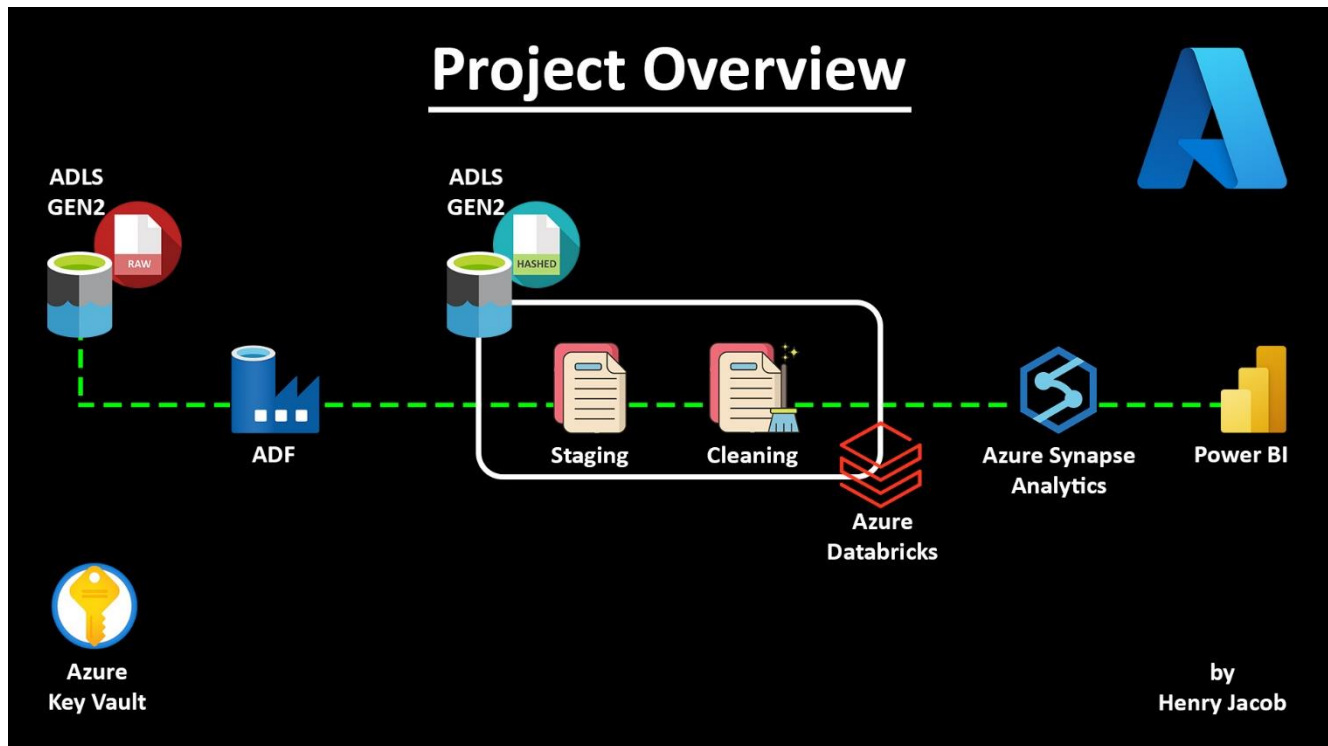# Azure end to end project – Healthcare



## Scenario:

- Sensitive patient information should only be stored in the raw storage.
- Access to sensitive patient data should be restricted after ingestion, and no one should have visibility to this data beyond the ingestion process.
