob storage account we'll create a container (studentdatablob) and upload our student dataset as a .csv format.



Blob Storage Account

## Azure Data Lake Storage Account

- We'll follow the same step as we perform in creating a blob storage account but this time we'll just **enable Hierarchical namespaces**
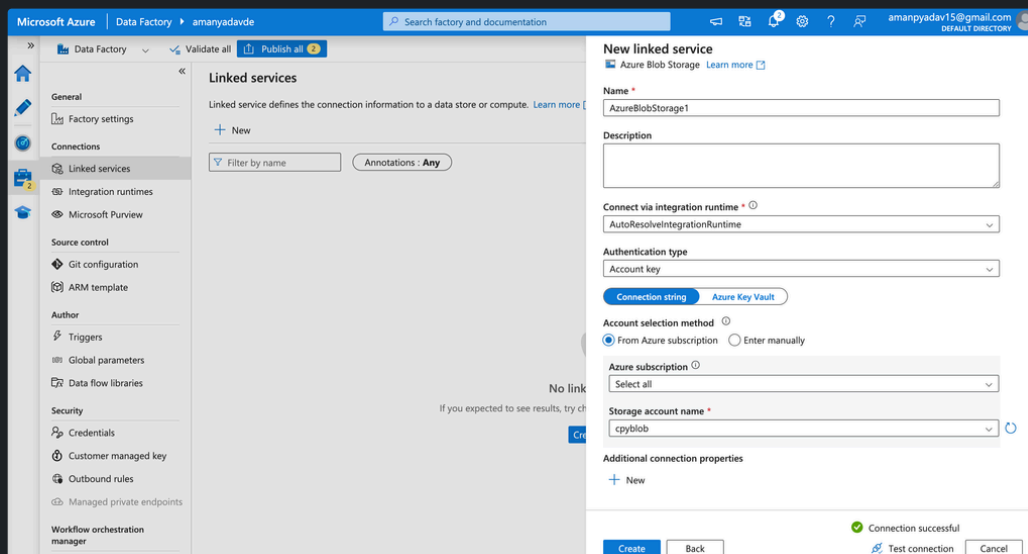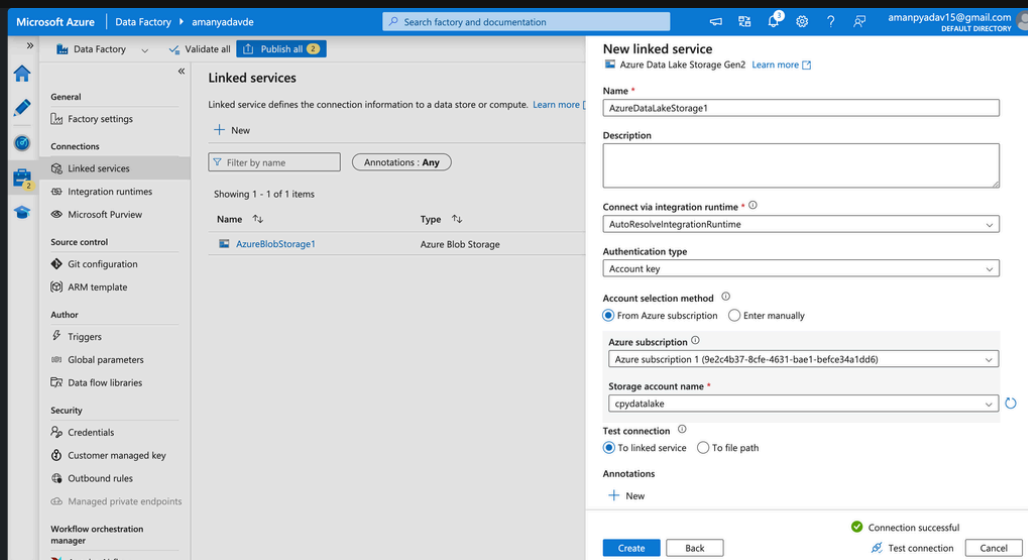


Azure Storage Account

# Pipeline 1

**Overview**

We have a student dataset and we'll be performing dataflow from Blob Storage to Azure Data Lake and perform some operations.

## Creating Linked Services

Step 1. After launching the Azure Data Factory. We'll navigate to Manage and create a Blob linked service & Data Lake Linked Services. While creating that we'll make sure to test connection before creating it to make sure it's been linked properly with our blob storage.
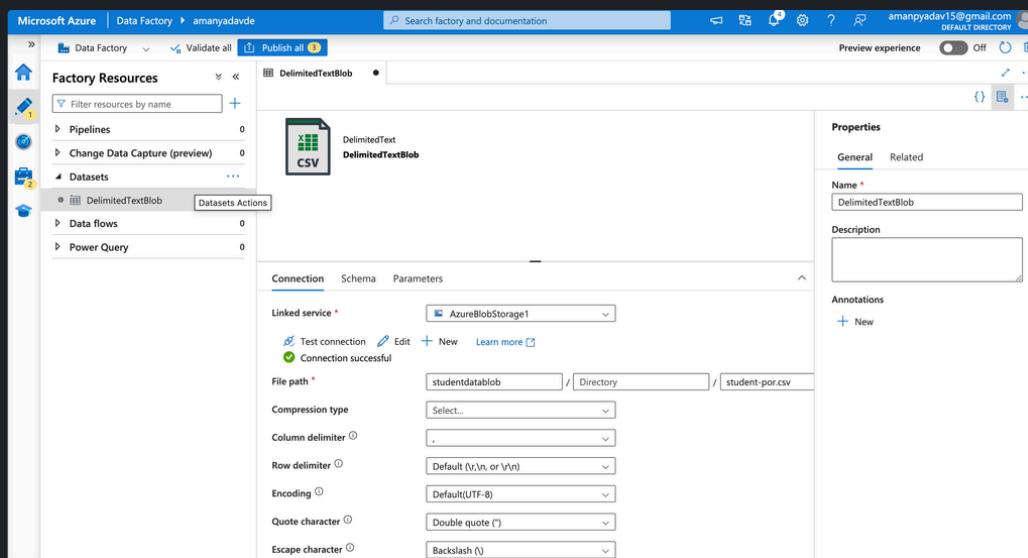


