



IST 615
Cloud Management

Final Project Report

E-COMMERCE ANALYTICS USING AZURE SERVICES FOR OLIST

Group 9:
M K Sowmeya
Dev Jindani
Raaj Mutreja
Sourabh Gavhane

Table of Contents

Sr No.	Topic	Page no.
1	Project Objective	1
2	Azure Cloud Services	2
3	Cloud Services Used	3
4	Integration of Cloud Services	4
5	Execution	5
6	Key Issues Faced	16
7	Visualization	18
8	Conclusion	19
9	References	20

PROJECT OBJECTIVE

This project is designed to leverage Azure cloud services, specifically Azure Blob Storage, Azure SQL Database, and Azure Data Factory, to create a robust solution for analyzing and predicting customer behavior. The primary goal is to enhance sales strategies for Olist, a prominent Brazilian e-commerce platform.

The core focus revolves around incorporating three key datasets—customer orders, order items, and customer details—into Azure's cloud environment. The datasets will undergo processing and transformation, including cleaning and structuring, to ensure they are well-prepared for analysis.

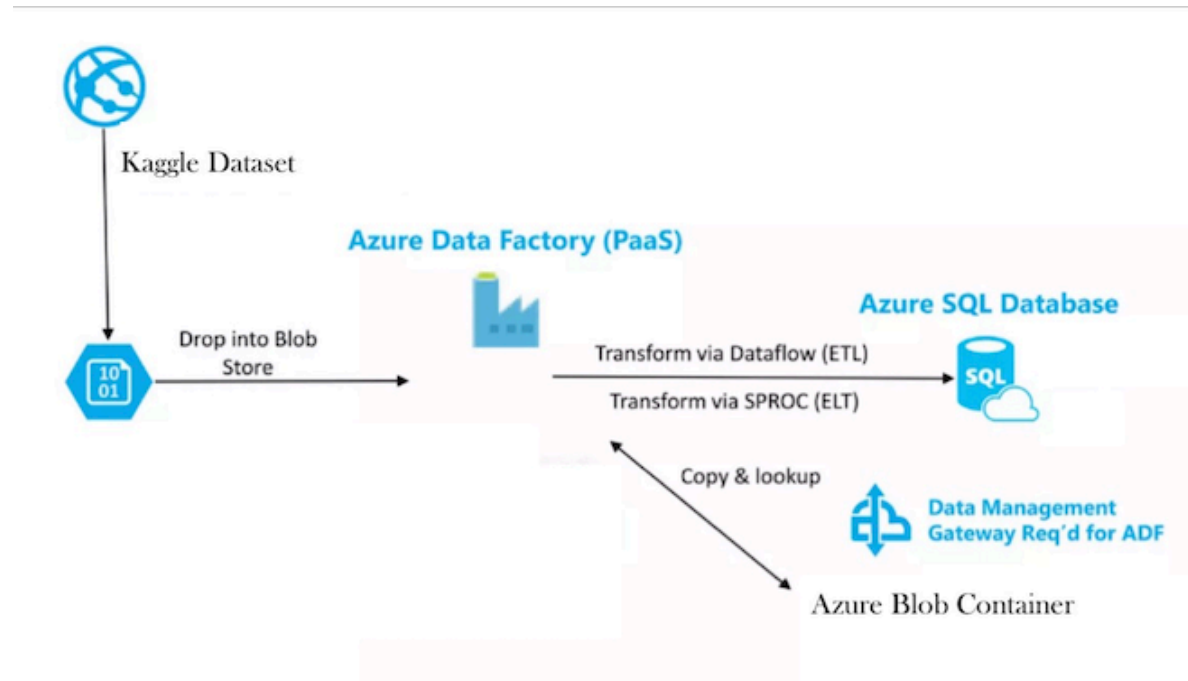
Within Azure's cloud ecosystem, these datasets will be stored and managed efficiently using Azure Blob Storage. Azure SQL Database will then come into play, providing a reliable and scalable relational database service for seamless integration and retrieval of structured data.

The innovative aspect of this project lies in the application of Azure Data Studio. Using this platform, analysis will be developed to identify trends and patterns in customer behavior. This analytical capability aims to provide valuable insights for Olist, enabling the platform to anticipate customer preferences and optimize sales strategies.

In summary, the project aims to harness Azure's cloud services to streamline data processing and develop predictive models for enhanced customer behavior analysis. The ultimate objective is to equip Olist with the tools needed to refine its sales strategies and stay ahead in the competitive e-commerce landscape.

AZURE CLOUD SERVICES

Flow diagram of Azure Service Integration:



INTEGRATION OF CLOUD SERVICES

Our project's integration strategy leverages the capabilities of Azure Blob Storage, Azure SQL Database, and Azure Machine Learning to create a cohesive and powerful data processing ecosystem:

1. Azure Blob Storage:

- Serves as the primary storage solution for datasets.
- Enables efficient data ingestion and staging.
- Provides secure and scalable storage for both raw and processed datasets.

2. Azure SQL Database:

- Facilitates seamless data integration from diverse sources, including customer orders, order items, and customer details.
- Acts as a fully managed relational database service for structured and organized data storage.
- Ensures optimal query performance, allowing for swift data retrieval during analysis.