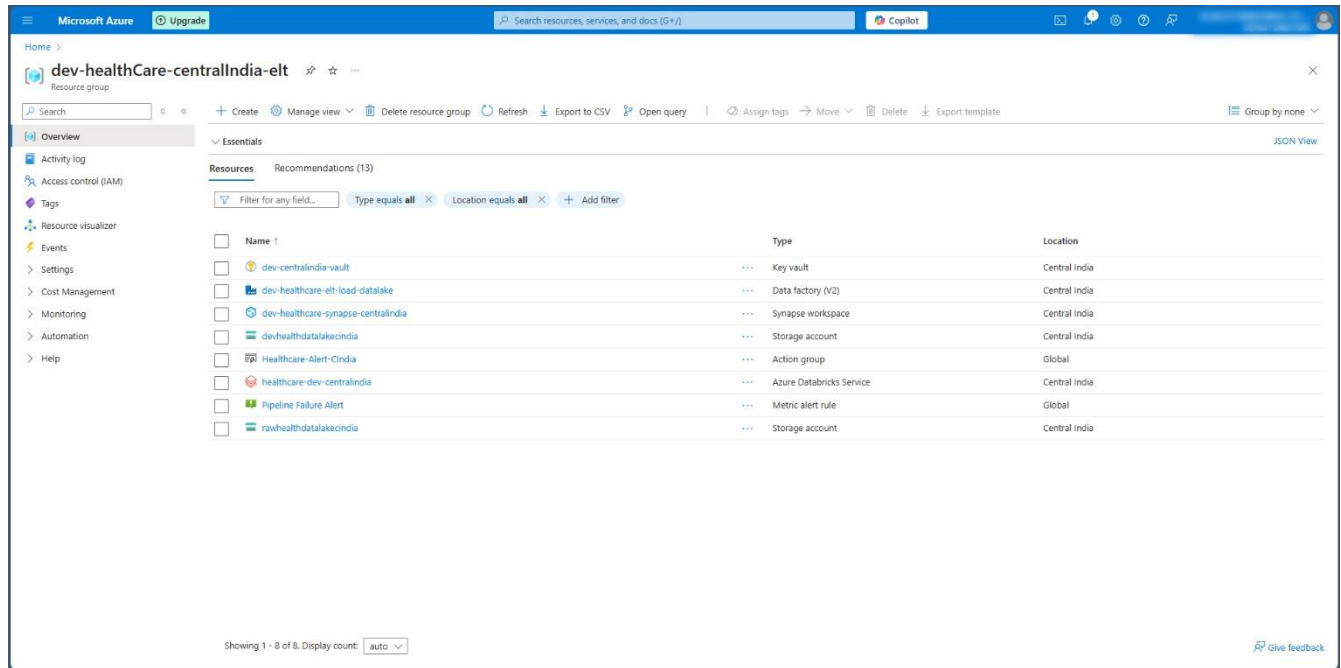
## Assumption:

Let's assume that the ADLS Gen2 raw data is stored outside of Azure.

## Services used:

- Azure Data Lake Gen2 (Raw & Staging and Cleaning)
- Azure Data Factory
- Azure Databricks
- Azure Synapse Analytics
- Power BI
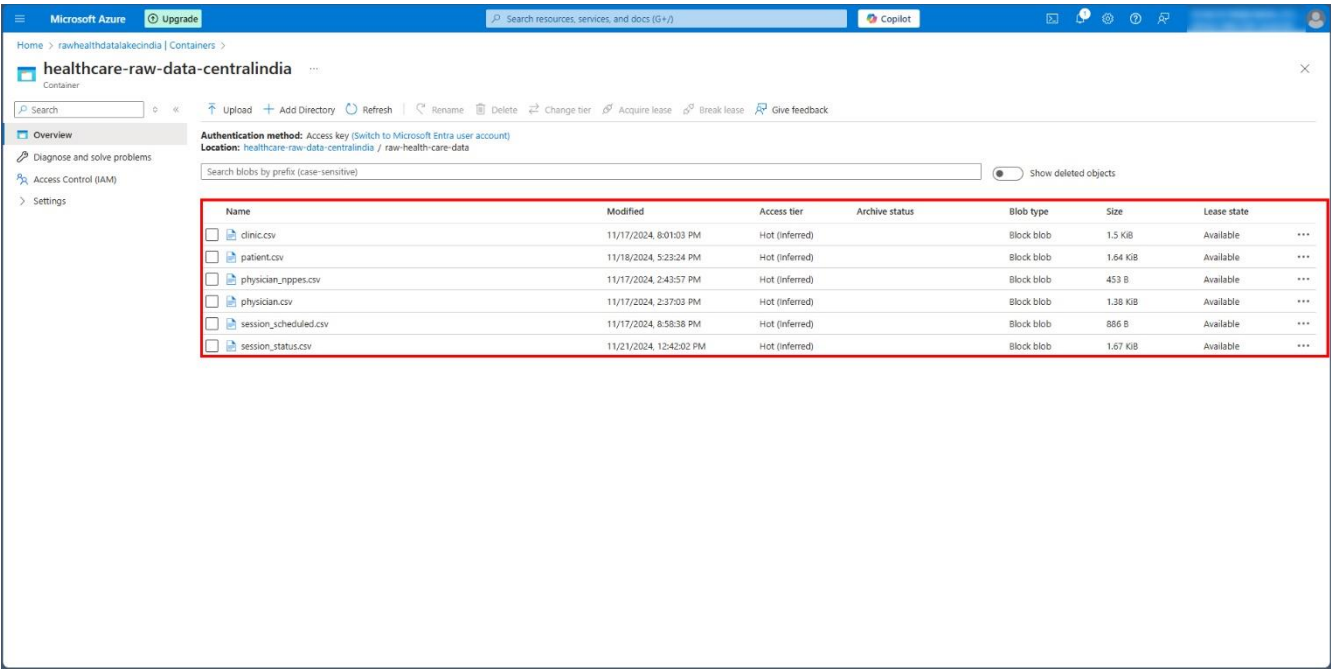- Azure Key Vault