Blob Linked Service
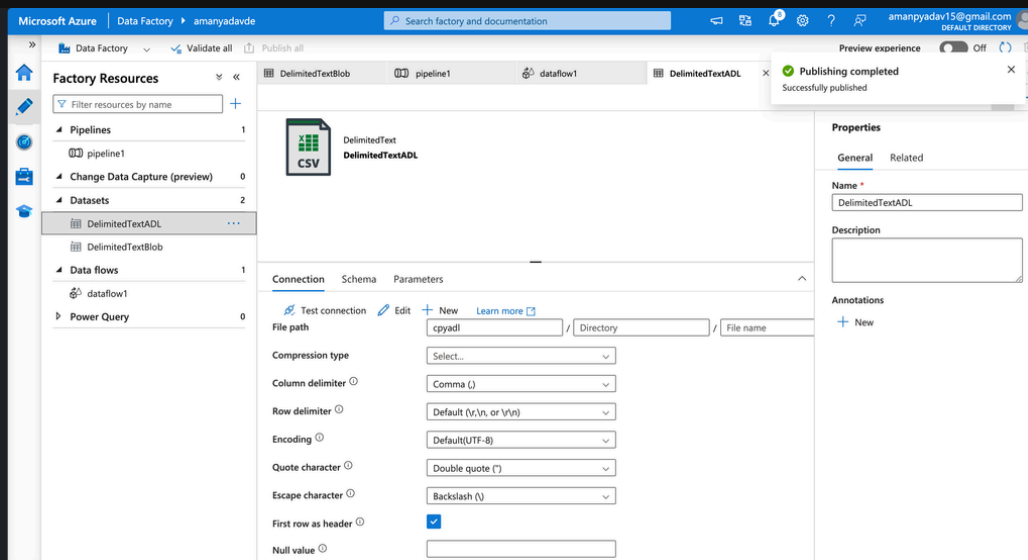
Data Lake Linked Service

## Creating Dataset

Step 2. We'll create Dataset for both Blob Storage and Azure Data Lake Storage. To Create Dataset navigate to Author → Datasets → Create New Datasets.
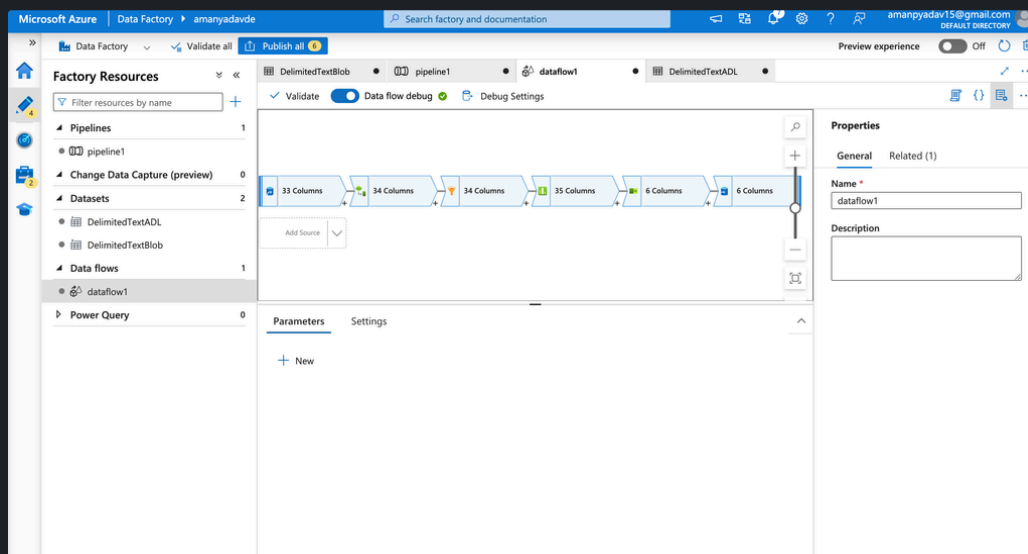

Blob Datasets

In Azure Data Lake Dataset we'll be using Deilmited Text as our format.
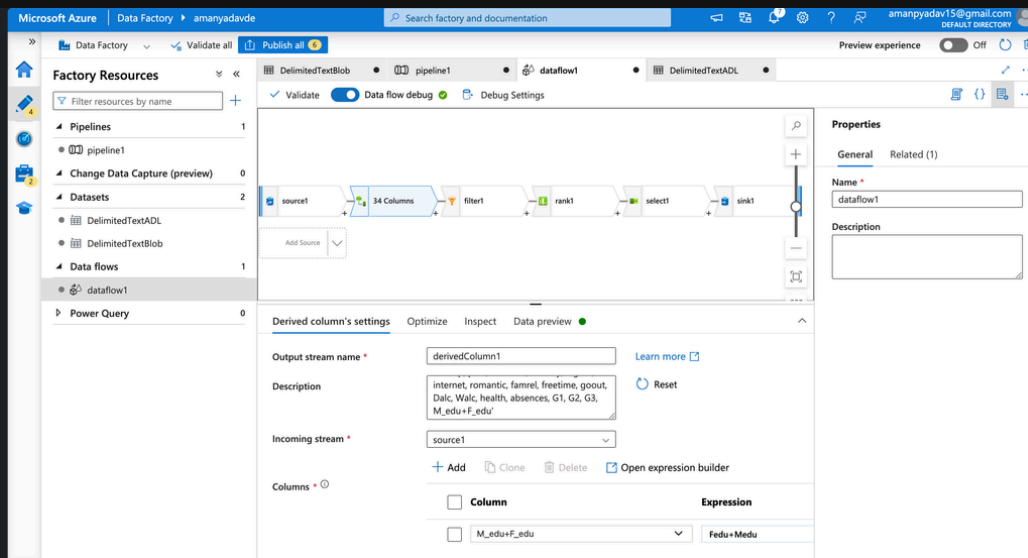
Azure Data Lake Dataset

