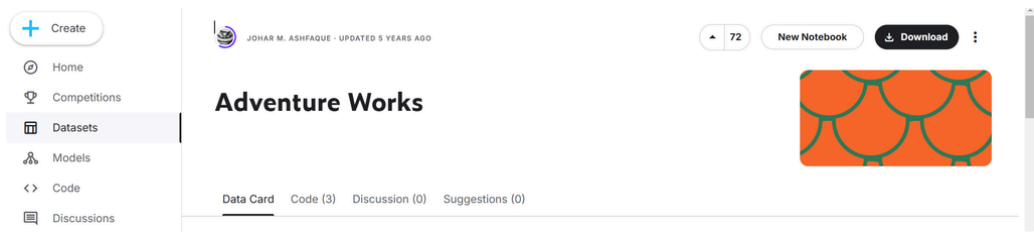


# Azure End-To-End Data Engineering Project Report

This report outlines the steps I followed to complete the Azure End-to-End Data Engineering Project. The project focuses on building a data pipeline using Azure services, starting from data ingestion, followed by transformation, storage, orchestration, and visualization. The goal is to create a scalable and automated data processing pipeline.

## DATASET

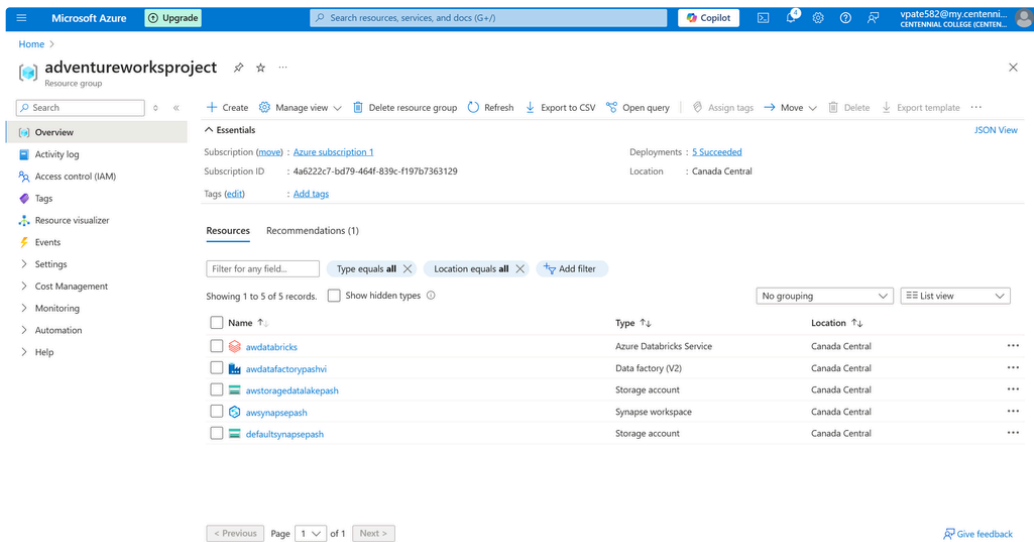
- The dataset used in this project is the **AdventureWorks** dataset from Kaggle. It is a widely-used sample database that contains transactional data for a fictitious multinational manufacturing company. It includes tables for sales, products, customers, and employees.
- The dataset is structured with various relational tables, such as:
  - **SalesOrderHeader**
  - **SalesOrderDetail**
  - **Product**
  - **Customer**
  - **Employee**



## PROJECT SETUP

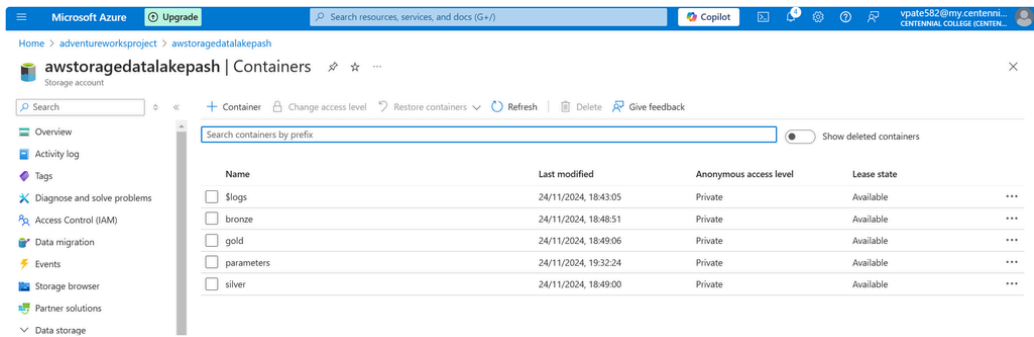
### 2.1 Environment Setup

- **Azure Subscription:** I started by using an existing Azure account to access the services required for the project.
- **Resource Group:** I created a new resource group on Azure to organize all the resources related to this project.



## Blob Storage Setup

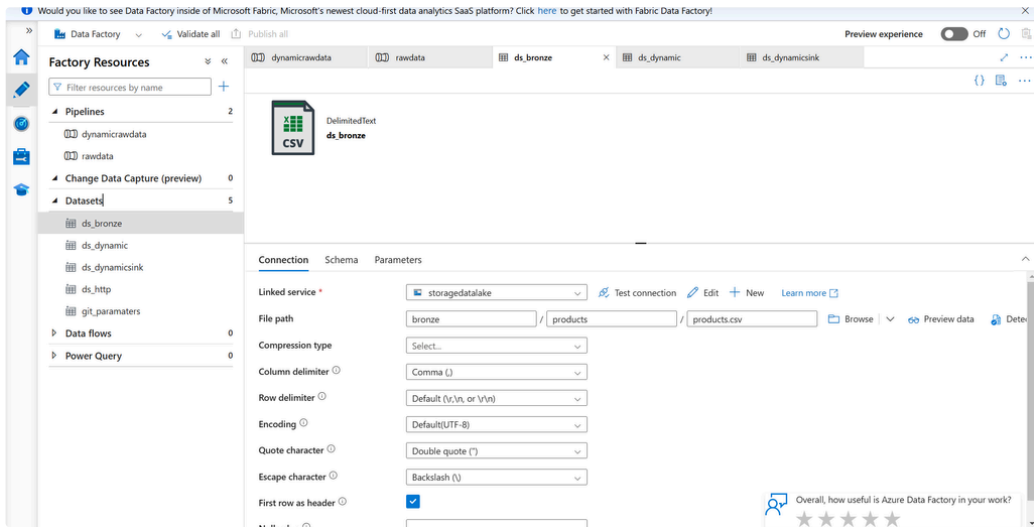
- I uploaded the **AdventureWorks** dataset to **Azure Blob Storage**, using CSV format for easy access and manipulation during the pipeline process.

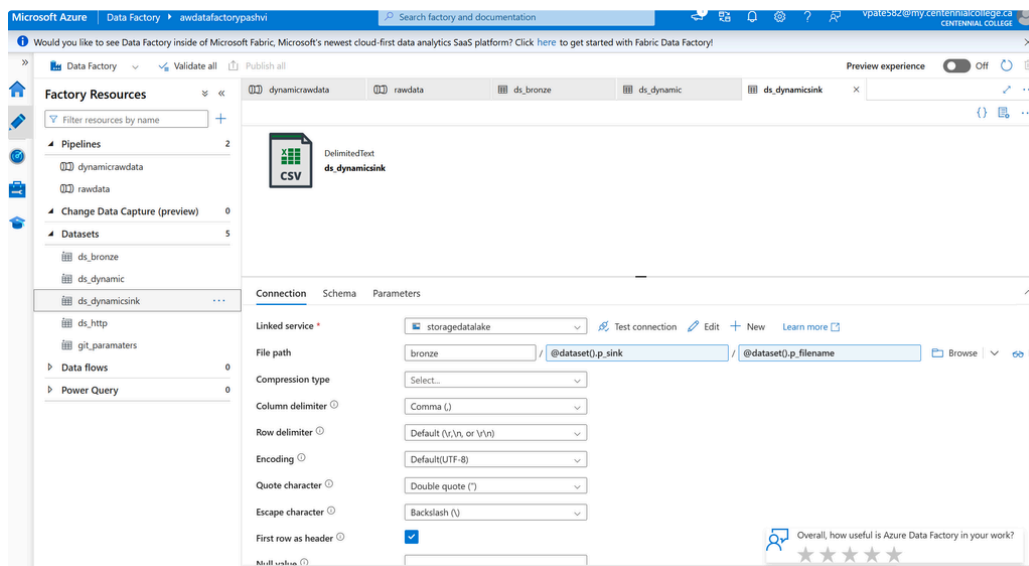
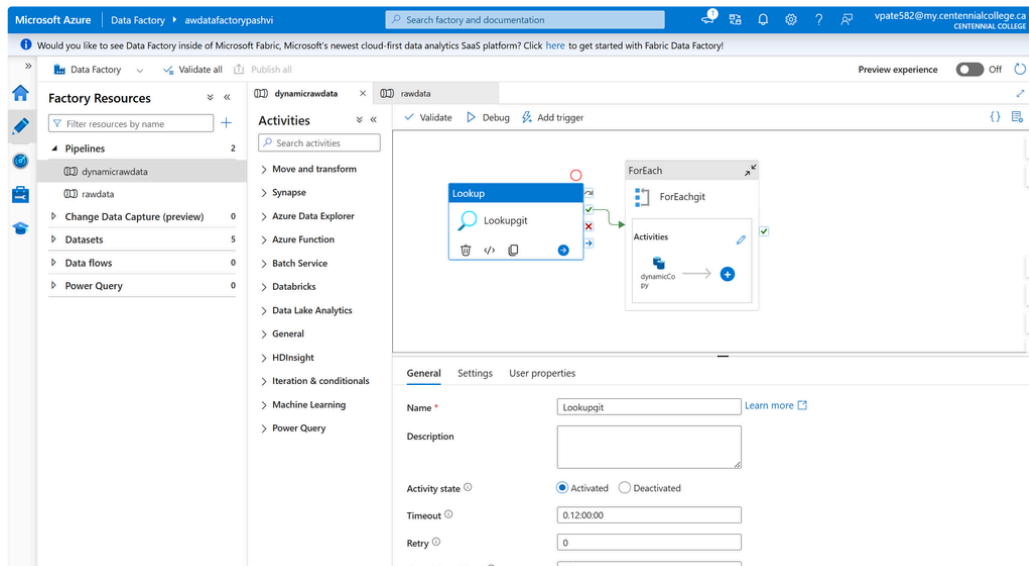
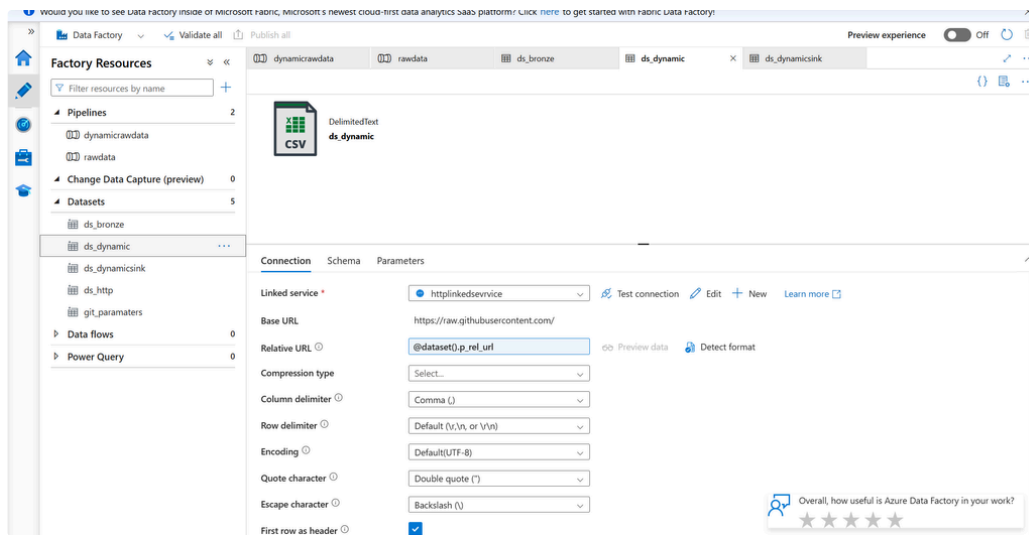


## Creating a Dynamic Pipeline in Azure Data Factory

To streamline the data ingestion process, I created a **dynamic pipeline** in **Azure Data Factory (ADF)**. This pipeline was designed to automatically ingest all files from the **Blob Storage** container without the need for hardcoding individual file names. The steps to achieve this dynamic ingestion were as follows:

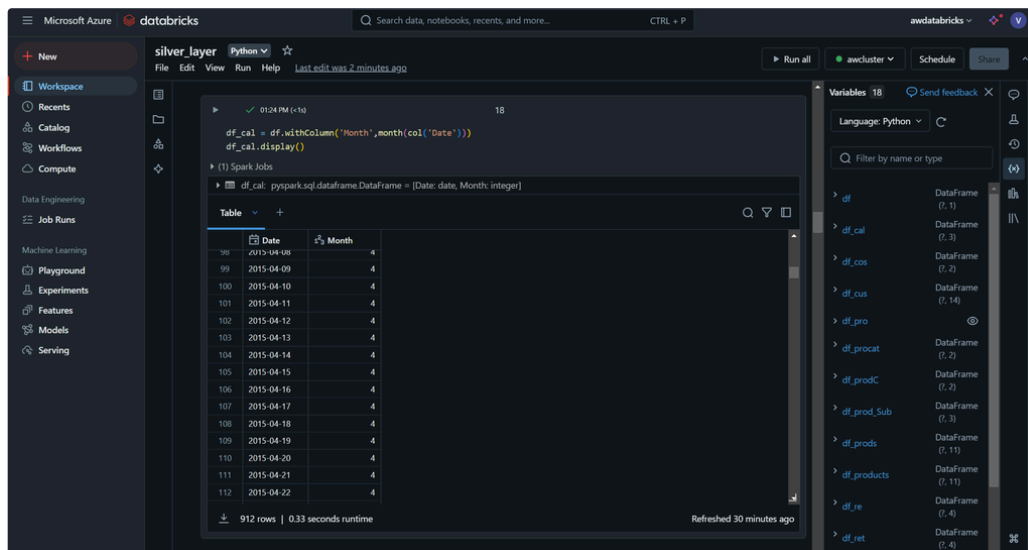
- Get List of Files:** The first step in the pipeline was to use the **Get Metadata** activity to dynamically list all the files within the Blob Storage container. This allowed the pipeline to detect any new or updated files in the container automatically.
- ForEach Activity:** After retrieving the list of files, I used the **ForEach** activity in ADF to iterate through each file. This loop processed each file dynamically, regardless of the number or name of the files. Inside the **ForEach** loop, I used the **Copy Data** activity to ingest each file into Azure Data Lake or Azure SQL Database.
- Dynamic Parameters:** To handle different file names and paths, I used **parameters** within the pipeline to pass dynamic values, such as file names and paths, to the copy activity. This ensured that the pipeline was flexible and scalable, accommodating new files added to Blob Storage without needing to modify the pipeline.





The **AdventureWorks** dataset was analyzed for its key tables and columns that would be relevant for transformation. The main transformation tasks involved:

- **Data Cleansing:** Removing any records with missing values or duplicates, particularly in critical fields like order IDs or customer details.
- **Data Aggregation:** Combining multiple related tables to get meaningful insights such as total sales by customer or product category.
- **Feature Engineering:** Creating new columns, by combining **SalesOrderDetail** line items and grouping by order.

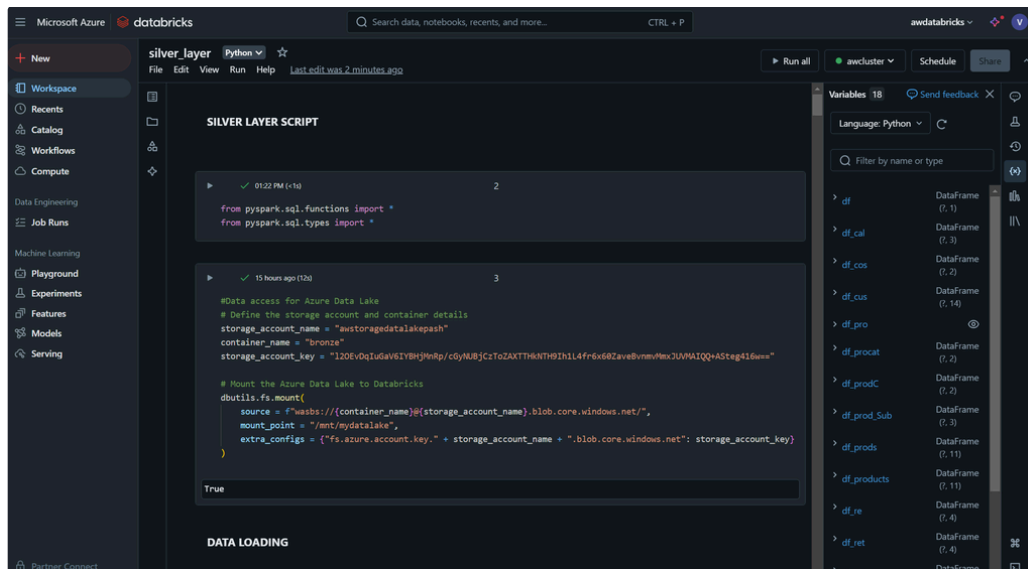


The screenshot shows the Databricks interface for a workspace named 'silver\_layer'. The notebook is written in Python and contains the following code:

```
df_cal = df.withColumn('Month', month(col('Date')))
df_cal.display()
```

The output of the code is a table with 912 rows and 2 columns: 'Date' and 'Month'. The 'Month' column is derived from the 'Date' column. The table is displayed in a table view, and the first 10 rows are visible. The table is refreshed 30 minutes ago.

Date	Month
2015-04-08	4
2015-04-09	4
2015-04-10	4
2015-04-11	4
2015-04-12	4
2015-04-13	4
2015-04-14	4
2015-04-15	4
2015-04-16	4
2015-04-17	4



The screenshot shows the Databricks interface for a workspace named 'silver\_layer'. The notebook is written in Python and contains the following code:

```
from pyspark.sql.functions import *
from pyspark.sql.types import *
```

The code defines the storage account name, container name, and storage account key:

```
#Data access for Azure Data Lake
# Define the storage account and container details
storage_account_name = "awstoragedatalakeash"
container_name = "bronze"
storage_account_key = "120EvDq1uGavSIvHjMhP/cyM8Jc1Z0ZAXTH9H9Ih1L4f6+68Zave@vneWfMcJUMAIQQ+A5teq46w=="
```

The code then mounts the Azure Data Lake to Databricks:

```
# Mount the Azure Data Lake to Databricks
dutils.fs.mount(
    source = f"wasbs://{container_name}@{storage_account_name}.blob.core.windows.net/",
    mount_point = "/mnt/mydatalake",
    extra_configs = {"fs.azure.account.key-" + storage_account_name + ".blob.core.windows.net": storage_account_key}
)
```

The code returns True.

```
True
```

Microsoft Azure databricks

silver\_layer Python

File Edit View Run Help Last edit was 2 minutes ago

Run all awclusterv Schedule Share

df\_cus = spark.read.format('csv')\n.option('header',True)\n.option('inferSchema',True)\n.load('abfss://bronze@awstoragedatalakeash.dfs.core.windows.net/Adventureworks\_Customers')

df\_cus: pyspark.sql.dataframe.DataFrame = [CustomerKey: integer, Prefix: string ... 11 more fields]

df\_procat = spark.read.format('csv')\n.option('header',True)\n.option('inferSchema',True)\n.load('abfss://bronze@awstoragedatalakeash.dfs.core.windows.net/Adventureworks\_Product\_Categories')

df\_procat: pyspark.sql.dataframe.DataFrame = [ProductCategoryKey: integer, CategoryName: string]

df\_pro = spark.read.format('csv')\n.option('header',True)\n.option('inferSchema',True)\n.load('abfss://bronze@awstoragedatalakeash.dfs.core.windows.net/Adventureworks\_Products')

df\_pro: pyspark.sql.dataframe.DataFrame = [ProductKey: integer, ProductSubcategoryKey: integer ... 9 more fields]

Variables 18 Send feedback

Language: Python

Filter by name or type

- df DataFrame (7, 1)
- df\_cal DataFrame (7, 3)
- df\_cos DataFrame (7, 2)
- df\_cus DataFrame (7, 14)
- df\_pro DataFrame (7, 14)
- df\_procat DataFrame (7, 2)
- df\_prodc DataFrame (7, 2)
- df\_prod\_Sub DataFrame (7, 3)
- df\_prods DataFrame (7, 11)
- df\_products DataFrame (7, 11)
- df\_re DataFrame (7, 4)
- df\_ret DataFrame (7, 4)

Microsoft Azure databricks

silver\_layer Python

File Edit View Run Help Last edit was 3 minutes ago

Run all awclusterv Schedule Share

df\_cal.write.format('parquet')\n.option('path', 'abfss://silver@awstoragedatalakeash.dfs.core.windows.net/Adventureworks\_Calendar')\n.save()

df\_cus = spark.read.format('csv')\n.option('header',True)\n.option('inferSchema',True)\n.load('abfss://bronze@awstoragedatalakeash.dfs.core.windows.net/Adventureworks\_Customers')

df\_cus.display()

CustomerKey	Prefix	FirstName	LastName	BirthDate	MaritalStatus
9	11009	MRL	SHANNON	1964-04-01	S
10	11010	MS	JACQUELYN	1964-02-06	S
11	11011	MRL	CURTIS	1963-11-04	M
12	11012	MRS	LAUREN	1963-05-15	M

Variables 18 Send feedback

Language: Python

Filter by name or type

- df DataFrame (7, 1)
- df\_cal DataFrame (7, 3)
- df\_cos DataFrame (7, 2)
- df\_cus DataFrame (7, 14)
- df\_pro DataFrame (7, 14)
- df\_procat DataFrame (7, 2)
- df\_prodc DataFrame (7, 2)
- df\_prod\_Sub DataFrame (7, 3)
- df\_prods DataFrame (7, 11)
- df\_products DataFrame (7, 11)
- df\_re DataFrame (7, 4)
- df\_ret DataFrame (7, 4)

Microsoft Azure databricks

silver\_layer Python

File Edit View Run Help Last edit was 2 minutes ago

Run all awclusterv Schedule Share

DATA LOADING

READ CALAENDER DATA

spark.conf.set('fs.azure.account.key.awstoragedatalakeash.dfs.core.windows.net', '120E4QJudaV6YBHjM8p/c6yN8Jc2ToZAXTHn7H91h1L4f6x6d2av8vnuh3UWMAIQ+5t6g16w==')

df\_cal = spark.read.format('csv').option('header', 'true').option('inferSchema', True).load('abfss://bronze@awstoragedatalakeash.dfs.core.windows.net/Adventureworks\_Calendar')

df\_cal: pyspark.sql.dataframe.DataFrame = [Date: date]

df\_cal.display()

Variables 18 Send feedback

Language: Python

Filter by name or type

- df DataFrame (7, 1)
- df\_cal DataFrame (7, 3)
- df\_cos DataFrame (7, 2)
- df\_cus DataFrame (7, 14)
- df\_pro DataFrame (7, 14)
- df\_procat DataFrame (7, 2)
- df\_prodc DataFrame (7, 2)
- df\_prod\_Sub DataFrame (7, 3)
- df\_prods DataFrame (7, 11)
- df\_products DataFrame (7, 11)
- df\_re DataFrame (7, 4)
- df\_ret DataFrame (7, 4)

Microsoft Azure databricks

silver\_layer Python

```
df_cal = df.withColumn('Month', month(col('Date')))\
            .withColumn('Year', year(col('Date')))\
df_cal.display()
```

(1) Spark Jobs

df\_cal: pyspark.sql.dataframe.DataFrame = [Date: date, Month: integer ... 1 more field]

	Date	Month	Year
1	2015-01-01	1	2015
2	2015-01-02	1	2015
3	2015-01-03	1	2015
4	2015-01-04	1	2015
5	2015-01-05	1	2015
6	2015-01-06	1	2015
7	2015-01-07	1	2015
8	2015-01-08	1	2015
9	2015-01-09	1	2015
10	2015-01-10	1	2015
11	2015-01-11	1	2015
12	2015-01-12	1	2015
13	2015-01-13	1	2015
14	2015-01-14	1	2015
15	2015-01-15	1	2015

912 rows | 0.28 seconds runtime

Refreshed 29 minutes ago

Variables 18

Language: Python

- df: DataFrame (7, 1)
- df\_cal: DataFrame (7, 3)
- df\_cos: DataFrame (7, 2)
- df\_cus: DataFrame (7, 14)
- df\_pro: DataFrame (7, 11)
- df\_procat: DataFrame (7, 2)
- df\_prodc: DataFrame (7, 2)
- df\_prod\_Sub: DataFrame (7, 3)
- df\_prods: DataFrame (7, 11)
- df\_products: DataFrame (7, 11)
- df\_re: DataFrame (7, 4)
- df\_ret: DataFrame (7, 4)

Microsoft Azure databricks

silver\_layer Python

```
df_pro = df_pro.withColumn('ProductSKU', split(col('ProductSKU'), '-')[0])\
            .withColumn('ProductName', split(col('ProductName'), '-')[0])\
df_pro: pyspark.sql.dataframe.DataFrame = [ProductKey: integer, ProductSubcategoryKey: integer ... 9 more fields]
```

(1) Spark Jobs

df\_pro: pyspark.sql.dataframe.DataFrame = [ProductKey: integer, ProductSubcategoryKey: integer ... 9 more fields]

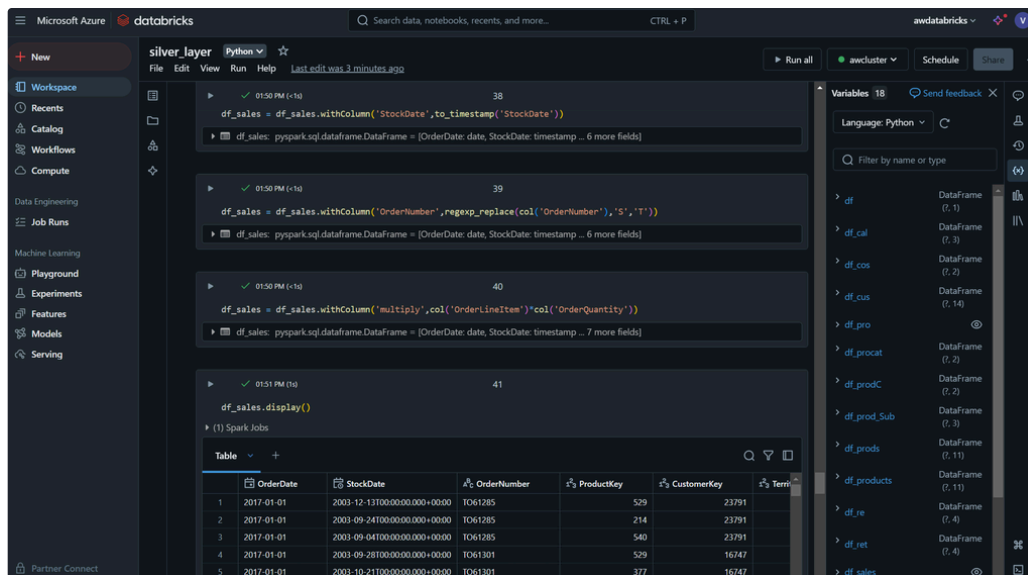
	ProductKey	ProductSubcategoryKey	ProductSKU	ProductName	ModelName	Category
1	214	31	HL	Sport-100	Sport-100	Uni
2	215	31	HL	Sport-100	Sport-100	Uni
3	218	23	SO	Mountain	Mountain Bike Socks	Cov
4	219	23	SO	Mountain	Mountain Bike Socks	Cov
5	220	31	HL	Sport-100	Sport-100	Uni
6	223	19	CA	AWC	Cycling Cap	Tra
7	226	21	LJ	Long-Sleeve	Long-Sleeve Logo Jers...	Uni
8	229	21	LJ	Long-Sleeve	Long-Sleeve Logo Jers...	Uni
9	232	21	LJ	Long-Sleeve	Long-Sleeve Logo Jers...	Uni
10	235	21	LJ	Long-Sleeve	Long-Sleeve Logo Jers...	Uni
11	238	14	FR	HL	HL Road Frame	> i
12	241	14	FR	HL	HL Road Frame	> i
13	244	14	FR	HL	HL Road Frame	> i
14	247	14	FR	HL	HL Road Frame	> i

Partner Connect

Variables 18

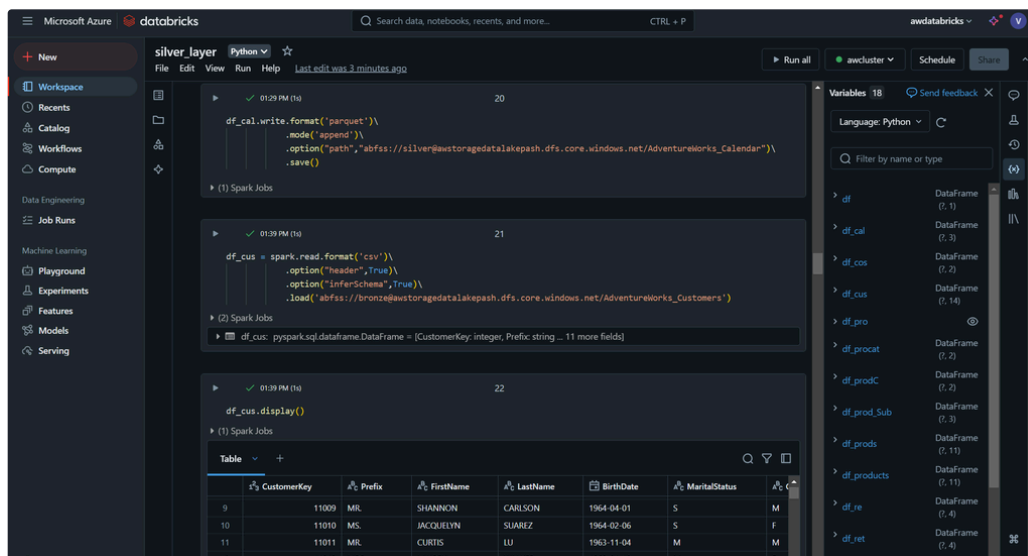
Language: Python

- df: DataFrame (7, 1)
- df\_cal: DataFrame (7, 3)
- df\_cos: DataFrame (7, 2)
- df\_cus: DataFrame (7, 14)
- df\_pro: DataFrame (7, 11)
- df\_procat: DataFrame (7, 2)
- df\_prodc: DataFrame (7, 2)
- df\_prod\_Sub: DataFrame (7, 3)
- df\_prods: DataFrame (7, 11)
- df\_products: DataFrame (7, 11)
- df\_re: DataFrame (7, 4)
- df\_ret: DataFrame (7, 4)



Using Azure Databricks, I applied PySpark to perform data transformations, including:

- **Join operations** to combine related tables
- **Data filtering** to clean out incomplete or incorrect records.
- **GroupBy operations** to summarize data by customer, region, or product.



Microsoft Azure databricks

silver\_layer Python

File Edit View Run Help Last edit was 3 minutes ago

Run all awdatabricks

Workspace

Recents

Catalog

Workflows

Compute

Data Engineering

Job Runs

Machine Learning

Playground

Experiments

Features

Models

Serving

```
df_cus = df_cus.withColumn('fullName', concat_ws(' ', col('Prefix'), col('FirstName'), col('LastName')))
```

```
df_cus.pyspark.sql.DataFrame = [CustomerKey: integer, Prefix: string ... 12 more fields]
```

```
df_cus.display()
```

(1) Spark Jobs

	CustomerKey	Prefix	FirstName	LastName	BirthDate	MaritalStatus	
1	11000	MIL	JON	YANG	1966-04-08	M	
2	11001	MIL	EUGENE	HUANG	1965-05-14	S	
3	11002	MIL	RUBEN	TORRES	1965-08-12	M	
4	11003	MS	CHRISTY	ZHU	1968-02-15	S	
5	11004	MRS.	ELIZABETH	JOHNSON	1968-08-08	S	
6	11005	MIL	JULIO	RUIZ	1965-08-05	S	
7	11007	MIL	MARCO	MEHTA	1964-05-09	M	
8	11008	MRS.	ROBIN	VERHOFF	1964-07-07	S	
9	11009	MIL	SHANNON	CARLSON	1964-04-01	S	
10	11010	MS	JACQUELYN	SUAREZ	1964-02-06	S	
11	11011	MIL	CURTIS	LU	1963-11-04	M	
12	11012	MRS.	LAUREN	WALKER	1968-01-18	M	
13	11013	MIL	IAN	JENKINS	1968-08-06	M	
14	11014	MRS.	SYDNEY	BENNETT	1968-05-09	S	

Variables 18

Language: Python

Filter by name or type

- df DataFrame (0, 1)
- df\_cal DataFrame (0, 3)
- df\_cos DataFrame (0, 2)
- df\_cus DataFrame (0, 14)
- df\_pro DataFrame (0, 2)
- df\_procat DataFrame (0, 2)
- df\_prodc DataFrame (0, 2)
- df\_prod\_Sub DataFrame (0, 3)
- df\_prods DataFrame (0, 11)
- df\_products DataFrame (0, 11)
- df\_re DataFrame (0, 4)
- df\_ret DataFrame (0, 4)

Microsoft Azure databricks

silver\_layer Python

File Edit View Run Help Last edit was 3 minutes ago

Run all awdatabricks

Workspace

Recents

Catalog

Workflows

Compute

Data Engineering

Job Runs

Machine Learning

Playground

Experiments

Features

Models

Serving

```
df_cus.withColumn("fullName", concat(col('Prefix'), lit(' '), col('FirstName'), lit(' '), col('LastName'))).display()
```

(1) Spark Jobs

	CustomerKey	Prefix	FirstName	LastName	BirthDate	MaritalStatus	
1	11000	MIL	JON	YANG	1966-04-08	M	
2	11001	MIL	EUGENE	HUANG	1965-05-14	S	
3	11002	MIL	RUBEN	TORRES	1965-08-12	M	
4	11003	MS	CHRISTY	ZHU	1968-02-15	S	
5	11004	MRS.	ELIZABETH	JOHNSON	1968-08-08	S	
6	11005	MIL	JULIO	RUIZ	1965-08-05	S	
7	11007	MIL	MARCO	MEHTA	1964-05-09	M	
8	11008	MRS.	ROBIN	VERHOFF	1964-07-07	S	
9	11009	MIL	SHANNON	CARLSON	1964-04-01	S	
10	11010	MS	JACQUELYN	SUAREZ	1964-02-06	S	
11	11011	MIL	CURTIS	LU	1963-11-04	M	
12	11012	MRS.	LAUREN	WALKER	1968-01-18	M	
13	11013	MIL	IAN	JENKINS	1968-08-06	M	
14	11014	MRS.	SYDNEY	BENNETT	1968-05-09	S	

10,000+ rows | Truncated data due to row limit | 0.70 seconds runtime

Refreshed 10 minutes ago

Variables 18

Language: Python

Filter by name or type

- df DataFrame (0, 1)
- df\_cal DataFrame (0, 3)
- df\_cos DataFrame (0, 2)
- df\_cus DataFrame (0, 14)
- df\_pro DataFrame (0, 2)
- df\_procat DataFrame (0, 2)
- df\_prodc DataFrame (0, 2)
- df\_prod\_Sub DataFrame (0, 3)
- df\_prods DataFrame (0, 11)
- df\_products DataFrame (0, 11)
- df\_re DataFrame (0, 4)
- df\_ret DataFrame (0, 4)

Microsoft Azure databricks

silver\_layer Python

File Edit View Run Help Last edit was now

Run all awdatabricks

Workspace

Recents

Catalog

Workflows

Compute

Data Engineering

Job Runs

Machine Learning

Playground

Experiments

Features

Models

Serving

10,000+ rows | Truncated data due to row limit | 0.58 seconds runtime

Refreshed 3 minutes ago

```
df_sales.groupby('OrderDate').agg(count('OrderNumber').alias('total_order')).display()
```

(2) Spark Jobs

Table Visualization 1

New charts: ON

400

total\_order

OrderDate

911 rows

Refreshed 1 minute ago

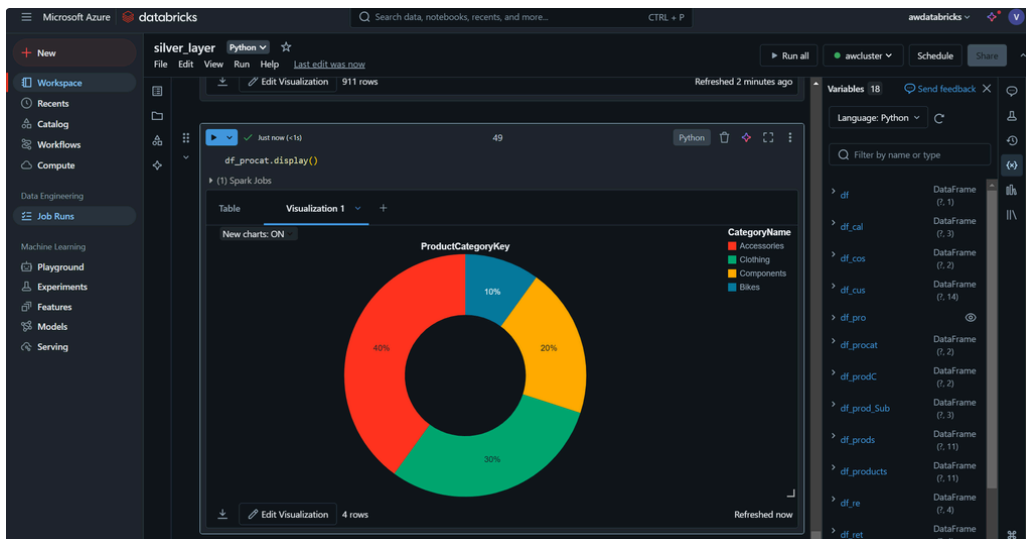
Variables 18

Language: Python

Filter by name or type

- df DataFrame (0, 1)
- df\_cal DataFrame (0, 3)
- df\_cos DataFrame (0, 2)
- df\_cus DataFrame (0, 14)
- df\_pro DataFrame (0, 2)
- df\_procat DataFrame (0, 2)
- df\_prodc DataFrame (0, 2)
- df\_prod\_Sub DataFrame (0, 3)
- df\_prods DataFrame (0, 11)
- df\_products DataFrame (0, 11)
- df\_re DataFrame (0, 4)
- df\_ret DataFrame (0, 4)

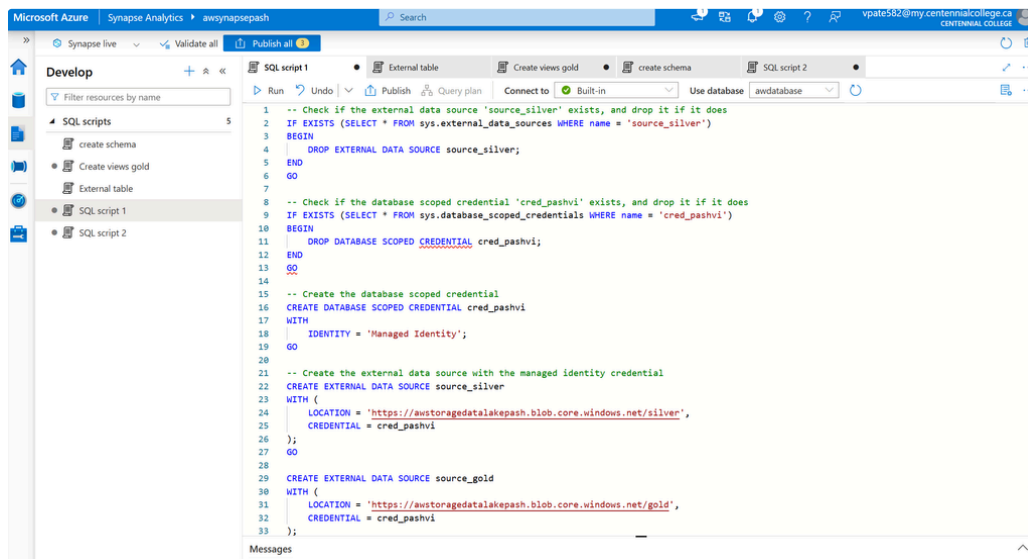
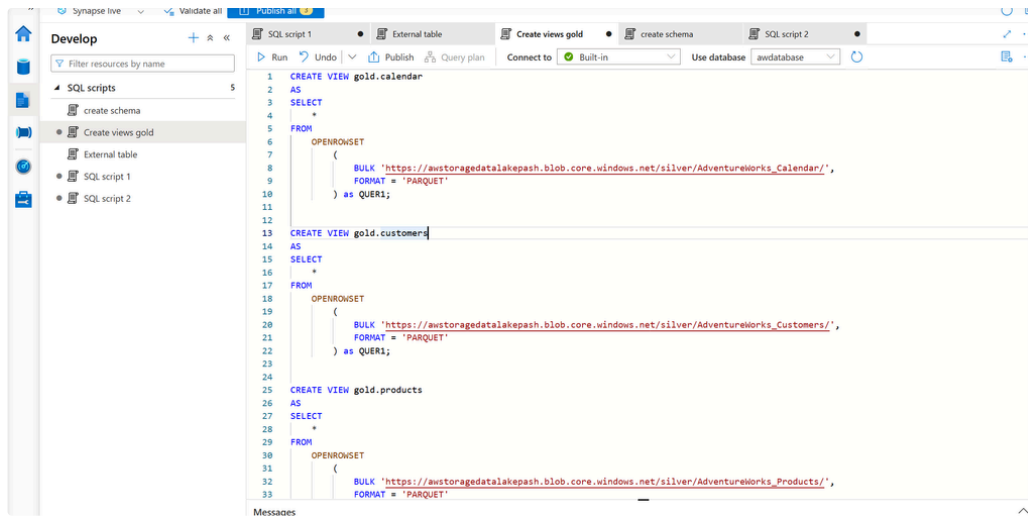




The screenshot shows the Microsoft Azure Databricks workspace. The main area displays a table titled "Table" with columns "OrderDate" and "total\_order". The table contains 15 rows of data. The code editor shows a Python script with `df_sales.groupby('OrderDate').agg(count('OrderNumber').alias('total_order')).display()`. The right sidebar lists variables including `df`, `df_cal`, `df_cos`, `df_cus`, `df_pro`, `df_procat`, `df_procatC`, `df_prod_Sub`, `df_prods`, `df_products`, `df_re`, `df_ret`, and `df_sales`.

	OrderDate	total_order
1	2017-01-06	151
2	2017-01-27	142
3	2017-02-26	119
4	2017-01-24	173
5	2017-06-29	172
6	2017-02-16	124
7	2017-04-09	140
8	2017-02-28	162
9	2017-03-28	149
10	2017-06-30	136
11	2017-01-30	145
12	2017-05-11	164
13	2017-02-10	168
14	2017-04-25	176
15	2017-03-19	152

For data analytics, I used **Azure Synapse Analytics** to run SQL-based queries and derive business insights from the transformed data. Synapse integrated with both Data Lake Storage and Azure SQL Database for seamless querying.



## Azure Logic Apps for Error Handling

- I configured **Azure Logic Apps** to send automated notifications in case of failures, ensuring that issues are quickly addressed.
- I set up **Azure Monitor** to track the performance of the entire data pipeline, including resource utilization and error tracking.

## Challenges & Learnings

Throughout this project, I faced the following challenges:

- **Data Quality Issues:** Encountered missing or inconsistent data, which required careful handling during the transformation phase.
- **Pipeline Orchestration:** Ensuring that the data ingestion, transformation, and loading steps were coordinated effectively was a key challenge, especially when dealing with large datasets.

### Learnings:

- Gained hands-on experience with core Azure services, such as Azure Data Factory, Databricks, and Synapse.
- Improved my skills in automation and monitoring within Azure.

The Azure End-to-End Data Engineering project successfully implemented an automated data pipeline using a range of Azure services. The project helped me build expertise in data ingestion, transformation, orchestration, and reporting, and I am confident that these skills will be valuable for future data engineering tasks.