**Resource Group:**



- **dev-centralindia-vault** – To store and use credentials.
- **dev-healthcare-elt-load-datalake** – To retrieve data from the raw storage account, hash, join, and load it into the staging container in the dev storage account.
- **dev-healthcare-synapse-centralindia** – To retrieve data from the cleaned container in the dev storage account and derive insights using SQL queries.
- **devhealthdatalakecindia** – To store the staging and cleaned data in two different containers.
- **Healthcare-Alert-CIndia** (auto-created by ADF) – To group notification channels.
- **healthcare-dev-centralindia** – To retrieve data from the staging container in the dev storage account, clean it, and load it into the cleaned container in the dev storage account.
- **Pipeline Failure Alert** (auto-created by ADF) – To send an email alert if the pipeline fails to run.
- **rawhealthdatalakecindia** – To store raw data only.

As we assumed the raw data is stored outside of Azure, two different storage accounts have been created:

- **rawhealthdatalakecindia** – to store raw data.
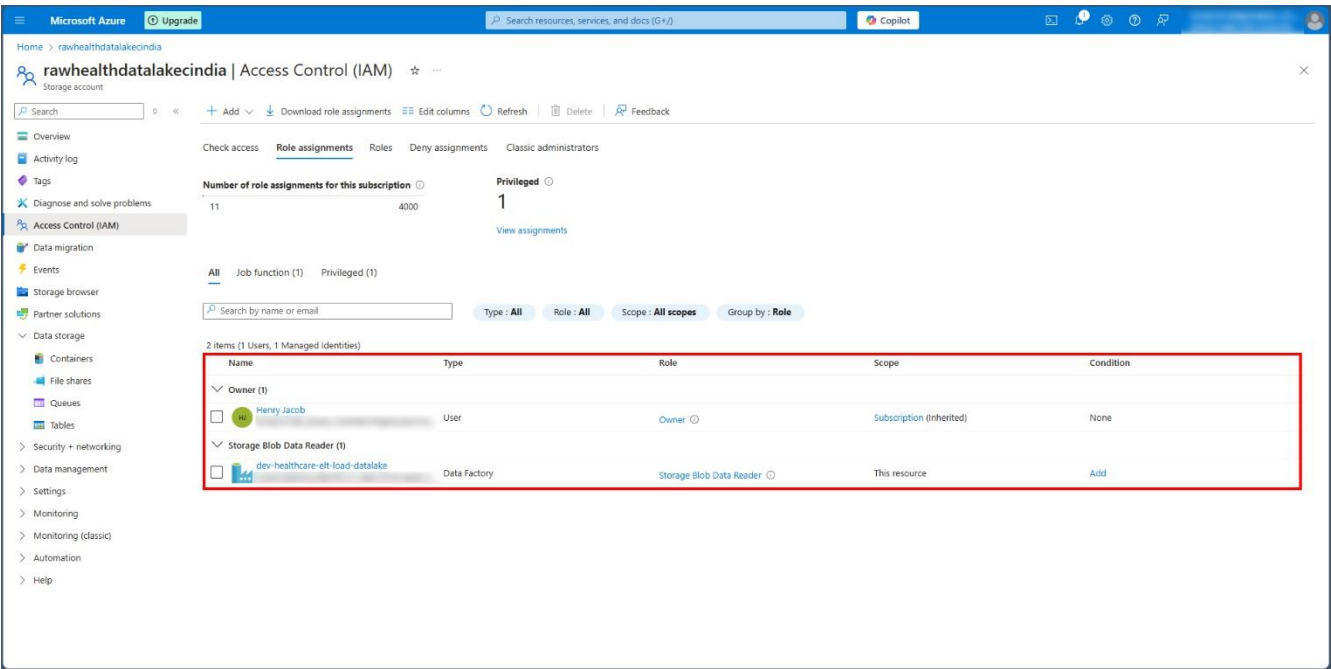- **devhealthdatalakecindia** – to store staged and cleaned data.