## Creating Pipeline

Step 3. We'll remain under author section and we'll navigate to Pipelines → Create New Pipeline.
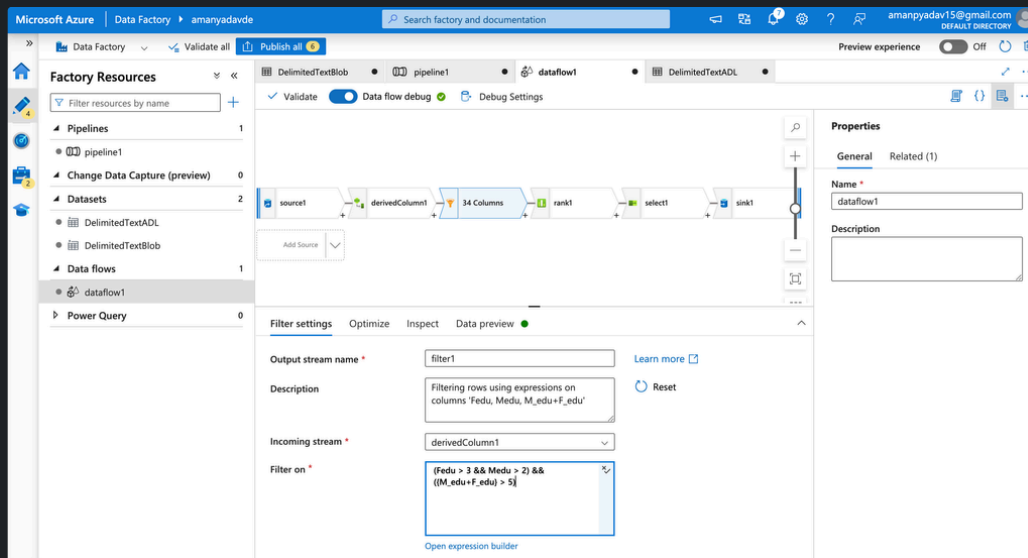


Pipeline 1

In Pipeline I have Defined my source as **DelimitedTextBlob**. I've linked then DerivedColumn in which I've used added the values of Fedu and Medu.
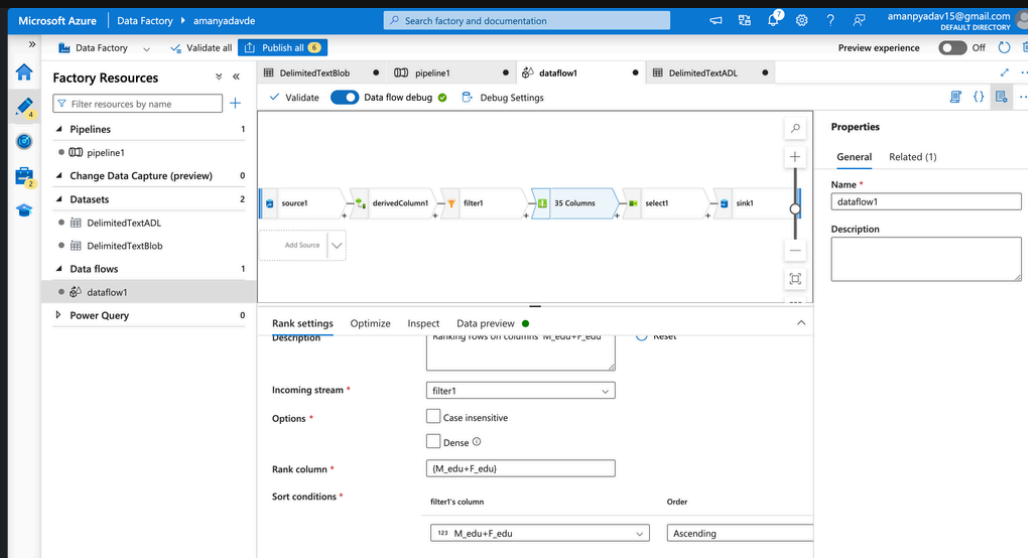
Derived Column - Pipeline 1

Then I've filtered the data based on the following conditions:

(Fedu > 3 && Medu > 2) && ({M_edu+F_edu} > 5)


