# Fraud Detection Using Azure Data Engineering

## Project Overview

This project is designed to detect potentially fraudulent transactions based on the distance between a customer's registered city and the city where the transaction occurred. Fraud criticality is categorized as **Low, Medium, or High** based on the calculated distance. The data pipeline is built using **Azure Data Factory (ADF), Azure Data Lake Storage (ADLS), Azure Databricks, and Azure SQL Database**.

## Requirement

Based on the city in which the transaction has happened, the fraud criticality needs to be calculated based on the below conditions,

- **Low:** Distance < 500 km
- **Medium:** Distance 5000-1000 km
- **High:** Distance > 1000 km

The following reports need to be developed based on the Fraud criticality,

- Bar chart displaying Fraud Alert Counts for past 12 months.
- Pie chart showing Fraud criticality trend
- Top 10 Customers and Cities of High Fraud Criticality
- Pie chart showing Age Group criticality trend

## Tech Stack

- **Azure Data Factory (ADF)** – Data ingestion and orchestration.

- **Azure Data Lake Storage (ADLS)** – Raw and landing data storage.

- **Azure Databricks** – Data transformation and fraud analysis.

- **Azure SQL Database** – Storing transformed data for reporting.

- **Power BI** – Visualization and reporting.

## Source Datasets Used

1. **Customer Metadata** – Contains customer details, including city.

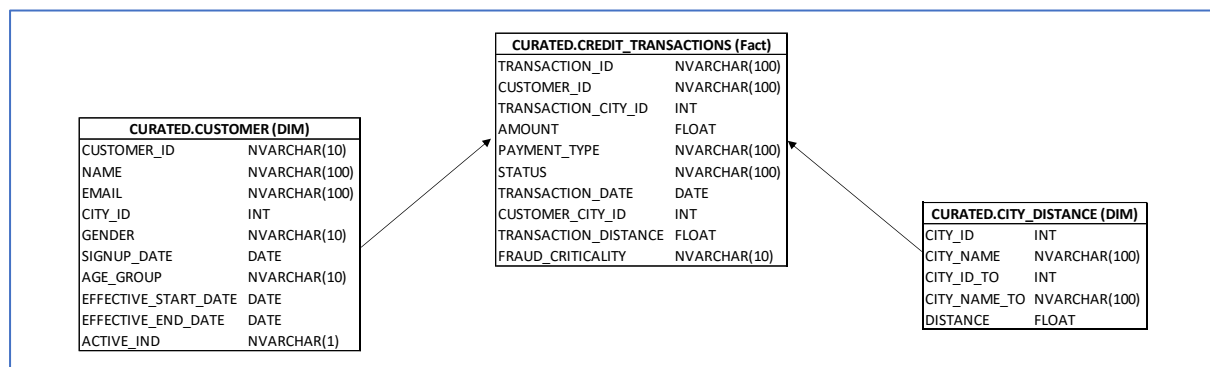| Customer | Name | Email | CityID | BirthYear | Gender | SignupDate |
|---|---|---|---|---|---|---|
| C001 | Umang Chokshi | inaaya-36@tak.com | 8 | 1998 | Female | 31-07-2021 |
| C002 | Kartik Dada | keyachaudhuri@hotmail.com | 10 | 1976 | Female | 06-04-2024 |
| C003 | Tushar Kar | pari98@yahoo.com | 9 | 1993 | Female | 24-07-2021 |

2. **City Distance Metadata** – Contains distance between every two cities in India.

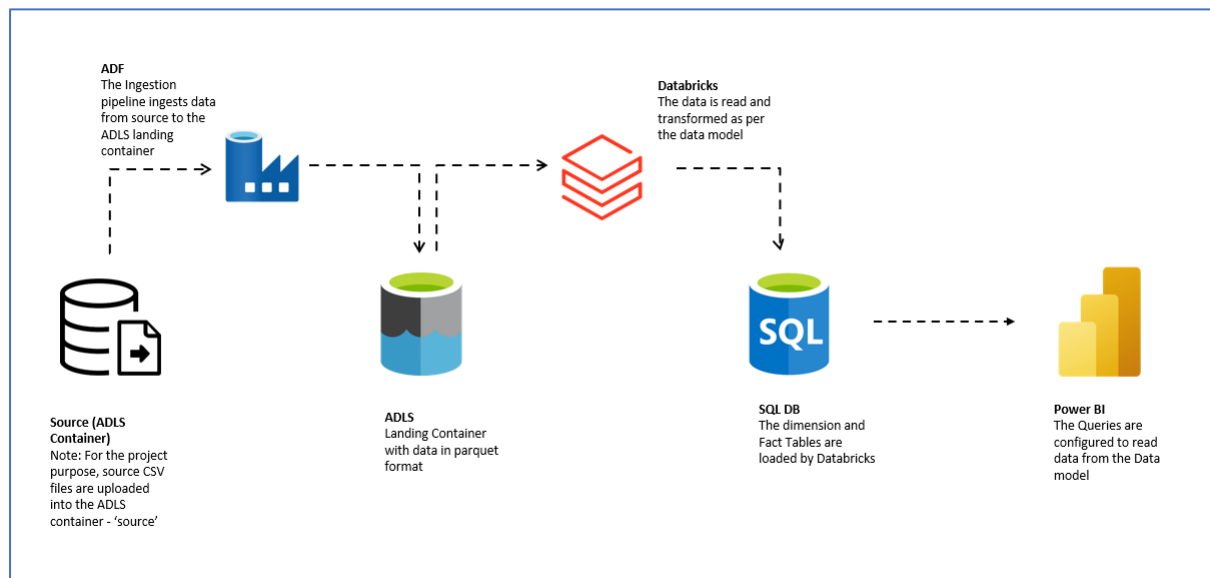| CityID_1 | CityName_1 | CityID_2 | CityName_2 | Distance_km |
|---|---|---|---|---|
| 1 | Mumbai | 1 | Mumbai | 0 |
| 1 | Mumbai | 2 | Delhi | 852 |
| 1 | Mumbai | 3 | Bangalore | 1940 |

3. **Transaction Data** – Contains transaction details, including transaction city.

| TransactionID | Customer | CityID | Date | Amount | PaymentType | Status |
|---------------|----------|--------|------|--------|-------------|--------|
| T0001 | C071 | 6 | 12-02-2025 | 8792.32 | Debit Card | Success |
| T0002 | C029 | 5 | 11-02-2025 | 40834.75 | UPI | Pending |
| T0003 | C007 | 7 | 08-02-2025 | 37359.97 | UPI | Failed |
| T0004 | C014 | 2 | 08-02-2025 | 17375.03 | Net Banking | Success |

## Data Model

**CURATED.CREDIT_TRANSACTIONS (Fact)**

| | |
|---|---|
| TRANSACTION_ID | NVARCHAR(100) |
| CUSTOMER_ID | NVARCHAR(100) |
| TRANSACTION_CITY_ID | INT |
| AMOUNT | FLOAT |
| PAYMENT_TYPE | NVARCHAR(100) |
| STATUS | NVARCHAR(100) |
| TRANSACTION_DATE | DATE |
| CUSTOMER_CITY_ID | INT |
| TRANSACTION_DISTANCE | FLOAT |
| FRAUD_CRITICALITY | NVARCHAR(10) |

**CURATED.CUSTOMER (DIM)**

| | |
|---|---|
| CUSTOMER_ID | NVARCHAR(10) |
| NAME | NVARCHAR(100) |
| EMAIL | NVARCHAR(100) |
| CITY_ID | INT |
| GENDER | NVARCHAR(10) |
| SIGNUP_DATE | DATE |
| AGE_GROUP | NVARCHAR(10) |
| EFFECTIVE_START_DATE | DATE |
| EFFECTIVE_END_DATE | DATE |
| ACTIVE_IND | NVARCHAR(1) |

**CURATED.CITY_DISTANCE (DIM)**

| | |
|---|---|
| CITY_ID | INT |
| CITY_NAME | NVARCHAR(100) |
| CITY_ID_TO | INT |
| CITY_NAME_TO | NVARCHAR(100) |
| DISTANCE | FLOAT |

## Architecture



**ADF**
The Ingestion pipeline ingests data from source to the ADLS landing container

**Databricks**
The data is read and transformed as per the data model

**Source (ADLS Container)**
Note: For the project purpose, source CSV files are uploaded into the ADLS container - 'source'

**ADLS**
Landing Container with data in parquet format

**SQL DB**
The dimension and Fact Tables are loaded by Databricks

**Power BI**
The Queries are configured to read data from the Data model

# Detailed Implementation:

**Azure Data Lake Storage:**

Two containers are created,

Source – to store the raw csv source files.

