

Portfolio Project:

Querying a Data Warehouse in Microsoft Fabric

Project Overview

This project focused on querying a data warehouse in Microsoft Fabric, which provides a relational database for large-scale analytics. The goal was to utilize SQL queries to extract insights from the data, ensuring data consistency and creating views for reporting purposes.

Objectives

- 1. Create a Workspace:**
 - Set up a workspace with the Fabric trial enabled.
- 2. Create a Sample Data Warehouse:**
 - Establish a sample data warehouse for analysis.
- 3. Query the Data Warehouse:**
 - Execute various SQL queries to analyze data.
- 4. Verify Data Consistency:**
 - Check for and handle inconsistent data.
- 5. Save as View:**
 - Create a view for filtered data reporting.

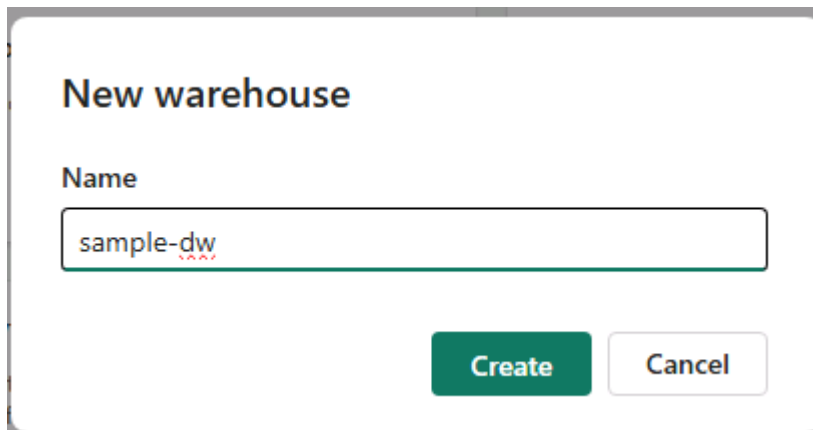
Experience

Create a Workspace

- Navigated to Microsoft Fabric Home and signed in.
- Selected Workspaces and created a new workspace.

Create a Sample Data Warehouse

- Selected Create, then Sample warehouse, and created a new data warehouse named sample-dw.



Query the Data Warehouse

- Opened a new SQL query and run the following code to get total trips and revenue by month:

```
1 SELECT
2   D.MonthName,
3   COUNT(*) AS TotalTrips,
4   SUM(T.TotalAmount) AS TotalRevenue
5 FROM dbo.Trip AS T
6 JOIN dbo.[Date] AS D
7   ON T.[DateID]=D.[DateID]
8 GROUP BY D.MonthName;
```



	MonthName	TotalTrips	TotalRevenue
1	September	231384	3562548
2	December	233767	3540532
3	November	236513	3547993
4	June	238432	3627186
5	May	249340	3794240
6	January	237842	3401334
7	April	246153	3679261
8	July	228327	3417355

- Ran another query to get average trip duration and distance by day of the week:

```

1 SELECT
2   D.DayName,
3   AVG(T.TripDurationSeconds) AS AvgDuration,
4   AVG(T.TripDistanceMiles) AS AvgDistance
5 FROM dbo.Trip AS T
6 JOIN dbo.[Date] AS D
7   ON T.[DateID]=D.[DateID]
8 GROUP BY D.DayName;

```

ABC	DayName	123 AvgDuration	12F AvgDistance
1	Tuesday	763	2.81239639559921
2	Thursday	796	19.3689186988126
3	Friday	790	2.84746039283699
4	Sunday	1147	3.21572794977795
5	Saturday	835	33.0886460668962
6	Monday	742	6.82623660465939
7	Wednesday	780	2.8184765031265

- Queried the top 10 most popular pickup and dropoff locations:

```

1 SELECT TOP 10
2   G.City,
3   COUNT(*) AS TotalTrips
4 FROM dbo.Trip AS T
5 JOIN dbo.Geography AS G
6   ON T.DropoffGeographyID=G.GeographyID
7 GROUP BY G.City
8 ORDER BY TotalTrips DESC;

```

ABC	City	123 TotalTrips
1	Manhattan	1523186
2	New York	595638
3	Brooklyn	153569
4	Prince	92711
5	Flushing	80533
6	Queens	58181
7	Planetarium	52907
8	Brooklyn Heights	47909

Verify Data Consistency

- Checked for trips with unusually long duration:

```

1 -- Check for trips with unusually long duration
2 SELECT COUNT(*) FROM dbo.Trip WHERE TripDurationSeconds > 86400; -- 24 hours

```

123	untitled1
1	49

- Checked for trips with negative trip duration:

```

1 -- Check for trips with negative trip duration
2 SELECT COUNT(*) FROM dbo.Trip WHERE TripDurationSeconds < 0;

```

123	untitled1
1	4

- Removed trips with negative trip duration:

```





-- Remove trips with negative trip duration
DELETE FROM dbo.Trip WHERE TripDurationSeconds < 0;



```

Save as View

- Created a view for January trips:

```
1 SELECT
2     D.DayName,
3     AVG(T.TripDurationSeconds) AS AvgDuration,
4     AVG(T.TripDistanceMiles) AS AvgDistance
5 FROM dbo.Trip AS T
6 JOIN dbo.[Date] AS D
7     ON T.[DateID]=D.[DateID]
8 WHERE D.Month = 1
9 GROUP BY D.DayName;
```

Messages Results    

 Search 

	ABC DayName	123 AvgDuration	12f AvgDistance
1	Tuesday	691	2.85065705430534
2	Thursday	715	2.79291961119472
3	Friday	713	2.62851462288349
4	Sunday	654	3.04907995357095
5	Saturday	663	2.69297393033589
6	Monday	667	2.87905740768423
7	Wednesday	716	2.86763962644205

- Saved the view as vw_JanTrip.

Save as view

This will save the text of your SQL query as a view. Make sure the SQL syntax for the view is correct below.

Warehouse

sample-dw

Schema

dbo

View name *

vw_JanTrip

SQL statement

```
CREATE VIEW [dbo].[vw_JanTrip]
AS
SELECT
    D.DayName,
    AVG(T.TripDurationSeconds) AS AvgDuration,
    AVG(T.TripDistanceMiles) AS AvgDistance
FROM dbo.Trip AS T
JOIN dbo.[Date] AS D
    ON T.[DateID]=D.[DateID]
WHERE D.Month = 1
GROUP BY D.DayName
```

Copy to clipboard

OK

Cancel

Results

- ✓ A workspace was successfully created in Microsoft Fabric.
- ✓ A sample data warehouse named sample-dw was established and populated with sample data for analysis.
- ✓ SQL queries were executed to analyze data, revealing:
 - Total trips and revenue by month.
 - Average trip duration and distance by day of the week.
 - The top 10 most popular pickup and dropoff locations.

- ✓ Data consistency checks were performed, identifying and handling trips with negative durations.
- ✓ A view named vw_JanTrip was created to filter and report on January trip data.

Conclusion

This project provided a practical introduction to querying a data warehouse in Microsoft Fabric. Key insights were gained into the use of SQL for data analysis and the importance of verifying data consistency. The ability to create views enhanced the usability of the data warehouse for reporting purposes. Overall, this exercise demonstrated the capabilities of Microsoft Fabric in managing and analyzing large-scale data efficiently.

Resources

GitHub: <https://github.com/ThatoMTNG/Microsoft-Fabric-Analytics-Engineer-DP-600->

Mentions

Project Author: Thato Metsing (<https://www.linkedin.com/in/thatometsing/>)

Project Mentor: Maureen Direro (<https://www.linkedin.com/in/maureen-direro-46a6b220/>)