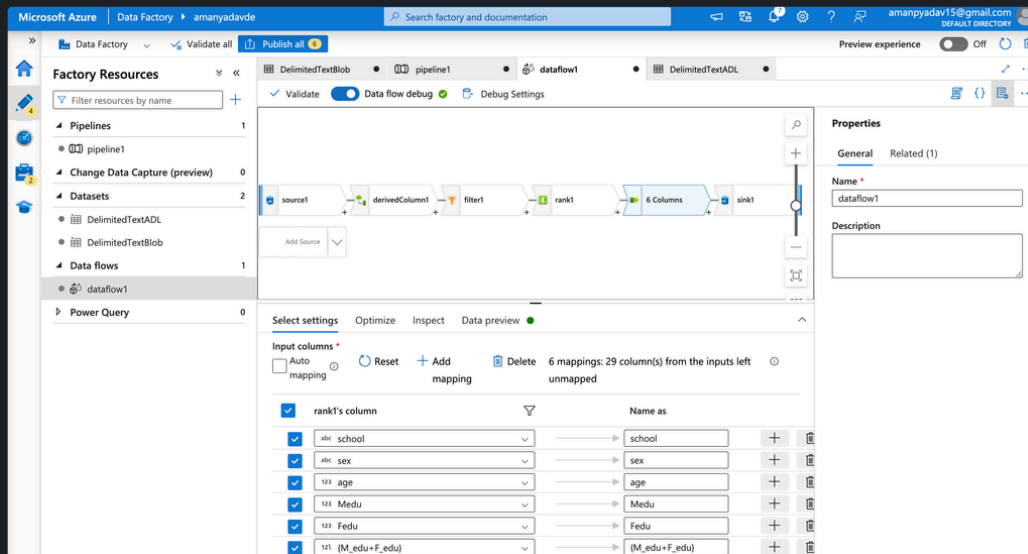Filtered - Pipeline 1
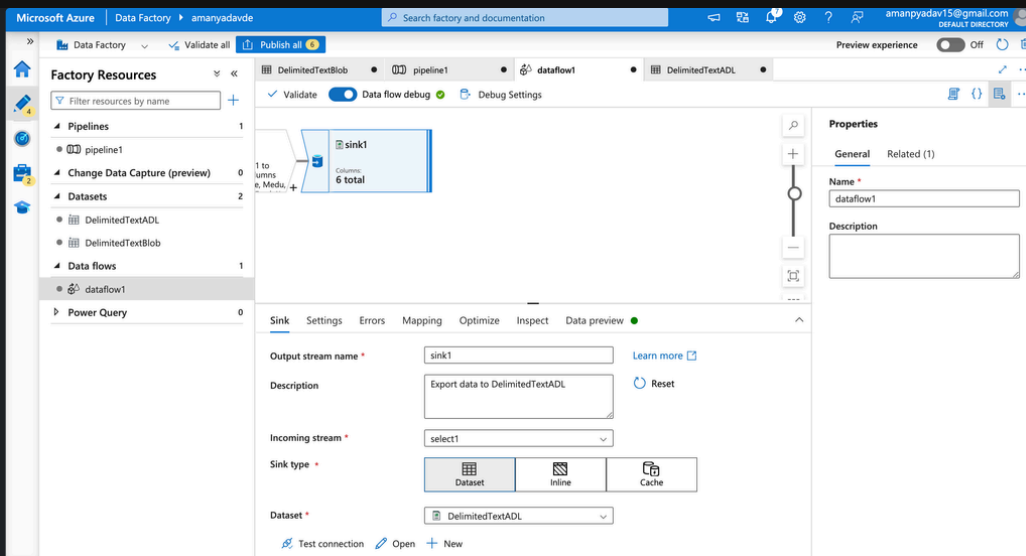
Then order the column based on M_Edu+F_edu

Ranking - Pipeline

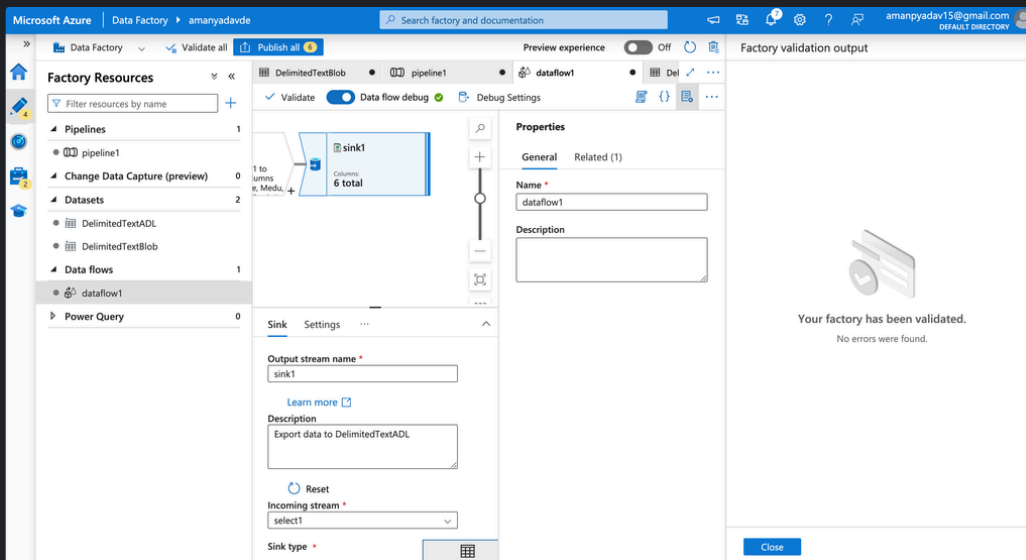Then I selected only School, Sex, Age, Medu, Fedu and {M_edu+F_edu}.


Selection - Pipeline 1

Then sinked into DelimetedTextADL
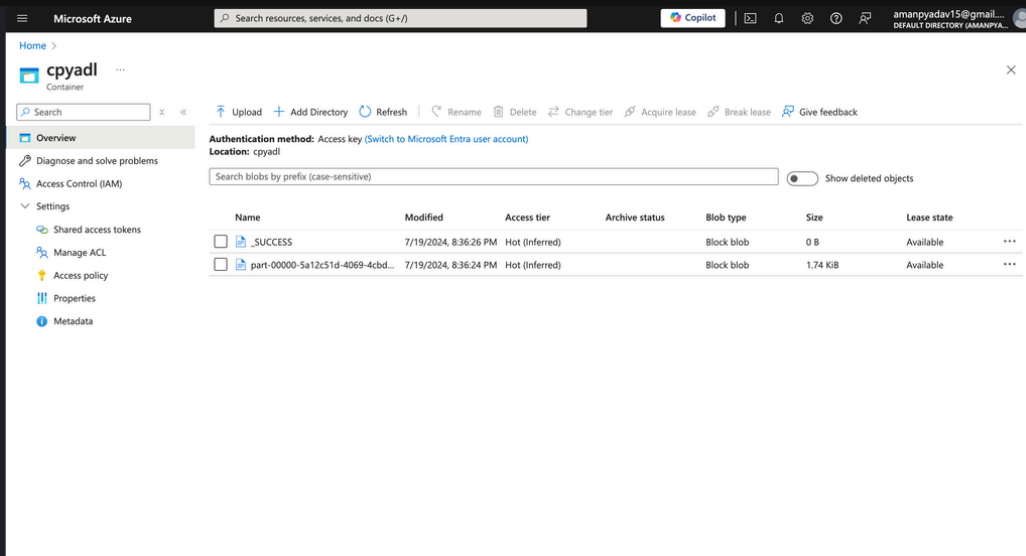
Sink - Pipeline 1