## Azure Data Lake Gen2 – Raw Data

The raw files can be accessed [here](here).



## Azure Data Lake Gen2 – Raw Data – IAM

- **Owner** – Since I am the owner of the account.
- **Storage Blob Data Reader** – Assigned to Azure Data Factory (since the files need to be read by Azure Data Factory for further processing).
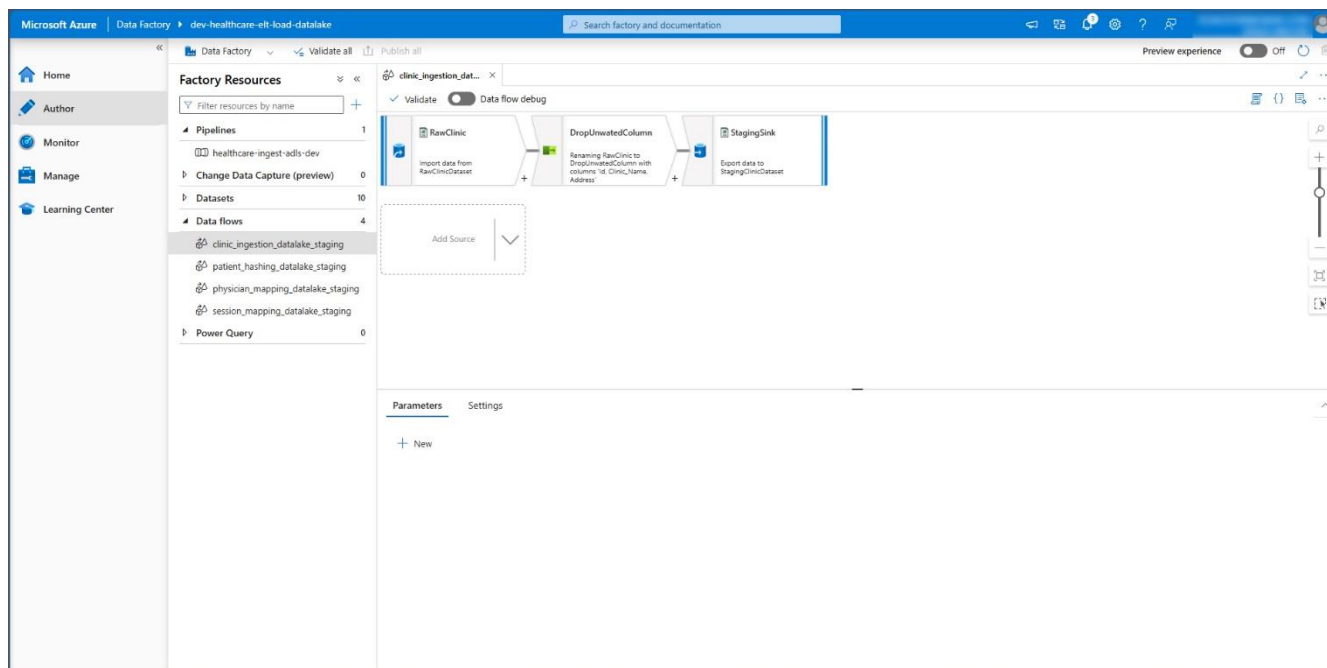
**Azure Data Factory**

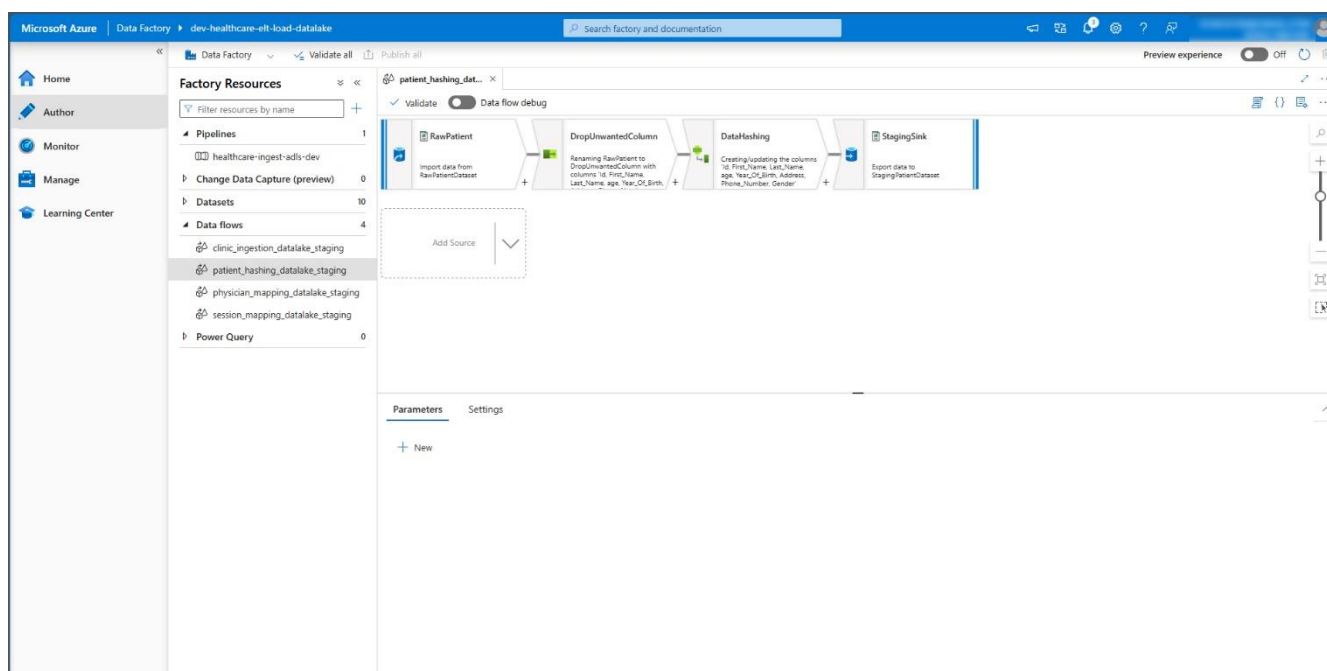Azure Data Factory is used only to perform data hashing, joining, and dropping unwanted columns.

**Data Flows**

**clinic_ingestion_datalake_staging**

Unwanted columns were removed from the clinic data and ingested into the staging layer.



**patient_hashing_datalake_staging**

The patients' data was hashed, and unwanted columns were removed.

## physician_mapping_datalake_staging

The physician data was joined with NPPES, names were filled in if missing, and unwanted columns were removed.



## session_mapping_datalake_staging

The scheduled sessions were joined with session status, and unwanted columns were dropped.

## Pipeline

All the data flow activities were added sequentially, and a notebook was included for the data cleaning process.



## Datasets

Source and sink datasets in Azure Data Factory.

**Linked Services**

- **AzureDatabricks1** – to connect with Azure Databricks.
- **AzureDataLakeStorage1RawDataLake** – to connect with the source (storage account containing raw data).
- **AzureDataLakeStorage2DevDataLake** – to connect with the sink (storage account which is ready to be loaded with staged and cleaned data).
- **AzureKeyVault1** – to securely access credentials (Databricks Access Token).



**Trigger – Daily Schedule**

A schedule-type trigger was added to run the pipeline on a daily basis.

# Alert – Pipeline Failure

Email notification if the pipeline failed to run.

**Azure Data Lake Gen2 – Staged and Cleaned Data**

**Staging Data Container**

Staged data loaded from Azure Data Factory can be accessed here.



**Cleaned Data Container**

Cleaned data loaded from Azure Databricks can be accessed here.

**Azure Data Lake Gen2 –Staged and Cleaned Data – IAM**

- **Owner** – Since I am the owner of the account.
- **Storage Blob Data Contributor** – Assigned to Azure Data Factory (since the files need to be written by Azure Data Factory).
- **Storage Blob Data Contributor** – Assigned to myself for Azure Databricks (As I enabled credential pass-through to read and write).
- **Storage Blob Data Reader** – Assigned to Azure Synapse Analytics (since the files need to be read by Azure Synapse Analytics).

## Azure Databricks

The notebooks can be accessed [here](#).

### For mounting

A separate notebook was created for the mounting purpose.



### For cleaning

Another notebook was created for the cleaning purpose.

## Azure Synapse Analytics

### Scripts

The scripts used in synapse analytics can be accessed [here](#).



### Views

The above scripts are used to create these views.

## Pipeline

All the scripts activities were added sequentially for reporting layer.



## Linked Services

- **AzureKeyVault1** – to securely access credentials (SQL password).
- **dev-healthcare-synapse-centralindia-WorkspaceDefaultSqlServer** – default linked service for Synapse Analytics SQL Server.
- **dev-healthcare-synapse-centralindia-WorkspaceDefaultStorage** – default linked service for Synapse Analytics Storage.
- **dev_healthcare_synapse_centralindia_serverlesssql** – to connect with the SQL database.

# Power BI

Power BI Desktop is used for visualizations.

The visualizations were checked by updating the raw session status. The visualizations are attached before and after updating the raw data.

## Before session status update
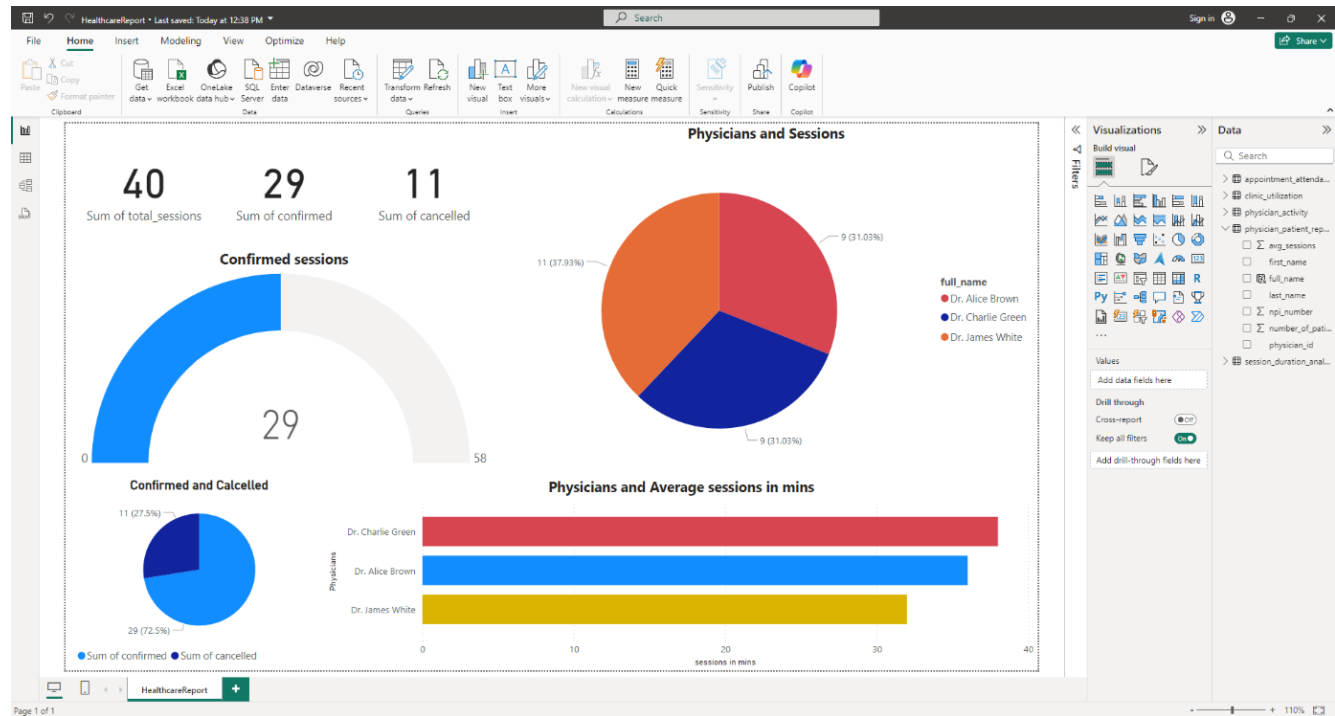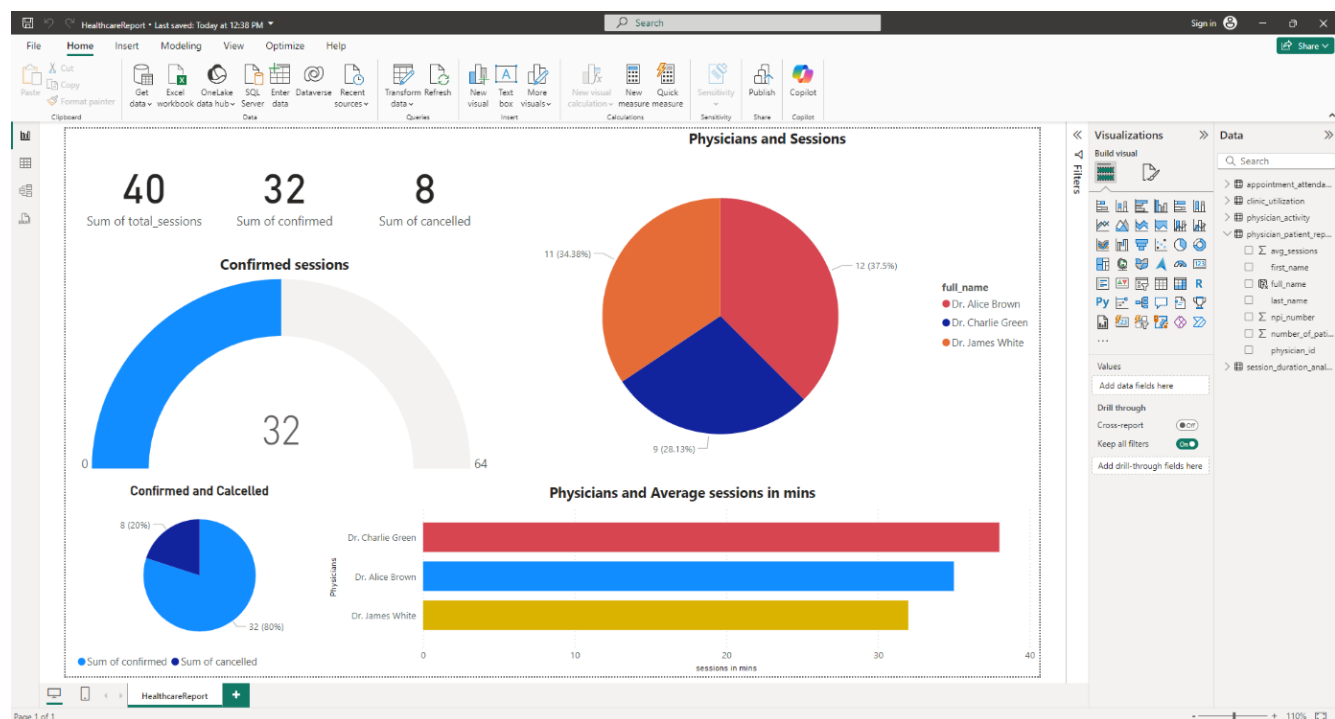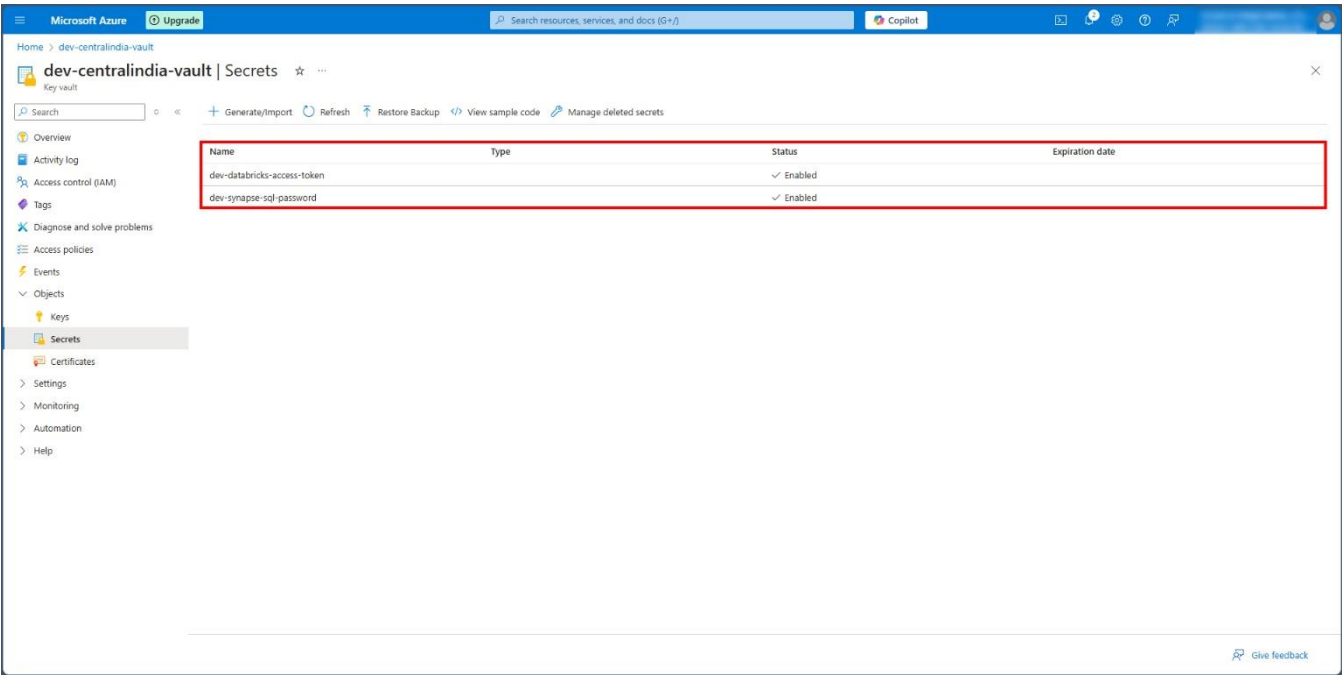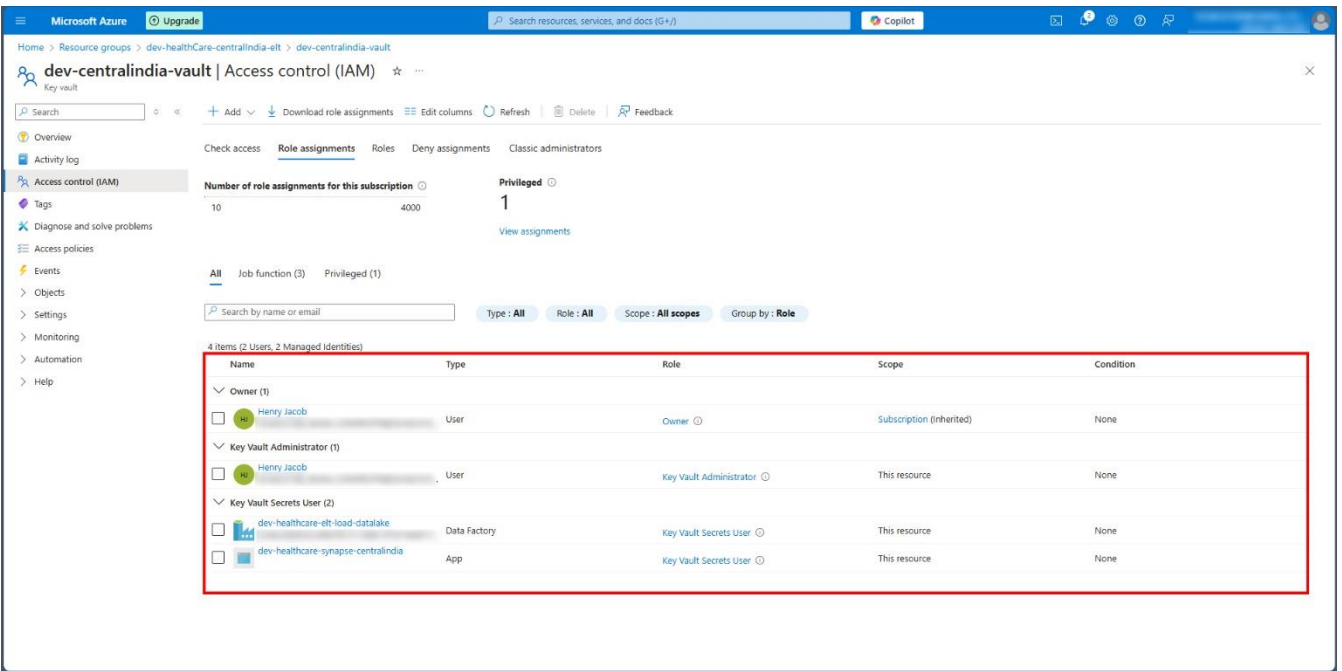


## After session status update

## Azure Key Vault

Credentials are stored inside the **Key Vault**.



## Key Vault – IAM

- **Owner** – Since I am the owner of the account.
- **Key Vault Administrator** – Assigned to myself (I need to be an administrator to assign permissions to others).
- **Key Vault Secret User** – Assigned to Azure Data Factory (since the credentials [Databricks Access Token] needs to be read by Azure Data Factory).
- **Key Vault Secret User** – Assigned to Azure Synapse Analytics (since the credentials [SQL Password] needs to be read by Azure Synapse Analytics).

*Don't forget to delete unused resources* after the project is completed

to avoid unnecessary charges.

---

It is **not recommended** to **share** your **storage account name**, **access keys**, or other sensitive information **in public**.

However, I am sharing the screenshot here to demonstrate the standard naming convention I maintained.

I deleted all my resources before publishing this document, and I do not plan to reuse or recover the resources.