# CRM Sales Analysis – Azure End-to-End Data Engineering Project

## 1. Project Overview

This project implements a complete **end-to-end Azure data engineering and analytics pipeline** using a CRM Sales dataset. The objective was to ingest raw CRM data, clean and transform it using Apache Spark, store it in a structured data lake, expose it through Azure Synapse Analytics, and finally build interactive dashboards in **Power BI** for business analysis and reporting.

The project closely follows real-world data engineering best practices, including layered data lake design, data quality checks, SQL-based analytics, and BI consumption.

Project link: https://github.com/KasturiDisale/CRM-dataset-azure-data-engineering-project

---

## 2. Dataset Description

The CRM dataset consists of multiple CSV files representing different business entities: - **Accounts** – company-level information such as sector, revenue, employees, and location - **Products** – product catalog with pricing and series information - **Sales Transactions** – opportunities, deal stages, sales values, and dates - **Sales Team** – sales agents, managers, and regional offices

These datasets together enable comprehensive analysis of sales performance, customer contribution, and product revenue.

---

## 3. High-Level Architecture

```
Source CSV Files
     ↓
Azure Data Factory (Ingestion)
     ↓
Azure Data Lake Storage Gen2 (Raw / Transformed Layers)
     ↓
Azure Databricks (Spark Data Processing)
     ↓
Azure Data Lake Storage Gen2 (Curated Data)
     ↓
Azure Synapse Analytics (Serverless SQL + Lake Database)
```
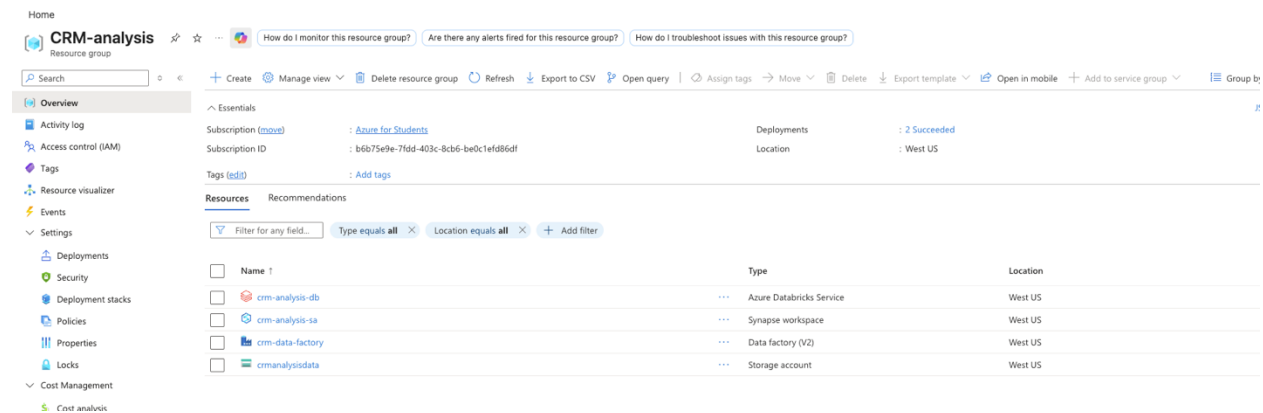
```
           ↓
Power BI (Semantic Model & Dashboards)
```



## 4. Azure Resource Setup

The following Azure services were created within a single resource group: - **Azure Data Factory** – for data ingestion pipelines - **Azure Data Lake Storage Gen2** – central data lake for raw and transformed data - **Azure Databricks** – Apache Spark environment for data processing - **Azure Synapse Analytics** – serverless SQL analytics and lake database - **Power BI** – reporting and visualization layer.

## 5. Data Ingestion using Azure Data Factory

### Steps Performed

1. Created **Linked Services** for:
   o Source CRM CSV files
   o Azure Data Lake Storage Gen2
2. Built an ADF pipeline with multiple **Copy Data** activities
3. Ingested each dataset (accounts, products, sales transactions, sales team) into ADLS Gen2
4. Stored ingested files in a **raw-data** folder structure

### Outcome

- All CRM source files were successfully ingested into the data lake
- Pipelines executed successfully and can be reused or scheduled

## Create Data Factory ...

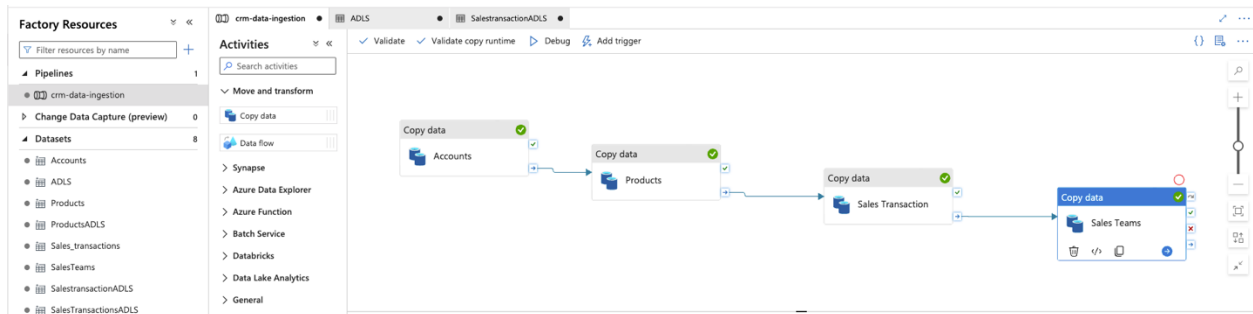Basics | Git configuration | Networking | Advanced | Tags | **Review + create**

👁 View automation template

**Basics**

| | |
|---|---|
| Subscription | Azure for Students |
| Resource group | CRM-analysis |
| Name | crm-data-factory |
| Region | West US |
| Version | V2 |

**Networking**

| | |
|---|---|
| Connect via | Public endpoint |



# 6. Data Lake Storage Design (ADLS Gen2)

The data lake was organized using a layered approach:

```
/raw-data
    ├── accounts
    ├── products
    ├── salestransaction
    └── salesteam

/transformed-data
    ├── accounts
```

```
├── products
├── salestransaction
└── salesteam
```

This design supports data traceability, reprocessing, and separation of raw and curated datasets.





# 7. Data Processing using Azure Databricks (Apache Spark)

Azure Databricks was used to preprocess and clean the CRM data before analytics.

## Key Processing Steps

1. Read raw CSV files from ADLS Gen2 using `abfss://` paths
2. Enabled schema inference and validated column data types
3. Performed **data quality checks**, including:
   - Column-wise null value analysis
   - Row count validation
4. Applied data cleaning logic:
   - Dropped records with null values in critical fields (e.g., account, close_value)
   - Retained or standardized optional fields where appropriate
5. Revalidated datasets to ensure no remaining nulls in required columns

6. Wrote cleaned and transformed data back to ADLS Gen2 in the **transformed-data** layer

## Outcome

- High-quality, analytics-ready datasets were produced
- Data consistency and integrity were ensured before downstream consumption

# 8. Azure Synapse Analytics – Lake Database & SQL

## Synapse Setup

- Created a **Synapse Workspace** with the built-in **Serverless SQL Pool**
- Created a **Lake Database** named `CRManalysis_DB`

## Tables Created

External tables were created over the transformed data for: - accounts - products - salestransaction - salesteam

These tables allow SQL-based querying directly on data lake files.

## Analytics Performed

Using Synapse SQL, multiple analytical queries were executed, including: - Top 10 accounts by total revenue - Revenue by product and product series - Sales performance by sales agent - Deal counts and average deal sizes

## Outcome

- Business-ready tables and queries exposed via SQL
- Serverless analytics enabled without infrastructure management

Microsoft Azure    Synapse Analytics ▸ crm-analysis-sa                    🔍 Search

ⓘ We use optional cookies to provide a better experience. Learn more ⧉

⟳ Synapse live ∨    ✓ Validate all    ⬆ Publish all ②

| CRMAnalysis_DB ● | SQL script 1 ● |

▷ Run  ↶ Undo ∨   ⬆ Publish  🔗 Query plan    Connect to  ✓ Built-in ∨   Use database  CRMAnalysis_DB ∨  ⟳

**Data**  + ∨ «

Workspace | Linked

▽ Filter resources by name

◢ Lake database                    1
  ▸ ◯ ▸ ⊟ CRManalysis_DB
    ▸ ⊟ Tables
      ▸ ⊞ accounts
        ▸ ◢ ▢ Columns
              account (string)
              sector (string)
              year_established (long)
              revenue (double)
              employees (long)
              office_location (string)
              subsidiary_of (string)
      ▸ ⊞ products
      ▸ ⊞ salesteam
      ▸ ⊞ salestransaction

```sql
1    -- Top 10 accounts by revenue
2    SELECT TOP 10
3      a.account,
4      a.sector,
5      a.office_location,
6      SUM(s.close_value) AS total_revenue,
7      COUNT(1) AS num_deals,
8      ROUND(SUM(s.close_value) / NULLIF(COUNT(1),0), 2) AS avg_deal_size
9    FROM CRManalysis_DB.dbo.salestransaction s
10   JOIN CRManalysis_DB.dbo.accounts a
11     ON s.account = a.account
12   GROUP BY a.account, a.sector, a.office_location
13   ORDER BY total_revenue DESC;
14
15
16   -- Sales performance metrics per sales agent
17   SELECT
18     s.sales_agent,
19     COALESCE(st.manager, 'Unknown') AS manager,
20     COALESCE(st.regional_office, 'Unknown') AS regional_office,
21     COUNT(1) AS deals_count,
22     SUM(s.close_value) AS total_revenue,
```

**Results**   Messages                                              ⌃

View  [ Table | Chart ]    ↦ Export results ∨

🔍 Search

| account | sector | office_location | total_revenue | num_deals | avg_deal_size |
|---------|--------|-----------------|---------------|-----------|---------------|
| Kan-code | software | United States | 341455 | 187 | 1825 |
| Konex | technolgy | United States | 269245 | 171 | 1574 |
| Condax | medical | United States | 206410 | 159 | 1298 |
| Cheers | entertainment | United States | 198020 | 90 | 2200 |
| Hottechi | technolgy | Korea | 194957 | 193 | 1010 |
| Condilean | marketing | United States | 182522 | 88 | 2074 |

✓ 00:00:02 Query executed successfully.

---

⬆ Publish all ②

| CRMAnalysis_DB ● | SQL script 1 ● |

▷ Run  ↶ Undo ∨   ⬆ Publish  🔗 Query plan    Connect to  ✓ Built-in ∨   Use database  CRMAnalysis_DB ∨  ⟳

```sql
39   -- Revenue by product and series (top products)
40   SELECT
41     p.series,
42     p.product,
43     COUNT(1) AS deals_count,
44     SUM(s.close_value) AS total_revenue,
45     ROUND(AVG(s.close_value),2) AS avg_deal_value
46   FROM CRManalysis_DB.dbo.salestransaction s
47   JOIN CRManalysis_DB.dbo.products p
48     ON s.product = p.product
49   GROUP BY p.series, p.product
50   ORDER BY total_revenue DESC;
51
52
```

**Results**   Messages

View  [ Table | Chart ]    ↦ Export results ∨

🔍 Search

| series | product | deals_count | total_revenue | avg_deal_value |
|--------|---------|-------------|---------------|----------------|
| GTX | GTX Plus Pro | 745 | 2629651 | 3529 |
| MG | MG Advanced | 1084 | 2216387 | 2044 |
| GTX | GTX Plus Basic | 1051 | 705275 | 671 |
| GTX | GTX Basic | 1436 | 499263 | 347 |
| GTK | GTK 500 | 25 | 400612 | 16024 |
| MG | MG Special | 1223 | 43768 | 35 |

```sql
-- Top 10 accounts by revenue
SELECT TOP 10
    a.account,
    a.sector,
    a.office_location,
    SUM(s.close_value) AS total_revenue,
    COUNT(1) AS num_deals,
    ROUND(SUM(s.close_value) / NULLIF(COUNT(1),0), 2) AS avg_deal_size
FROM CRManalysis_DB.dbo.salestransaction s
JOIN CRManalysis_DB.dbo.accounts a
    ON s.account = a.account
GROUP BY a.account, a.sector, a.office_location
ORDER BY total_revenue DESC;
```

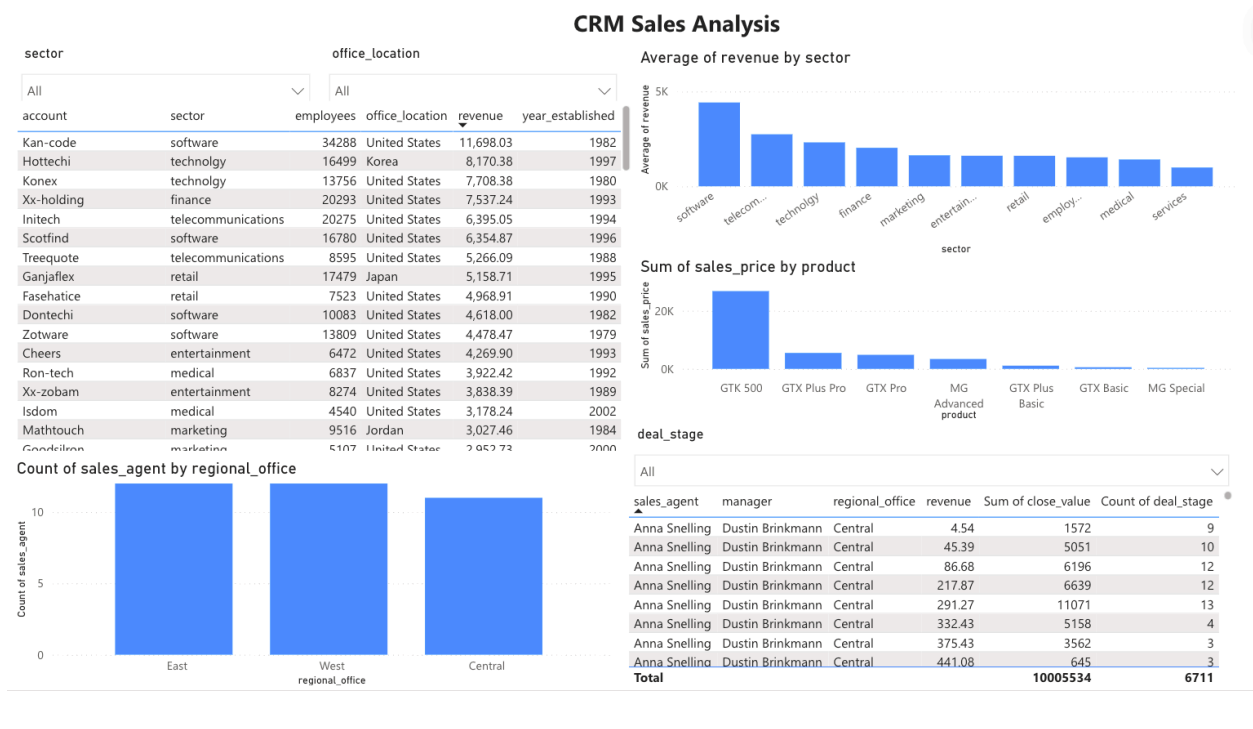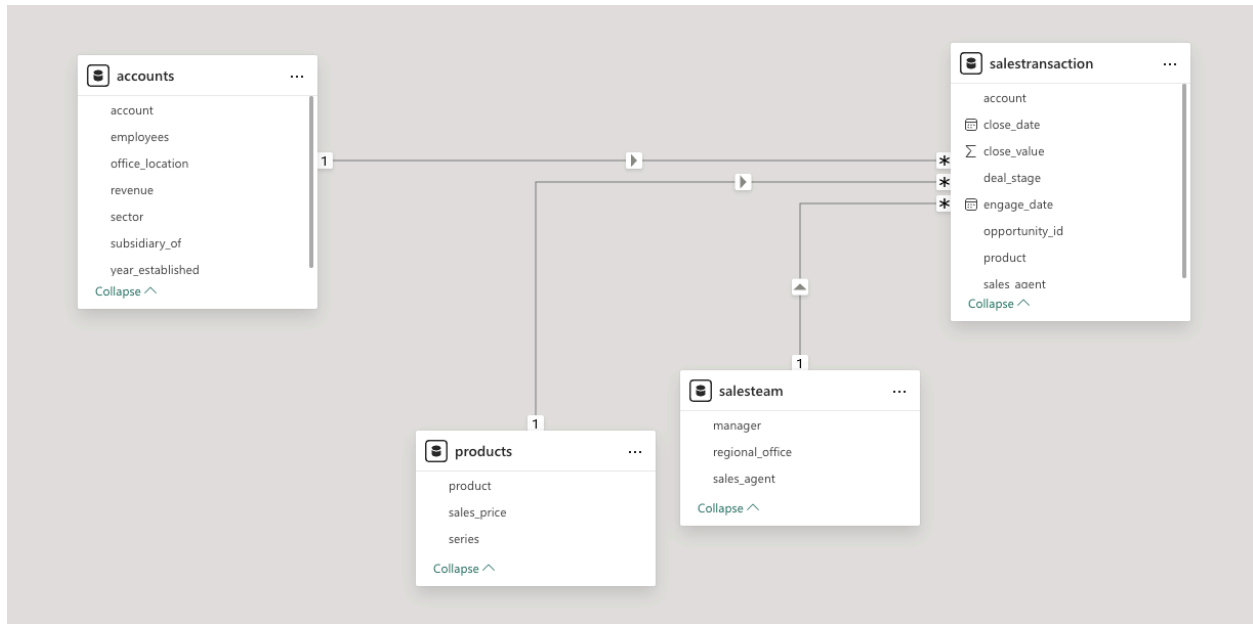# 9. Power BI Integration & Dashboarding

## Semantic Model

1. Connected Power BI to Synapse using the **Azure Synapse Analytics connector**
2. Created a **semantic model** from Synapse tables
3. Defined relationships in a **star schema**:
   o Fact table: salestransaction
   o Dimension tables: accounts, products, salesteam

## Dashboards Built

The Power BI dashboard includes: - Average revenue by sector - Revenue by product and product series - Sales agent distribution by regional office - Interactive tables with account-level details - Slicers for sector, office location, and deal stage

## Outcome

- Interactive and business-friendly dashboards
- End users can explore sales performance and trends dynamically

## 10. Tools & Technologies Used

- Azure Data Factory
- Azure Data Lake Storage Gen2
- Azure Databricks (Apache Spark, PySpark)
- Azure Synapse Analytics (Serverless SQL)
- Power BI

- SQL
- Python

## 11. Key Learnings

- Designing and implementing an end-to-end Azure data pipeline
- Applying data quality checks and transformations using Spark
- Using Synapse Serverless SQL for lakehouse analytics
- Building semantic models and dashboards in Power BI
- Applying star schema concepts for analytics

## 12. Conclusion

This project demonstrates a **real-world CRM sales analytics pipeline** built entirely on Azure. It showcases skills across data ingestion, big data processing, cloud analytics, and business intelligence. The architecture is scalable, cost-efficient, and aligned with industry best practices, making it suitable for portfolio demonstration and interview discussions.