Once done with it I clicked on "Validate All"


Validating & Publishing

Once there were no errors found I've published it.

Checking on our destination dataset.

**cpyadl** ...
Container

Search

- Overview
- Diagnose and solve problems
- Access Control (IAM)
- Settings
  - Shared access tokens
  - Manage ACL
  - Access policy
  - Properties
  - Metadata

⬆ Upload  + Add Directory  🔄 Refresh  |  Rename  🗑 Delete  Change tier  Acquire lease  Break lease  Give feedback

**Authentication method:** Access key (Switch to Microsoft Entra user account)
**Location:** cpyadl

Search blobs by prefix (case-sensitive)   ⬤ Show deleted objects

| Name | Modified | Access tier | Archive status | Blob type | Size | Lease state | |
|------|----------|-------------|----------------|-----------|------|-------------|---|
| _SUCCESS | 7/19/2024, 8:36:26 PM | Hot (Inferred) | | Block blob | 0 B | Available | ... |
| part-00000-5a12c51d-4069-4cbd... | 7/19/2024, 8:36:24 PM | Hot (Inferred) | | Block blob | 1.74 KiB | Available | ... |

Desitnation Directory

Output Data:

select1_data_preview

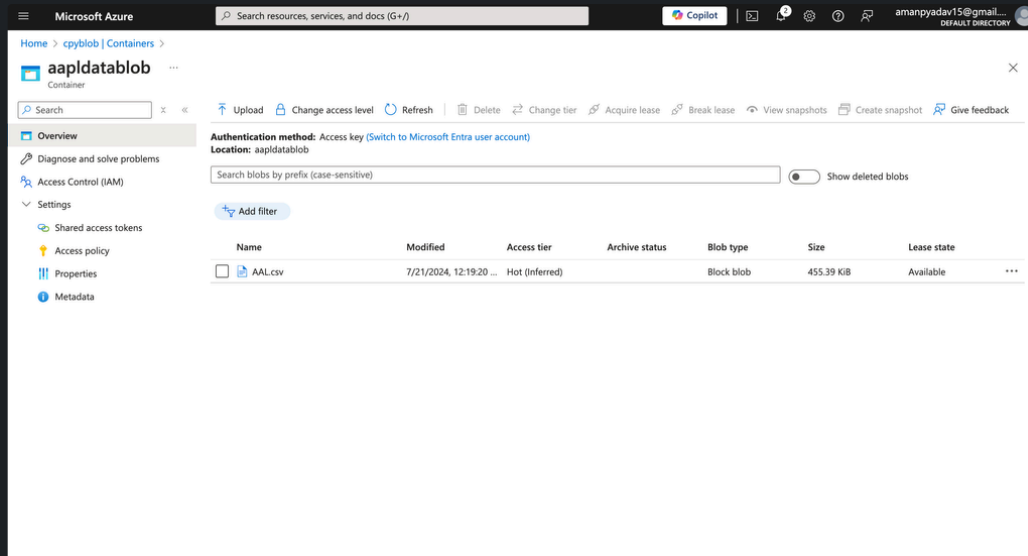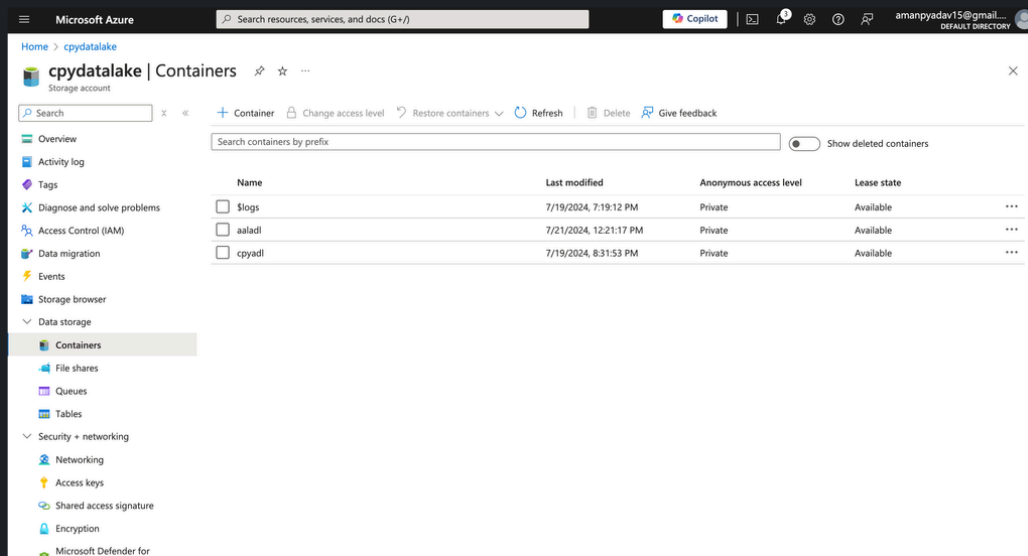| school | sex | age | Medu | Fedu | M_edu+F_edu |
|--------|-----|-----|------|------|-------------|
| GP | M | 15 | 3 | 4 | 7 |
| GP | M | 16 | 3 | 4 | 7 |
| GP | F | 15 | 3 | 4 | 7 |
| GP | F | 16 | 3 | 4 | 7 |
| GP | M | 15 | 3 | 4 | 7 |
| GP | M | 15 | 3 | 4 | 7 |
| GP | M | 15 | 3 | 4 | 7 |
| GP | F | 15 | 3 | 4 | 7 |
| GP | F | 15 | 3 | 4 | 7 |
| GP | F | 16 | 3 | 4 | 7 |
| GP | M | 15 | 3 | 4 | 7 |
| GP | M | 17 | 3 | 4 | 7 |
| GP | M | 16 | 3 | 4 | 7 |
| GP | F | 17 | 3 | 4 | 7 |
| GP | M | 18 | 3 | 4 | 7 |
| GP | F | 17 | 3 | 4 | 7 |
| GP | F | 18 | 3 | 4 | 7 |
| GP | F | 17 | 3 | 4 | 7 |
| GP | M | 18 | 3 | 4 | 7 |
| GP | M | 18 | 3 | 4 | 7 |
| MS | M | 16 | 3 | 4 | 7 |
| MS | F | 16 | 3 | 4 | 7 |
| MS | M | 17 | 3 | 4 | 7 |
| GP | F | 18 | 4 | 4 | 8 |
| GP | F | 17 | 4 | 4 | 8 |
| GP | F | 15 | 4 | 4 | 8 |
| GP | M | 15 | 4 | 4 | 8 |