Landing – to store the parquet file ingested by ADF

NOTE: This action is implemented to mimic the data ingestion from source to ADLS



The transaction data files are stored in YYYY/MM/DD date folder format.

**Azure Data Factory:**

The ADF contains a main pipeline which starts with setting the load date (current date -1 day). It then executes two ingestion pipeline (Master data & Transaction data) and one curation pipeline.

PL_FRAUD_TRANSACTION_REPORT (main pipeline):



PL_MASTER_INGESTION: Ingests Customer and Cities data into landing container on full load basis

PL_CREDIT_TRANSACTION_INGESTION: For ingesting the transaction data on an incremental basis



Incremental logic: The Load date is passed as parameter from the main pipeline which is used in the wild card path

PL_CURATION: The curation pipeline contains two notebook activity. One for executing the master curation and another on for executing the transaction curation



All pipelines are implemented with logging logic in which the stored procedure is called. This stored procedure INSERTs and UPDATEs the necessary logs at the appropriate steps.

Logging Table:



## Databricks:

### Connection and Helper Function:

- This notebook contains the necessary properties declared for establishing connections from Databricks notebook to ADLS and SQL DB.
- This notebook also contains the helped function that will accept the query string. This query string is executed in the SQL DB through ODBC Driver 17 for SQL Server.
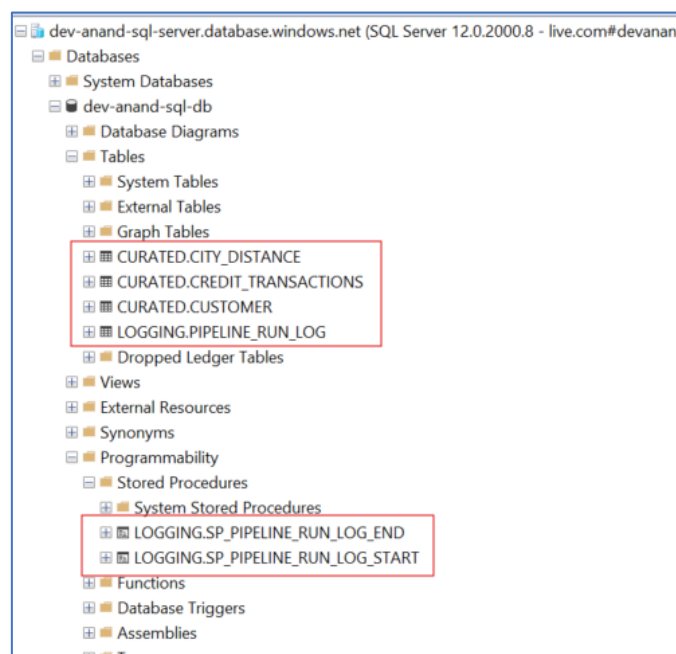
### Master Table Curation

- The customers and cities files from landing container are read, processed and loaded to the CURATED.CUSTOMER and CURATED.CITY_DISTANCE tables.
- The cities data are overwritten every time the notebook is executed.
- From the customer data, the age group of each customer is calculated. Every time when the notebook is executed, the data is compared against existing data in the table and instead of rewriting the updated data, the old data is marked as Inactive and the new updated data is inserted as Active. This helps in maintaining the customer history (SCD2) for future analysis.

**Transaction Curation:**

- This notebook reads the incremental transactions on a (current date – 1 day) basis.
- Using the customer and city curated data, the distance between the customer city and the city where the transaction has happened is calculated. Based on this distance, the criticality for fraud is calculated.
    - **Low:** Distance < 500 km
    - **Medium:** Distance 5000-1000 km
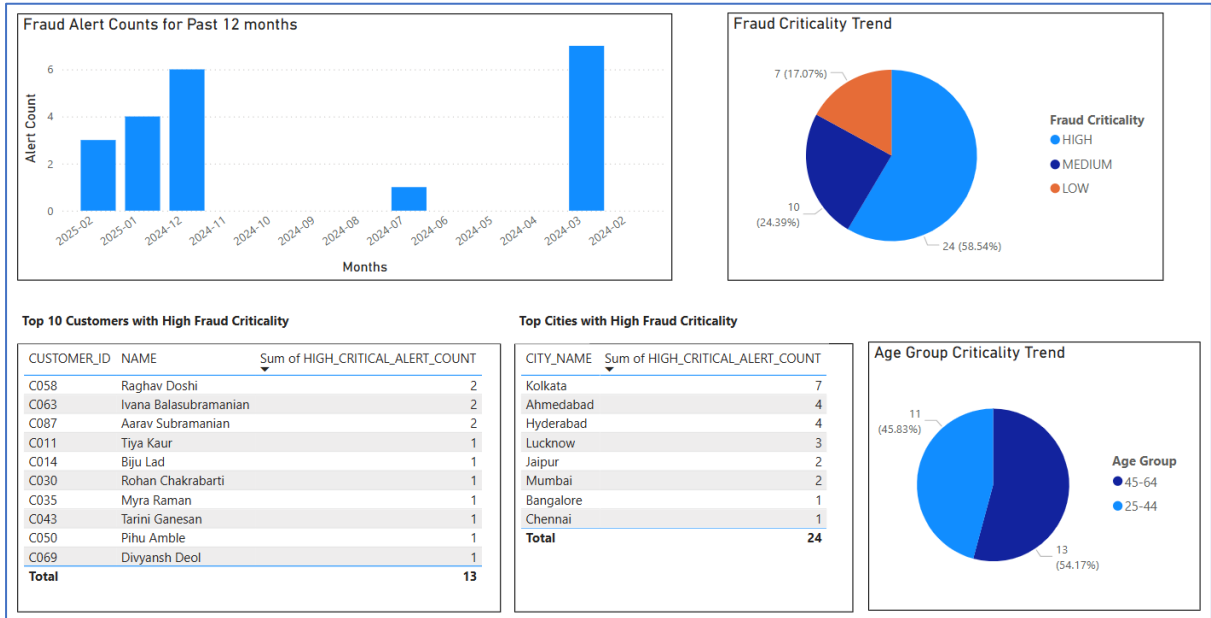    - **High:** Distance > 1000 km

**SQL Database:**

- CURATED.CREDIT_TRANSACTIONS – This is the fact tables loaded with transaction containing the calculated fraud criticality column. It gets loaded on an incremental basis.
- CURATED.CUSTOMER – Contains customer master data. Follows SCD2.
- CURATED.CITY_DISTANCE – Contains city master data and is loaded on full load basis.
- LOGGING.SP_PIPELINE_RUN_LOG_START, LOGGING.SP_PIPELINE_RUN_LOG_END – Stored Procedures used to INSERT and UPDATE logs respectively
- LOGGING.PIPELINE_RUN_LOG – Contains pipeline execution logs

**Power BI:**

Finally, in the Power BI, 5 different queries are created to generate reports as per the requirement.



# Security:

### ADF <-> ADLS:

The Connection for the ADF to read and write files from and to ADLS is performed using managed identity with below access roles.
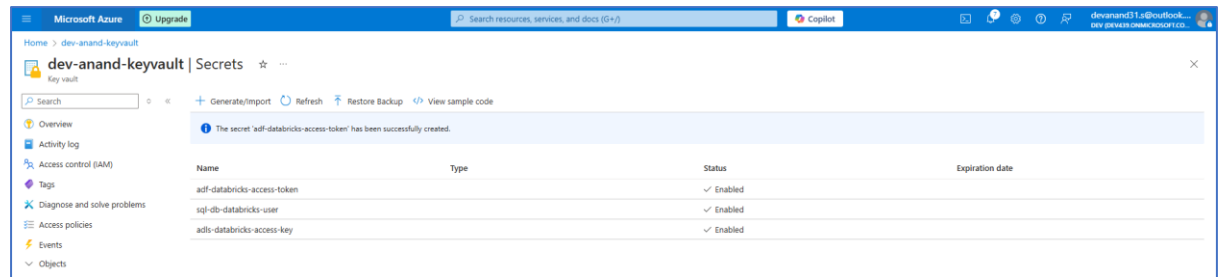


### ADF <-> SQL DB:

The Connection for the ADF to read and write files from and to SQL DB is performed using managed identity with below access roles.



### Databricks with ADF, ADLS and SQL DB:

Since this is a personal account, only the Standard edition of Databricks workspace is created. This version of Databricks does not allow Managed identity. So, all the Databricks related access are established used access tokens/access keys via Key Vault.

**Power BI <-> SQL DB**

A SQL user – powerbi_user is created for the Power BI to access SQL DB.

```sql
USE MASTER
GO
CREATE LOGIN powerbi_user WITH PASSWORD = '<<encrypted>>'
GO
CREATE USER powerbi_user FOR LOGIN powerbi_user
GO
ALTER ROLE db_datareader ADD MEMBER powerbi_user
GO
```