# End-to-End Sales Data Analytics Platform on Azure
# Trainee Name: Swathi Baskaran

## Table of Contents

## Project Statement

The goal of this project is to design and implement an end-to-end cloud-based data analytics platform that ingests raw sales, customer, and product data from CSV files, processes and transforms it into a structured data warehouse using Azure Data Factory and Azure Synapse Analytics, orchestrates the workflows and deployments with Azure DevOps, and delivers interactive business intelligence dashboards through Power BI. This solution provides real-time and historical insights into sales performance, customer behavior, and product trends, enabling data-driven decision-making.

## Project Overview

This project demonstrates the implementation of a modern data engineering and analytics solution on Microsoft Azure.

- **Data Sources:** Sales, Products, and Customers data in CSV format.
- **Data Ingestion:**
  - Azure Data Factory (ADF) pipelines were developed to extract raw data from the Azure Storage Account.
  - Incremental loads and scheduled triggers ensure automated data refresh.
- **Data Transformation and Storage:**
  - Data is cleaned, standardized, and transformed using Azure Synapse Analytics.
  - Fact and dimension tables (FactSales, DimCustomer, DimProduct) form a star schema optimized for analytics.
  - SQL scripts and stored procedures manage transformations.
- **Orchestration & CI/CD:**
  - Azure DevOps manages source control, versioning, and deployment of ADF pipelines, Synapse scripts, and ARM templates.
  - CI/CD pipelines (YAML) ensure seamless deployment across environments.

- **Visualization & Reporting**:
  - Power BI dashboards provide interactive analytics with KPIs such as total sales, sales by product, customer segmentation, and time-based sales trends.
  - DAX measures ensure accurate aggregation and business logic representation.
- **Business Value:**
  - Enables management to analyze customer purchasing patterns, product performance, and revenue trends.
  - Provides a scalable, automated, and reusable cloud data analytics solution.

## Prerequisities

1. **Azure Subscription** with permission to create and manage resources.
2. Basic knowledge of **SQL, ETL concepts**, and **data warehousing.**
3. Familiarity with **Azure Portal, Azure Storage** and **Power BI.**
4. **Power BI** installed and configured for building and publishing reports.
5. Access to **Azure DevOps** for source control and deployment.
6. Source CSV files available:
   - sales.csv (transactional sales data)
   - products.csv (product catalog)
   - customers.csv (customer details)

## Azure Resources Used for this Project

- **Azure Storage Account** – Stores raw CSV data files
- **Azure Data Factory (ADF)** – Ingests and transforms raw data through pipelines.
- **Azure Synapse Analytics (Synapse Studio)** – Acts as the data warehouse, hosts star schema (FactSales, DimCustomer, DimProduct)
- **Azure DevOps** – Provides version control, CI/CD pipelines, and deployment automation.
- **Power BI** – Used for building and publishing dashboards for reporting.

## Project Objectives

- Ingest CSV files from Azure Storage into a cloud-based system.
- Automate ETL processes using Azure Data Factory pipelines.
- Build a star schema data warehouse in Synapse (FactSales, DimCustomer, DimProduct).
- Implement CI/CD using Azure DevOps for managing ADF pipelines, Synapse scripts, and templates.
- Develop interactive dashboards in Power BI for sales, customer, and product analytics.
- Enable automation and scalability for handling larger datasets and scheduled refreshes.

## Tools Used

- **Azure Data Factory (ADF)**
  - Cloud-based ETL tool
  - Extracts data from Azure Storage
  - Transforms and loads into Synapse Analytics
- **Azure Synapse Analytics**
  - Cloud data warehouse
  - Stores transformed data in star schema (FactSales, DimCustomer, DimProduct)
  - Used for running SQL queries and managing tables
- **Azure Storage Account**
  - Landing zone for raw CSV files (sales.csv, products.csv, customers.csv)
  - Source for ADF pipelines
- **Azure DevOps**
  - Source control for JSON pipeline definitions, SQL scripts, and YAML files
  - CI/CD pipeline automation for deployments
- **Power BI**
  - Used for data modeling, visualization, and creating interactive dashboards from Synapse data.

