

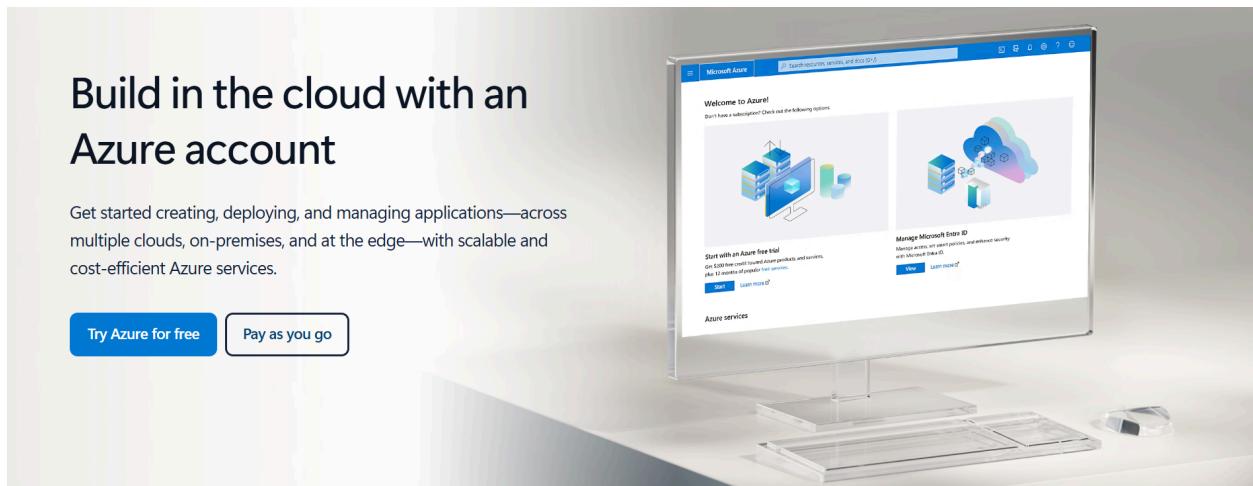
Project overview:

2 sources of data are chosen here. One is an HTTP (saved on GitHub and getting the same via HTTPS requests). But in industry, we will be connecting to a sophisticated data source so that the data won't be lying on GitHub or anywhere; it will primarily be in a database. Second is the SQL table source.

So, once we will have the data in different sources, we will be ingesting it into Azure Data Factory and we will be moving it into a Data Lake solution. **Azure Data Factory** is a tool which helps us to extract data from multiple sources and then we can do some minor transformations as well. Once, we have the data, we will be sending the raw data to a datalake ADLS (**Azure Data Lake Storage**).

Databricks is something which we will be using to transform the data before sending to ADLS. Now, this transformed data will be reading via Azure Synapse.

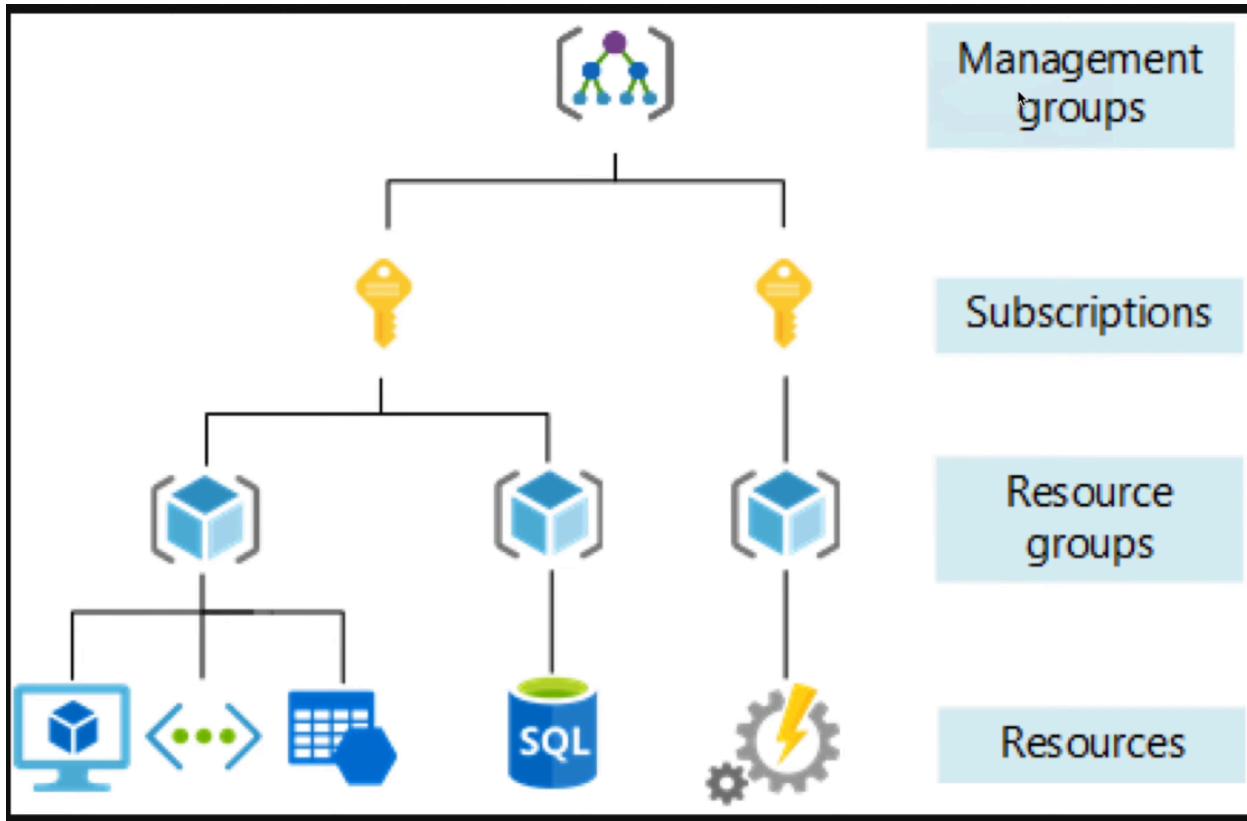
Medallion Architecture is used here which is a data lakehouse design pattern that organizes data into three increasingly refined layers—**Bronze**, **Silver**, and **Gold**—to improve data quality and consistency for analytics and business intelligence. Raw data lands in the Bronze layer, is cleaned and validated in the Silver layer, and finally transformed into enriched, aggregated data ready for reporting and machine learning in the Gold layer. We will also be reading the data from a NoSQL database inside the Azure Databricks. This is just to make sure the complexity of project increases.



Why Cloud?

So, MNCs like Amazon, GCP, Microsoft have their own data centres. The basic funda is that they will give us cloud service and they are going to charge us for it.

Azure account structure:



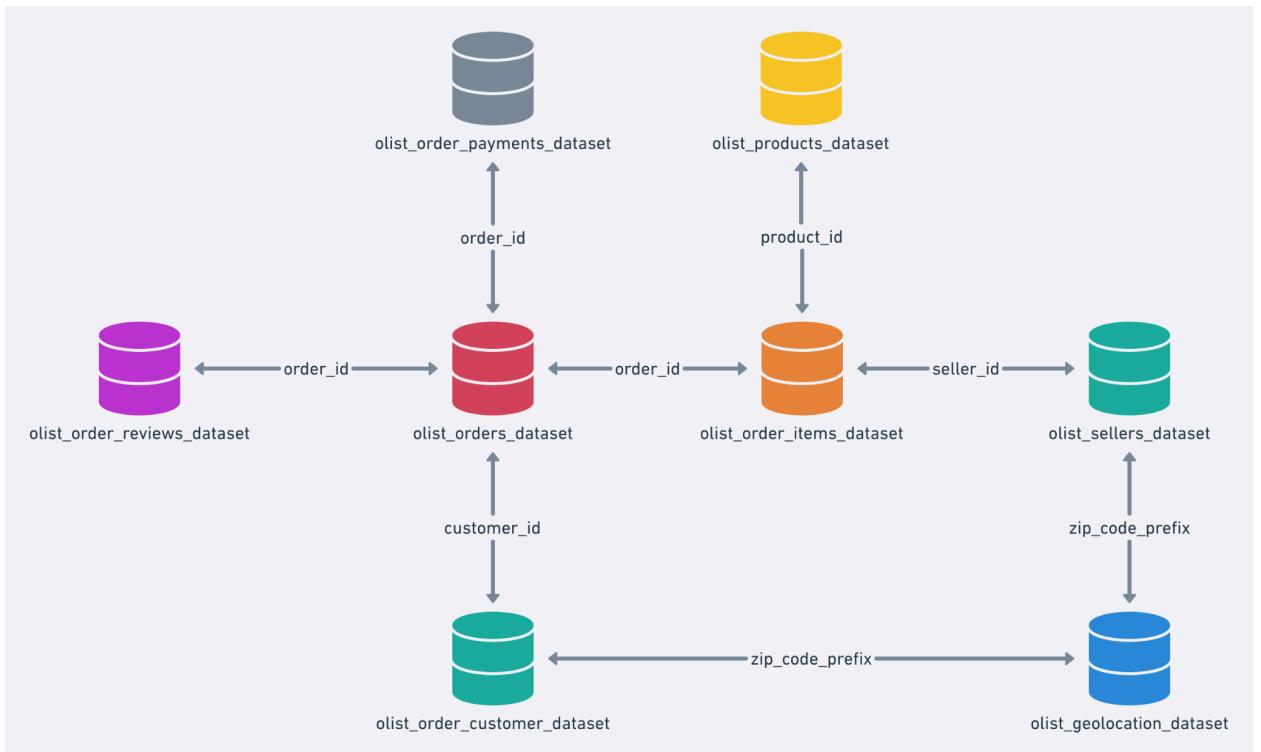
Management groups are like Principal of a school and **Subscriptions** can be like for different classes. And, classes has some **Resource Groups**. And, under it, we have **Resources**.

Steps:

1. Created a new subscription “Azure for Students”
2. Under it, created a new “Resource group”
3. Nearer the Azure Data Centre region, better will be the ping and project will work in a more better way
4. Resource group is created

The screenshot shows the Microsoft Azure Resource Groups page. At the top, there's a search bar and a Copilot button. Below the header, it says "Resource groups" and "Northwestern University". There are buttons for "Create", "Manage view", "Refresh", "Export to CSV", "Open query", and "Assign tags". A message at the top says "You are viewing a new version of Browse experience. Click here to access the old experience." Below that, there's a filter bar with "Subscription equals all", "Location equals all", and "Add filter". The main table has columns for "Name", "Subscription", and "Location". One row is shown: "ecommerce_live" (Subscription: Azure for Students, Location: East US). At the bottom, it says "Showing 1 - 1 of 1. Display count: auto".

- 5.
6. Here, I am referring Azure Data Factory, Azure Data Lake, Azure Databricks, Azure Synapse Analytics as Resources: ecommerce_live
7. **Olist Dataset Schema:**



The reason of choosing this dataset is somewhat analogous to the user and have some complexities which is somehow similar to the stuff that happens in the industry or real life. Dataset link:
<https://www.kaggle.com/datasets/olistbr/brazilian-e-commerce>

<https://filess.io/>

What is [files.io](#)?

Filess.io: A Free & Scalable Database Hosting Solution for Developers! 100% Free Database Hosting Forever. DBaaS – We offer plan-based database services for developers. Free Free plan Forever Availability 99 % Availability

While working on a data pipeline project, I needed a reliable database host for data ingestion. That's when I discovered Filess.io, an excellent platform offering free forever database hosting with support for MySQL, MariaDB, PostgreSQL, and MongoDB!

- 💡 Why Choose Filess.io?
- ✓ Free Forever Plan – Get 10MB per database, support for up to 2 databases, and weekly backups.
- ✓ Global Infrastructure – Deploy databases in Europe, Asia, and America for faster access.
- ✓ Scalable & Flexible – Upgrade seamlessly with more storage, frequent backups, and security enhancements.
- ✓ Developer-Friendly – Simple UI for easy setup and management.

Here, we have created a SQL database in which we have to send the data.

The screenshot shows the Filess.io database management interface. On the left, there's a sidebar with icons for connections, tables, views, functions, and settings. The main area is titled 'CONNECTIONS' and shows a connection named 'olistproject_proudsolar on injaxe.h.files'. Below this, under 'TABLES, VIEWS, FUNCTIONS', it says 'Database olistproject_proudsolar is empty or structure is not loaded, press Refresh button to reload structure'. There are buttons for 'Refresh' and 'New table'. The central part of the screen is titled 'Table #1' and shows a table structure with 'Columns (0)' and a 'Add new' button. It also includes sections for 'Primary key', 'Indexes (0)', 'Unique constraints (0)', and 'Foreign keys (0)'. At the bottom, there are buttons for 'Save', 'Reset changes', 'Add column', 'Add index', and 'Open data'. The status bar at the bottom indicates the connection is 'Connected' to 'MySQL 5.7.38-41' and was 'a few seconds ago'.

Home > Data factories >

Create Data Factory

One-click to create Data factory with sample pipeline and datasets. [Try it](#)

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *	Azure for Students
Resource group *	ecommerce_live
Create new	

Instance details

Name *	olist-ecom-data-factory-sid
Region *	East US
Version *	V2

8.

Home >

Microsoft.DataFactory-20250915013157 | Overview

Deployment

Deployment is in progress

Deployment name : Microsoft.DataFactory-2025091... Start time : 9/15/2025, 1:36:25 AM
 Subscription : Azure for Students Correlation ID : 1bd37a77-a131-4171-844a-db...

Resource group : ecommerce_live

[Deployment details](#)

Resource	Type	Status	Operator
There are no resources to display.			

[Give feedback](#)

[Tell us about your experience with deployment](#)

... Deployment in progress...
 Deployment to resource group 'ecommerce_live' is in progress.