3. Azure Data Factory:

- **Primary Storage Solution:** Serves as the central repository for datasets.
- **Efficient Data Ingestion:** Facilitates streamlined data ingestion and staging processes.
- **Secure and Scalable Storage:** Provides a secure and scalable environment for both raw and processed datasets.

This integrated approach establishes a harmonized cloud environment, where Azure Blob Storage acts as the central repository, Azure SQL Database serves as the relational foundation, and Azure data factory drives advanced analytics. Through this synergy, our project aims to provide Olist with valuable insights, enhancing its sales strategies through a data-driven approach.

EXECUTION

Azure Blob storage containers:

Home > Recent > ProjectStorage >

projectteam9

Storage account

Search

Overview

Activity log

Tags

Diagnose and solve problems

Access Control (IAM)

Data migration

Events

Storage browser

Storage Mover

Data storage

Containers

File shares

Queues

Tables

Security + networking

Networking

Front Door and CDN

Access keys

Shared access signature

Encryption

Upload

Open in Explorer

Delete

Move

Refresh

Open in mobile

CLI / PS

Feedback

Essentials

JSON View

Resource group (move)

Location

Primary/Secondary Location

Subscription (move)

Subscription ID

Disk state

Performance

Replication

Account kind

Provisioning state

Created

Tags (edit)

Add tags

Properties

Monitoring

Capabilities (7)

Recommendations (0)

Tutorials

Tools + SDKs

Blob service

Hierarchical namespace

Default access tier

Blob anonymous access

Blob soft delete

Container soft delete

Versioning

Change feed

NFS v3

Allow cross-tenant replication

Disabled

Hot

Disabled

Enabled (7 days)

Enabled (7 days)

Disabled

Disabled

Disabled

Disabled

File service

Security

Require secure transfer for REST API operations

Storage account key access

Minimum TLS version

Infrastructure encryption

Enabled

Enabled

Version 1.2

Disabled

Networking

Allow access from

Number of private endpoint connections

Network routing

Access for trusted Microsoft services

All networks

0

Microsoft network routing

Yes

Home > Recent > projectteam9 | Containers >

projectcontainer

Container

Search

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

Metadata

Upload

Change access level

Refresh

Delete

Change tier

Acquire lease

Break lease

View snapshots

Create snapshot

Give feedback

Authentication method: Access key (Switch to Microsoft Entra user account)

Location: projectcontainer

Search blobs by prefix (case-sensitive)

Show deleted blobs

Add filter

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state	
<input type="checkbox"/> olist_customers_dataset.csv	11/6/2023, 6:39:39 PM	Hot (Inferred)		Block blob	8.62 MiB	Available	***
<input type="checkbox"/> olist_geolocation_dataset.csv	11/6/2023, 6:39:46 PM	Hot (Inferred)		Block blob	58.44 MiB	Available	***
<input type="checkbox"/> olist_order_items_dataset.csv	11/6/2023, 6:39:37 PM	Hot (Inferred)		Block blob	14.72 MiB	Available	***
<input type="checkbox"/> olist_order_payments_dataset.csv	11/6/2023, 6:39:37 PM	Hot (Inferred)		Block blob	5.51 MiB	Available	***
<input type="checkbox"/> olist_order_reviews_dataset.csv	11/6/2023, 6:39:36 PM	Hot (Inferred)		Block blob	13.78 MiB	Available	***
<input type="checkbox"/> olist_orders_dataset.csv	11/6/2023, 6:39:39 PM	Hot (Inferred)		Block blob	16.84 MiB	Available	***
<input type="checkbox"/> olist_products_dataset.csv	11/6/2023, 6:39:37 PM	Hot (Inferred)		Block blob	2.27 MiB	Available	***
<input type="checkbox"/> olist_sellers_dataset.csv	11/13/2023, 12:54:40...	Hot (Inferred)		Block blob	170.61 KiB	Available	***
<input type="checkbox"/> product_category_name_translati...	11/13/2023, 12:54:40...	Hot (Inferred)		Block blob	2.55 KiB	Available	***

Microsoft Azure | Data Factory | ProjectDF9

Search factory and documentation

Microsoft recently announced the public preview of Microsoft Fabric, a brand new and exciting way to build cloud-first data analytics. Click [here](#) to get started with Fabric Data Factory!

Factory Resources

- Pipelines: 1
 - pl_customer
- Change Data Capture (preview): 0
- Datasets: 5
 - ds_customer
 - ds_orderItems
 - ds_orders
 - ds_products
 - ds_sellers
- Data flows: 0
- Power Query: 0

Activities

- Move and transform
- Synapse
- Azure Data Explorer
- Azure Function
- Batch Service
- Databricks
- Data Lake Analytics
- General
- HDInsight
- Iteration & conditionals
- Machine Learning
- Power Query

Validate | Debug | Add trigger

Properties

General

Name: pl_customer

Description:

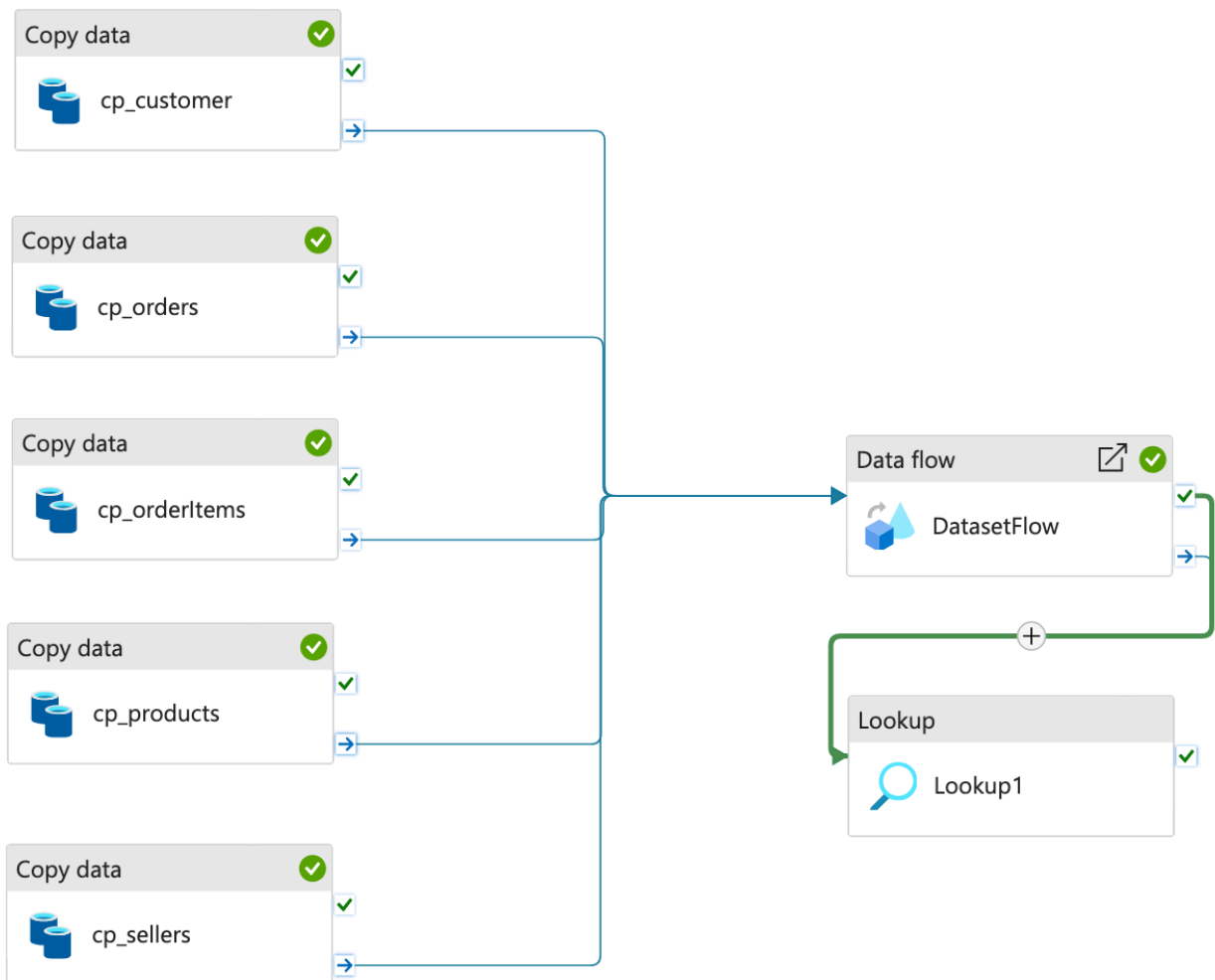
Annotations

+ New

Parameters | Variables | Settings | Output

Showing 1 - 5 of 5 items

Activity name	Activity status	Activity type	Rt
cp_orderItems	Succeeded	Copy data	1'
cp_customer	Succeeded	Copy data	1'
cp_sellers	Succeeded	Copy data	1'
cp_products	Succeeded	Copy data	1'
cp_orders	Succeeded	Copy data	1'



- In ADF, a pipeline is a logical grouping of activities that together perform a task. The activities in a pipeline define actions to perform on your data.
- The pipeline shown in the image used for copying data from various sources, likely to aggregate them into a common data store or to move them into a transformation process (Dataset Flow).
- The "Copy Data" activities indicate that data is being moved from the source datasets to a destination, which could be a database, a data lake, or some other storage service.
- The "Data flow" activity would then apply transformations to the data that has been copied, such as joining, cleaning, enriching, or aggregating the data before it is loaded into its destination.
- To put it all together, this ADF pipeline is being used to automate the extraction and transformation of data from multiple sources.
- Once the data is prepared through the "Data flow", it can then be utilized for analytics, reporting, or further data processing tasks.
- To perform SQL operations as part of your pipeline, we added a "Lookup" activity if we need to retrieve data.

Home > projectteam9 | Containers >

projectsink ...

Container

Search <<

Upload Change access level Refresh Delete Change tier Acquire lease Break lease View snapshots Create snapshot Give feedback

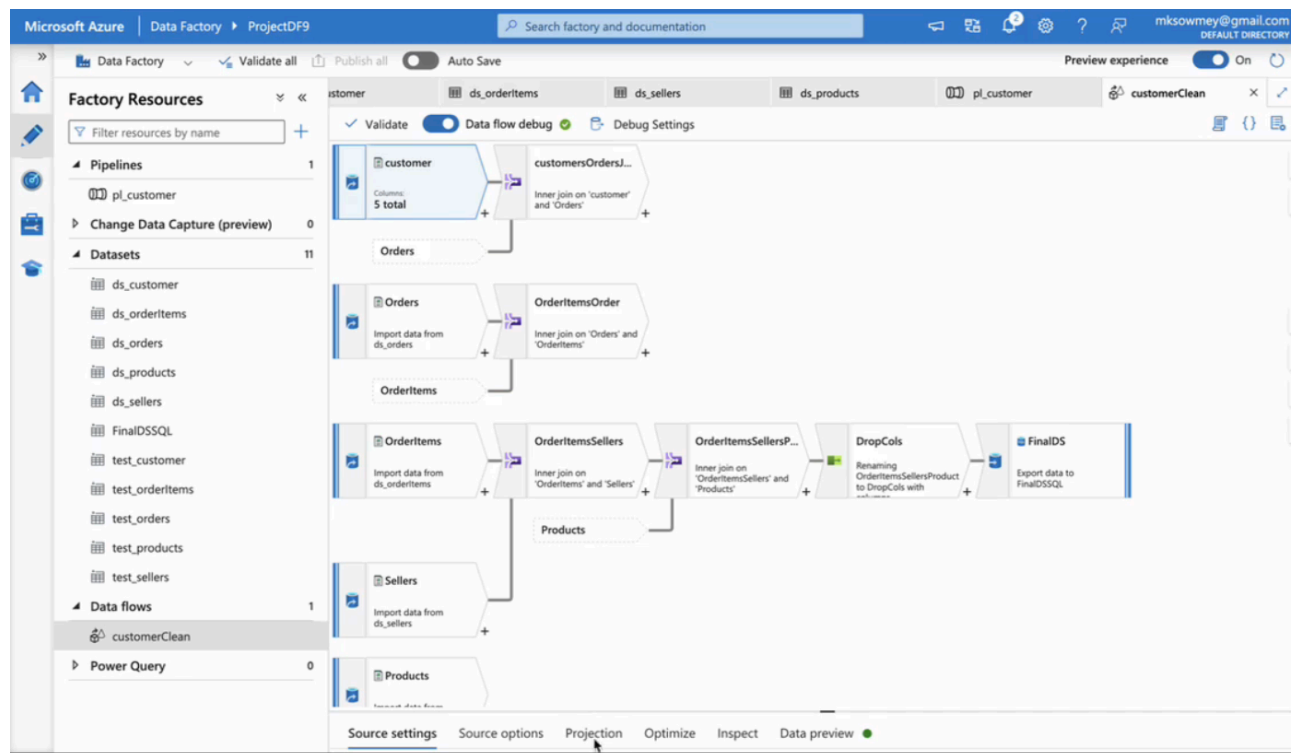
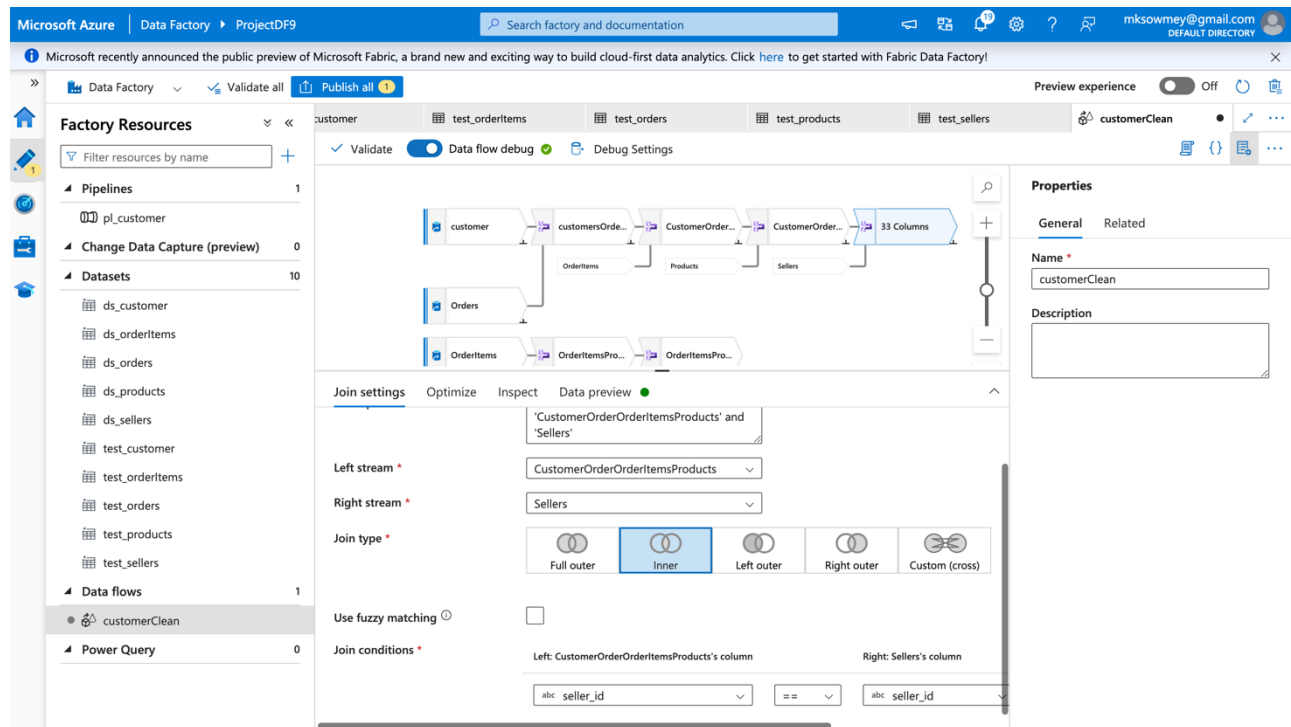
Authentication method: Access key (Switch to Microsoft Entra user account)

Location: projectsink / output

Search blobs by prefix (case-sensitive) Show deleted blobs

Add filter

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
<input type="checkbox"/> [.]						...
<input type="checkbox"/> op_customer.csv	11/20/2023, 2:06:25 ...	Hot (Inferred)		Block blob	8.62 MiB	Available ...
<input type="checkbox"/> op_orderItems.csv	11/20/2023, 2:06:23 ...	Hot (Inferred)		Block blob	14.72 MiB	Available ...
<input type="checkbox"/> op_orders.csv	11/20/2023, 2:06:24 ...	Hot (Inferred)		Block blob	16.84 MiB	Available ...
<input type="checkbox"/> op_products.csv	11/20/2023, 2:06:22 ...	Hot (Inferred)		Block blob	2.27 MiB	Available ...
<input type="checkbox"/> op_sellers.csv	11/20/2023, 2:06:23 ...	Hot (Inferred)		Block blob	170.61 KiB	Available ...

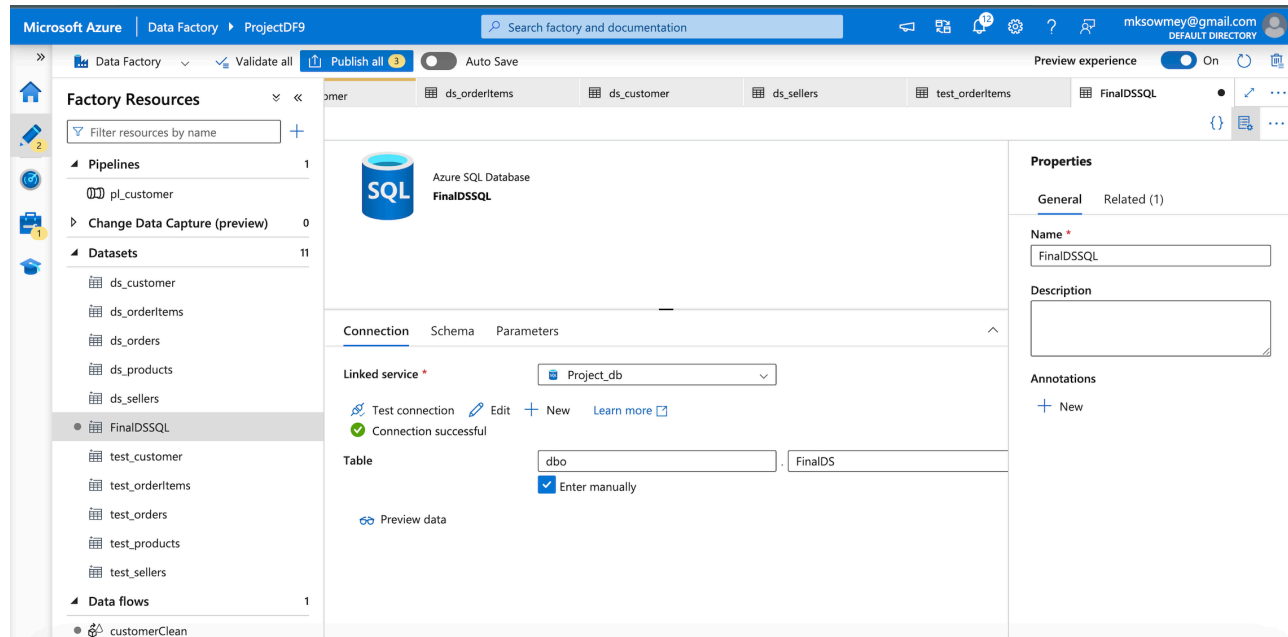


- In Azure Data Factory, a data flow represents a series of data transformation steps that are visually designed using a drag-and-drop interface. The use of data flow within Azure Data Factory is multifaceted and includes the following key aspects as per our scenario.
- Data Transformation: Data flows allow users to perform data transformation processes such as sorting, filtering, aggregating, and merging data without writing code. Transformations

are defined visually and executed as part of the Azure Data Factory pipeline.

- **ETL Processes:** They are essential for creating ETL (Extract, Transform, Load) processes, which enable businesses to consolidate data from various sources, transform it into a usable format, and load it into a destination data store for analytics or other purposes.

The above data flow is likely part of a larger data pipeline that includes several steps like joining different datasets, transforming datasets through various operations, and eventually preparing a final dataset for reporting, analytics, & operational use.



- **Sink Configuration:** Once the data is processed within the data flow, the final output, known as the "sink," is configured to be stored or loaded into a SQL Server database, commonly referred to as SQL Studio.
- **SQL Server Integration:** The data is sent to the SQL Server database, where it can be stored, analyzed, or used for various purposes like reporting, querying, or serving as a source for other applications.
- **Utilization in SQL Studio:** Within SQL Studio, the received data can be further analyzed, joined with other datasets, queried using SQL queries, or used in various business intelligence and reporting tools to derive insights or support decision-making processes.
- **This seamless flow from data transformation in the data flow to storing the finalized data in SQL Studio** ensures the availability of processed and structured data in a database environment for efficient storage, analysis, and utilization within an organization's data ecosystem.

SQLQuery_2 - porjec...sowmey) 7 X SQLQuery_3 - porjec...sowmey)

```

1 SELECT TOP (1000) [shipping_limit_date]
2     , [price]
3     , [seller_zip_code_prefix]
4     , [seller_city]
5     , [seller_state]
6     , [product_category_name]
7 FROM [dbo].[FinalDS]

```

Results Messages

	shipping_limit_date	price	seller_zip_code_prefix	seller_city	seller_state	product_category_name
1	2018-08-07 02:15:06	39.90	07194	guarulhos	SP	beleza_saude
2	2018-08-07 02:15:06	39.90	07194	guarulhos	SP	beleza_saude
3	2018-08-07 02:15:06	39.90	07194	guarulhos	SP	beleza_saude
4	2018-08-07 02:15:06	39.90	07194	guarulhos	SP	beleza_saude
5	2018-06-22 17:31:11	349.97	03694	sao paulo	SP	perfumaria
6	2018-02-19 16:27:35	28.50	14709	bebedouro	SP	automotivo
7	2018-02-19 16:27:35	28.50	14709	bebedouro	SP	automotivo
8	2018-07-03 18:30:18	37.99	04160	sao paulo	SP	telefonica
9	2018-07-03 18:30:18	37.99	04160	sao paulo	SP	telefonica
10	2018-06-12 12:31:07	46.20	11704	praia grande	SP	eletrodomesticos
11	2018-04-24 08:51:40	8.99	85960	marechal candido rondon	PR	sinalizacao_e_seguranca
12	2018-04-24 08:51:40	8.99	85960	marechal candido rondon	PR	sinalizacao_e_seguranca
13	2018-04-24 08:51:40	8.99	85960	marechal candido rondon	PR	sinalizacao_e_seguranca
14	2018-04-17 11:10:52	362.45	89820	xanxere	SC	automotivo
15	2017-11-23 16:10:23	69.90	14940	ibitinga	SP	cama_mesa_banho
16	2018-08-06 05:05:11	29.00	13060	campinas	SP	utilidades_domesticas
17	2018-03-12 03:40:11	99.90	71691	brasilia	DF	eletronicos
18	2018-08-17 11:45:21	33.90	74230	goiania	GO	utilidades_domesticas
19	2018-05-22 04:54:43	399.00	04782	sao paulo	SP	utilidades_domesticas
20	2018-03-07 21:00:50	18.00	14092	ribeirao preto	SP	informatica_acessorios
21	2018-07-24 10:10:17	79.90	11704	praia grande	SP	eletrodomesticos
22	2018-03-29 02:15:49	190.00	14075	ribeirao preto	SP	climatizacao
23	2018-07-25 15:00:22	65.00	80310	curitiba	PR	moveis_decoracao
24	2018-07-25 15:00:22	65.00	80310	curitiba	PR	moveis_decoracao
25	2018-07-25 15:00:22	65.00	80310	curitiba	PR	moveis_decoracao
26	2018-07-25 15:00:22	65.00	80310	curitiba	PR	moveis_decoracao
27	2018-02-13 17:55:30	39.90	80230	curitiba	PR	papelaria
28	2017-07-04 02:44:18	15.00	11701	praia grande	SP	esporte_lazer
29	2017-09-13 11:35:13	279.99	22210	rio de janeiro	RJ	utilidades_domesticas

Ln 7, Col 23 Spaces: 4 UTF-8 CRLF 1,000 rows MSSQL 00:00:00 projectserver.database.windows.net : Project

```

10
11 SELECT AVG(price) AS average_price FROM [dbo].[FinalDS];
12

```

Results Messages

	average_price
1	120.653739

```

1  -- SELECT DISTINCT product_category_name FROM [dbo].[FinalDS];
2
3  SELECT product_category_name, COUNT(*) AS product_count
4  FROM [dbo].[FinalDS]
5  GROUP BY product_category_name
6  ORDER BY product_count DESC;
7  |

```

Results Messages

	product_category_name	product_count
1	cama_mesa_banho	22230
2	beleza_saude	19340
3	esporte_lazer	17282
4	moveis_decoracao	16668
5	informatica_acessorios	15654
6	utilidades_domesticas	13928
7	relogios_presentes	11982
8	telefonica	9090
9	ferramentas_jardim	8694
10	automotivo	8470
11	brinquedos	8234
12	cool_stuff	7592
13	perfumaria	6838
14	bebes	6130
15	eletronicos	5534
16	papelaria	5034
17	fashion_bolsas_e_acessorios	4062
18	pet_shop	3894
19	moveis_escritorio	3382
20	NULL	3206
21	consoles_games	2274
22	malas_acessorios	2184
23	construcao_ferramentas_construcao	1858
24	eletrodomesticos	1542
25	instrumentos_musicais	1360
26	eletroportateis	1358
27	casa_construcao	1208
28	livros_interesse_geral	1106
29	alimentos	1020
30	moveis_sala	1006
31	casa_conforto	868
32	bebidas	758
33	audio	728

SQLQuery_2 - porjec...sowmey) 7

```
7
8 -- ALTER TABLE [dbo].[FinalDS]
9 -- ALTER COLUMN price DECIMAL(10, 2); -- Assuming a decimal data type with 10 total digits and 2 decimal places
10
11 -- SELECT AVG(price) AS average_price FROM [dbo].[FinalDS];
12
13 -- SELECT COUNT(DISTINCT seller_zip_code_prefix) AS unique_seller_count
14 -- FROM [dbo].[FinalDS];
15
16 SELECT seller_city, COUNT(*) AS product_count
17 FROM [dbo].[FinalDS]
18 GROUP BY seller_city
19 ORDER BY product_count DESC;
20
21
```

Results Messages Chart

	seller_city	product_count
1	sao paulo	55966
2	ibitinga	15500
3	curitiba	6032
4	santo andre	5928
5	belo horizonte	5186
6	sao jose do rio preto	5158
7	rio de janeiro	4884
8	guarulhos	4724
9	ribeirao preto	4538
10	maringa	4440
11	piracicaba	3806
12	itaguapecetuba	3306
13	campinas	2860
14	salto	2692
15	praia grande	2666
16	campo limpo paulista	2390
17	guariba	2312
18	sao bernardo do campo	2250
19	jacarei	1908
20	limeira	1866
21	brasilia	1722
22	petropolis	1680
23	sao jose dos campos	1606
24	porto alegre	1604
25	ilicinea	1540
26	pedreira	1416
27	joinville	1340

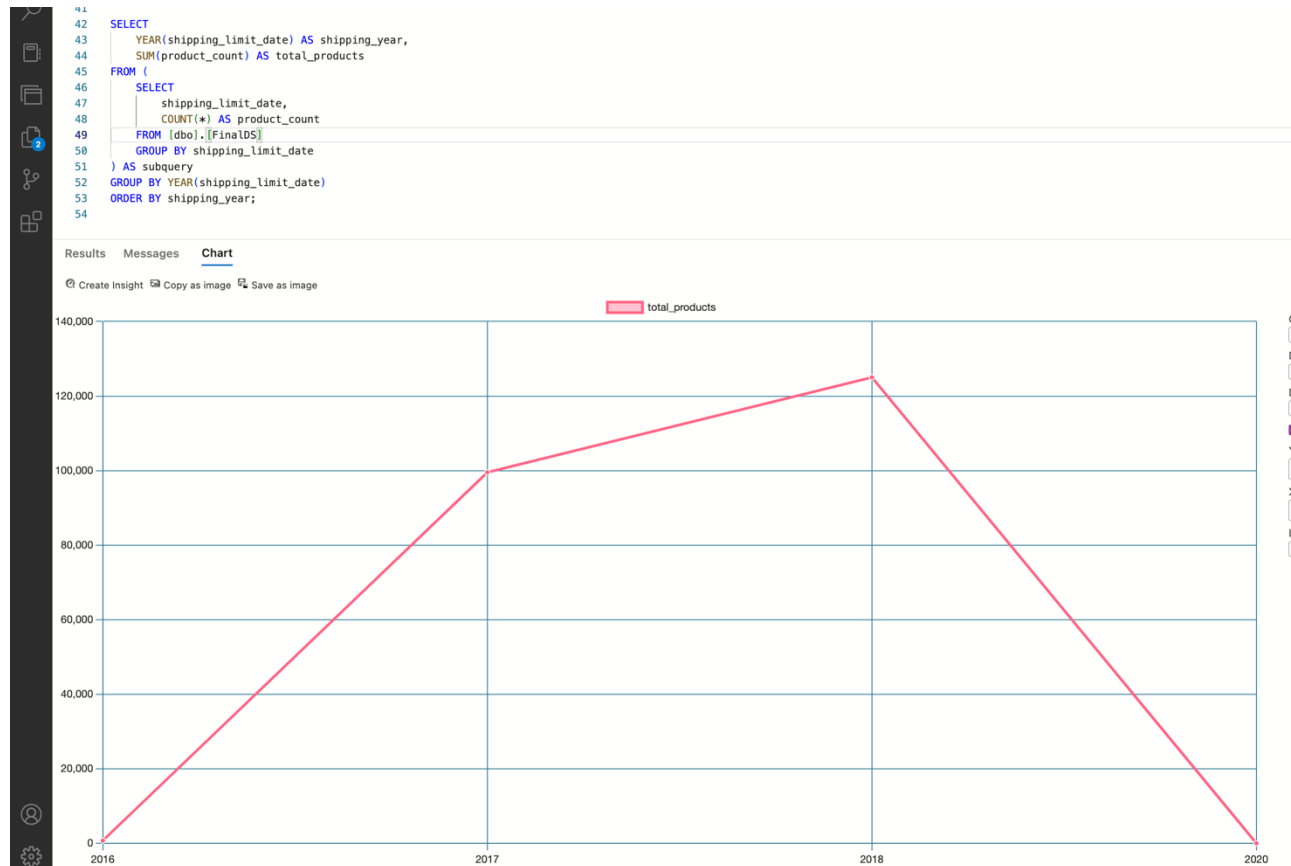
Ln 21, Col 1 Spaces: 4 UTF-8 CRLF 611 rows MSSQL Average: 15620.333 Count: 12 Sum: 93770 00:00:00 porjectserver.database.windows.net : Project

```

42  SELECT
43      YEAR(shipping_limit_date) AS shipping_year,
44      SUM(product_count) AS total_products
45  FROM (
46      SELECT
47          shipping_limit_date,
48          COUNT(*) AS product_count
49      FROM [dbo].[FinalDS]
50      GROUP BY shipping_limit_date
51  ) AS subquery
52  GROUP BY YEAR(shipping_limit_date)
53  ORDER BY shipping_year;
54

```

	Results	Messages	Chart
	shipping_year	total_products	
1	2016	740	
2	2017	99530	
3	2018	125022	
4	2020	8	



KEY CHALLENGES FACED

1. The challenge we faced was with datetime zone in Power Query within Azure Data Factory (ADF) which is a common hurdle when integrating various cloud services. Handling datetime zones, especially in data transformation processes, can be complex due to different time zone formats or inconsistencies across data sources.
2. In cases where the dataset lacks certain qualities required for predictive modeling, such as insufficient features, data quality issues, or inadequate sample size, predictive modeling might not yield reliable results. When faced with a dataset that isn't suitable for predictive modeling, leveraging SQL queries for analytics purposes becomes a valuable alternative enabling the extraction of insights critical for informed decision-making and strategy formulation.

Queries [2]

ADFRResource [1]

UserQuery

UserQuery : Expression.Error: The Power Query Spark Runtime does not support the function DateTimeZone.From.

✕

✓

fx

{{"Count", each Table.RowCount(_), Int64.Type}}

	1.2 Year	1 ² 3 Count
1	2017	49765
2	2018	62511
3	2016	370
4	2020	4

Queries [2]

UserQuery : Expression.Error: The Power Query Spark Runtime does not support the function DateTimeZone.From.

↶

↷

📄

📖

— 100% +

↶ ↷

^

^

✕

✓

fx

Table.TransformColumnTypes(Table.AddColumn

	AB shipping_limit_date	AB price	AB freight_value	Parse
1	2017-09-19 09:45:35	58.90	13.29	9/19/2017

Query setting >

Properties

Name

UserQuery

Applied steps

S...

In...

1 warning

Completed (8.42 s)

Columns: 8

Rows: 99+

Step

Settings

Parameters

CONCLUSION

In conclusion, this project marks a strategic initiative to harness the power of Azure cloud services—Azure Blob Storage, Azure SQL Database, and Azure Data Factory—to revolutionize the analysis and prediction of customer behavior for Olist, a leading Brazilian e-commerce platform. By seamlessly integrating and transforming key datasets within Azure's cloud

pg. 15

environment, we have laid the foundation for a robust and scalable solution.

Azure Blob Storage facilitates efficient data storage and management, while Azure SQL Database ensures reliable integration and retrieval of structured data. The innovative use of Azure Data Studio adds a layer of sophistication, enabling us to discover trends and patterns in customer behavior, thereby providing Olist with actionable insights to refine sales strategies.

This project not only demonstrates the prowess of Azure's cloud ecosystem but also underscores its practical applications in enhancing e-commerce analytics. By empowering Olist with advanced predictive models and analytical capabilities, the project aligns with the objective of staying ahead in the competitive e-commerce landscape. In essence, the successful implementation of this project equips Olist with the tools needed to optimize customer engagement and elevate sales strategies to new height.

REFERENCES

- <https://learn.microsoft.com/en-us/azure/?product=popular>
- <https://olist.com/>
- <https://azure.microsoft.com/en-us/blog/>
- <https://azure.microsoft.com/en-us/products/azure-sql/database>