# Pipeline 2

Dataset: AAL stocks

**Objective**

To analyze the market trend of Apple.

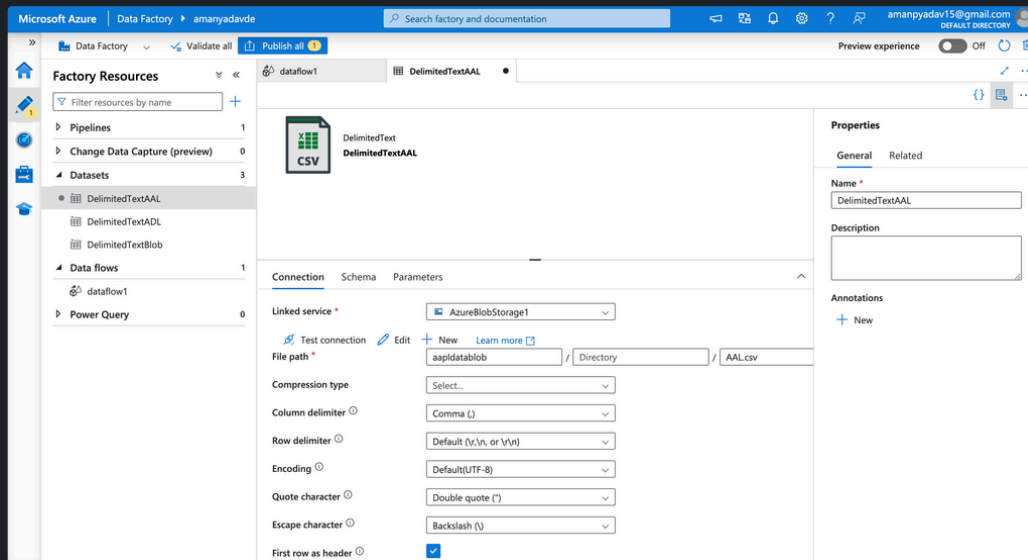We will first upload our csv files in our blob storage account which we previously created.
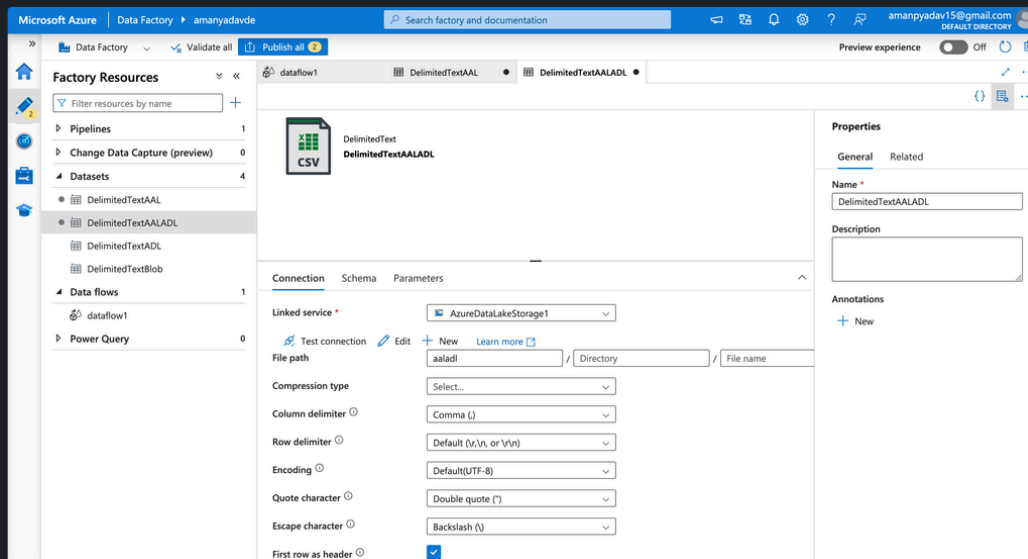