## Execution Overview

### 1. Data Ingestion

- Raw CSV files (sales.csv, products.csv, customers.csv) are uploaded to the Azure Storage Account.
- Azure Data Factory (ADF) pipelines are triggered to pick up these files.

### 2. Data Transformation & Loading

- ADF Pipelines transform the raw data (Data type conversions, formatting, cleaning).
- Transformed data is loaded into Azure Synapse Analytics tables.
- The star schema is created with:
  - FactSales (Transactional Sales Data)
  - DimCustomer (Customer Details)
  - DimProduct (Product Details)

### 3. Data Storage & Querying

- Synapse Analytics acts as the central data warehouse
- SQL scripts define and manage schema, relationships, and business rules.

### 4. Source Control & Deployment

- All ADF pipeline definitions, Synapse SQL scripts, and configurations are stored in Azure DevOps.
- YAML Pipelines in DevOps handle CI/CD automation, ensuring consistent deployments across environments.

### 5. Data Visualization

- Power BI connects directly to Synapse Analytics.
- Data models are built linking FactSales with DimCustomer and DimProduct.
- Interactive dashboards are created to show:
  - Sales by product
  - Sales by customer
  - Overall sales performance trends

**6. End-to-End Flow**

- New CSV files land in Storage → ADF ingests and loads → Synapse stores in warehouse → Power BI visualizes → DevOps manages version control & automation.

## Implementation – Tasks Performed

1. **Data Preparation**
   - Collected and organized source data files: sales.csv, products.csv, and customers.csv.
   - Uploaded raw data files into the Azure Storage Account (Blob container).

2. **Pipeline Development in ADF**
   - Created Linked Services to connect ADF with Azure Storage and Synapse Analytics.
   - Designed Datasets representing input CSV files and output tables.
   - Built Data Pipelines to:
     o Extract CSV data from storage
     o Transform the data into structured format
     o Load the transformed data into Synapse Analytics

3. **Data Warehouse Setup in Synapse Analytics**
   - Designed and created a star schema with:
     o FactSales table (Transactional sales data)
     o DimCustomer table (customer attributes)
     o DimProduct table (product details)
   - Executed SQL scripts for table creation, constraints, and data loading.

4. **Integration with Azure DevOps**
   - Configured Git repository in Azure DevOps to store JSON pipeline definitions, SQL scripts, and ARM templates.
   - Implemented CI/CD pipelines using YAML for automated deployments.

5. **Data Visualization in Power BI**

- Connected Power BI to Synapse Analytics as the data source.
- Built relationships between FactSales, DimCustomer, and DimProduct.
- Designed dashboards and reports showing key insights such as:
  - Total sales by Product
  - Total sales by Customer
  - Sales Performance trends

6. **Testing and Validation**
- Verified correctness of ETL pipelines by comparing source and target data.
- Validated star schema relationships and Power BI reports for accuracy.

7. **Final Deployment**
- Published final dashboards in Power BI for reporting.
- Ensured automation of pipeline runs for regular data refresh.

## Practical Implementation on Azure Portal

## Step 1: Create Azure Storage Account

- Provision a Storage Account for raw and processed data

- Create a container called **'datalake'** in the storage account for the raw CSV files in the Azure Storage Account



- Upload CSV files into the source container

## Step 2: Create an Azure Synapse Analytics Workspace

- Creating an Azure Synapse Studio



- Creating a SQL Pool in Azure Synapse Studio

- Monitoring the SQL pool created



# Step 3: Create Azure Data Factory Studio

- Creating datasets within Azure Data Factory (ADF)



- Creating linked service for Azure Blob Storage

- Create linked service for Azure Synapse Analytics



- Connecting Azure Blob Storage Linked Service to the raw dataset



- Connecting Azure Synapse Analytics Linked Service to the data warehouse dataset

## Step 4: Creating SQL table in Synapse Studio

- Connecting to the SQL Database



- Creating staging tables in Synapse Studio



- Creating ETL Logs in Synapse Studio

- Creating Fact table for Sales



# Step 5: Creation of Pipeline

- Creating a pipeline to copy data from raw files into staging tables in Azure Data Factory (ADF)

- Setting the source for the pipeline



- Setting the sink for the pipeline



- Executing the Pipeline

- Successful execution of the pipeline



- Monitoring the Pipeline



## Step 6: Executing SQL Queries

- Upsert Dimensions

- Querying the Facts table



- Creating Report Views



- Business Queries

## Step 7: Implementing Source Control and CI/CD with Azure DevOps

- Git Integration with Azure Data Factory (ADF)



- Storing SQL Scripts and ARM Templates



- CI/CD Pipeline setup

**Step 8: Verify Data and Visualize**

- Open Power BI > Get Data from SQL Server Database



- Connect to the SQL Server and Database and select the respective tables



- Perform some analysis and create a report

**Successful Output Generated**

**ADF Pipeline Execution**

- Pipeline executed successfully, extracting data from Azure Storage (CSV files).
- Data was transformed and loaded into Azure Synapse Analytics tables.
- ADF run history shows all activities completed without failure.

## Query Validation in Azure Synapse Analytics

- SQL queries executed successfully on FactSales, DimCustomer, and DimProduct.
- Verified data correctness and relationships between dimension and fact tables.
- Confirmed the data model supports analytical use cases.

**Screenshot 1 — Creating_Fact_Tables (SQL script 2)**

```
50   GO
51
52   SELECT TOP 10 * FROM dw.FactSales;
53
54   -- To verify relationships (Star Schema)
55   SELECT TOP 10
56       f.SaleID,
57       f.SaleDate,
58       dc.CustomerName,
59       dp.ProductName,
60       f.Quantity,
61       f.Amount
62   FROM dw.FactSales f
63   JOIN dw.DimCustomer dc ON f.CustomerSK = dc.CustomerSK
64   JOIN dw.DimProduct dp ON f.ProductSK = dp.ProductSK;
```

Results

| SaleID | SaleDate | CustomerName | ProductName | Quantity | Amount |
|--------|----------|--------------|-------------|----------|--------|
| 1 | 2024-03-14T00:... | Customer 020 | Product 021 | 2 | 386.32 |
| 52 | 2024-09-28T00:... | Customer 015 | Product 010 | 4 | 211.48 |
| 39 | 2025-08-07T00:... | Customer 019 | Product 016 | 2 | 886.60 |

00:00:04 Query executed successfully.

Properties — General — Name *: Creating_Fact_Tables — Type: .sql script — Size: 0 bytes

---

**Screenshot 2 — Create_Report_View**

```
1    -- Drop the view if it already exists
2    IF OBJECT_ID('dw.vw_SalesAnalysis', 'V') IS NOT NULL
3        DROP VIEW dw.vw_SalesAnalysis;
4    GO
5
6    -- Now create the view
7    CREATE VIEW dw.vw_SalesAnalysis
8    AS
9    SELECT
10       f.SaleID,
11       f.SaleDate,
12       dc.CustomerName,
13       dc.City,
14       dc.State,
15       dp.ProductName,
```

Results

| CustomerName | TotalSpent |
|--------------|------------|
| Customer 009 | 449514.60 |
| Customer 002 | 432028.20 |
| Customer 008 | 393765.00 |

00:00:02 Query executed successfully.

Properties — General — Name *: Create_Report_View — Type: .sql script — Size: 0 bytes

---

**Screenshot 3 — BusinessQueries**

Other users in your workspace may have access to modify this item. Do not use this item unless you trust all users who may have access to the workspace.

```
1    -- Build Some Business Queries
2    -- Total Sales by Product Category
3    SELECT dp.Category, SUM(f.Amount) AS TotalSales
4    FROM dw.FactSales f
5    JOIN dw.DimProduct dp ON f.ProductSK = dp.ProductSK
6    GROUP BY dp.Category
7    ORDER BY TotalSales DESC;
```
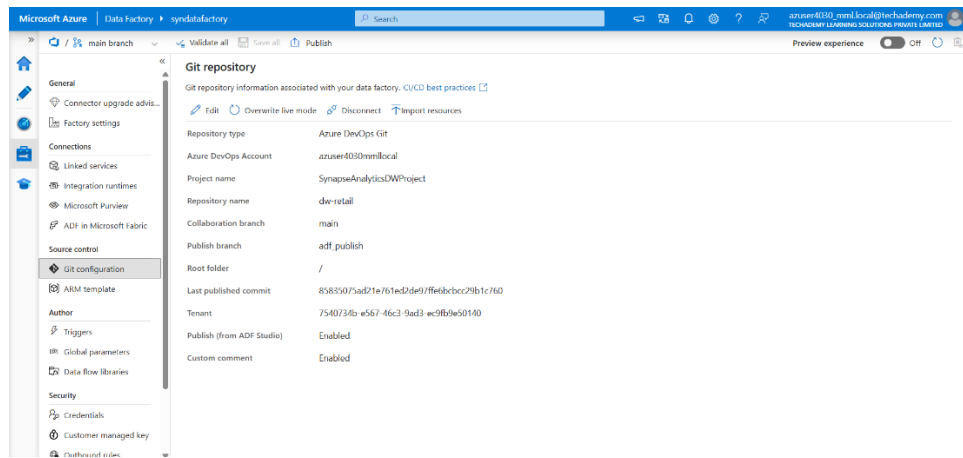
Results

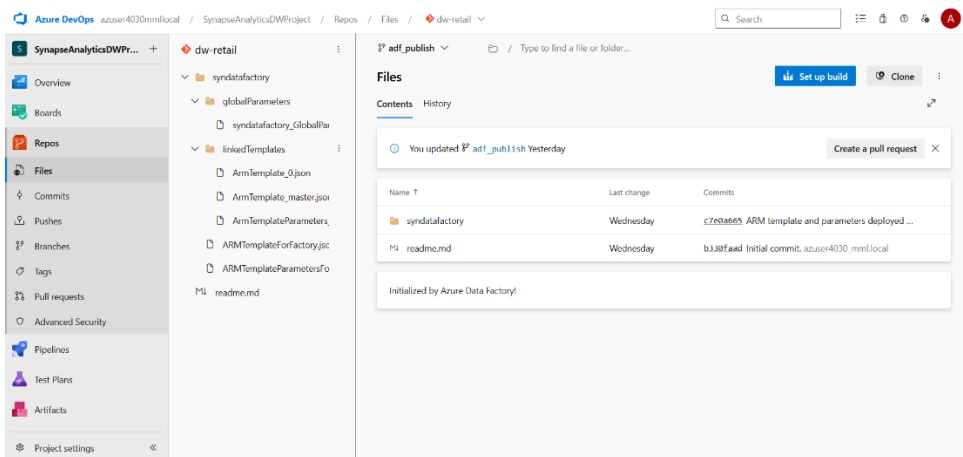| Category | TotalSales |
|----------|------------|
| Electronics | 1400867.40 |
| Grocery | 1187857.20 |
| Home | 882835.80 |
| Clothing | 507745.20 |
| Sports | 208106.40 |

00:00:03 Query executed successfully.

**Azure DevOps Pipeline Execution**

- CI/CD pipeline deployed ADF ARM templates and Synapse SQL scripts from the repo.
- Build stage executed successfully with no errors





**Power BI Visualization**

- Power BI connected to Synapse and consumed validated data
- Reports built showing "Sum of Quantity by Year, Quarter, Month and Day", "Count of CustomerID by CustomerName and City", "Sum of Amount by ProductName", "Sum of UnitPrice by ProductName and Category"

- Dashboards confirmed end-to-end functionality of the project.



## Strategies for Optimizing the Process

1. **Perform Data Cleaning and Standardization** – Handle nulls, duplicates, inconsistent formats, and ensure schema alignment during ingestion.
2. **Adopt Incremental Data Loads** – Instead of full refreshes, ingest only new or changed data in ADF to save time and resources.
3. **Leverage Staging Layers** – Use Azure Storage as a staging area before transformations in Synapse for smoother workflows.
4. **Design Efficient Data Models** – Implement star schema with Fact and Dimension tables for optimal query and reporting performance.
5. **Use Partitioning and Indexing** – Partition large fact tables and add indexes in Synapse to accelerate query execution.
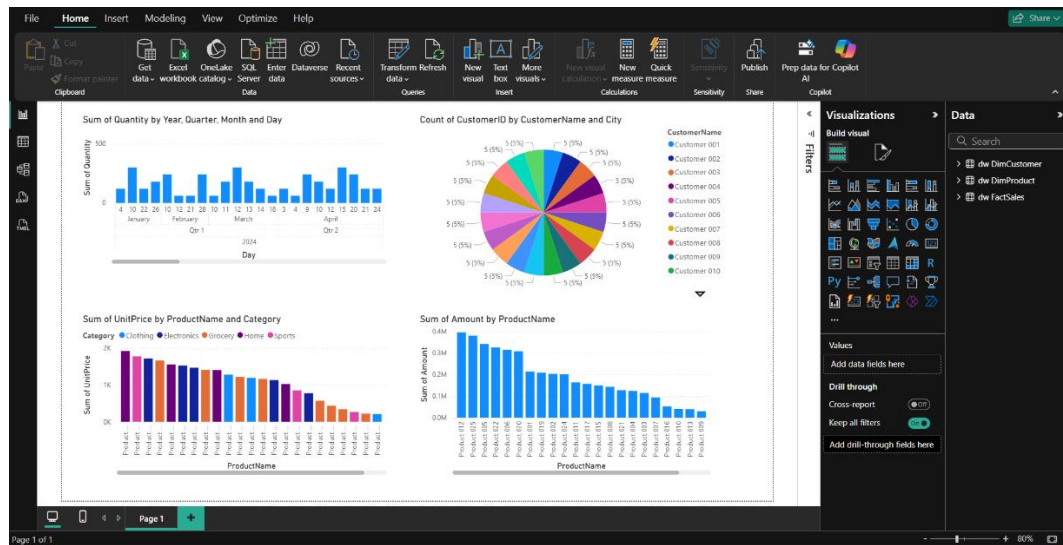6. **Automate Deployments with CI/CD** – Utilize Azure DevOps pipelines for automated deployments of ADF, Synapse, and Power BI assets.
7. **Optimize Power BI Models** – Import only necessary fields, use measures instead of calculated columns, and apply row-level security.
8. **Monitor and Scale Resources** – Use Azure Monitor to track performance, and scale Synapse or ADF integration runtimes dynamically based on workload.

9. **Apply Cost Management Practices** – Move cold data to cheaper storage tiers and pause Synapse compute when not in use.
10. **Maintain Governance and Documentation** – Enforce naming conventions, document pipelines, and review processes periodically for continuous optimization.

## Conclusion

This project successfully implemented a comprehensive and scalable sales data processing pipeline using Azure Data Factory (ADF), Azure Synapse Analytics, Azure Storage, and Power BI, converting raw CSV files into structured fact and dimension tables to support advanced reporting and analytics.

### Successful Implementation of the Data Pipeline

- Azure Data Factory orchestrated the end-to-end workflow, from ingesting raw CSV files (customers, products, sales) from Azure Storage to loading them into Synapse Analytics.
- The pipeline supported both scheduled and on-demand executions, enabling flexible and reliable data movement.

### Efficient Data Cleaning, Transformation, and Modeling

- Raw CSV data was validated and cleaned to remove duplicates, handle nulls, and enforce consistency.
- Transformations such as joins, aggregations, and schema alignment were applied before loading data into Synapse.
- A star schema was designed with FactSales and dimension tables (DimCustomer, DimProduct), enabling optimized querying and analysis.

### Exploration, Optimization, and Querying in Synapse Analytics

- Synapse enabled execution of analytical SQL queries across large datasets with improved performance.

- Optimization techniques such as partitioning, indexing, and efficient table design were applied.
- The structured model supported quick retrieval of insights like total sales by product, region, or customer segment.

**Integration and Reporting**

- Azure DevOps pipelines ensured CI/CD, with version-controlled ADF ARM templates, Synapse SQL scripts, and deployment automation.
- Power BI dashboards were developed, connecting directly to Synapse tables, to provide interactive reports and visualizations.
- Business insights such as customer purchasing patterns, top-selling products, and revenue trends were delivered through intuitive visualizations.

As a whole, this project demonstrates a scalable, end-to-end data engineering and analytics solution for sales data, leveraging the synergy between ADF, Synapse, Azure Storage, DevOps, and Power BI. The pipeline ensures reliable ingestion, efficient transformation, governed storage, automated deployments, and business intelligence dashboards enabling data-driven decision-making with accuracy and efficiency.