# AZURE DATA FACTORY

# CAPSTONE – COVID USE CASE

**Emp id : 2320074**

**Cohort id : CSDAIA24AZ003**

# Introduction

The purpose of the Covid use case exercise is to learn how to build a real-world data pipeline in Azure Data Factory (ADF) to analyze the covid trend across the regions using Azure cloud data services. By performing this case study, you will learn.
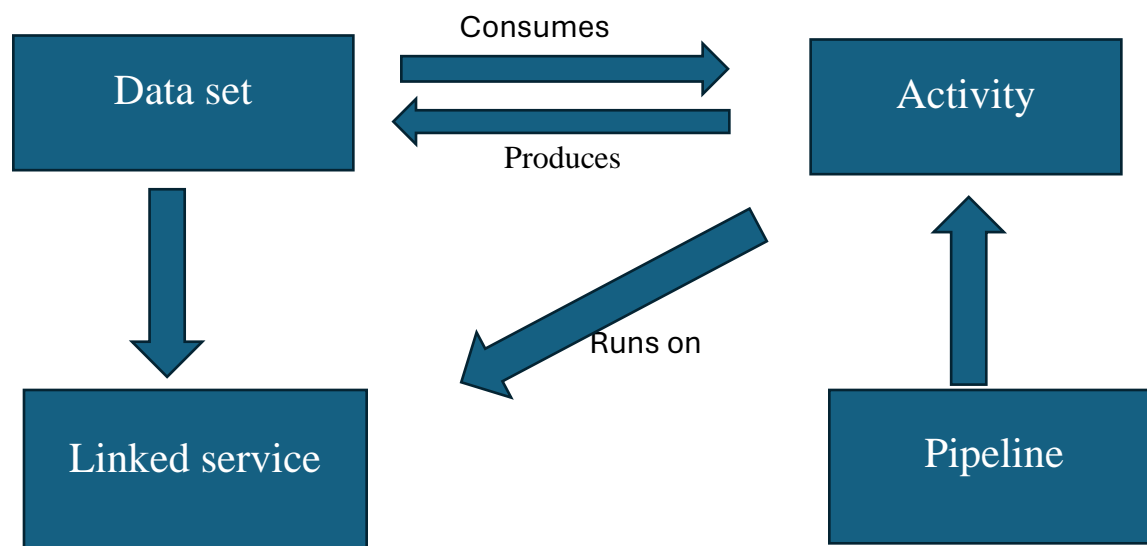
- How to ingest data from flat files into Azure Data Lake Gen2 and Azure Synapse using Azure Data Factory (ADF)
- How to transform data using Data Flows in Azure Data Factory (ADF) and load into Azure Synapse

Through this exercise, you will be having a hands-on experience on Storage, ADF Pipeline, Mapping Dataflow, Azure Synapse.
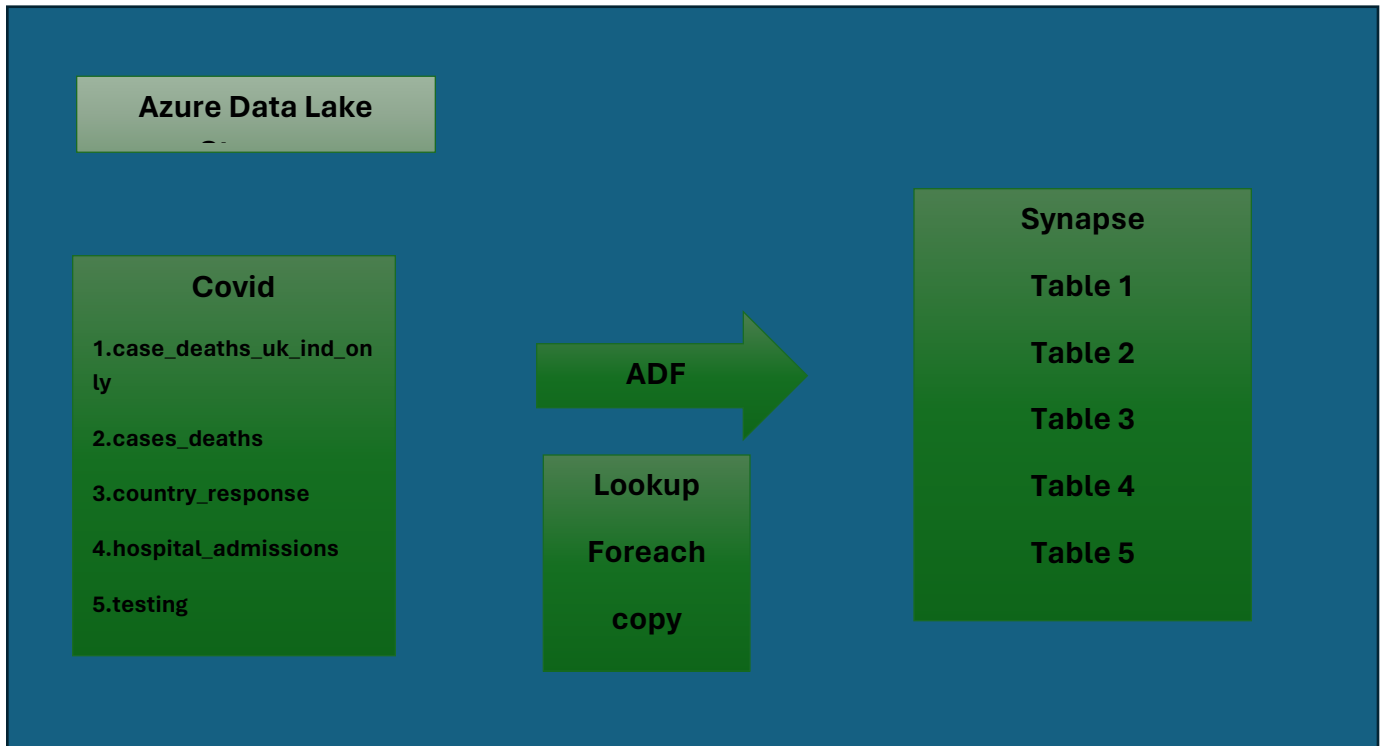
## Purpose of the project

The main aim of this project (Covid Use Case Exercise) is that we will be having a hands-on experience on Storage, ADF Pipeline, Mapping Dataflow, Azure Synapse along with getting to know how to ingest data from flat files into Azure Data Lake Gen2 and Azure Synapse using Azure Data Factory (ADF) and also knowing how to transform data using Data Flows in Azure Data Factory (ADF) and load into Azure Synapse. This report gives a summary of the entire project making us realize and interpret the use case scenario of Azure and its applications.
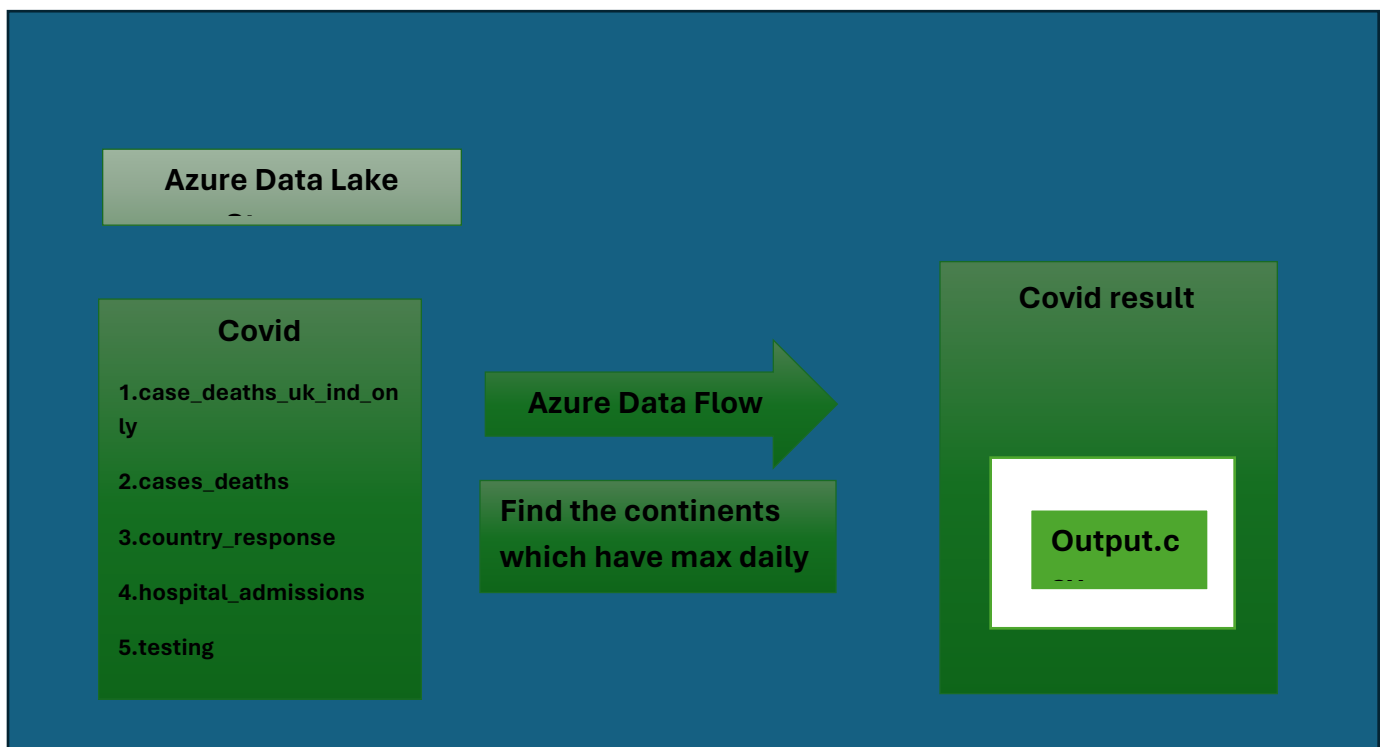
## Architecture of the project