Blob Storage Account

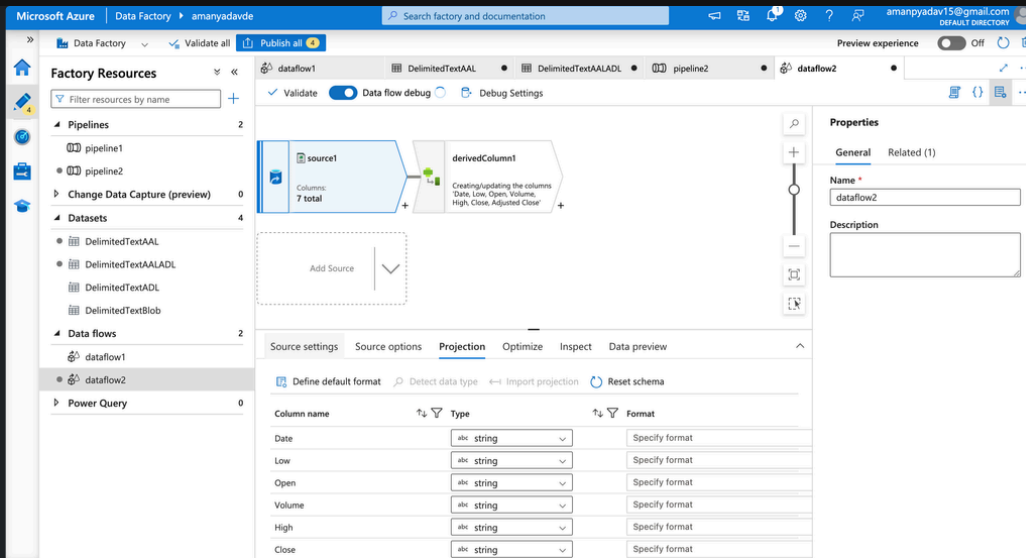Then will just create a new container at our Data lake Storage Account as **aaladl** .



Then in our azure data factory studio we will create a new dataset for our blob storage and we will just be using the linked services which we created previously.
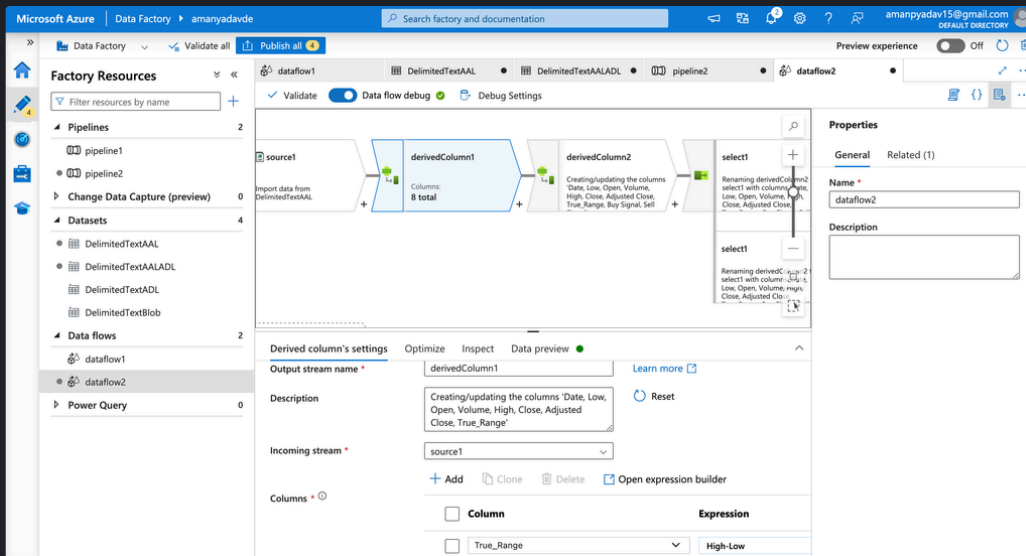
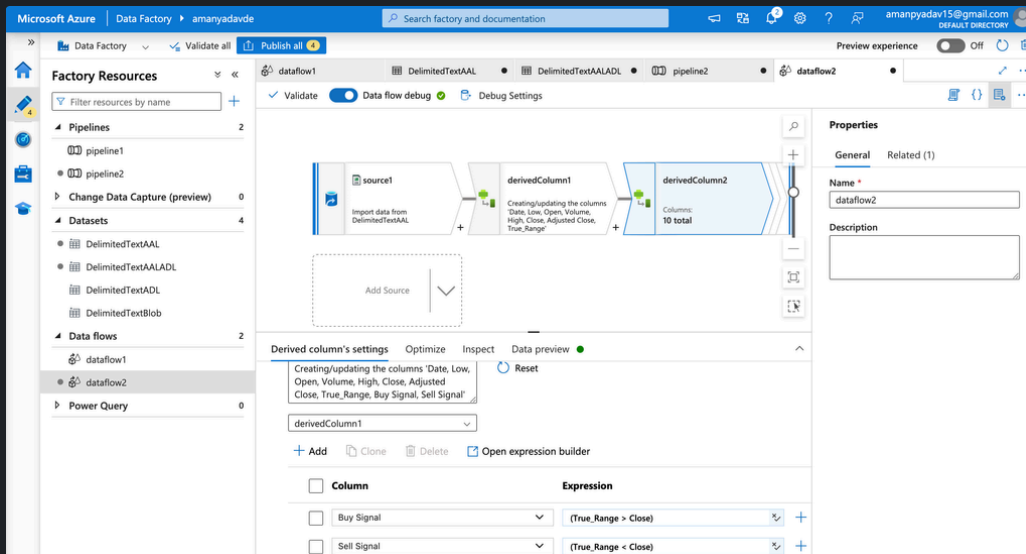Creating a dataset for our data lake storage as well.
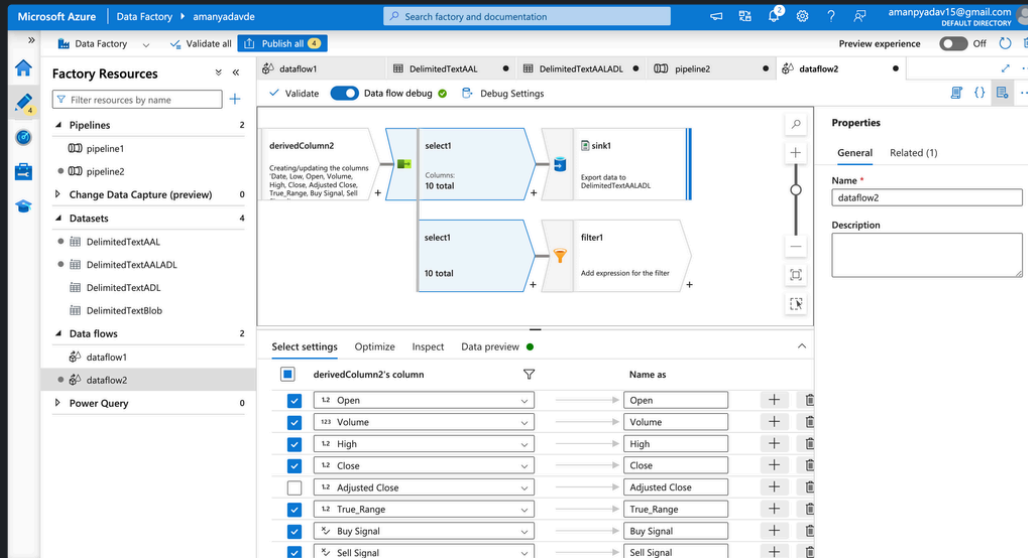