Microsoft Defender for Cloud
 Secure your apps and infrastructure
[Go to Microsoft Defender for Cloud >](#)

Free Microsoft tutorials
[Start learning today >](#)

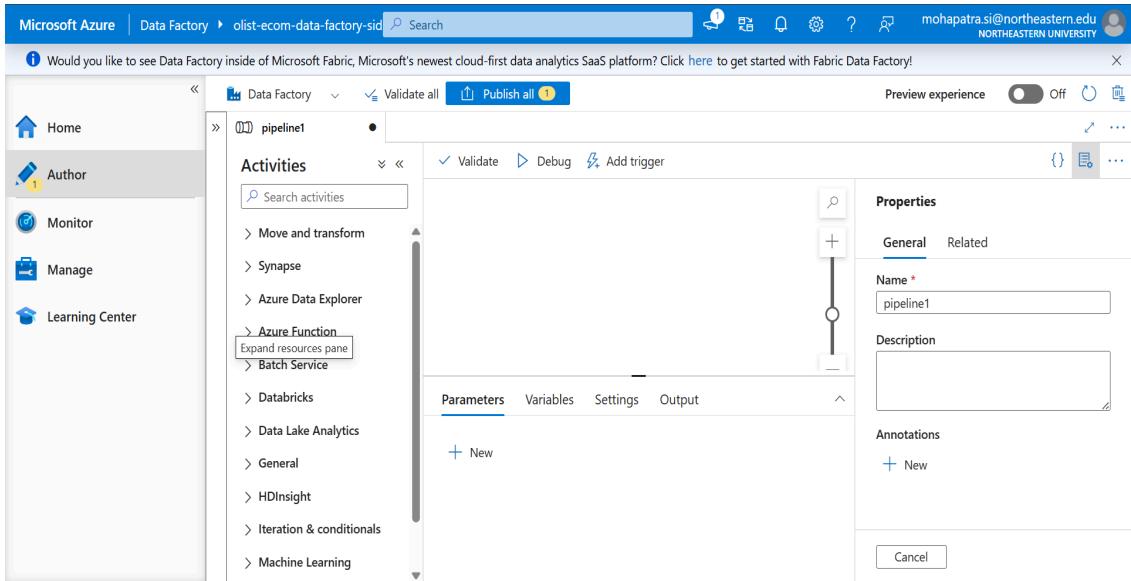
Work with an expert
 Azure experts are service provider partners who can help manage your assets on Azure and be your first line of support.
[Find an Azure expert >](#)

9.

Azure Data Factory is a cloud tool that collects data from different sources, organizes it and moves it where we need it: Data Collection (Ingestion), Data Transformation (Cleaning and organizing), and Data movement. (**Why do we use data movement?** Because it helps us to build perfect pipelines.

Click on Azure Data Factory (which we created) -> Launch Data Studio -> Go to Author -> Click on '+' icon -> Create a pipeline

It will look something like this:



Azure Data Factory is like a data delivery and transformation service. It helps us move data from one place to another and prepare it for use, often combining multiple sources along the way. We can think of it as a factory assembly line for data.

Azure Data Factory is a managed service that automates the movement, transformation, and integration of data, much like a factory assembly line prepares raw materials into a finished product.

While creating a storage account, there are few things to keep in mind and consider. There is something called **Hierarchical Namespace**, which we have to manually enable it. Just because ADF is a very big data lake where we can store any kind of data, be it structured or unstructured. But **Microsoft** uses a **Blob** format. Whereas **Amazon** uses **S3** format and Hadoop uses **HDFS** format. So, basically, if it is disabled, it means we cannot create folders inside the container, which means that the container will be of no use. It's just like we have a new brand laptop, but we are unable to create a folder. Obviously, nobody would like to buy that laptop.

Azure Data Lake Storage (ADLS) Gen 2: It's a secure cloud platform that provides scalable, cost-effective storage for big data analytics. It stores a massive amount of data, keeps data organized, fast and efficient as well. And, it is also built on **Azure Blob Storage**. We can store files, logs, images, IoT data, or structured/unstructured data. Think of it like **Google Drive** for huge amounts of data, but made for businesses that deal with terabytes (or even petabytes) of information. Azure uses (Advanced Encryption Standard) AES-256 encryption, one of the strongest security standards used by banks and governments.

10. Creating a container in ADLS:

Name	Last modified	Anonymous access level	Lease state
\$logs	9/25/2025, 1:19:37 AM	Private	Available
olistdata	9/25/2025, 2:00:12 AM	Private	Available

I created a directory inside olistdata storage named bronze followed by silver followed by gold.

The screenshot shows the Microsoft Azure Storage Explorer interface. The left sidebar has a 'Container' view selected. The main area displays a table of blobs in the 'olistdata' container. The table columns are Name, Last modified, Access tier, Blob type, Size, and Lease state. There are three items listed:

Name	Last modified	Access tier	Blob type	Size	Lease state
bronze	9/25/2025, 2:02:33 AM				...
gold	9/25/2025, 2:03:31 AM				...
silver	9/25/2025, 2:03:24 AM				...

Here, basically I am trying to follow the Medallion architecture which is followed in most of the industries.

So, instead of doing manually, if we have 100s or 1000s of csv data in github, we can put the base url as <https://raw.githubusercontent.com/SIDDHARTH107/>

And relative url as

Real-Time-E-Commerce-Data-Processing-and-Analytics-using-Medallion-architecture-on-Azure-Cloud
[/refs/heads/main/olist_customers_dataset.csv](https://raw.githubusercontent.com/SIDDHARTH107/main/refs/heads/main/olist_customers_dataset.csv)

This relative url is something that we can configure again and again when we have tons of csv data in github.

11. Now after Data Ingestion from source to ADF, I worked on Azure Data Lake Studio Gen 2 to transfer Raw Data from ADF to ADLS Gen 2. Sometimes, there will be a question why are we not transferring data directly from Github to ADLS Gen 2? Why **ADF** is used in between?(**Azure Data Factory** is like a cloud-based “data logistics & control center” that securely connects, moves, and orchestrates data from many sources into your storage or analytics systems—on time, at scale, with monitoring built-in.)

Key Benefits of ADF

Feature	Why It's Important
Many Connectors	Easily connect to SQL, APIs, SaaS, on-prem
Scheduling	Automate runs (hourly, daily, etc.)
Monitoring	Dashboard for pipeline status, alerts, logs
Security	Secure credentials, managed identities, encryption
Scalability	Serverless – no infra to manage
Integration	Works smoothly with ADLS, Synapse, Databricks

ADF isn't just "copying" — it's "coordinating, scheduling, securing, and monitoring" our data pipelines at scale. Without it, we would be writing and maintaining a lot of custom scripts. **ADF** is basically the **ETL** engine at the front door.

■ What Azure Data Factory Does That Direct Copy Doesn't

1. Orchestration / Scheduling

- ADF lets you **schedule** when and how data moves.
- Ex: Pull GitHub data at 2 a.m. daily, load new SQL table snapshots every hour.
- Without ADF, you'd have to write and maintain your own code/scripts to do this.

2. Data from Multiple Sources

- ADF can ingest from dozens of sources at once — APIs, FTP, on-prem SQL, SaaS apps.
- You don't have to manually connect each one to ADLS Gen2.

3. Pre-processing (Light Transformations)

- ADF can filter, map columns, join simple files before landing into ADLS Gen2.
- Example: only pull rows from GitHub CSV where "status=active" before storage.

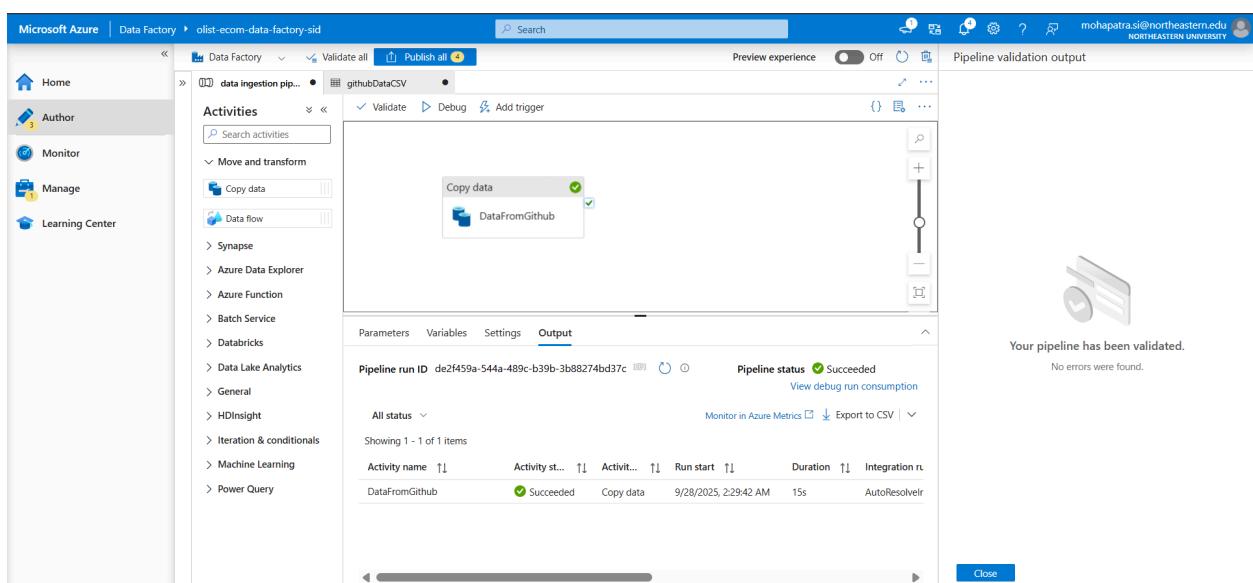
4. Monitoring & Error Handling

- ADF has built-in dashboards, retry logic, and alerts.
- Without it, if a job fails, you'd only know after the fact (or not at all).

5. Security & Governance

- Secure credential management, data lineage, and compliance baked in.

12.



Now, we can see, the pipeline has been created connecting from Azure Data Factory to Azure Data Lake Studio Gen 2. And, just because both service belongs to Microsoft Azure, the connection took few time and was done seamlessly. While building the pipeline, I preferred "Validate" and "Validate copy runtime" and after clicking "Debug", it shows like this:

The screenshot shows the Microsoft Azure Data Factory interface. On the left, the navigation menu includes Home, Author, Monitor, Manage, and Learning Center. The main workspace displays a pipeline named 'data ingestion pip...' with a single activity: 'githubDataCSV'. Below the pipeline, the 'Activities' section lists various options like Move and transform, Copy data, and Data flow. The 'Output' tab is selected, showing details of the most recent pipeline run. The run ID is 'de2f459a-544a-489c-b39b-3b88274bd37c', status is 'Succeeded', and it occurred on 9/28/2025 at 2:29:42 AM. The activity 'DataFromGitHub' is listed with a duration of 15s and ran on the 'AutoResolveIntegrationRuntime (East US)'.

The screenshot shows the Microsoft Azure Storage Explorer interface. It displays a container named 'olistdata' under 'Recent'. The 'Overview' tab is selected, showing 1 item named 'bronze'. A 'Change tier' dialog box is open over the storage interface, prompting the user to optimize storage costs by selecting an access tier. The current selection is 'Hot (Inferred)'. Other options include 'Hot (Inferred)', 'Cool', 'Cold', and 'Archive'. At the bottom of the dialog are 'Save', 'Cancel', and 'Give feedback' buttons.

Tier	What it is (plain English)	How fast can you get data	Cost to store	Cost to read	Best for...
Hot	“Always in use” area	Fast (milliseconds)	Highest	Lowest	Files/data you use or update very often (daily/hourly)

Cool	“Sometimes in use” area	Fast (milliseconds)	Cheaper than Hot	A bit more expensive to read	Files/data you use once or twice a month
Cold	“Rarely in use” area (newer)	Fast (milliseconds) but higher transaction cost	Even cheaper than Cool	Even more expensive to read	Files/data you might check only a few times a year
Archive	“Deep storage / vault” area	Slow (can take hours to rehydrate before you can read)	Cheapest	Highest retrieval cost	Old data you must keep but almost never use (e.g., legal or compliance records)

13. Now, I will create a **link** service so that I don't keep on entering this everytime. And, this should be done **parametrically**, (which means we will provide the base url as base url is same for all the datasets in github and the relative url will be provide as a parameter which changes from dataset to dataset) so that there are 1000s of csv file, we can do it smartly.

The screenshot shows the Microsoft Azure Data Factory interface. The left sidebar has 'Manage' selected. The main area is titled 'Linked services' and shows one item: 'httpGitHubLinkedService' (Type: HTTP). The URL is listed as 'http://github.com'.

Name	Type	Related	Annotations
httpGitHubLinkedService	HTTP	0	

(Above figure: We can see, I have added 1 linked service i.e. github as a data source)

Now, in order to connect my local postgres with cloud based Microsoft server which is Azure Data Factory, I need **Self Hosted Integration Runtime**:

The screenshot shows the 'Integration runtimes' configuration page in the Azure Data Factory portal. On the left, a navigation sidebar lists various settings like General, Connections, and Security. The main area displays a table of existing integration runtimes, with one entry named 'AutoResolveIntegrationRuntime'. To the right, a panel titled 'Integration runtime setup' lists three options: 'Azure, Self-Hosted', 'Azure-SSIS', and 'Airflow (Preview)'. The 'Azure, Self-Hosted' option is highlighted.

Here, after clicking **Azure Self-Hosted**, I got this:

This screenshot shows the 'Integration runtime setup' page for the 'Azure Self-Hosted' runtime. It includes sections for 'Network environment' (describing Azure as a managed serverless compute) and 'External Resources' (listing 'Linked Self-Hosted'). At the bottom, there are 'Continue', 'Back', and 'Cancel' buttons.

Then **Self-Hosted** is clicked,

The screenshot shows the Microsoft Azure Data Factory interface. On the left, there's a navigation menu with options like Home, Author, Monitor, Manage, and Learning Center. The main area is titled 'Integration runtimes' and contains a table showing one item: 'AutoResolveIntegrationRuntime' of type 'Azure' and sub-type 'Public'. A right-hand panel is titled 'Integration runtime setup' and shows a form for creating a new runtime. The 'Name' field is populated with 'MyLocal-Integration-Runtime-PostgreSQL'. There's a 'Description' field with placeholder text 'Enter description here...' and a 'Type' field set to 'Self-Hosted'. At the bottom of this panel are 'Create', 'Back', and 'Cancel' buttons.

Then, I clicked **Express Setup** to install the Integration Runtime.

The screenshot shows the 'Microsoft Integration Runtime Express Setup' window. The title bar says 'Microsoft Integration Runtime Express Setup'. The main content area has a progress bar at the top. Below it, the text reads 'Installing and registering the Integration Runtime (Self-hosted) node.' A green checkmark icon followed by the text 'Loading configuration' is displayed. The progress bar is at 23%. Below the progress bar, the text 'Downloading Integration Runtime (Self-hosted)' is shown with a progress bar indicating 23% completion. Underneath, it says 'File of size: 1.08 GB downloaded: 246.00 MB at download speed: 14.37 MB/s'. The next step is 'Installing Integration Runtime (Self-hosted)'. The final step is 'Registering Integration Runtime (Self-hosted)'. In the bottom right corner, there is a 'Close' button.

The screenshot shows the Microsoft Azure Data Factory interface. On the left, there's a navigation sidebar with options like Home, Author, Monitor, Manage, and Learning Center. The main area is titled 'Linked services' under the 'Data Factory' section. It lists two entries: 'filesSQLDB' (PostgreSQL) and 'httpGitHubLinkedService' (HTTP). A success message box in the top right corner says 'Successfully created' and 'Successfully created filesSQLDB (Linked service)'. The URL in the browser is 'Microsoft Azure | Data Factory > olist-ecom-data-factory-sid'.

Now, we can see my **PostgreSQL** is connected to my **Azure Data Factory**.

We can see there are options like sequential and batch count:

The screenshot shows the Microsoft Azure Data Factory interface. The left sidebar has 'Activities' selected. In the center, a pipeline named 'data ingestion pipeline' is shown. It contains a 'ForEach' activity which contains a 'Copy data' activity. The 'Properties' panel on the right shows the 'General' tab with 'Name' set to 'data ingestion pipeline' and 'Description' empty. The 'Settings' tab is selected, showing 'Sequential' checked and 'Items' set to 'This property should be parameterized.'. The URL in the browser is 'Microsoft Azure | Data Factory > olist-ecom-data-factory-sid'.

The **Sequential** checkbox in a **ForEach** activity dictates how the loop runs. When we check this box, the **ForEach** activity will process each item in our input list one at a time, in order. The second iteration will not start until the first one has completely finished.

By default, if we leave this box unchecked, the **ForEach** activity runs in parallel, which is what "in batches" refers to. It processes multiple items simultaneously to speed up the total execution time.

Sequential Processing

Basically, it is used for control and safety. When we check the Sequential box, we are forcing the loop to operate in a strict, single-file line.

When to Use It?

Order is Critical: Need to use this when the outcome of one iteration depends on the completion of the previous one. For example, if we are processing daily sales files and day_2.csv can only be processed after day_1.csv is successfully loaded.

Avoiding Race Conditions: When multiple iterations writing to the same file or database table could corrupt the data if they run at the same time. Sequential processing ensures that only one operation is happening at any given moment.

Batches (Parallel Processing)

Basically, it is used for speed and efficiency. When the order doesn't matter, parallel processing leverages the power of the cloud to get the job done much faster. When we uncheck the Sequential box, we are allowing the loop to run on a multi-lane highway. By default, it can run up to 20 iterations at once. We can control this by setting a Batch count (e.g., setting it to 5 would mean it runs 5 iterations at a time).

When to Use It:

Maximum Performance: This is the default and most common use case. Use it when your goal is to complete the entire job as fast as possible.

Independent Tasks: It's perfect when the items in your list have no dependency on each other. For example, copying 1,000 independent product images from one storage account to another. Copying product_A.jpg has no impact on the process for product_B.jpg.

Bulk Data Ingestion: When you need to copy hundreds of separate tables from a source database. Each table copy is an independent task, and running them in parallel drastically reduces the total pipeline run time.

14. Now we can see, all the data has been ingested successfully to the bronze layer:

RECAP:

Phase 1: Setup & Connection

Established a Hybrid Cloud Connection: You first installed and configured a Self-Hosted Integration Runtime (SHIR) on your local machine. This created a secure bridge to allow Azure to access your on-premises PostgreSQL database without exposing it to the internet.

Created Linked Services: You connected Azure Data Factory to your data stores by creating two key Linked Services: one for your local PostgreSQL database (using the SHIR) and another for your Azure Data Lake Storage (the "olistdata" container).

Phase 2: Building the Automated Pipeline

Created Reusable Data Templates: You designed two generic, parameterized datasets. One for the PostgreSQL source (using a tableName parameter) and one for the CSV sink in the data lake (using a fileName parameter).

Engineered a Dynamic Pipeline: You built a master pipeline (`Ingest_All_PostgreSQL_Tables`) to automate the entire workflow:

Discovered Tables Automatically: You used a **Lookup** activity with a SQL query to dynamically fetch a list of all table names from your PostgreSQL database that matched the "`olist_%`" pattern.

Looped Through Each Table: You added a **ForEach** activity that takes the list of tables from the **Lookup** activity as its input, preparing to process each one individually.

Configured the Data Copy: Inside the **ForEach** loop, you placed a **Copy Data** activity and used the parameterized datasets to make it generic. You passed the current table name from the loop (`@item().table_name`) to both the source and sink, dynamically creating the correct file name for each table.

Phase 3: Execution & Finalization

Tested and Validated the Pipeline: You ran the entire pipeline using the **Debug** button. This executed a live test run, proving that the logic worked by successfully copying all 9 tables and their data to the "bronze" folder in your data lake.

Troubleshoot and Cleaned Up: You encountered a publishing error caused by old, incomplete datasets (`GithubDataCSV` and `DelimitedText1`). You successfully identified and deleted these unused assets to resolve the validation issue.

Saved the Final Pipeline: You successfully used the "Publish all" button to permanently save the entire pipeline structure—including the linked services, datasets, and pipeline logic—to the Azure Data Factory service.

15.

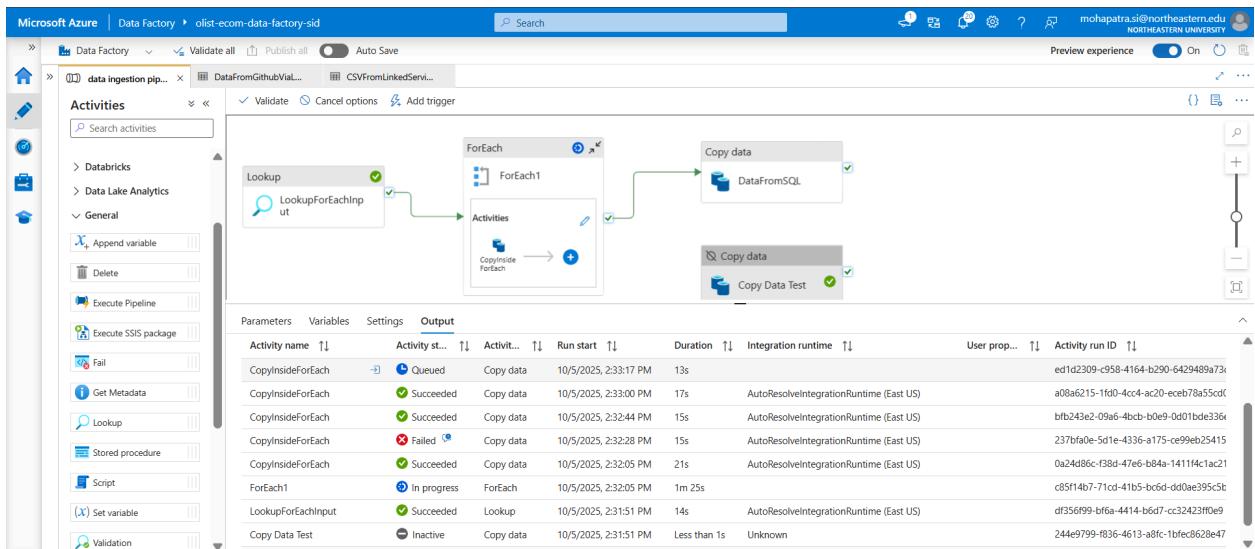
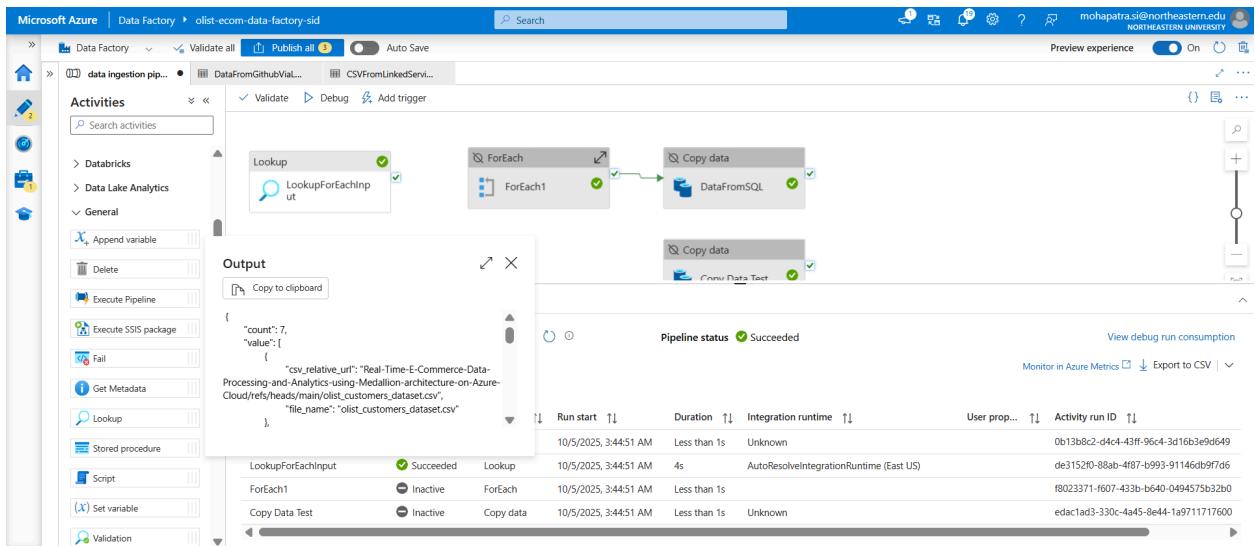
In Azure Data Factory, the **Lookup** activity is used to retrieve a small amount of data from any data source. Think of it as asking a question and getting a small piece of information back that you can use in the next steps of your pipeline.

Instead of us manually telling the pipeline which tables or files to copy one by one, the **Lookup** activity does it for you. Here is the exact process:

Get a To-Do List: The **Lookup** activity runs first. It connects to a source (like a database query or a configuration file) to get a list of items that need to be processed. In your case, it's likely getting a list of all the table names (e.g., '`olist_customers_dataset`', '`olist_orders_dataset`', etc.) from your PostgreSQL database.

Pass the List to the Next Step: The output of the **Lookup** activity is a list (an array) of these table names.

Feed the Loop: This list is then passed directly into the **ForEach** activity. The **ForEach** loop takes that list and runs the **Copy Data** activity inside it once for every single item on the list.

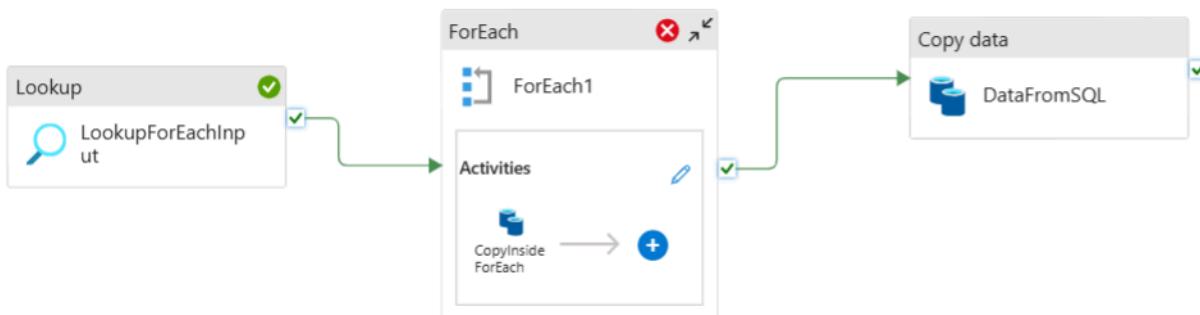


The screenshot shows the Microsoft Azure Storage Explorer interface. The left sidebar shows the 'Home' section with the 'Container' tab selected. The main area displays the contents of the 'olistdata' container, which contains six CSV files: 'olist_customers_dataset.csv', 'olist_order_items_dataset.csv', 'olist_order_reviews_dataset.csv', 'olist_orders_dataset.csv', 'olist_products_dataset.csv', and 'olist_sellers_dataset.csv'. These files are categorized under the 'bronze' blob type. The table includes columns for Name, Last modified, Access tier, Blob type, Size, and Lease state.

The screenshot shows the Microsoft Azure Data Factory Pipeline expression builder. It displays a pipeline diagram with a 'Lookup' activity followed by a 'ForEach' activity. Inside the 'ForEach' activity, there is an 'Activities' section containing a 'CopyInsideForEach' activity. A 'Copy data' activity is connected to the 'CopyInsideForEach' activity. The 'Copy data' activity is configured to copy data from 'DataFromSQL' to 'DataFromStorage'. On the right side of the screen, the 'Pipeline expression builder' pane is open, showing the expression `@activity('LookupForEachInput').output.value`. Below the expression, there is a list of available variables and functions, and buttons for OK and Cancel.

The value represents the JSON structure that we got from the raw link of GitHub, not the count.

What is LOOKUP, ForEach, Copy Data?



This combination of activities is a powerful and common pattern in data engineering. Think of it like a smart system for completing a list of chores.

Lookup (LookupForEachInput):

What it is: The "List Maker." The Lookup activity's job is to read a small amount of data from a source.

Its Role Here: It runs first to create a "to-do list." It queries a source (like a database control table or a configuration file) to get a list of all the tables or files that need to be copied. Its output is simply this list.

ForEach (ForEach1):

What it is: The "Worker." The ForEach is a loop. It takes the "to-do list" generated by the Lookup activity and performs a set of actions for every single item on that list.

Its Role Here: It examines the list of tables provided by the Lookup and states, "Okay, for table 1, I will do these tasks. Now for table 2, I will do the same tasks..." and so on, until the list is complete.

Copy Data (CopyInsideForEach):

What it is: The "Specific Chore." This is the workhorse activity in Azure Data Factory that moves data from a source to a sink.

Its Role Here: It is the task inside the ForEach loop. For every table name the ForEach loop processes, the Copy Data activity is the specific action that actually connects to that table and moves its data to the destination.

In summary: The Lookup finds out what to do, the ForEach manages doing it for everything on the list, and the Copy Data is the actual work being done for each item.

Can we assume this as a Production-Grade Pipeline? If yes, then why?

Yes, absolutely. The architectural pattern you have built (Lookup -> ForEach -> Copy) is not just a simple pipeline; it is a metadata-driven framework, which is a best practice and definitely considered production-grade.

While a full production deployment would also include robust error handling, logging, and monitoring, the core design pattern itself is excellent for the following reasons:

It is Dynamic and Data-Driven: The pipeline is not static or hardcoded. It intelligently discovers the work it needs to do at runtime by using the Lookup activity. If a new table is added to your source database that matches the Lookup's query, the pipeline will automatically include it in the next run without you ever needing to modify or republish the pipeline.

It is Highly Scalable: This pattern works just as well for 10 tables as it does for 1,000. Because you don't have to manually create a Copy data activity for each table, the solution scales effortlessly as your data sources grow.

It is Easy to Maintain: Imagine you needed to change how the data is copied (e.g., add a new mapping rule or change a setting). You only have to modify the single Copy Data activity inside the ForEach loop, and that change is instantly applied to all 1,000 tables. In a static design, you would have to manually edit 1,000 separate activities, which is slow and prone to errors.

It is Efficient: The ForEach activity can be configured to run in parallel. This means it can copy multiple tables simultaneously (e.g., 20 at a time), drastically reducing the total time it takes to ingest all your data compared to processing them one by one.

NOW THE BRONZE LEVEL OF WORK IS COMPLETED IN MEDALLION ARCHITECTURE.

NEXT STEPS: We need to read the bronze layer data to **Azure Databricks** and read some data from **MongoDB** for enrichment & transformations, and we will see how we can transform the data in the **Silver** layer.

DATABRICKS:

Databricks is actually a third-party software. All major cloud services provide Databricks on top of their platform.

The screenshot shows the Azure portal's deployment overview for a workspace named "ecommerce_live_olist-spark-workspace". The deployment status is "Deployment is in progress". Key details include:

- Deployment name: ecommerce_live_olist-spark-workspace
- Subscription: Azure for Students
- Resource group: ecommerce_live
- Start time: 10/5/2025, 5:40:12 PM
- Correlation ID: 6bd891ba-2373-479e-9a22-e8c583ef5a3e

On the left, there's a navigation bar with "Overview", "Inputs", "Outputs", and "Template" options. A "Deployment details" section is expanded, showing a table with columns: Resource, Type, Status, and Operation details. The table body contains the message: "There are no resources to display." On the right side, there are promotional links for Microsoft Defender for Cloud, Microsoft tutorials, and working with experts.

The screenshot shows the Microsoft Azure portal's deployment overview page. The deployment name is 'ecommerce_live_olist-spark-workspace'. The status is 'Your deployment is complete'. Deployment details include: Deployment name: ecommerce_live_olist-spark-workspace, Subscription: Azure for Students, Resource group: ecommerce_live. The start time was 10/5/2025, 5:40:13 PM, and the Correlation ID is 6bd891ba-2373-479e-9a22-e8c583ef5a3e. On the right side, there are promotional cards for Cost management, Microsoft Defender for Cloud, Free Microsoft tutorials, and Work with an expert.

You can see I have deployed my Databricks workspace in the Azure cloud service.

AZURE DATABRICKS:

- Spark-powered cloud service
- Integrated with Azure
- Handles Big Data easily
- Great for Machine Learning

AZURE DATABRICKS WORKFLOW:

- Reading the data from ADLS Gen 2
- Performing Basic Transformation like cleaning, renaming, & filtering
- Using the JOIN operation to integrate multiple datasets
- Enriching the data via MongoDB
- Performing aggregation and deriving insights
- Writing final data back to ADLS Gen 2

What is a Compute in Databricks?

Compute is just a hardware that we will ask Azure to give to us. In **Databricks** (and most cloud platforms), “compute” means the actual processing power — the virtual machines (VMs), CPUs, RAM, and sometimes GPUs — that do the work. In **Databricks**, **compute** refers to the clusters of virtual machines that actually process your data and run your notebooks or jobs. Storage just holds the data, while compute does the heavy lifting. **Databricks** lets us scale this compute up or down automatically, so we only pay for what we use.

Azure Databricks Workflow which I am planning to implement:

- Reading the data from ADLS Gen 2
- Performing basic transformations like cleaning, renaming, & filtering
- Using the JOIN operation to integrate multiple datasets
- Enriching or transferring the data via MongoDB

- Performing aggregations & deriving insights
- Writing final data back to ADLS Gen 2

Step 1: Here, 1st we need to connect our Azure Databricks to Azure Data Lake Studio Gen2. ([Link](#)) So, to make the connection, we need some key or permission here to maintain the equality I mean, both are Microsoft Azure services so both of them get connected inside Microsoft and no third party can enter.

Step 7: Connect to Azure Data Lake Storage using python

You can now securely access data in the Azure storage account using OAuth 2.0 with your Microsoft Entra ID application service principal for authentication from an Azure Databricks notebook.

1. Navigate to your Azure Databricks workspace and create a new python notebook.
2. Run the following python code, with the replacements below, to connect to Azure Data Lake Storage.

```
Python Copy
service_credential = dbutils.secrets.get(scope="<scope>",key="<service-credential-key>")

spark.conf.set("fs.azure.account.auth.type.<storage-account>.dfs.core.windows.net", "OAuth")
spark.conf.set("fs.azure.account.oauth.provider.type.<storage-account>.dfs.core.windows.net", "org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider")
spark.conf.set("fs.azure.account.oauth2.client.id.<storage-account>.dfs.core.windows.net", "<client-id>")
spark.conf.set("fs.azure.account.oauth2.client.secret.<storage-account>.dfs.core.windows.net", "<client-secret>")
spark.conf.set("fs.azure.account.oauth2.client.endpoint.<storage-account>.dfs.core.windows.net", "<client-endpoint>")
```

After selecting App Registrations in Azure search, This was the error I was getting when trying to connect Databricks with my ADLS Gen2.

What This Error Means

The error "You do not have access" with Error code: 401 means "Unauthorized." It's a permissions issue. I was successfully logged into Azure, but my specific account (mohepatra.si@northeastern.edu) does not have the administrative rights to view or manage the "App registrations" section.

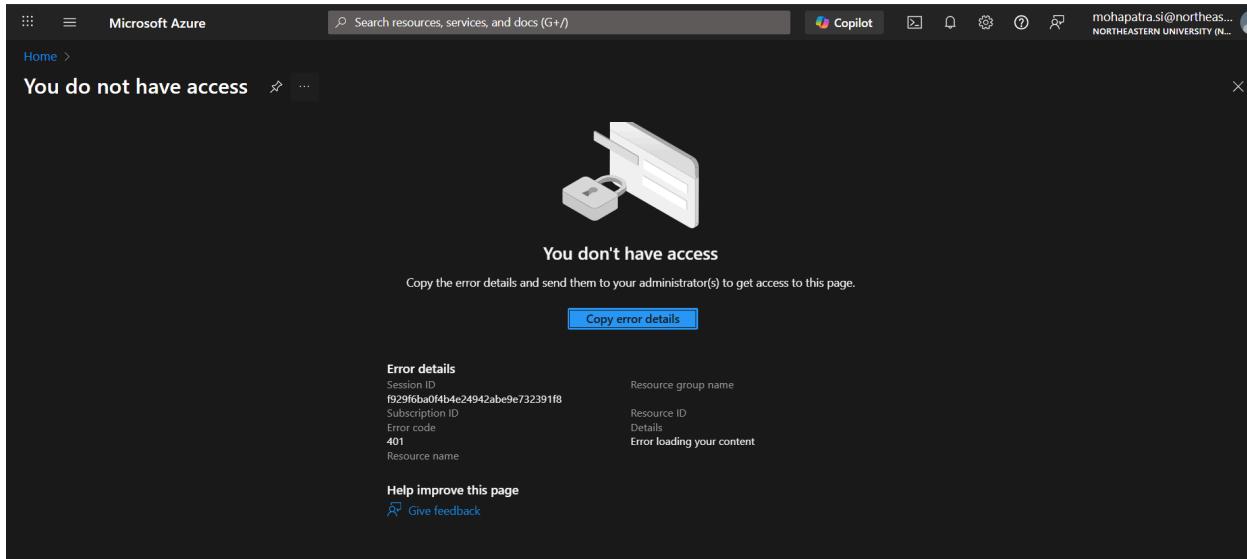
Why I didn't Have Access (The University Analogy)

Think of your university's entire Microsoft Azure setup as a large office building.

The Building (Azure Active Directory/Tenant): The entire building is managed by your university's IT administrators. They control the main entrance, the security office, and who gets what keys. The "App registrations" page is like the building's main security office. It's a central, high-level area.

Your Office (Your Azure for Students Subscription): The university has given you your own private office inside the building—this is your "Azure for Students" subscription. Inside your office, you are the owner. You can set up furniture (like Azure Data Factory, Databricks, etc.) and do whatever you want.

The error we are seeing is because our keycard (our student account) lets us into our own office, but it doesn't grant us access to the main security office for the entire building. This is a security measure to prevent students from viewing or changing administrative settings that affect the whole university.



However, the tutorial I was following demonstrated the most common and robust method for production environments: creating a **Service Principal**. A **Service Principal** is like creating a robot user with its own specific permissions. However, as you discovered, creating these "robot users" requires high-level administrative access to the "security office" (App Registrations), which I don't have with a student account.

There is another, simpler method that is perfect for this situation and still keeps the connection secure within the Microsoft cloud. I will use a method called **mounting with an Account Key**.

Instead of creating a new "**robot user**," you will temporarily use your own powerful key (the Storage Account Access Key) to establish a permanent link (a mount point) between Databricks and your Data Lake.

What I did here?

Think of it like giving your workshop (Databricks) a secure key to your storage unit (ADLS).

We got the Master Key: I went to my Azure storage account and copied its Access Key (key1). This key is like a master password that grants full access to all the data inside the service.

We Used a **Secure Lockbox**: Instead of pasting this sensitive key directly into your code (which is insecure), we created a **Databricks Widget** outside of the cell. This widget is like a secure lockbox at the top of my notebook. I pasted the key into this box, keeping it separate from my saved code.

We Built a **Secure Bridge**: Finally, I ran the **dbutils.fs.mount()** command. This command did two things:

- It securely took the key from the lockbox.
- It used the key to build a permanent, direct link—a mount point—from Databricks to my storage container.

The result is that my olistdata container now appears as a simple folder inside Databricks at the path /mnt/olistdata. Anyone can now read and write files to our storage as if it were a local directory. 

For Widget:

```
dbutils.widgets.text("storage_key", "Paste Your Key Here", "Storage Account Key")
```

For Connection:

```
# 1. Defining my storage account and container names
storage_account_name = "olistdatastorageaccount4"
container_name = "olistdata"
mount_point = f"/mnt/{container_name}"

# 2. Securely getting the key from the text box
storage_key_from_widget = dbutils.widgets.get("storage_key")

# 3. Unmounting if it already exists
try:
    dbutils.fs.unmount(mount_point)
except:
    pass

# 4. Creating the connection using a consistent protocol
# We will use 'wasbs' and '.blob.core.windows.net' for both the source and the configuration.
source_url = f"wasbs://{container_name}@{storage_account_name}.blob.core.windows.net/"
config_key = f"fs.azure.account.key.{storage_account_name}.blob.core.windows.net" # <-- This key is now
consistent

dbutils.fs.mount(
    source = source_url,
    mount_point = mount_point,
    extra_configs = {config_key: storage_key_from_widget}
)

print(f"Success! The '{container_name}' container is now available at {mount_point}")
```

```
# 5. Verifying the connection by listing the files INSIDE the 'bronze' folder
display(dbutils.fs.ls(f"{mount_point}/bronze/"))
```

The screenshot shows a Databricks notebook interface with a table view. At the top, a message says "Success! The 'olistdata' container is now available at /mnt/olistdata". The table has four columns: path, name, size, and modificationTime. The data consists of 8 rows, each representing a CSV file from the 'bronze' folder:

	path	name	size	modificationTime
1	dbfs:/mnt/olistdata/bronze/olist_customers_dataset.csv	olist_customers_dataset.csv	9033957	1759689747000
2	dbfs:/mnt/olistdata/bronze/olist_geolocation_dataset.csv	olist_geolocation_dataset.csv	51229443	1759692571000
3	dbfs:/mnt/olistdata/bronze/olist_order_items_dataset.csv	olist_order_items_dataset.csv	15438671	1759689792000
4	dbfs:/mnt/olistdata/bronze/olist_order_payments_dataset.csv	olist_order_payments_dataset.csv	6354773	1759688872000
5	dbfs:/mnt/olistdata/bronze/olist_order_reviews_dataset.csv	olist_order_reviews_dataset.csv	14451670	1759689810000
6	dbfs:/mnt/olistdata/bronze/olist_orders_dataset.csv	olist_orders_dataset.csv	17654914	1759689825000
7	dbfs:/mnt/olistdata/bronze/olist_products_dataset.csv	olist_products_dataset.csv	2379446	1759689839000
8	dbfs:/mnt/olistdata/bronze/olist_sellers_dataset.csv	olist_sellers_dataset.csv	174703	1759689856000

8 rows | 21.43s runtime Refreshed 46 minutes ago

```
# Reading data in pyspark (all 8 csv files)
# Defining the base path
base_path = "/mnt/olistdata/bronze"

# Reading each file into its own variable
orders_df = spark.read.option("header", "true").csv(f"{base_path}/olist_orders_dataset.csv")
payments_df = spark.read.option("header", "true").csv(f"{base_path}/olist_order_payments_dataset.csv")
reviews_df = spark.read.option("header", "true").csv(f"{base_path}/olist_order_reviews_dataset.csv")
items_df = spark.read.option("header", "true").csv(f"{base_path}/olist_order_items_dataset.csv")
customers_df = spark.read.option("header", "true").csv(f"{base_path}/olist_customers_dataset.csv")
sellers_df = spark.read.option("header", "true").csv(f"{base_path}/olist_sellers_dataset.csv")
geolocation_df = spark.read.option("header", "true").csv(f"{base_path}/olist_geolocation_dataset.csv")
products_df = spark.read.option("header", "true").csv(f"{base_path}/olist_products_dataset.csv")

# Now we can display any of them
display(customers_df)
```

Step 2: # Need to deal with mongodb inside pyspark.

To analyze the product_category_translation data, we first had to move it from a private, local database to a cloud database where Databricks could access it. First, we used MongoDB Compass to export the local product_category_translation collection into a JSON file. Then, we went to MongoDB Atlas (our cloud database service) and set up a new free cluster. This involved creating a database user (with a secure password) and whitelisting all IP addresses to allow our Databricks Factory and Databricks services to connect. Next, we used Compass again to connect to our new cloud cluster and imported the JSON file into a new database called olist_data. Finally, in our Databricks notebook, we wrote Python code using the pymongo

library. This code securely connected to our Atlas cluster, read all 71 documents from the collection, and loaded them into a new Spark DataFrame named `translation_df` for analysis.

Code:

```
# 1. Importing the necessary library
from pymongo import MongoClient
import pandas as pd
```

```
# 2. Setting up my connection string
# IMPORTANT: Replacing <username>, <password>, and the cluster URL with your
actual credentials from MongoDB Atlas.

username = "abc"
password = "abc"
cluster_url = "abc" # Atlas connection string
```

```
# These are the names of the database and collection I created
db_name = "olist_data"
collection_name = "product_category_translation"
```

```
# Building the connection string
uri =
f"mongodb+srv://{{username}}:{{password}}@{{cluster_url}}/?retryWrites=true&w=majorit
y"
```

```
print("Variables are set. Ready to connect in the next cell.")
```

```
# This uses the 'uri' variable from Cell 1

print("Attempting to connect to MongoDB Atlas...")
try:
    client = MongoClient(uri)
    db = client[db_name]
    collection = db[collection_name]

    # Pinging the server to confirm the connection
    client.admin.command('ping')
    print("✅ MongoDB connection successful!")
```

```
except Exception as e:  
    print("X Connection failed. Please check your username, password, or IP  
whitelist settings.")  
    print(f"Error details: {e}")
```

```
# Cell 3: Reading the Data from the Collection  
# This uses the 'collection' variable from Cell 2  
  
print(f"Reading data from the '{collection_name}' collection...")  
  
# Finding all documents ({})) and convert them to a list  
data = list(collection.find({}))  
  
if data:  
    # Your file had 71 documents or records, so this should match i guess  
    print(f"✓ Successfully found {len(data)} documents.")  
else:  
    print("X No data found. Make sure you imported the data to your Atlas  
collection.")
```

```
# Cell 4: Converting Data to a Spark DataFrame  
# This uses the 'data' variable from Cell 3  
  
if 'data' in locals() and data:  
    print("Converting data to a Spark DataFrame...")  
  
    # Converting the list of data to a Pandas DataFrame  
    pandas_df = pd.DataFrame(data)  
  
    # Dropping the '_id' column that MongoDB automatically adds  
    if '_id' in pandas_df.columns:  
        pandas_df = pandas_df.drop(columns=['_id'])  
  
    # Converting the Pandas DataFrame into a Spark DataFrame  
    translation_df = spark.createDataFrame(pandas_df)  
  
    print("✓ Data successfully loaded into a Spark DataFrame named  
'translation_df'.")
```

```
else:  
    print("No data to convert. Please run Cell 3 first.")
```

```
# Cell 5: Checking Shape and Display the Data  
# This uses the 'translation_df' variable from Cell 4  
  
if 'translation_df' in locals():  
    # Getting the row and column counts  
    row_count = translation_df.count()  
    column_count = len(translation_df.columns)  
  
    print(f"--- Data Shape ---")  
    print(f"Number of rows: {row_count}")  
    print(f"Number of columns: {column_count}")  
    print(f"-----")  
  
    # Displaying the final DataFrame  
    display(translation_df)  
else:  
    print("No DataFrame to display. Please run Cell 4 first.")
```

Step 3: Joining the datasets.

```
# Joining Orders and Customers  
orders_customers_df = orders_df.join(  
    customers_df,  
    orders_df.customer_id == customers_df.customer_id,  
    "left"  
)  
  
# Joining the above result with Payments  
orders_payments_df = orders_customers_df.join(  
    payments_df,  
    orders_customers_df.order_id == payments_df.order_id,  
    "left"  
)  
  
# Joining Orders + Payments with Items using order_id  
orders_items_df = orders_payments_df.join(items_df, "order_id", "left")  
  
# Joining Orders + Items with Products
```

```

orders_items_products_df = orders_items_df.join(
    products_df,
    orders_items_df.product_id == products_df.product_id,
    "left"
)

final_df = orders_items_products_df.join(
    sellers_df,
    orders_items_products_df.seller_id == sellers_df.seller_id,
    "left"
)

```

(NOTE: Spark Transformations & Actions: **Transformations** are instructions or recipes you give Spark to describe what you want to do with data, but no work is actually done yet.

✿ Examples of Transformations in PySpark		
Transformation	What it does	Example
.select()	Pick specific columns	<code>df.select("name", "age")</code>
.filter()	Keep only rows matching a condition	<code>df.filter(df.age > 30)</code>
.withColumn()	Add or modify a column	<code>df.withColumn("bonus", df.salary * 0.1)</code>
.groupBy()	Group data before aggregating	<code>df.groupBy("department").sum("salary")</code>
.join()	Combine two DataFrames	<code>df1.join(df2, "id", "left")</code>

Actions are commands that instruct Spark to execute the transformations and compute the actual results.

Examples of Actions in PySpark

Action	What it does	Example
<code>.show()</code>	Displays data on screen	<code>df.show(5)</code>
<code>.count()</code>	Counts total rows	<code>df.count()</code>
<code>.collect()</code>	Returns all rows to the driver	<code>df.collect()</code>
<code>.first()</code>	Returns the first row	<code>df.first()</code>
<code>.save() / .write.csv()</code>	Saves results to storage	<code>df.write.csv("/path/output")</code>

)

Step 4: Loading the transformed data into the Silver Layer of olistdata container.

```
# Dropping duplicate columns in the final_df
# This function scans through all the columns in your DataFrame, finds any
column names that appear more than once, and removes the duplicate ones –
keeping only the first occurrence.

def remove_duplicate_columns(df):
    columns = df.columns
    seen_columns = set()
    columns_to_drop = []

    for column in columns:
        if column in seen_columns:
            columns_to_drop.append(column)
        else:
            seen_columns.add(column)

    df_cleaned = df.drop(*columns_to_drop)
    return df_cleaned

final_df = remove_duplicate_columns(final_df)

# This line is saving the final Spark DataFrame (final_df) as a Parquet file
into the Azure Databricks storage location /mnt/olistdata/silver (Silver)
folder. It's like a read_csv in pandas.
```

```
final_df.write.mode("overwrite").parquet("/mnt/olistdata/silver")
```

In the above container (olistdata), we can see that the transformed data has been stored in the **SILVER** layer.

NOW THE SILVER LEVEL OF WORK IS COMPLETED IN MEDALLION ARCHITECTURE.

NEXT STEPS: We need to work on **Synapse Analytics**. Here, we will be building the **GOLD** layer of the Medallion architecture.

Step 1: Setting up the Azure Synapse Analytics. **Azure Synapse** is a cloud-based data warehouse and analytics service that combines **data integration, enterprise data warehousing, & big data analytics**.

Sometimes, we get these kind of errors in Azure while deploying the Synapse Analytics. There are reasons for it as well.

Why do Synapse deployments fail?

Think of a Synapse workspace like opening a new mall. We're not just creating the mall—we also need parking (SQL server), security (managed identity/permissions), roads (networking), and a storage warehouse (Data Lake). If any of those can't be created in the location you selected, the entire "open the mall" action fails.

One reason for this can be:

[Region blocks / capacity limits \(my exact error\)](#)

Message: SqlServerRegionDoesNotAllowProvisioning ... Location 'eastus' is not accepting creation of new ... SQL Database servers...

What it means: In that region (e.g., East US), my subscription isn't allowed to create new Azure SQL Server resources right now—often due to capacity or subscription-type limits (e.g., Free Trial/Azure for Students).

Real example: We pick East US for Synapse, but Synapse needs to spin up a SQL server there; the region says "no new SQL servers allowed," so deployment fails.

Instead we can pick another region (e.g., East US 2, Central US, West US 2) to solve this, which is what happened in my case.

The screenshot shows the Microsoft Azure portal interface. The top navigation bar includes 'Microsoft Azure', a search bar, and user information. Below the navigation is a breadcrumb trail: 'Home > Microsoft.Azure.SynapseAnalytics-20251024011122 | Overview'. The main content area displays a deployment status message: 'Your deployment failed' with a red exclamation mark icon. It provides deployment details: Deployment name: Microsoft.Azure.SynapseAnalytics-202510... Start time: 10/24/2025, Subscription: Azure for Students, Resource group: ecommerce_live. A red warning box states: 'Workspace request validation failed, check error details for more information. Click here for details'. To the right, a 'Errors' panel is open under the 'Summary' tab. It contains an error message: 'Location 'eastus' is not accepting creation of new Windows Azure SQL Database servers for the subscription 'f6234497-0192-4d89-b2ea-a962d46f29d3' at this time. (Code: SqlServerRegionDoesNotAllowProvisioning)' and a link to 'Copy Full Error'. Below the error message is a 'WAS THIS HELPFUL?' section with a thumbs up/down icon. The 'Troubleshooting Options' section lists links for 'Common Deployment Errors', 'Check Usage and Quotas', and 'New Support Request'. At the bottom of the page, there are 'Give feedback' and 'Give Feedback' sections with links to 'Tell us about your experience with deployment' and 'Tell us about your experience with the Errors page'.

The screenshot shows the Microsoft Azure Synapse Analytics deployment overview page. It displays a success message: "Your deployment is complete". Deployment details include: Deployment name: Microsoft.Azure.SynapseAnalytics-20251024012735, Start time: 10/24/2025, 1:31:49 AM, Subscription: Azure for Students, Resource group: ecommerce_live. There are links for "Deployment details" and "Next steps" (including "Go to resource group"). A feedback section asks for deployment experience. On the right, there are promotional cards for "Cost management", "Microsoft Defender for Cloud", "Free Microsoft tutorials", and "Work with an expert".

Step 2: Now, we will work on Azure Synapse Analytics. Below is the UI of Synapse Analytics. Just like the way we gave a permission to Azure Databricks to read data from ADLS Gen2, similarly, we need to give permission to Azure Synapse Analytics to read data from ADLS Gen2.

The screenshot shows the Azure Synapse Analytics workspace interface for the "olist-synapse-workspace". The workspace name is displayed prominently. The interface includes a sidebar with icons for Home, Ingest, Explore and analyze, Visualize, and Discover more. A large central area features a 3D bar chart and network graph visualization. Three main buttons are visible: "Ingest" (Perform a one-time or scheduled data load), "Explore and analyze" (Learn how to get insights from your data), and "Visualize" (Build interactive reports with Power BI capabilities).

It enables users to ingest, prepare, manage, and serve data for business intelligence (BI) and machine learning — all within a single, serverless environment. In other words, it's like a one-stop platform where we can bring in data (ingest), clean and analyze it (transform), and build dashboards (visualize)—without leaving Azure.

Component	Description	Example
Synapse Studio	The web interface for managing all Synapse tasks — similar to a "control center" for your data.	You access it via the Azure portal or a direct URL.

SQL Pools	Provide powerful querying for structured data using SQL. Includes Dedicated SQL Pool (provisioned) and Serverless SQL Pool (on-demand).	Running <code>SELECT * FROM sales_data</code> on large datasets.
Spark Pools	Allow data exploration and transformation using PySpark, Scala, or .NET .	Cleaning and aggregating big data stored in Data Lake.
Data Integration Pipelines	Automate movement and transformation of data from multiple sources (like Azure Data Lake, Databricks, or APIs).	A daily job that loads sales data into Synapse tables.
Linked Services	Secure connections to external data sources such as Azure Storage, SQL DB, Databricks, or Power BI .	Connecting Synapse to ADLS Gen2 to read Parquet files.
Data Flows	Visual, drag-and-drop transformations (no coding).	Converting CSV data into normalized tables visually.

Type	Description	Real-Life Example
SQL Database	Uses SQL language for querying and analytics. Best for structured, relational data (tables, columns).	Storing customer transactions, order records, or product details.
Lake Database	Schema-on-read database created directly over data in ADLS. Works well for raw or semi-structured files like CSV, JSON, or Parquet.	Creating a database view directly over millions of sales CSV files stored in your data lake.

Category	Description	Example Use Case
Synapse	Run SQL scripts or notebooks directly inside Synapse.	Execute a SQL query to aggregate sales data.
Move and Transform	Copy, transform, or clean data between sources.	Copy data from ADLS → Synapse SQL database.

Azure Data Explorer	Integrate with Kusto Query Language (KQL) for log analytics.	Analyze IoT or telemetry data.
Azure Function	Trigger serverless functions during workflow.	Validate data quality before loading.
Batch Service	Execute batch compute jobs on Azure Batch.	Process multiple files in parallel.
Databricks	Run Databricks notebooks or jobs directly from Synapse.	Execute a PySpark job for feature engineering.
Data Lake Analytics	Run U-SQL jobs on large datasets.	Perform distributed queries on raw logs.
General	Simple tasks like "Wait," "Execute Web," or "Set Variable."	Pause execution for 30 seconds.
HDInsight	Integrate with Hadoop or Spark-based clusters.	Run MapReduce or Hive queries.
Iteration & Conditionals	Control logic like loops or conditional branching.	"If file exists, then load; else skip."
Machine Learning	Connect to Azure ML models for predictions.	Run a trained model to predict customer churn.

We will be reading data from **the SILVER** layer and creating our own database within Synapse Analytics.

Step 3: We will give **permissions** to **Synapse**, so we will go to **App Registration** like this.

The screenshot shows the Microsoft Azure Storage Account Access Control (IAM) page for 'olistdatastorageaccount4'. The left sidebar lists various options like Overview, Activity log, Tags, and Access Control (IAM). Under IAM, 'Data migration', 'Events', and 'Storage browser' are listed under 'Data storage', while 'Containers' and 'File shares' are under 'Partner solutions'. The main content area has tabs for 'Check access', 'Role assignments', 'Roles', 'Deny assignments', and 'Classic administrators'. The 'Check access' tab is selected, showing sections for 'My access' (with a 'View my access' button), 'Check access' (with a 'Check access' button), 'Grant access to this resource' (with a 'Learn more' link), 'View access to this resource' (with a 'View' button), and 'View deny assignments' (with a 'View' button).

Go to IAM (Identity Access Management)

The screenshot shows the 'Add role assignment' page for 'blob'. The top navigation bar includes 'Copilot can help pick a role'. The main section is titled 'Job function roles' and shows 'Privileged administrator roles'. It says 'Grant access to Azure resources based on job function, such as the ability to create virtual machines.' A search bar contains 'blob'. Below it is a table with columns: Name, Description, Type, Category, and Details. The table lists several roles:

Name	Description	Type	Category	Details
Defender CSPM Storage Data Scanner	Grants access to read blobs and files. This role is used by the data scanner of Defender CSPM.	BuiltInRole	None	View
Defender for Storage Data Scanner	Grants access to read blobs and update index tags. This role is used by the data scanner of Defender for Stor...	BuiltInRole	None	View
Storage Actions Blob Data Operator	Used by the Storage Actions - Storage Task to list & perform operations on the Storage Account blobs	BuiltInRole	None	View
Storage Blob Data Contributor	Allows for read, write and delete access to Azure Storage blob containers and data	BuiltInRole	Storage	View
Storage Blob Data Owner	Allows for full access to Azure Storage blob containers and data, including assigning POSIX access control.	BuiltInRole	Storage	View
Storage Blob Data Reader	Allows for read access to Azure Storage blob containers and data	BuiltInRole	Storage	View
Storage Blob Delegator	Allows for generation of a user delegation key which can be used to sign SAS tokens	BuiltInRole	Storage	View

At the bottom are buttons for 'Review + assign', 'Previous', 'Next', and 'Feedback'.

Search 'Blob' and click on 'Storage Blob Data Contributor.' It gives both **write** and **read** access.

Microsoft Azure

Home > olistdatastorageaccount4 | Access Control (IAM) > Add role assignment

Role Members* Conditions Review + assign

Selected role Storage Blob Data Contributor

Assign access to User, group, or service principal Managed identity

Members + Select members

Description Optional

Review + assign Previous Next

Select managed identities

Some results might be hidden due to your ABAC condition.

Subscription * Azure for Students

Managed identity Select

Search by resource type

User-assigned managed identity (1)

System-assigned managed identity

All system-assigned managed identities (3)

Access Connector for Azure Databricks (1)

Data factory (V2) (1)

Synapse workspace (1)

Learn more about RBAC

Select Close Feedback

Microsoft Azure

Home > olistdatastorageaccount4 | Access Control (IAM) > Add role assignment

Role Members* Conditions Review + assign

Selected role Storage Blob Data Contributor

Assign access to User, group, or service principal Managed identity

Members + Select members

Description Optional

Review + assign Previous Next

Select managed identities

Some results might be hidden due to your ABAC condition.

Subscription * Azure for Students

Managed identity Synapse workspace (1)

Search by name

Selected members:

olist-synapse-workspace /subscriptions/f6234497-0192-4d89-b2ea-a902d46f29d3/resourceGroups/... Remove

Select Close Feedback

Click Next. Then Click on ‘Managed Identity’. Click on ‘Select Members’. Select ‘Synapse Workspace’

Step 4: This is how Azure Synapse Analytics looks like.

```

1 This is auto-generated code
2
3 SELECT
4     TOP 100 *
5 FROM
6     OPENROWSET(
7         BULK 'https://synapselfstoredefault1.dfs.core.windows.net/synapsefilesys/olist_customers_dataset.csv',
8         FORMAT = 'CSV',
9         PARSE_VERSION = '2.0'
10    ) AS [result]

```

When I am trying to create a new SQL Database, it is showing me two types of SQL Pool Type.

Create SQL database

Create database to organize your workload into databases and database objects.

Select SQL pool type *

Serverless Dedicated

Database *

There are 2 types of SQL Pool in Azure:

Factor	Dedicated SQL Pool	Serverless SQL Pool
Cost Model	Pay-per-capacity, charges apply whether resources are used or not	Pay-per-query, charged based on data processed during execution
Tech Stack	Uses MPP (Massively Parallel Processing) architecture	Uses distributed query processing, optimized for data lake storage
Implementation Effort	Requires more setup and configuration	Minimal setup, often ready-to-use out of the box
Performance	Consistent, optimized for high-performance workloads	Performance can vary based on the complexity and size of the query
Storage and Compute	Storage and computing are decoupled but need to be managed independently	Storage is handled via a data lake, and compute is allocated on demand
Use Cases	Ideal for data warehousing, large scale ETL processes, and fixed workload patterns	Best for ad-hoc queries, data exploration, and dynamic workloads

⚙️ 1. Serverless SQL Pool (Pay-as-you-go)

💡 Simple Explanation:

- You don't have to set up or manage any server.
- You just run a query directly on files in your **Data Lake (ADLS)** — like CSV, JSON, or Parquet files.
- Azure will automatically provide the computing power needed to run that query.
- Once done, it shuts down automatically.
- You only **pay for the data processed** by your queries (not for idle time).

💡 Real-life Analogy:

Think of it like **using Uber or Ola 🚗**

- You **don't own the car** (no servers to maintain).
- You **call a ride only when you need it** (on-demand query).
- You **pay only for that trip** (pay-per-query).
- When your ride ends, you stop paying — no ongoing cost.

So, if you just need to **occasionally explore data**, check reports, or run analytics — this model is perfect.

⚙️ 2. Dedicated SQL Pool (Provisioned)

💡 Simple Explanation:

- You **reserve computing resources (servers)** ahead of time.
- The servers are **always running** and ready to process queries.
- You **load data into Synapse's internal storage** (not directly from Data Lake).
- You **pay a fixed cost**, even if the system is idle.

💡 Real-life Analogy:

Think of it like **buying your own car 🚗**

- You **own the car** (fixed resources).
- It's **always available** whenever you want to drive.
- You **pay maintenance, fuel, insurance** — even if it sits in the garage unused.
- But it's **much faster and predictable** than waiting for an Uber.

This is ideal for companies with **regular heavy workloads**, like daily ETL pipelines, dashboards, or reports that must run at high speed.

Feature	Serverless SQL Pool (Pay-as-you-go)	Dedicated SQL Pool (Provisioned)
How it works	On-demand, no servers to manage	Fixed compute resources reserved
Payment model	Pay only for queries run	Pay fixed cost for reserved compute
Data location	Data stays in Data Lake (ADLS)	Data loaded into Synapse database
Use case	Great for exploring data occasionally	Best for large, frequent workloads
Setup	Easy — no setup or maintenance	Needs setup, configuration, and monitoring
Speed	Slower for big workloads	Very fast and consistent performance

```

-- This is auto-generated code
SELECT
    TOP 100 *
FROM
    OPENROWSET(
        BULK 'https://olistdatastorageaccount4.dfs.core.windows.net/olistdata/silver/',
        FORMAT = 'PARQUET'
    ) AS result1
10 -- The data is in the Lake only and we are able to query it thus we will save a lots of cost because we are able to read the data only
  
```

The screenshot shows the Microsoft Azure Synapse Analytics workspace interface. In the center, there's a code editor window displaying a T-SQL query. The query reads data from a Parquet file located in Azure Data Lake Storage Gen2. The results pane below the editor shows a table with several columns, including product_category, order_status, and various dates. The status bar at the bottom indicates that the query was executed successfully.

Now after running the query (with little modifications accordingly), we can see the data from the **SILVER** layer of the Azure Data Lake Studio Gen 2.

Parquet File Format:

Parquet is a **columnar storage format** designed for efficient storage and processing of large datasets in distributed environments like **Azure Synapse, Databricks, Hadoop, and Spark**. Unlike traditional row-based formats (e.g., CSV), Parquet stores data **column by column**, allowing queries to read only the specific columns they need. This structure significantly reduces I/O operations, speeds up query performance, and minimizes storage space. It's especially suited for **analytical workloads**, where large datasets are scanned but only a few columns are used.

Snappy Compression Algorithm:

Snappy is a **fast data compression and decompression algorithm** developed by Google. Its primary goal is **speed over maximum compression ratio** — it compresses data efficiently enough to save space while ensuring very fast read/write performance. When combined with Parquet, Snappy makes file storage both **compact and quick to process**, which is crucial for big data systems that handle terabytes of information daily. Remember **OPENROWSET** is like a temporary connector between SQL tables and it is **Microsoft Specific**. But, it is only applicable **Azure Synapse Analytics & Microsoft SQL Server**.

 **BUT... They have similar features under different names:**

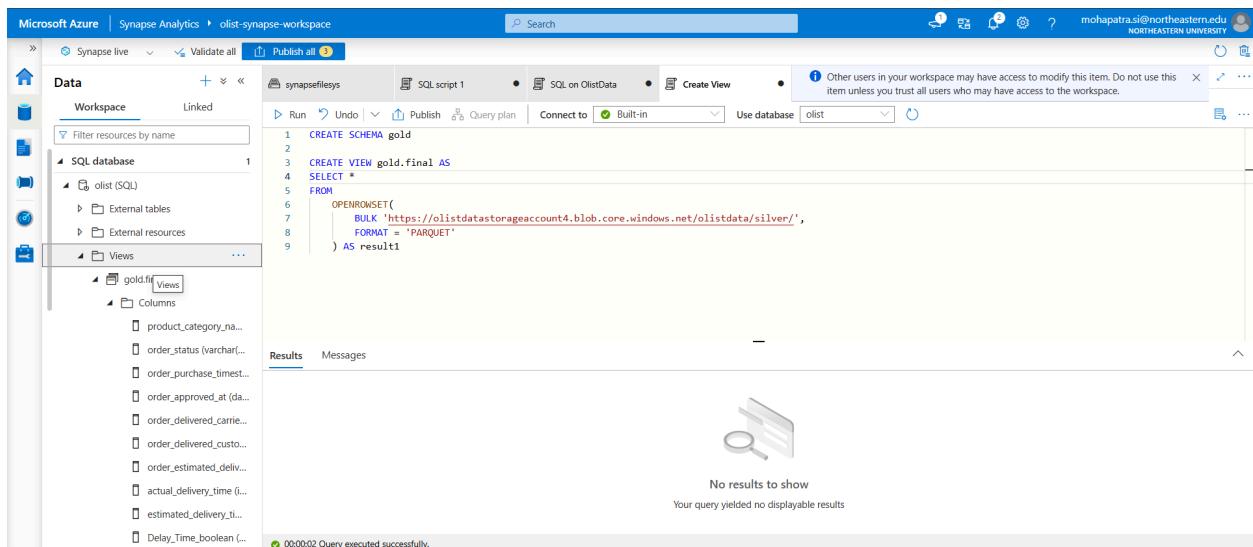
Snowflake
Equivalent: `EXTERNAL TABLE`, `STAGE`, `COPY INTO`

BigQuery
Equivalent: `LOAD DATA`, external tables using Cloud Storage

Redshift
Equivalent: `COPY`, `UNLOAD`, `SPECTRUM`

Databricks
Equivalent: `spark.read.parquet()`, `CREATE TABLE USING delta LOCATION '...'`

Step 5:



The screenshot shows the Microsoft Azure Synapse Analytics workspace interface. The left sidebar displays the Data section with a workspace named "synapsefilessys". Under "SQL database", there is a "olist" database containing "External tables" and "External resources". A "Views" folder contains a view named "gold.mf". The main pane shows a SQL script being run:

```
1 CREATE SCHEMA gold
2
3 CREATE VIEW gold.final AS
4 SELECT *
5 FROM
6 OPENROWSET(
7     BULK 'https://olistdatastorageaccount4.blob.core.windows.net/olistdata/silver/',
8     FORMAT = 'PARQUET'
9 ) AS result1
```

The results pane at the bottom indicates "No results to show" and "Your query yielded no displayable results". A status bar at the bottom right shows "00:00:02 Query executed successfully."

Here, I have created a view to store the data (that I retrieved from the **SILVER** layer) I created the View from the Silver layer because I think it's easy to fetch and see data from the SQL editor itself I mean Azure Synapse Analytics.

Step 5: Now, I am going to send those view's data to **GOLD** layer so that it can be used for serving. So, I will be creating the physical tables will be the one which anyone can now read from my data lake and this will be stored inside the **GOLD** layer. So, here we will use **CETAS** command in SQL (**Create External Table as Select**) (<https://learn.microsoft.com/en-us/azure/synapse-analytics/sql/develop-tables-cetas>)

```
SQL Copy

CREATE DATABASE [<mydatabase>];
GO

USE [<mydatabase>];
GO

CREATE MASTER KEY ENCRYPTION BY PASSWORD = '<strong password>';

CREATE DATABASE SCOPED CREDENTIAL [WorkspaceIdentity] WITH IDENTITY = 'Managed Identity';
GO

CREATE EXTERNAL FILE FORMAT [ParquetFF] WITH (
    FORMAT_TYPE = PARQUET,
    DATA_COMPRESSION = 'org.apache.hadoop.io.compress.SnappyCodec'
);
GO

CREATE EXTERNAL DATA SOURCE [SynapseSQLwriteable] WITH (
    LOCATION = 'https://<mystorageaccount>.dfs.core.windows.net/<mycontainer>/<mybaseoutputfold>',
    CREDENTIAL = [WorkspaceIdentity]
);
GO

CREATE EXTERNAL TABLE [dbo].[<myexternaltable>] WITH (
    LOCATION = '<myoutputsubfolder>/',
    DATA_SOURCE = [SynapseSQLwriteable],
    FILE_FORMAT = [ParquetFF]
) AS
SELECT * FROM [<myview>];
GO
```

AZURE SYNAPSE WORKFLOW:

1. Medallion Architecture and Lake House. (Pre-requisite)
2. Creating a Serverless SQL Pool. (For just reading & querying)
3. Creating a schema & then user/password

```

CREATE MASTER KEY ENCRYPTION BY PASSWORD = '<strong password>';

CREATE DATABASE SCOPED CREDENTIAL [WorkspaceIdentity] WITH IDENTITY = 'Managed Identity';

```

4. Using OPENROWSET to read data from the **SILVER** layer. (Here, data is not stored, just referenced to **SILVER** layer)
5. Creating a **VIEW** to easily query the data.

Now, I will be making sure, I am able to create external physical tables which will be stored in the **GOLD** layer which can be used further by Data Analysts and Data Scientists.

```

CREATE EXTERNAL FILE FORMAT [ParquetFF] WITH (
    FORMAT_TYPE = PARQUET,
    DATA_COMPRESSION = 'org.apache.hadoop.io.compress.SnappyCodec'
);

```

```

CREATE EXTERNAL DATA SOURCE [SynapseSQLwriteable] WITH (
    LOCATION =
    'https://<mystorageaccount>.dfs.core.windows.net/<mycontainer>/<mybaseoutputfolderpath>',
    CREDENTIAL = [WorkspaceIdentity]
);

```

```

Microsoft Azure | Synapse Analytics > olist-synapse-workspace
Develop + <> Run Undo Publish... Query plan Connect to Built-in Use database olist Publishing Deploying changes to the workspace
Filter resources by name
SQL scripts 6
Create View
SQL on OlistData
SQL script 1
SQL script 2
3
4 -- SELECT *
5 -- FROM sys.database_credentials
6
7 CREATE EXTERNAL FILE FORMAT extfileformat WITH (
8     FORMAT_TYPE = PARQUET,
9     DATA_COMPRESSION = 'org.apache.hadoop.io.compress.SnappyCodec'
10 );
11

```

After this, I published my final data to **Gold Layer** in the Serving Folder:

Name	Last modified	Access tier	Blob type	Size	Lease state
...					
BD3DEF45-BFC7-4934-8ECE-73D16F45FC...	10/26/2025, 5:23:48 PM	Hot (Inferred)	Block blob	2.08 MiB	Available
BD3DEF45-BFC7-4934-8ECE-73D16F45FC...	10/26/2025, 5:23:48 PM	Hot (Inferred)	Block blob	2.17 MiB	Available
BD3DEF45-BFC7-4934-8ECE-73D16F45FC...	10/26/2025, 5:23:48 PM	Hot (Inferred)	Block blob	2.07 MiB	Available
BD3DEF45-BFC7-4934-8ECE-73D16F45FC...	10/26/2025, 5:23:48 PM	Hot (Inferred)	Block blob	2.09 MiB	Available
...	10/26/2025, 5:23:40 PM	Hot (Inferred)	Block blob	0	Available

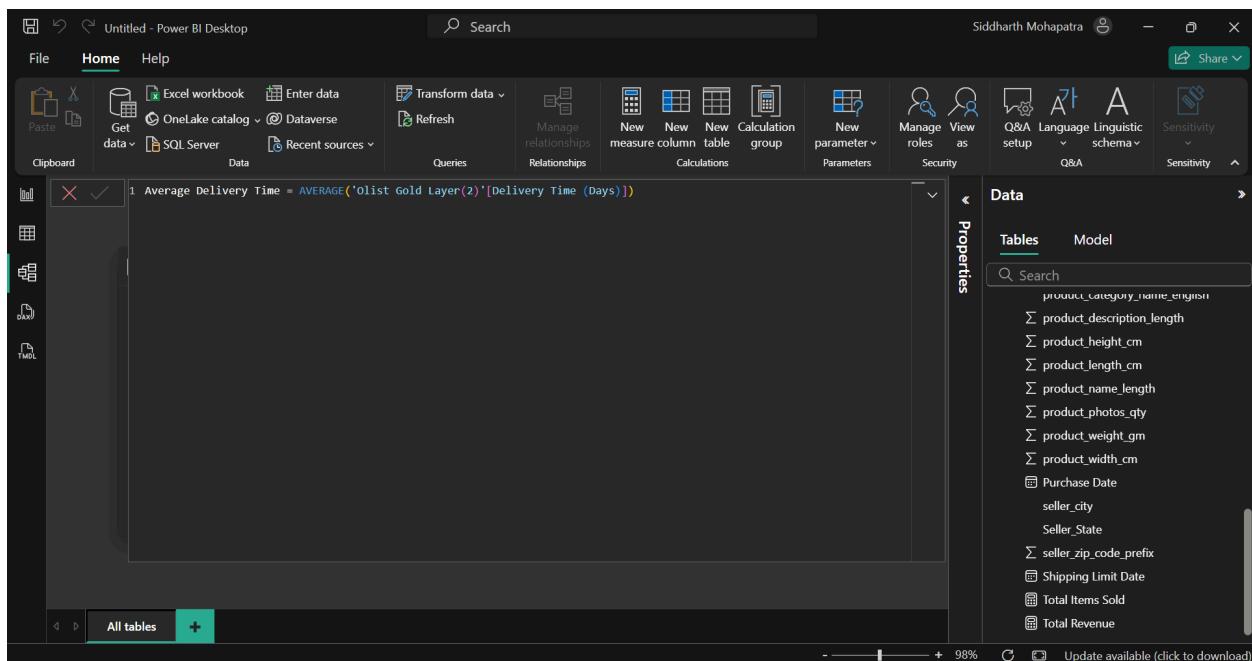
Final Cleaning of Gold Layer Data before dashboarding in Power BI

- Removed duplicate rows
- Removed null value - rows in **product category name** column
- Collected a State - State code data of Brazilian states and used a **LOOKUP** function to replace all the customer and seller state codes with actual full state names for better understanding for the users and audience

POWER BI DASHBOARD GOAL:

Now, I will show how my **Azure Medallion architecture** (Bronze → Silver → Gold) powers real-time insights on e-commerce performance, focusing on **customer behavior**, **sales performance**, and **delivery efficiency**.

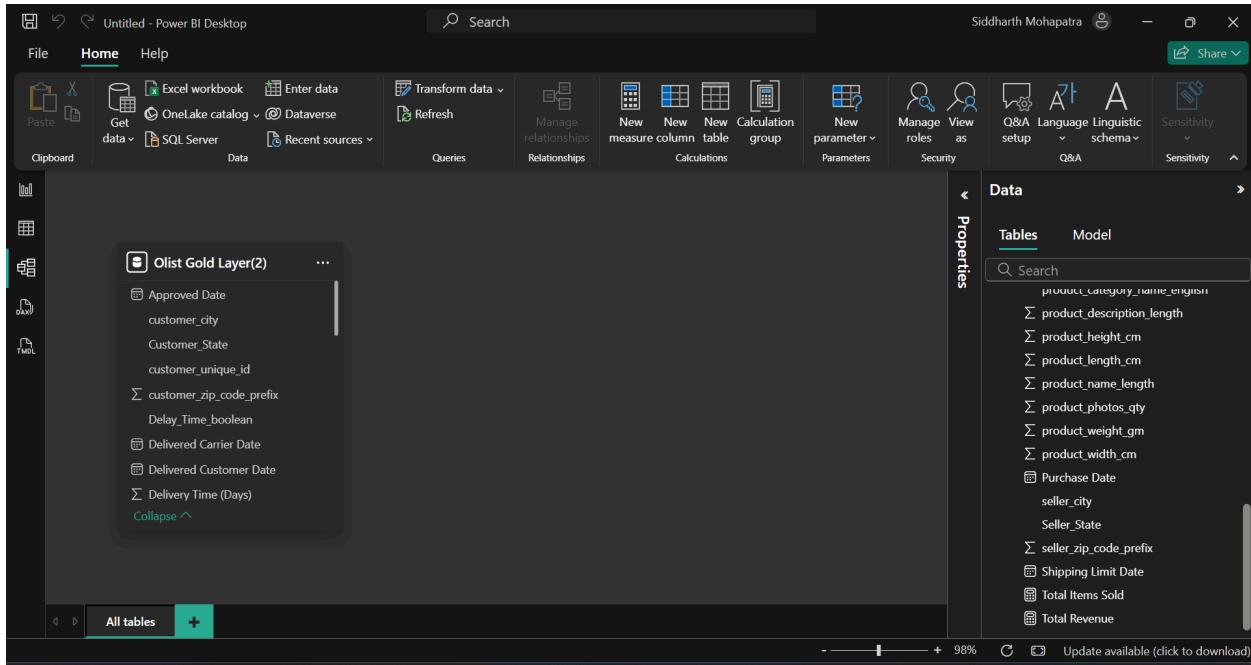
POWER BI INTERFACE



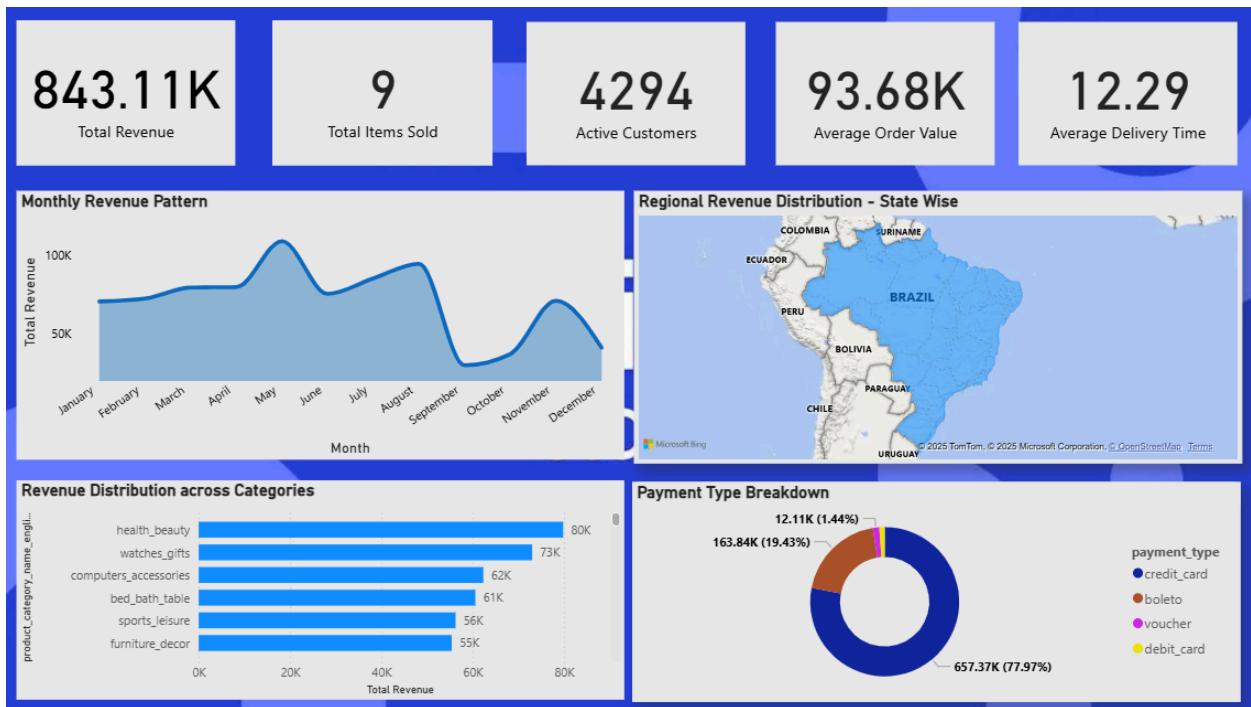
DAX (Data Analytics Expressions) for KPIs

- Total Revenue = SUM('Olist Gold Layer(2)'[payment_value (Brazilian Real)])
- Total Items Sold = DISTINCTCOUNT('Olist Gold Layer(2)'[order_item_id])
- Active Customers = DISTINCTCOUNT('Olist Gold Layer(2)'[customer_unique_id])
- Average Order Value = DIVIDE([Total Revenue], [Total Items Sold])

- Average Delivery Time = AVERAGE('Olist Gold Layer(2)'[Delivery Time (Days)])



MY POWER BI DASHBOARD:



This dashboard provides a high-level view of Olist's e-commerce performance across Brazil — showing revenue trends, regional distribution, product performance, and customer insights. It helps management and analysts monitor performance, identify growth opportunities, and improve operational efficiency.

1KPI: Total Revenue (843.11K BRL)

- **What it shows:** The total value of all payments received across all orders.
- **Why it matters:** This is the most direct measure of overall business performance and sales growth.
- **What we infer:** Revenue is healthy and concentrated in a few regions and categories. Any significant monthly drop could indicate seasonal trends or logistics issues.
- **Improvement:** Drill down by region or product category to find where sales are underperforming and launch promotions there.

2KPI: Total Items Sold (9)

- **What it shows:** Number of product categories or unique SKUs (Stock Keeping Unit) sold.
- **Why it matters:** Reflects product variety and inventory turnover.
- **What we infer:** Olist focuses on a smaller curated set of products, which simplifies supply chain but limits category diversification.
- **Improvement:** Expand high-demand categories or add complementary products to increase basket size.

3KPI: Active Customers (4,294)

- **What it shows:** Total number of unique customers who made at least one purchase.
- **Why it matters:** Indicates customer base strength and market reach.
- **What we infer:** The platform has a good number of active buyers, but retention trends can be explored further.
- **Improvement:** Build a loyalty or referral program to retain customers and boost repeat purchases.

4KPI: Average Order Value (93.68K BRL)

- **What it shows:** Average amount spent per customer order.

- **Why it matters:** Higher AOV means customers are buying more expensive or multiple items per transaction.
- **What we infer:** The pricing strategy or bundling seems effective.
- **Improvement:** Introduce cross-selling or up-selling strategies to lift AOV further.

5 KPI: Average Delivery Time (12.29 Days)

- **What it shows:** Average duration between order placement and delivery completion.
- **Why it matters:** Affects customer satisfaction and repeat business.
- **What we infer:** Delivery speed is moderate; some regions likely experience delays.
- **Improvement:** Optimize last-mile logistics and vendor coordination to reduce delivery time.

6 Chart: Monthly Revenue Pattern (Line Chart)

- **What it shows:** Revenue trends from January to December.
- **Why it matters:** Helps identify seasonality, peak months, and slow periods.
- **What we infer:** Revenue spikes mid-year (May–August), possibly due to sales events or festivals, and dips toward year-end.
- **Improvement:** Introduce marketing campaigns during low-sales months and increase inventory before peak months.

7 Chart: Regional Revenue Distribution (Map Visual)

- **What it shows:** State-wise contribution to total revenue across Brazil.
- **Why it matters:** Helps pinpoint high-performing and low-performing regions.
- **What we infer:** Most sales are concentrated in southeastern Brazil (SP, RJ), showing strong urban demand.
- **Improvement:** Expand operations and seller partnerships in underperforming states like northern or northeastern Brazil.

8 Chart: Revenue Distribution Across Categories (Bar Chart)

- **What it shows:** Top-performing product categories by total revenue.
- **Why it matters:** Guides inventory planning, marketing, and product focus.
- **What we infer:** *Health & Beauty* and *Watches & Gifts* generate the most revenue, suggesting strong lifestyle-driven demand.
- **Improvement:** Increase marketing for high-margin categories and run discount campaigns for weaker categories like *Furniture Décor*.

9| Chart: Payment Type Breakdown (Donut Chart)

- **What it shows:** Share of different payment modes (Credit Card, Boleto, Voucher, Debit Card).
- **Why it matters:** Reveals customer payment preferences and financial inclusivity.
- **What we infer:** Over 75% of transactions are made via credit cards — customers rely heavily on credit-based purchases.
- **Improvement:** Offer more flexible payment options (e.g., UPI, installment plans) to cater to different customer segments.

Key Takeaways

- **Top Categories:** Health & Beauty and Watches & Gifts are the main revenue drivers.
- **Geographic Focus:** São Paulo and Rio de Janeiro dominate the market share.
- **Customer Experience:** Improving delivery time could increase satisfaction and reduce churn.
- **Strategic Focus:** Marketing and logistics optimization in lower-performing regions and months will balance overall revenue.

<u>Visual / KPI</u>	<u>Purpose</u>	<u>Insight Derived</u>	<u>Business Impact</u>	<u>Improvement Suggestions</u>
---------------------	----------------	------------------------	------------------------	--------------------------------

Total Revenue	Measure total payments received from all transactions	Indicates overall business performance and sales health	Reveals financial strength and helps evaluate marketing or pricing effectiveness	Drill down by region and product to identify underperforming areas
Total Items Sold	Shows total unique products/categories sold	Reflects product diversity and SKU performance	Helps in inventory planning and demand forecasting	Add more products in high-demand segments to increase variety and sales
Active Customers	Tracks the number of unique customers who purchased	Indicates customer base size and engagement	Helps identify customer growth or retention rate	Launch loyalty/referral programs to boost repeat customers
Average Order Value (AOV)	Measures average spend per order	High AOV means effective pricing and cross-selling	Correlates with profitability and marketing effectiveness	Encourage product bundling and targeted promotions to lift AOV
Average Delivery Time	Tracks mean time from order to delivery	Measures efficiency of logistics and fulfillment	Impacts customer satisfaction and brand trust	Optimize delivery routes and improve warehouse coordination
Monthly Revenue Pattern (Line Chart)	Visualizes monthly revenue trend	Shows seasonality — peaks and dips in sales	Aids budgeting and forecasting	Introduce offers in low-sales months and prepare inventory for peaks
Regional Revenue Distribution (Map)	Highlights state-wise revenue performance	SE region (SP, RJ) dominates; others underperform	Guides regional marketing and logistics	Expand seller network and delivery centers in low-performing states
Revenue by Category (Bar Chart)	Compares total revenue across product categories	Health & Beauty and Watches lead in revenue	Helps focus investment on high-performing categories	Rebalance stock and create campaigns for slow-selling items

Payment Type Breakdown (Donut Chart)	Shows customer payment preferences	78% of users use credit cards; few use vouchers/debit	Helps tailor financial partnerships and promotions	Introduce flexible payment options (installments, wallets)
---	------------------------------------	---	--	--

STRATEGIC INSIGHTS & RECOMMENDATION

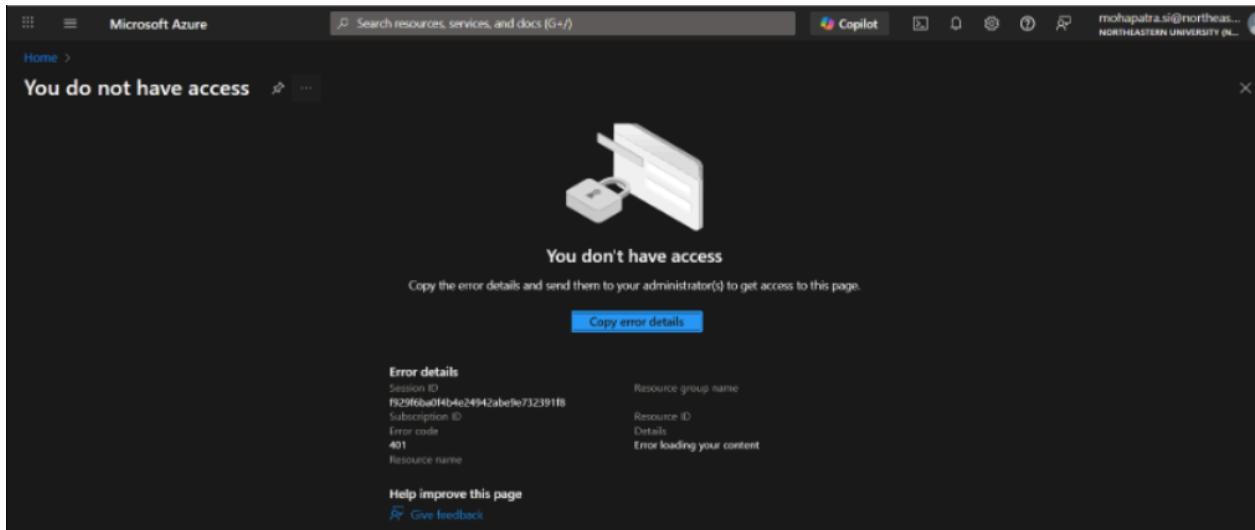
<u>Insight</u>	<u>Business Interpretation</u>	<u>Recommended Action</u>
Revenue is heavily concentrated in SE Brazil (SP, RJ)	Indicates high dependency on a few markets	Expand seller onboarding and delivery infrastructure in underpenetrated states (North & Northeast)
Health & Beauty and Watches & Gifts dominate revenue	Product mix is skewed toward lifestyle categories	Diversify inventory — promote tech/home products through discounts or cross-sell campaigns
Average Delivery Time is relatively high (12.29 days)	Slower deliveries may hurt customer satisfaction	Partner with regional carriers or introduce local warehouses for faster fulfillment
Majority of payments are via credit cards (~78%)	High credit dependency may limit financially constrained customers	Add more payment options (Pix, digital wallets, UPI equivalents) to expand reach
Revenue dips in certain months	Possible seasonality or supply issues	Use predictive analytics to forecast demand and plan marketing campaigns during off-peak months

SOME OF THE MAJOR & MINOR CHALLENGES I HAVE FACED DURING THE PROJECT

Challenges Faced

1. Built a separate pipeline for olist_geolocation from PostgreSQL to Azure bronze layer (ADLS Gen2) because of the timeout issue due to the large volume of data (around 10 lakhs records).
2. Used Self Hosted Integrated Runtime (SHIR), which acts as a gateway from the Local PostgreSQL server to the Azure cloud, as local to cloud is not possible without any gateway.
3. Used the Publish All button for permanent save otherwise I would have taken around 2-3 weeks more to start building the pipeline from scratch, which took me a lot of time.

4. To make the connection, we need some key or permission here to maintain the equality. Both are Microsoft Azure services, so they are connected within Microsoft, and no third party can access them. So, here in this case when I tried to connect Azure Databricks and ADLS, I got this error.



What This Error Means

The error "You do not have access" with Error code: 401 means "Unauthorized." It's a permissions issue. I was successfully logged into Azure, but my specific account (mohapatra.si@northeastern.edu) does not have the administrative rights to view or manage the "App registrations" section.

Why I didn't Have Access (The University Analogy)

Think of your university's entire Microsoft Azure setup as a large office building.

The Building (Azure Active Directory/Tenant): The entire building is managed by your university's IT administrators. They control the main entrance, the security office, and who gets what keys. The "App registrations" page is like the building's main security office. It's a central, high-level area.

Your Office (Your Azure for Students Subscription): The university has given you your own private office inside the building—this is your "Azure for Students" subscription. Inside your office, you are the owner. You can set up furniture (like Azure Data Factory, Databricks, etc.) and do whatever you want.

The error we are seeing is because our keycard (our student account) lets us into our own office, but it doesn't grant us access to the main security office for the entire building. This is a security measure to prevent students from viewing or changing administrative settings that affect the whole university.

However, the tutorial I was following demonstrated the most common and robust method for production environments: creating a Service Principal. A Service Principal is like creating a robot user with its own specific permissions. However, as you discovered, creating these "robot users" requires high-level administrative access to the "security office" (App Registrations), which I don't have with a student account.

There is another, simpler method that is perfect for this situation and still keeps the connection secure within the Microsoft cloud. I will use a method called mounting with an Account Key.

Instead of creating a new "robot user," you will temporarily use your own powerful key (the Storage Account Access Key) to establish a permanent link (a mount point) between Databricks and your Data Lake.

What I did here?

Think of it like giving your workshop (Databricks) a secure key to your storage unit (ADLS).

We got the Master Key: I went to my Azure storage account and copied its Access Key (key1). This key is like a master password that grants full access to all the data inside the service.

We Used a Secure Lockbox: Instead of pasting this sensitive key directly into your code (which is insecure), we created a Databricks Widget outside of the cell. This widget is like a secure lockbox at the top of my notebook. I pasted the key into this box, keeping it separate from my saved code.

We Built a Secure Bridge: Finally, I ran the dbutils.fs.mount() command. This command did two things:

It securely took the key from the lockbox.

It used the key to build a permanent, direct link—a mount point—from Databricks to my storage container.

The result is that my olistdata container now appears as a simple folder inside Databricks at the path /mnt/olistdata. Anyone can now read and write files to our storage as if it were a local directory.

5. Why do Synapse deployments fail?

Think of a Synapse workspace like opening a new mall. We're not just creating the mall—we also need parking (SQL server), security (managed identity/permissions), roads (networking), and a storage warehouse (Data Lake). If any of those can't be created in the location you selected, the entire "open the mall" action fails.

One reason for this can be:

Region blocks / capacity limits (my exact error)

Message: SqlServerRegionDoesNotAllowProvisioning ... Location 'eastus' is not accepting creation of new ... SQL Database servers...

What it means: In that region (e.g., East US), my subscription isn't allowed to create new Azure SQL Server resources right now—often due to capacity or subscription-type limits (e.g., Free Trial/Azure for Students).

Real example: We pick East US for Synapse, but Synapse needs to spin up a SQL server there; the region says "no new SQL servers allowed," so deployment fails.

Instead we can pick another region (e.g., East US 2, Central US, West US 2) to solve this, which is what happened in my case.

Dataset Link: <https://www.kaggle.com/datasets/olistbr/brazilian-e-commerce>