# Requirements of project

## Requirement 1



## Requirement 2

# Procedure of Requirement 1



**Step 1:** created one Resource group and required resources for project like storage account, Synapse workspace (data warehouse), Azure Data Factory.
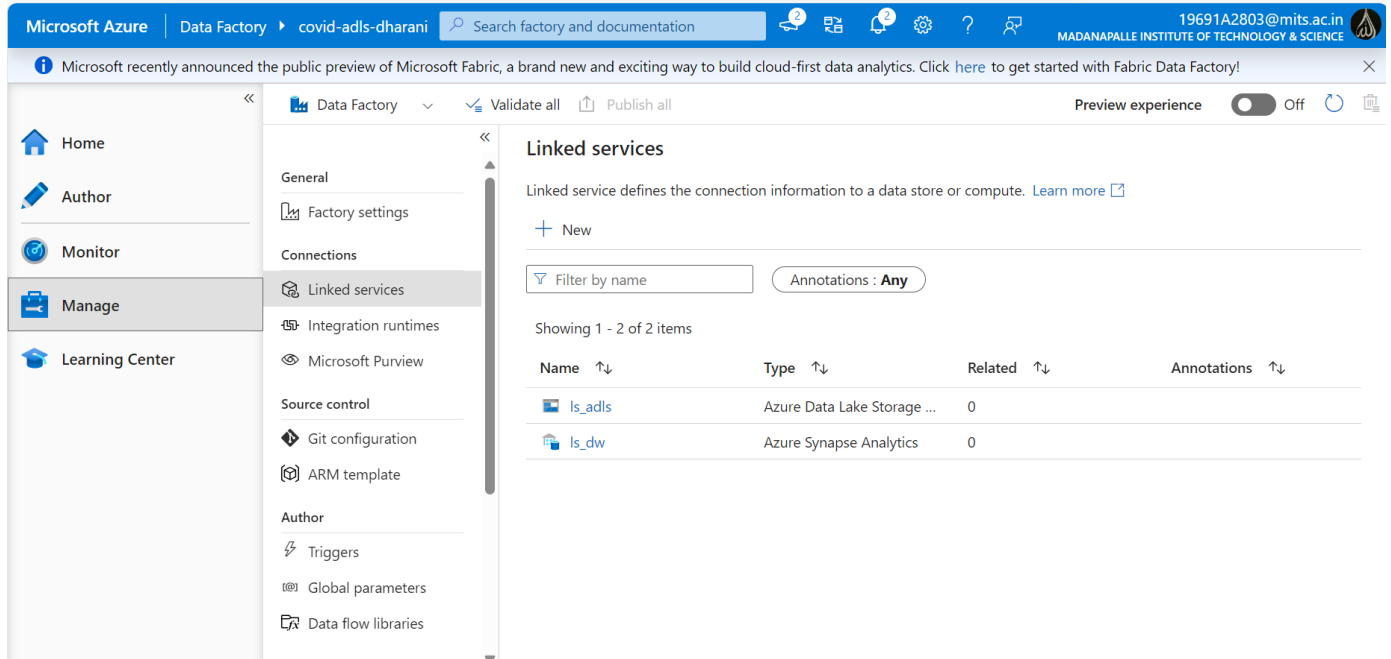


**Step 2 :** Created Container with name "**covid**" in Storage account(data lake) for holding folder which contains flat files.