Creating a new pipeline 2 with dataflow2.

Added Derived column in which the column name is **True_Range** and expression is High - Low
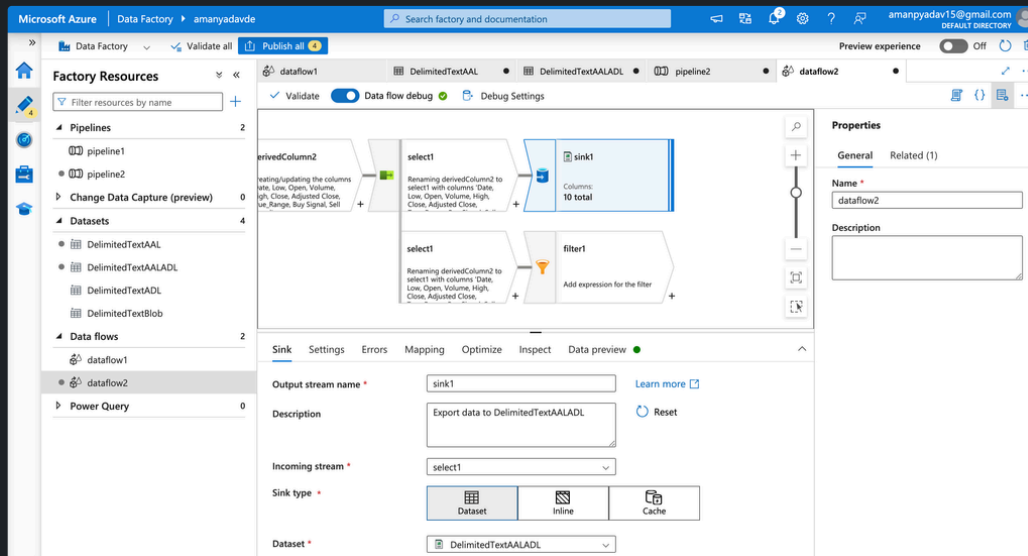


Added another derived column in which I've added two more columns name **BUY SIGNAL** and **SELL SIGNAL**
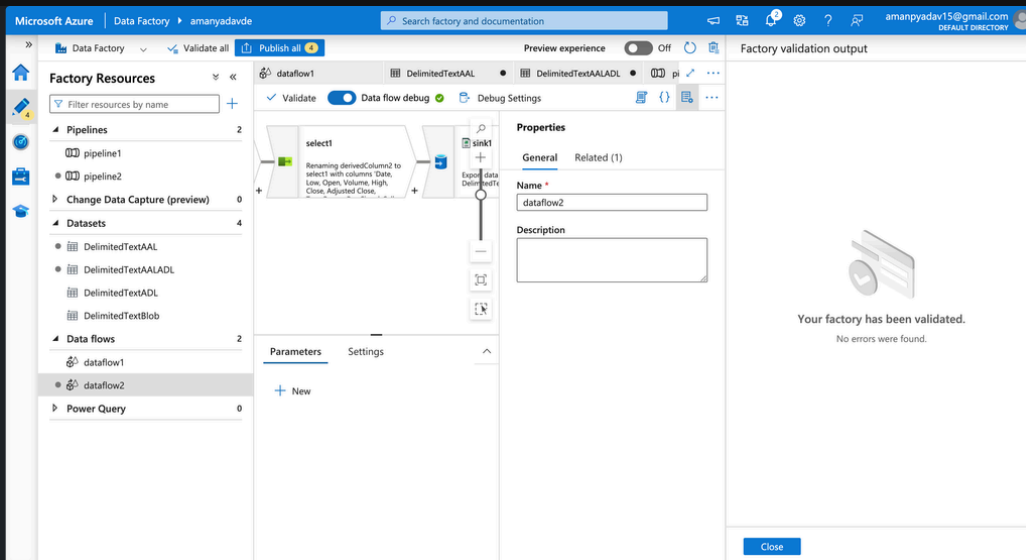
Then selected the column below:



Then sinked in to our destination storage account i.e. Data Lake Storage.



Validated everything to see if there's any error.

Published it, and let's check in our destination folder whether we got the file or not.

Output Data: