# Explore Data Analytics in Azure

## Create a Synapse workspace



## Create a Lakehouse



## Ingest Data

- A simple way to ingest data is to use a **Copy Data** activity in a pipeline to extract the data from a source and copy it to a file in the lakehouse.



1. Connecting to a data source



2. Choosing an existing data destination

**Copy data into Lakehouse**

Choose data source ✓
Connect to data source ✓
**Choose data destination**
Define the data store as destination.
Connect to data destination
Review + save

**Choose data destination**

Lakehouse
Learn more

○ Existing Lakehouse   ○ Create new Lakehouse

Lakehouse *
sn_lakehouse    ↻ Refresh

Back   Next

---



**Copy data into Lakehouse**

Choose data source ✓
Connect to data source ✓
Choose data destination ✓
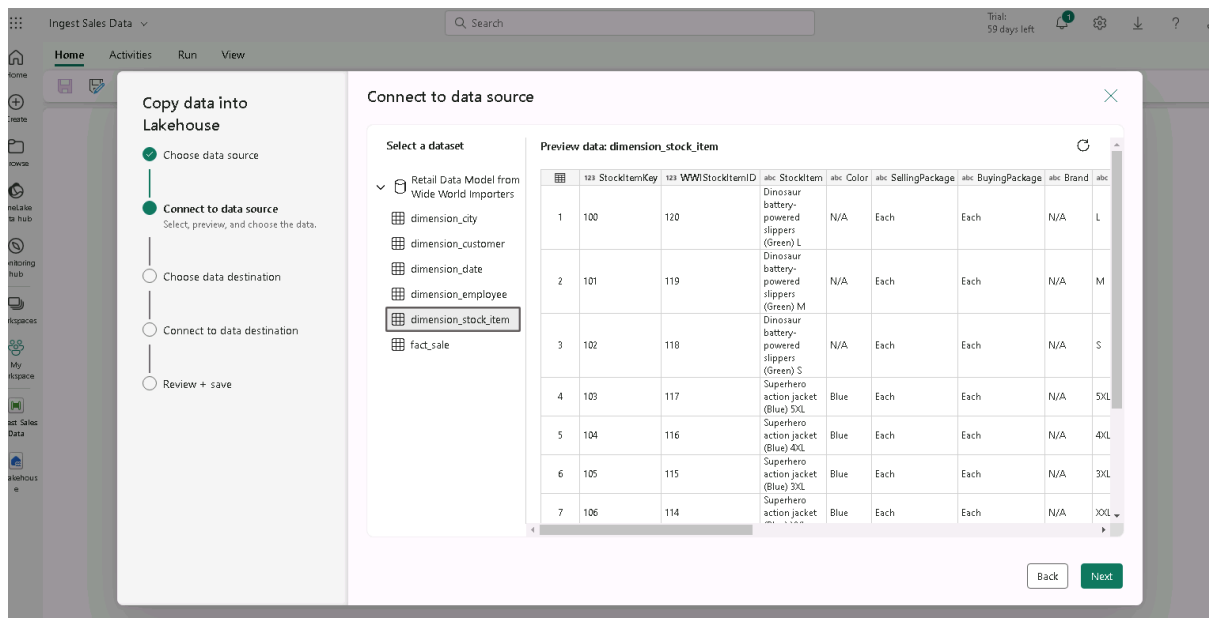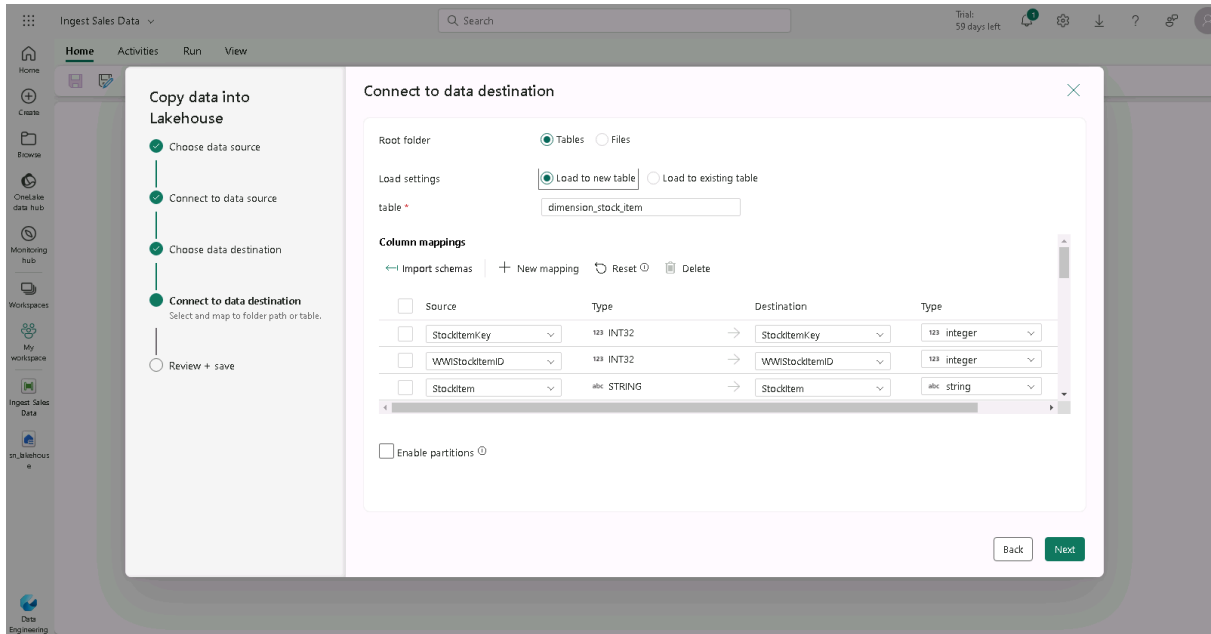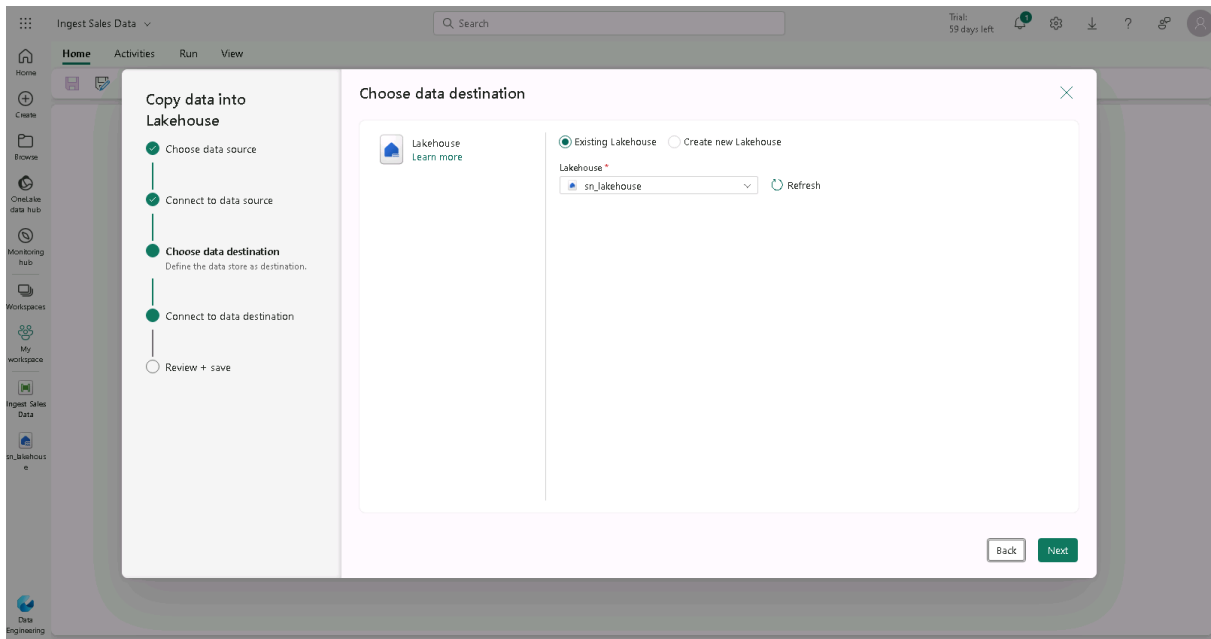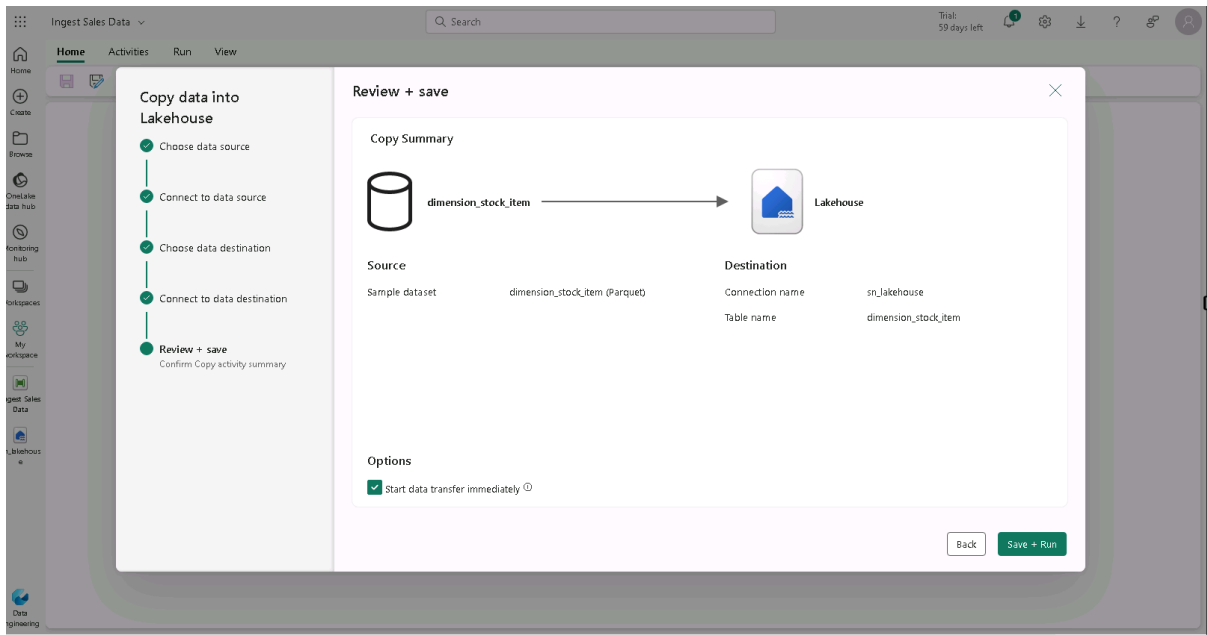**Connect to data destination**
Select and map to folder path or table.
Review + save

**Connect to data destination**

Root folder          ○ Tables   ○ Files

Load settings        ○ Load to new table   ○ Load to existing table

table *              dimension_stock_item

**Column mappings**

← Import schemas   + New mapping   ↺ Reset ⓘ   🗑 Delete

| | Source | Type | | Destination | Type |
|---|---|---|---|---|---|
| ☐ | StockItemKey | 123 INT32 | → | StockItemKey | 123 integer |
| ☐ | WWIStockItemID | 123 INT32 | → | WWIStockItemID | 123 integer |
| ☐ | StockItem | abc STRING | → | StockItem | abc string |

☐ Enable partitions ⓘ

Back   Next

A new pipeline containing a **Copy Data** activity is created, as shown here



Results:

Query data in a lakehouse:

1. Select SQL analytical endpoint



2. Select new sql query and entered the sql code into the query editor. Ran the query and it revealed that there are two brand values - N/A and Northwide

and shows the number of products in each



Visualise you data in a lakehouse:

1. Select model tab. This will let you see the data for the tables in the lakehouse.

2. Selected stacked barchart and input brand and stockitemkey fields. Then the aggregation in the y-axis and x-axis was changed to count. Next, added colours and resized the chart to make it more visualsing.



3.

4. Saved it as brand quantity report



Create a KQL Database

1. Switched to Real Time Analytics. The real time analytics home page includes tiles to create commonly used assets for real-time data analyis. Named the database ive created ' my_kql_database'.

# Create an eventstream

1. Eventstreams provide a scalable and flexible way to ingest real-time data from a streaming source.

## KQL Database

**Destination name** *
taxi-data

**Workspace** *
sn-kqldatabase

**KQL Database** *
my_kql_db

**Destination table** *
(New) taxi-data

✓ Schema validation successful.
Create new

**Input data format** * ⓘ
Json

### Event processing

Event processor enables you to transform and preview the data that is being processed for the destination.

Open event processor

Add

**Data preview**     Data insights              ⟳ Refresh  ⌄

Data format   Json
Last refreshed  01/18/24 06:49:38 AM
Last event time  01/18/24 06:49:31 AM          ▦ Show details

| VendorID | tpep_pickup_datetime | tpep_dropoff_datetime |
|----------|----------------------|------------------------|
| 2 | 2022-06-01 08:10:50 | 2022-06-01 08:50:24 |
| 2 | 2022-06-01 08:00:27 | 2022-06-01 08:56:17 |
| 2 | 2022-06-01 08:55:24 | 2022-06-01 09:17:44 |
| 2 | 2022-06-01 08:52:06 | 2022-06-01 08:59:46 |
| 2 | 2022-06-01 08:30:06 | 2022-06-01 08:50:27 |

---

Synapse Real-Time Analytics      🔍 Search          Trial: 59 days left  🔔8  ⚙  ⤓  ?  ⛶  👤

**Home**

→ New source ⌄    → New destination ⌄    ⟳ Refresh

**Data**  «

∨ Sources
   📊 taxis

∨ Destinations
   📊 taxi-data

taxis
● Streaming          →    My_eventstream    →    taxi-data
                                                 ● Ingesting

**Details**    Data preview    Data insights    Runtime logs          ⌄

Name          taxi-data 📋

Type          KQL Database

Status        ● Ingesting

Related item  Kusto item: my_kql_db   Open item ⎋
              Kusto table: taxi-data

Activate Windows
Go to Settings to activate Windows.

Query real time data in in a KQL data base:

1. select **Query table > Records ingested in the last 24 hours**.



2. This KQL query shows all taxi records ingested from the streaming source in the last 24 hours.

Explore your data

▷ Run    🖫 Save as KQL queryset                                        🗋 Copy query    ▥◦ Build Power BI repor

```
1
2  ['taxi-data']
3  | where ingestion_time() between (now(-1d) .. now())
4
5
```

⊞ Table 1    ◎ Stats                           🔍 Search    ✔ Done (2.828 s)    🔢 59,900 records    👁 🗑

| VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | passenger_count | trip_distance | RatecodeID | store_and_fwd_flag |
|---|---|---|---|---|---|---|
| 2 | 2022-06-01 11:25:42 | 2022-06-01 11:41:44 | 1.0 | 1.63 | 1.0 | N |
| 2 | 2022-06-01 11:48:46 | 2022-06-01 12:10:20 | 1.0 | 2.56 | 1.0 | N |
| 1 | 2022-06-01 11:38:30 | 2022-06-01 11:48:56 | 2.0 | 1.6 | 1.0 | N |
| 1 | 2022-06-01 11:51:13 | 2022-06-01 11:56:06 | 3.0 | 0.5 | 1.0 | N |
| 2 | 2022-06-01 11:37:07 | 2022-06-01 12:07:35 | 1.0 | 3.23 | 1.0 | N |
| 2 | 2022-06-01 11:46:26 | 2022-06-01 12:49:57 | 2.0 | 18.79 | 2.0 | N |
| 2 | 2022-06-01 11:18:08 | 2022-06-01 11:25:00 | 1.0 | 1.0 | 1.0 | N |
| 2 | 2022-06-01 11:23:49 | 2022-06-01 11:38:04 | 1.0 | 2.14 | 1.0 | N |

Activate Windows
Go to Settings to activate Windows.

3. This KQL query code shows the number of taxi pickups for each hour.

Search

⌕ Search

▷ Get data ⌄    ⊞ Query

Explorer

⌕ Search

🗄 my_kql_db
  ⌄  ⊞ Tables
    >  ⊞ taxi-data
  >  🗂 Shortcuts
  >  ⊞ Materialized view
  >  ƒx Functions
  >  ⟲ Data streams

# Explore your data

▷ Run    💾 Save as KQL queryset    📋 Copy query    ▥ Build Power BI report

```kql
1   // This query returns the number of taxi pickups per hour
2   ['taxi-data']
3   | summarize PickupCount = count() by bin(todatetime(tpep_pickup_datetime), 1h)
```

⊞ Table 1    ⊛ Stats    ⌕ Search    ✓ Done (0.233 s)    🔢 15 records    👁    🗑

| tpep_pickup_datetime ☰ | PickupCount ☰ |
|---|---|
| > 2022-06-01 00:00:00.0000 | 1 |
| > 2022-06-01 05:00:00.0000 | 1 |
| > 2022-06-01 06:00:00.0000 | 9 |
| > 2022-06-01 09:00:00.0000 | 1 |
| > 2022-06-01 10:00:00.0000 | 73 |
| > 2022-06-01 11:00:00.0000 | 4,663 |
| > 2022-06-01 12:00:00.0000 | 6,669 |
| > 2022-06-01 13:00:00.0000 | 6,761 |
| 2022-06-01 14:00:00.0000 | 7,177 |

Extension: Fabric Exercise

## Use delta tables in Apache Spark

Module: Work with Delta Lake tables in Microsoft Fabric

1. Create a lakehouse and upload data
   a. Named new lakehouse
   b. Download the data file for this exercise from
      https://github.com/MicrosoftLearning/dp-data/raw/main/products.csv, saving it as products.csv on your local compute



   c. select New subfolder and create a folder named products.
   d. Upload saved file from downloads.



2. Explore data in a dataframe

   a. On the Home page while viewing the contents of the products folder in your datalake, in the Open notebook menu, select New notebook.
      After a few seconds, a new notebook containing a single *cell* will open.

Notebooks are made up of one or more cells that can contain *code* or *markdown* (formatted text).

b. Select the existing cell in the notebook, which contains some simple code, and then use its 🗑 (*Delete*) icon at its top-right to remove it - you will not need this code.



c. In the Lakehouse explorer pane on the left, expand Files and select products to reveal a new pane showing the products.csv file you uploaded previously:

d.  In the … menu for products.csv, select Load data > Spark. A new code cell containing the following code should be added to the notebook: Run it.



Create delta tables

Create a *managed* table; *Managed* tables are tables for which both the schema metadata and the data files are managed by Fabric. The data files for the table are created in the Tables folder.

a.  Under the results returned by the first code cell, use the + Code button to add a new code cell if one doesn't already exist. Then enter the following code in the new cell and run it:

```
1   df.write.format("delta").saveAsTable("managed_products")
```

[3]   ✓   23 sec -Command executed in 23 sec 172 ms by User1-37054891 on 9:24:46 AM, 1/18/24   PySpark (Python) ∨

∨  ☰ Spark jobs (6 of 6 succeeded)   ▤ Log   ...

🔍 Search          ☐  Only show errors and warnings

```
∨2024-01-18 17:24:41,543 INFO DAGScheduler [dag-scheduler-event-loop]: Register:
 2024-01-18 17:24:41,543 INFO DAGScheduler [dag-scheduler-event-loop]: Got map :
 2024-01-18 17:24:41,543 INFO DAGScheduler [dag-scheduler-event-loop]: Final st.
 2024-01-18 17:24:41,543 INFO DAGScheduler [dag-scheduler-event-loop]: Parents
 2024-01-18 17:24:41,544 INFO DAGScheduler [dag-scheduler-event-loop]: Missing
 2024-01-18 17:24:41,545 INFO DAGScheduler [dag-scheduler-event-loop]: Submittii
 2024-01-18 17:24:41,546 INFO ExecutorMonitor [spark-listener-group-executorMan.
 2024-01-18 17:24:41,607 INFO MemoryStore [dag-scheduler-event-loop]: Block bro.
∨2024-01-18 17:24:41,610 INFO MemoryStore [dag-scheduler-event-loop]: Block bro.
 2024-01-18 17:24:41,610 INFO BlockManagerInfo [dispatcher-BlockManagerMaster]:
 2024-01-18 17:24:41,610 INFO SparkContext [dag-scheduler-event-loop]: Created I
∨2024-01-18 17:24:41,611 INFO DAGScheduler [dag-scheduler-event-loop]: Submittii
 2024-01-18 17:24:41,611 INFO YarnClusterScheduler [dag-scheduler-event-loop]: /
 2024-01-18 17:24:41,612 INFO FairSchedulableBuilder [dag-scheduler-event-loop]
∨2024-01-18 17:24:41,613 INFO TaskSetManager [dispatcher-CoarseGrainedScheduler
```
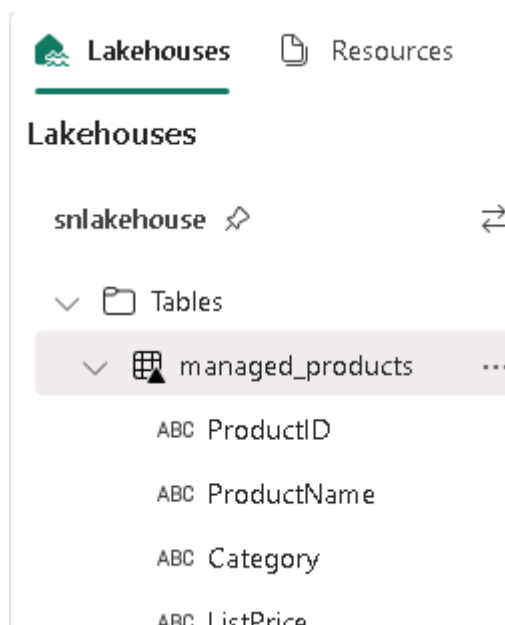
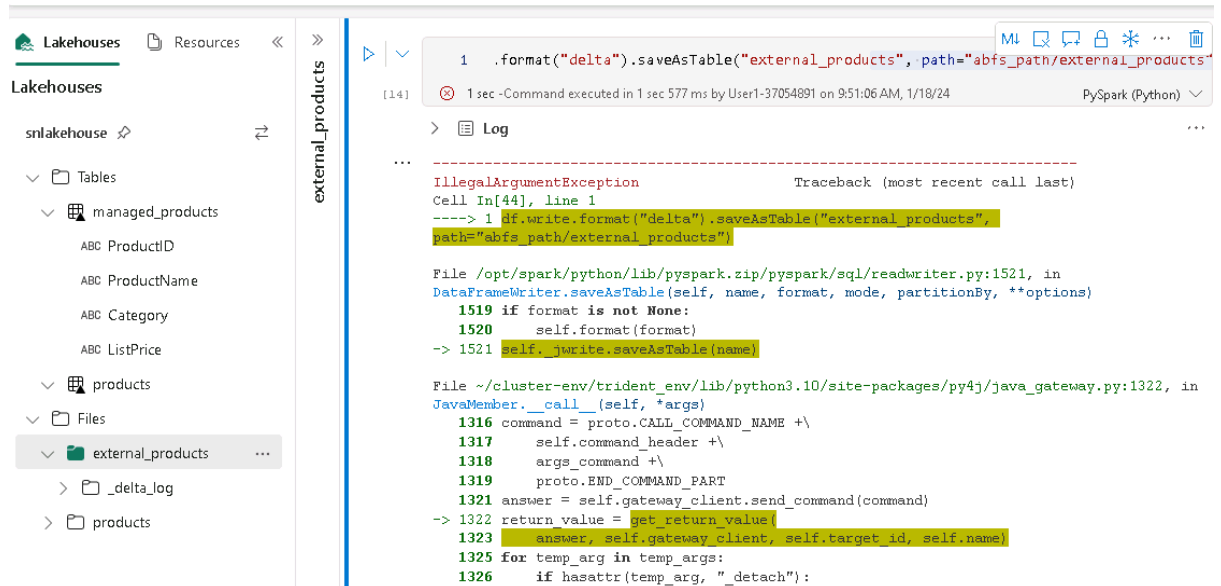b.  In the Lakehouse explorer pane, in the … menu for the Tables folder,
    select Refresh. Then expand the Tables node and verify that the
    managed_products table has been created.



🏠 Lakehouses     ▱ Resources

Lakehouses

snlakehouse  ⚲              ⇄

∨  ▱  Tables

  ∨  ▦  managed_products    ...

     ABC ProductID

     ABC ProductName

     ABC Category

     ABC ListPrice

Create an *external* table

You can also create *external* tables for which the schema metadata is defined in the metastore for the lakehouse, but the data files are stored in an external location.

     a.  Add another new code cell, and add the following code to it:



     b.  In the code you entered into the code cell, replace abfs_path with the path you copied to the clipboard so that the code saves the dataframe as an external table with data files in a folder named external_products in your Files folder location

# Create a medallion architecture in a Microsoft Fabric lakehouse

## Module: Organize a Fabric lakehouse using medallion architecture design

- Microsoft Fabric is an end-to-end analytics platform that provides a single, integrated environment for data professionals and the business to collaborate on data projects. Fabric provides a set of integrated services that enable you to ingest, store, process, and analyse data in a single environment.

1. Create a workspace: Select synapse data engineering, create a new workspace, navigate to the workspace setting and enable 'data modelling editing'.

Synapse Data Engineering  snworkspace2

# Workspace settings

About

Premium

Azure connections

System storage

Git integration

OneLake

Other

**Power BI**

General

Data connections

Embed codes

Organization apps

**Secure update**

Allow contributors to update the app for this workspace

☐ Allow contributors to update the app

Template apps

**Template apps**

Template apps are developed for sharing outside your organization. A template app workspace will be created for developing and releasing the app. Learn more about template apps ↗

☐ Develop template apps

Data model settings

**Data model settings**

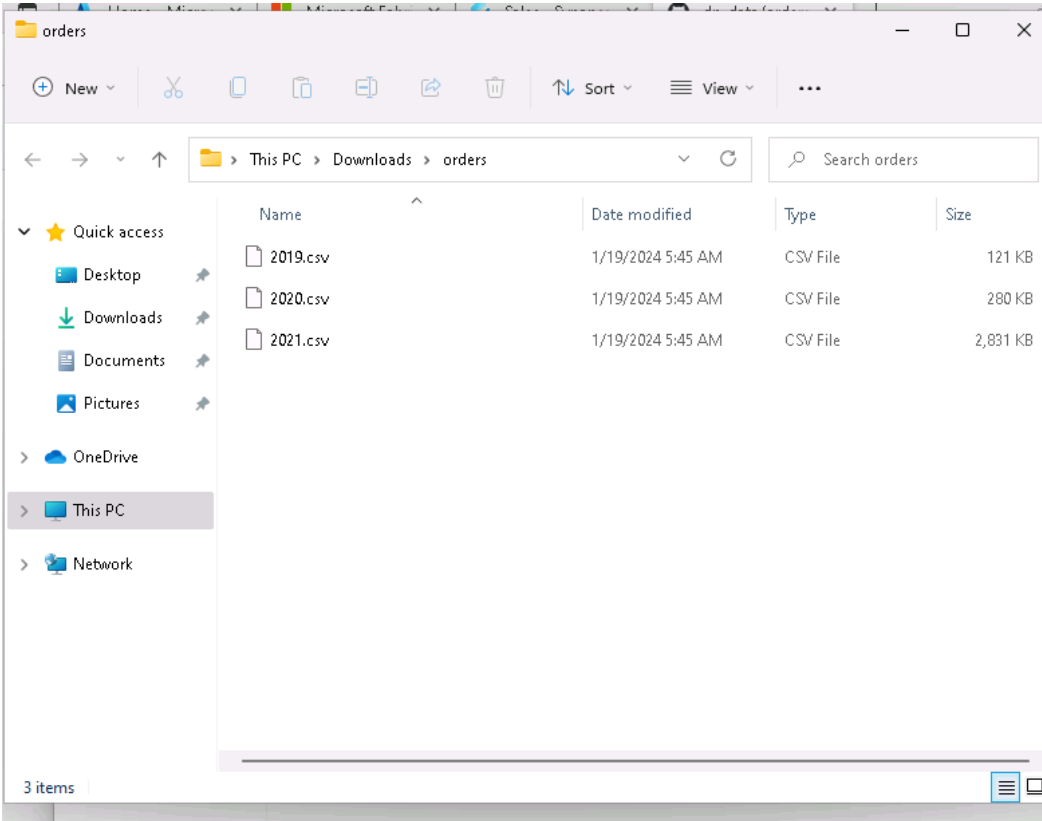Allow workspace members to edit data models in the service. Edits are permanent and automatically saved in this feature preview, and version history isn't saved. This setting doesn't apply to Direct Lake datasets or editing a dataset through an API or XMLA endpoint. Learn more ↗
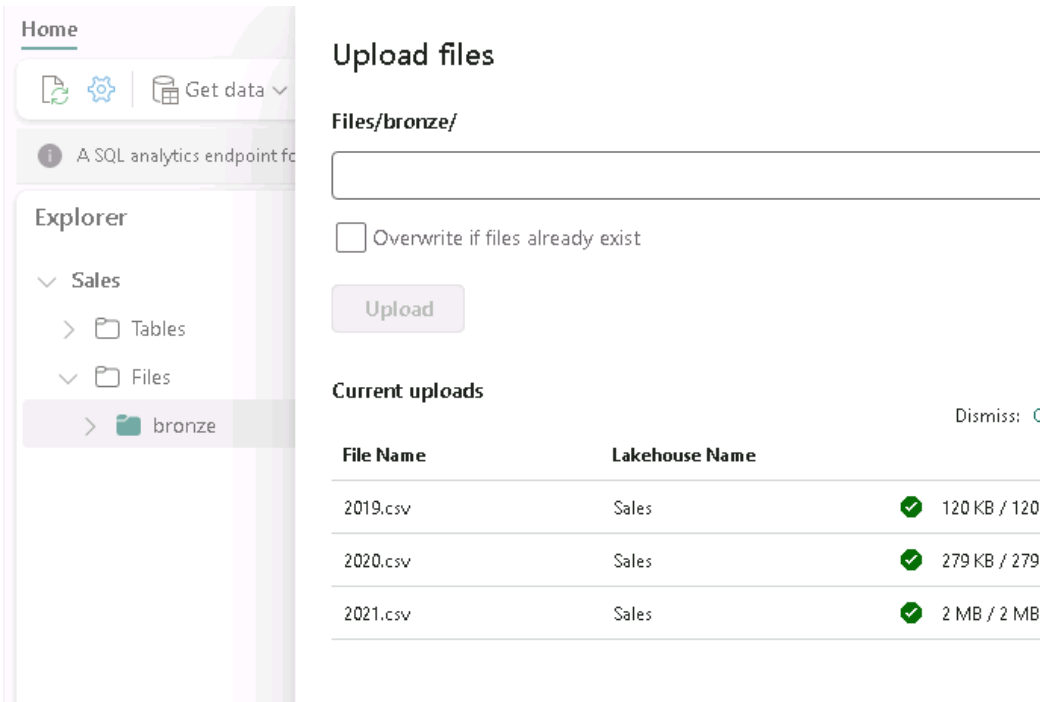
☑ Users can edit data models in the Power BI service (preview)

2. Create a lakehouse and upload data to bronze layer

a. created a new lakehouse named sales, extracted folder from github containing three sales data.



b. Created a new subfolder names bronze under files and upload the extracted files.



3. Transform data and load to silver delta table

a. Created a new notebook named transform data for silver and inserted a new sparks code into the code cell and ran it. The code loaded the data from the CSV files in the bronze folder into a Spark dataframe, and then displayed the first 10 rows of the dataframe.

The code:

```
1   from pyspark.sql.types import *
2
3   # Create the schema for the table
4   orderSchema = StructType([
5       StructField("SalesOrderNumber", StringType()),
6       StructField("SalesOrderLineNumber", IntegerType()),
7       StructField("OrderDate", DateType()),
8       StructField("CustomerName", StringType()),
9       StructField("Email", StringType()),
10      StructField("Item", StringType()),
11      StructField("Quantity", IntegerType()),
12      StructField("UnitPrice", FloatType()),
13      StructField("Tax", FloatType())
14      ])
15
16  # Import all files from bronze folder of lakehouse
17  df = spark.read.format("csv").option("header", "true").schema(orde
18
19  # Display the first 10 rows of the dataframe to preview your data
20  display(df.head(10))
21
```

✓ 14 sec -Apache Spark session ready in 9 sec 788 ms. Command executed in 3 sec 833   PySpark (Python) ⌄

The output:



b. Now adding columns for data validation and cleanup. The first line of the code imports the necessary functions from PySpark.  then adding new columns to the dataframe to track the source file name, whether the order was flagged as

being a before the fiscal year of interest, and when the row was created and modified.

```python
from pyspark.sql.functions import when, lit, col, current_timestamp, input_file_name

# Add columns IsFlagged, CreatedTS and ModifiedTS
df = df.withColumn("FileName", input_file_name()) \
    .withColumn("IsFlagged", when(col("OrderDate") < '2019-08-01',True).otherwise(False)) \
    .withColumn("CreatedTS", current_timestamp()).withColumn("ModifiedTS", current_timestamp())

# Update CustomerName to "Unknown" if CustomerName null or empty
df = df.withColumn("CustomerName", when((col("CustomerName").isNull() | (col("CustomerName")=="")),lit("Unknown")).otherwise(col("CustomerN
```
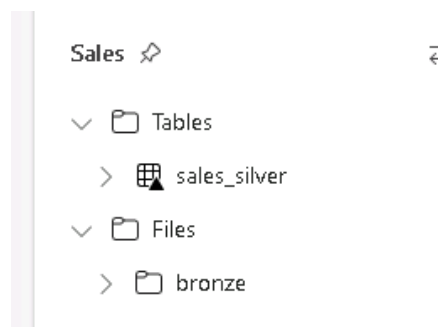
c. Now define the scheme for the sales silver table.

```python
# Define the schema for the sales_silver table

from pyspark.sql.types import *
from delta.tables import *

DeltaTable.createIfNotExists(spark) \
    .tableName("sales.sales_silver") \
    .addColumn("SalesOrderNumber", StringType()) \
    .addColumn("SalesOrderLineNumber", IntegerType()) \
    .addColumn("OrderDate", DateType()) \
    .addColumn("CustomerName", StringType()) \
    .addColumn("Email", StringType()) \
    .addColumn("Item", StringType()) \
    .addColumn("Quantity", IntegerType()) \
    .addColumn("UnitPrice", FloatType()) \
    .addColumn("Tax", FloatType()) \
    .addColumn("FileName", StringType()) \
    .addColumn("IsFlagged", BooleanType()) \
    .addColumn("CreatedTS", DateType()) \
    .addColumn("ModifiedTS", DateType()) \
    .execute()
```

PySpark (Python) ∨

d. New sales silver table was created. The triangle indicates that it's a delta table.

Sales ⌕

∨ ⊟ Tables
> ⊞ sales_silver
∨ ⊟ Files
> ⊟ bronze

e.  perform an upsert operation on a Delta table, updating existing records based on specific conditions and inserting new records when no match is found. Add a new code block. This operation is important because it enables you to update existing

records in the table based on the values of specific columns, and insert new records when no match is found.

```python
1   # Update existing records
2
3   from delta.tables import *
4
5   deltaTable = DeltaTable.forPath(spark, 'Tables/sale
6
7   dfUpdates = df
8
9   deltaTable.alias('silver') \
10      .merge(
11          dfUpdates.alias('updates'),
12          'silver.SalesOrderNumber = updates.SalesOrderNu
13      ) \
14      .whenMatchedUpdate(set =
15      {
16
17      }
18      ) \
19      .whenNotMatchedInsert(values =
20      {
21          "SalesOrderNumber": "updates.SalesOrderNumber
22          "SalesOrderLineNumber": "updates.SalesOrderLi
23          "OrderDate": "updates.OrderDate",
24          "CustomerName": "updates.CustomerName",
25          "Email": "updates.Email",
26          "Item": "updates.Item",
27          "Quantity": "updates.Quantity",
28          "UnitPrice": "updates.UnitPrice",
29          "Tax": "updates.Tax",
30          "FileName": "updates.FileName",
31          "IsFlagged": "updates.IsFlagged",
32          "CreatedTS": "updates.CreatedTS",
33          "ModifiedTS": "updates.ModifiedTS"
34      }
35      ) \
36      .execute()
```

4. Explore data in the silver layer using the SQL endpoint

    a.  This query calculates the total sales for each year in the sales_silver table.

b. This query calculates the total quantity of items purchased by each customer in the sales_silver table, and then returns the top 10 customers in terms of quantity.

```sql
SELECT TOP 10 CustomerName, SUM(Quantity) AS TotalQuantity
FROM sales_silver
GROUP BY CustomerName
ORDER BY TotalQuantity DESC
```

| | ABC CustomerName | 123 TotalQuantity |
|---|---|---|
| 1 | Samantha Jenkins | 41 |
| 2 | Henry Garcia | 39 |
| 3 | Charles Jackson | 38 |
| 4 | April Shan | 35 |
| 5 | Ryan Thompson | 34 |
| 6 | Hailey Patterson | 34 |
| 7 | Mason Roberts | 33 |
| 8 | Dalton Perez | 32 |
| 9 | Jason Griffin | 31 |
| 10 | Fernando Barnes | 28 |

4. Transform Data for Gold Layer

a. Create a new notebook named transform data for gold and add Sales Lakehouse, this will add sales_silver table under tables.
b. Load schema and then add a new code to define the schema

```python
from pyspark.sql.types import *
from delta.tables import*

# Define the schema for the dimdate_gold table
DeltaTable.createIfNotExists(spark) \
    .tableName("sales.dimdate_gold") \
    .addColumn("OrderDate", DateType()) \
    .addColumn("Day", IntegerType()) \
    .addColumn("Month", IntegerType()) \
    .addColumn("Year", IntegerType()) \
    .addColumn("mmmyyyy", StringType()) \
    .addColumn("yyyymm", StringType()) \
    .execute()
```

[2] ✓ 5 sec -Command executed in 4 sec 924 ms by User1-37089922 on 7:15:26 AM, 1/19/24          PySpark (Python) ∨

> Spark jobs (4 of 4 succeeded)   ▤ Log

<delta.tables.DeltaTable at 0x7f9cac2694e0>

```python
# Load data to the dataframe as a starting point to create the gold layer
df = spark.read.table("Sales.sales_silver")
```

[1] ✓ 19 sec -Apache Spark session ready in 9 sec 340 ms. Command executed in 10 sec 94 ms by User1-37089  PySpark (Python) ∨

> Spark jobs (1 of 1 succeeded)

c. Created 4 tables which was curated, modelled gold layer that can be used for reporting and analysis.

## 5. Create a database

    a.  create a new dataset that includes the gold tables.

B. set relationship.

The common relationship is sharing order date and it is many to 1.