**Step 3 :** Create Folder with name **"Ingest"** inside the container "covid" and uploaded csv files(data sets) from personal computer which are there in zip file given in project document.



**Step 4 :** Create Azure Synapse resource and one dedicated pool inside the azure synapse for data warehouse creation and it should be turn on.

**Step 5 :** Creating **two linked services** as per the project requirement in Azure data factory

° **Storage account (data lake) to Azure Data factory.**

°**Azure synapse workspace to Azure Data factory.**



**Step 6 :** Created dataset **(ds_adls)** for fetching flat files from storage account(data lake) which are present in ingest folder inside covid container.

**Step 7 :** Created required SQL Tables in Synapse SQL database by writing create table queries in Synapse workspace (SQL Script).



**Step 8 :** Created **dataset(ds_configdataset)** for inserting into SQL Tables created in synapse(data warehouse).

**Step 9 :** Created **stored procedure(sp3final)** for fetching records from parameters table.



**Step 10 :** Created dataset(ds_configdataset ) for fetching parameters table from synapse(data warehouse).

**Step 12 :** Before creating the pipeline in the data factory, we need to turn on the dedicated pool. we need to check these two to three times while moving on to creation of pipeline.



**Step 13 :** Drag and Drop the LOOK UP activity into pipeline workspace and set the source dataset (ds_configdataset) for lookup activity and choose the option **Stored procedure** and give the stored procedure name created in synapse(data warehouse).

**Step 14 :** After clicking on debug, lookup will be run successfully.



**Step 15 :** Drag and Drop **For each activity** in pipeline workspace and configure the for each activity settings like **Items** with output of lookup activity (**@activity('Lookup1').Output.value**).

**Step 16 :** Click on add activity symbol present on foreach activity and inside foreach activity add a **copy activity** for copy data from CSV file into SQL Table.

Configure settings at source side in copy activity by giving dataset**(ds_adls)** and giving folder name and file name dynamically by taking from foreach activity by **item**.

**Folder Name(@{item().FolderName})Name(@item().FileName})**



**Step 17 :** Configure setting in copy activity at sink by giving dataset(ds_configdataset) and giving sqltableName dynamically by taking from foreach activity by item->sqltableName({@item().sqltableName})

**Step 18 :** After setting whole pipeline by using lookup and foreach activity recheck all parameters given in each configuration setting, check the dedicated pool is turn on and then turn on the debug option in pipeline. Finally, all the activities are successfully ruined.



**Step 19 :** After successfully run of pipeline ,now we need to check the data inserted into tables in data warehouse by performing two SQL queries operation given in project documentation.

**Step 20 :** After running first SQL query written in SQL script and it is successfully has ran and given output as per the query.



**Step 21 :** Wrote the second SQL query as per the question given in the project documentation and click on run.

**Step 22 :** Successfully query has ran and given output as per the query.
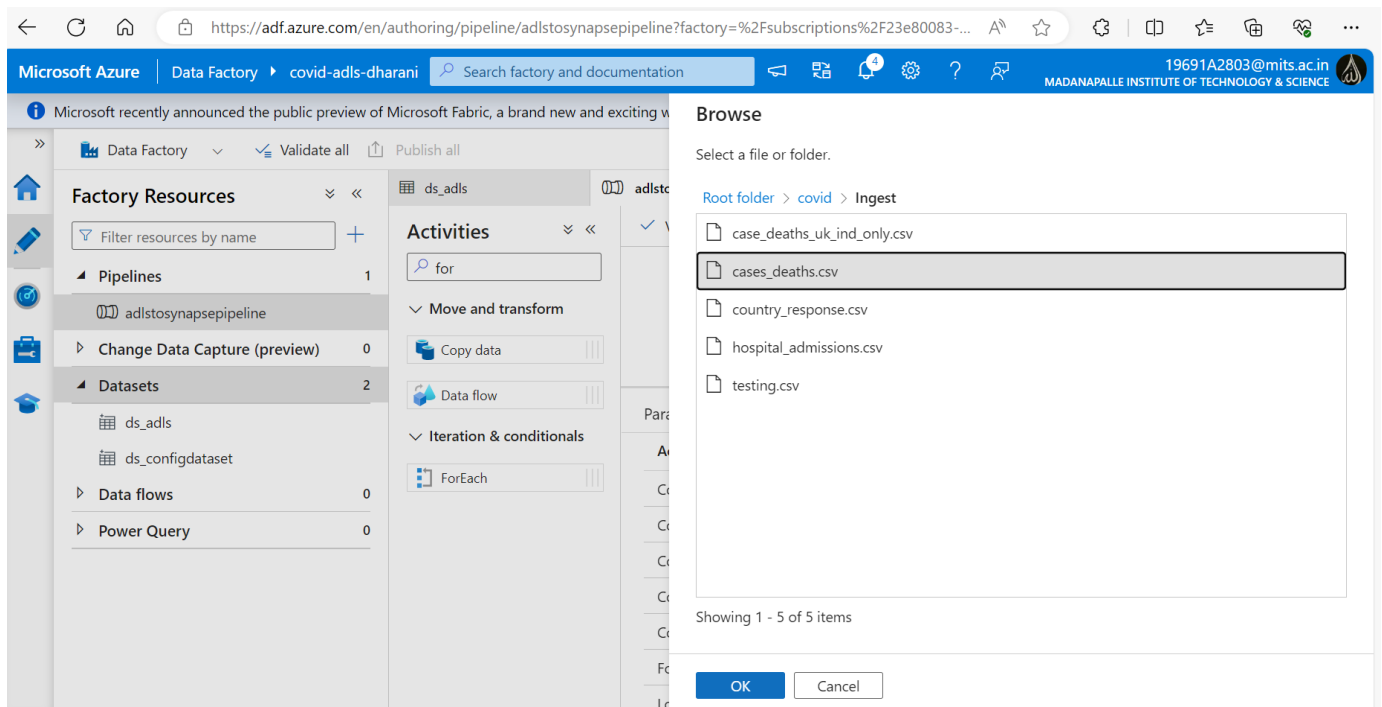


**Step 23 :** Compare the above SQL query output with data which is appear from csv file which is open in excel for reference check.
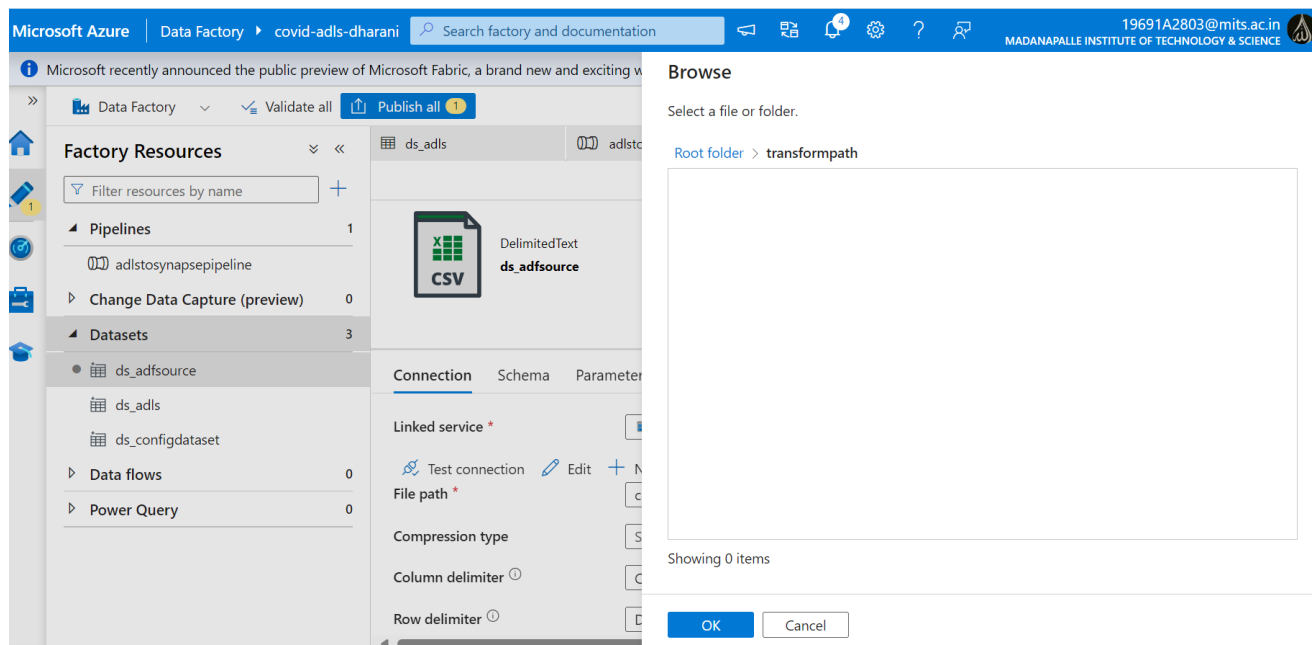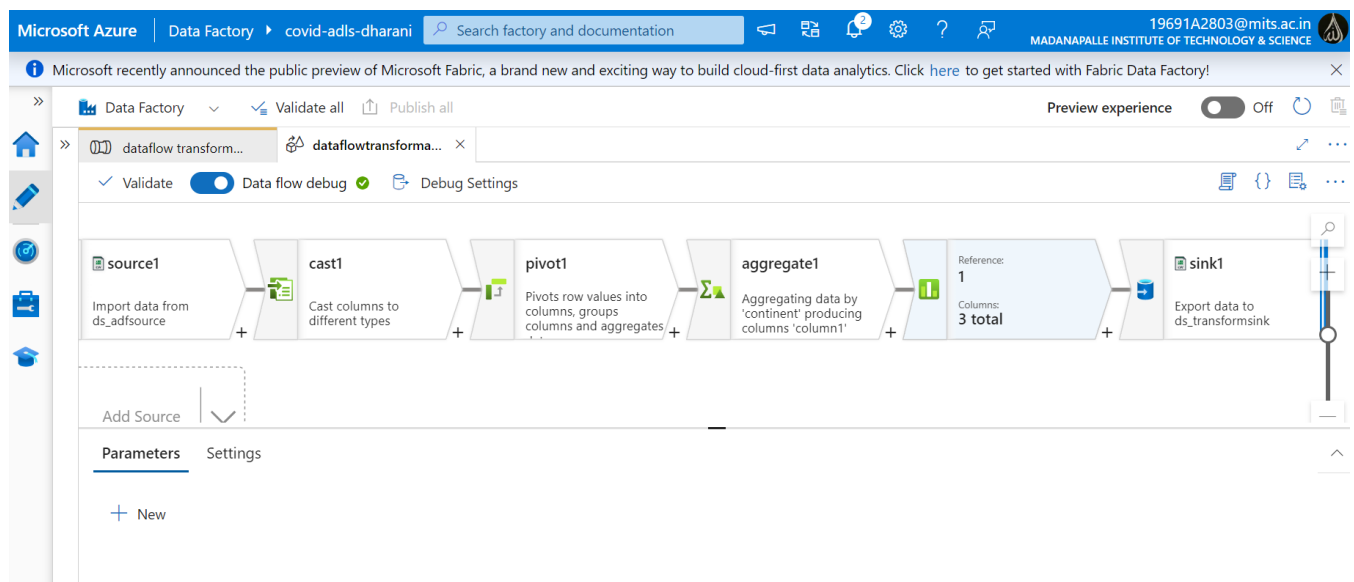
# Procedure for Project Requirement 2 :



**Step 23 :** Create another container with name **"transformpath"** for second requirement given in project for storing transformed data file by using data flow.



**Step 24 :** Create source dataset (ds_adls) for dataflow by giving a file specific file name on which data transformation need to be taken place as per project requirement.

**Step 25 :** Create target dataset (ds_adfsource) for dataflow to keep that data transformed file in specific place for further use.



**Step 26 :** After successfully creating dataflow I click on dataflow debug.

**Step 27 :** After clicking debug, drag the dataflow in pipeline and run the pipeline. you will get the output once it is run successfully.

v

**Step 28 :** After that we need to check file appear in the transformpath container in storage account (data lake). I successfully got that file in my container as per the question given in the project documentation.

# MY DETAILS: