# MICROSOFT AZURE CAPSTONE PROJECT

## BY BHAVANA JAISWAL

# INTRODUCTION

- Azure Synapse Analytics is a cloud-based analytics service offered by Microsoft. It provides a workspace for data professionals to ingest, prepare, manage, and serve data for immediate business intelligence and machine learning needs. With Synapse, you can explore, analyze, and create reports on your data using a unified experience, leveraging a powerful combination of big data and data warehousing technologies. It allows you to connect to various data sources, including on-premises data stores, cloud data stores, and other external sources, and analyze them with a wide range of tools such as Apache Spark, SQL, and machine learning models. Synapse Analytics is a powerful platform that can help you scale up  our data analytics capabilities, improve decision-making, and gain valuable insights into our business.

# PROBLEM FIRST

- **Problem statement 1:** The task is to explore data analytics workspace by using Azure Synapse Analytics. You will create ADLS Gen 2 accounts and define the pipelines in Azure Synapse Analytics to transfer data from various data sources into the workspace for analysis. The data will be ingested in Azure Synapse with Built-in copy task option, and you will query the uploaded data.

# INTEGRATED WITH AZURE SYNAPSE ANALYTICS WORKSPACE

AZURE SYNAPSE ANALYTICS: A CLOUD-BASED ANALYTICS SERVICE THAT BRINGS TOGETHER BIG DATA AND DATA WAREHOUSING. IT USES APACHE SPARK AND PROVIDES INTEGRATION WITH OTHER AZURE SERVICES SUCH AS AZURE DATA FACTORY, AZURE DATA LAKE STORAGE, AND POWER BI.

PRE-REQUISITES:  MICROSOFT ACCOUNT, RESOURCE GROUP, SYNAPSE WORKSPACE

# syn-rg
Resource group

Search

- Overview
- Activity log
- Access control (IAM)
- Tags
- Resource visualizer
- Events

**Settings**

- Deployments
- Security
- Policies
- Properties
- Locks

**Cost Management**

- Cost analysis
- Cost alerts (preview)
- Budgets
- Advisor recommendations

---

+ Create    Manage view ⌄    🗑 Delete resource group    ↻ Refresh    ⬇ Export to CSV    ⚙ Open query    |    Assign tags    → Move ⌄    🗑 Delete    ⬇ Export template    •••

⌄ Essentials                                                                    JSON View

**Resources**    Recommendations

Filter for any field...    Type equals **all** ✕    Location equals **all** ✕    + Add filter

Showing 1 to 1 of 1 records.    ☐ Show hidden types ⓘ    No grouping ⌄    ☰ List view ⌄

| ☐ Name ↑↓ | Type ↑↓ | Location ↑↓ |
|---|---|---|
| ☐ 🔷 projectworkk | Synapse workspace | South Central US  ••• |

This is my synapse workspace.

< Previous    Page 1 ⌄ of 1    Next >

👤 Give feedback

**Microsoft Azure**

Search resources, services, and docs (G+/)

bhavana.66@outlook.co...
DEFAULT DIRECTORY

Home > syn-rg >

## projectworkk
Synapse workspace

Search

- Overview
- Activity log
- Access control (IAM)
- Tags
- Diagnose and solve problems

**Settings**

- Azure Active Directory
- Properties
- Locks

**Analytics pools**

- SQL pools
- Apache Spark pools
- Data Explorer pools (preview)

**Security**

- Encryption
- Networking
- Identity

+ New dedicated SQL pool    + New Apache Spark pool    + New Data Explorer pool (preview)    ↻ Refresh    ✎ Reset SQL admin password    | 🗑 Delete

⌃ Essentials                                                                                                    JSON View

| | | | |
|---|---|---|---|
| Resource group (move) | : syn-rg | Networking | : Show firewall settings |
| Status | : Succeeded | Primary ADLS Gen2 acco... | : https://storagedestinat.dfs.core.windows.net |
| Location | : South Central US | Primary ADLS Gen2 file s... | : startgen2 |
| Subscription (move) | : Azure for Students | SQL admin username | : sqladminuser |
| Subscription ID | : 5448780b-5af5-45de-812a-a121aa797fbd | SQL Active Directory ad... | : live.com#bhavana.66@outlook.com |
| Managed virtual network | : No | Dedicated SQL endpoint | : projectworkk.sql.azuresynapse.net |
| Managed Identity object ... | : a0ac1929-574d-4e22-9d10-ff6757b464cf | Serverless SQL endpoint | : projectworkk-ondemand.sql.azuresynapse.net |
| Workspace web URL | : https://web.azuresynapse.net?workspace=%2fsubscriptions%2f54... | Development endpoint | : https://projectworkk.dev.azuresynapse.net |
| Tags (edit) | : Click here to add tags | | |

### Getting started

**Open Synapse Studio**
Start building your fully-integrated analytics solution and unlock new insights.

Open ⧉

**Read documentation**
Learn how to be productive quickly. Explore concepts, tutorials, and samples.

Learn more ⧉

This is my Synapse studio ready to lauch.

### Analytics pools

Search to filter items...

Name                                Type                                Size

Microsoft Azure | Synapse Analytics ▶ projectworkk

bhavana.66@outlook.com
DEFAULT DIRECTORY

# Copy Data tool

✓ Properties

② Source

    Dataset

    Configuration

③ Destination

④ Settings

⑤ Review and finish

File form

**File format**

Delimited

**Column del**

Comma (,)

☐ Edit

**Row delimi**

Default (\r

☐ Edit

☑ First row

> Advance

**Compressio**

None

**Additional**

+ New

## Preview data

Linked service: Http_to_CSV

Object: https://raw.githubusercontent.com/MicrosoftLearning/DP-900T00A-Azure-Data-Fundamentals/master/Azure-Syna...

**Preview**     Schema

| ProductID | ProductName | Category | ListPrice |
|---|---|---|---|
| 771 | Mountain-100 Silver, 38 | Mountain Bikes | 3399.9900 |
| 772 | Mountain-100 Silver, 42 | Mountain Bikes | 3399.9900 |
| 773 | Mountain-100 Silver, 44 | Mountain Bikes | 3399.9900 |
| 774 | Mountain-100 Silver, 48 | Mountain Bikes | 3399.9900 |
| 775 | Mountain-100 Black, 38 | Mountain Bikes | 3374.9900 |
| 776 | Mountain-100 Black, 42 | Mountain Bikes | 3374.9900 |
| 777 | Mountain-100 Black, 44 | Mountain Bikes | 3374.9900 |
| 778 | Mountain-100 Black, 48 | Mountain Bikes | 3374.9900 |
| 779 | Mountain-200 Silver, 38 | Mountain Bikes | 2319.9900 |

Preview of my data which is successfully transform to analyze further.

< Previous    Next >      Cancel

Here my data is fully ingest into ADLS GEN 2 STORAGE ACCOUNT FOR FURTHER PROCESS.

Here is the source and sink location.

Microsoft Azure | Synapse Analytics ▶ projectworkk

Search

bhavana.66@outlook.com
DEFAULT DIRECTORY

Synapse live ⌄    Validate all    ⬆ Publish all 2

**Integrate**                    +  ⌄  «

🔍 Filter resources by name

▲ Pipelines                              1

    ⬚⬚ CopyPipeline_product

startgen2    SQL script 1    Notebook 2    Notebook 3    ⬚⬚ CopyPipeline_prod... ✕

**Activities**              ⌄  «
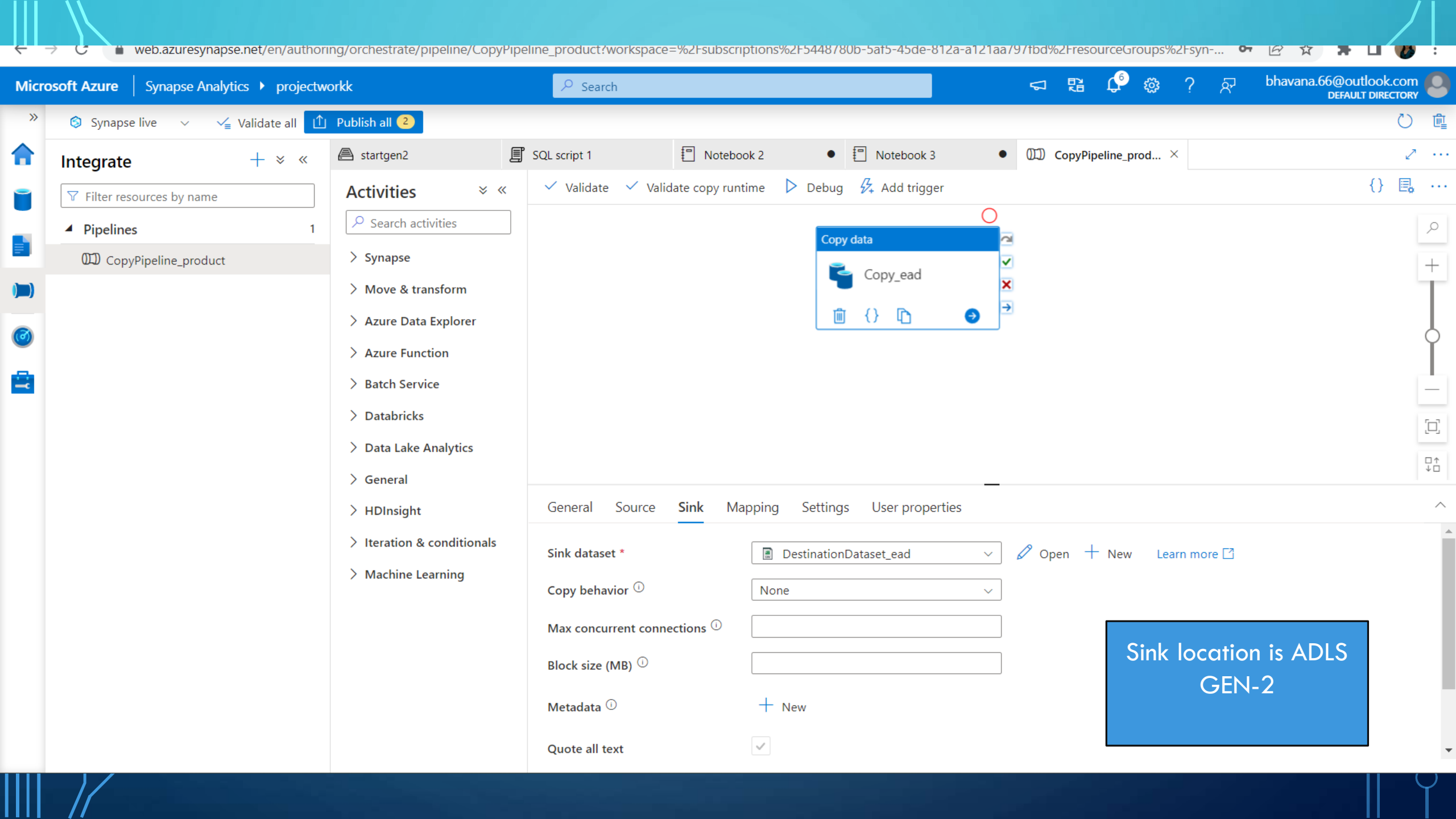
🔍 Search activities

> Synapse
> Move & transform
> Azure Data Explorer
> Azure Function
> Batch Service
> Databricks
> Data Lake Analytics
> General
> HDInsight
> Iteration & conditionals
> Machine Learning

✓ Validate    ✓ Validate copy runtime    ▶ Debug    ⚡ Add trigger

Copy data
  🗄 Copy_ead
  🗑 {} 📋 →

General    Source    **Sink**    Mapping    Settings    User properties

Sink dataset *            📄 DestinationDataset_ead  ⌄        ✏ Open   + New   Learn more ⧉

Copy behavior ⓘ          None                                 ⌄

Max concurrent connections ⓘ

Block size (MB) ⓘ

Metadata ⓘ              + New

Quote all text          ☑

Sink location is ADLS
GEN-2

bhavana.66@outlook.com
DEFAULT DIRECTORY

## Analytics pools

- SQL pools
- Apache Spark pools
- Data Explorer pools (preview)

## Activities

- SQL requests
- KQL requests
- Apache Spark applications
- Data flow debug

## Integration

- Pipeline runs
- Trigger runs
- Integration runtimes
- Link connections

# Pipeline runs

Triggered    Debug          ▶ Rerun        ⊘ Cancel options ⌄        ↻ Refresh        ☰ Edit columns        [ List    Gantt ]

▽ Filter by run ID or name          Chennai, Kolkata, Mu... : **Last 24 hours**          Pipeline name : **CopyPipeline_product**          Status : **All**

Runs : Latest runs          Triggered by : **All**          ▽ Add filter          ✕

🗐 Copy filters    ⬇ Export to CSV   | ⌄

Showing 1 - 1 items                                                                                                          Last refreshed 0 minutes ago

| ☐ | Pipeline name ↑↓ | Run start ↑↓ | Run end ↑↓ | Duration | Triggered by | Status ↑↓ | Run | Parameters |
|---|---|---|---|---|---|---|---|---|
| ☐ | CopyPipeline_product | 6/9/2023, 5:38:39 PM | 6/9/2023, 5:38:54 PM | 00:00:15 | Manual trigger | ✅ Succeeded | Original | |

Here is my pipeline is successfully running.

Microsoft Azure | Synapse Analytics ▸ projectworkk

bhavana.66@outlook.com
DEFAULT DIRECTORY

# Copy Data tool

- ✓ Properties
- ✓ Source
- ✓ Destination
- ✓ Settings
- 5 **Review and finish**
  - Review
  - **Deployment**



HTTP ──────────────▶ Azure Data Lake Storage Gen2

In this place the data is fully stored into ADLS gen2

## Deployment complete

| Deployment step | Status |
|---|---|
| Validating copy runtime environment | ✓ Succeeded |
| ❯ Creating datasets | ✓ Succeeded |
| ❯ Creating pipelines | ✓ Succeeded |
| ❯ Running pipelines | ✓ Succeeded |

Datasets and pipelines have been created. You can now monitor and edit the copy pipelines or click finish to close Copy Data Tool.

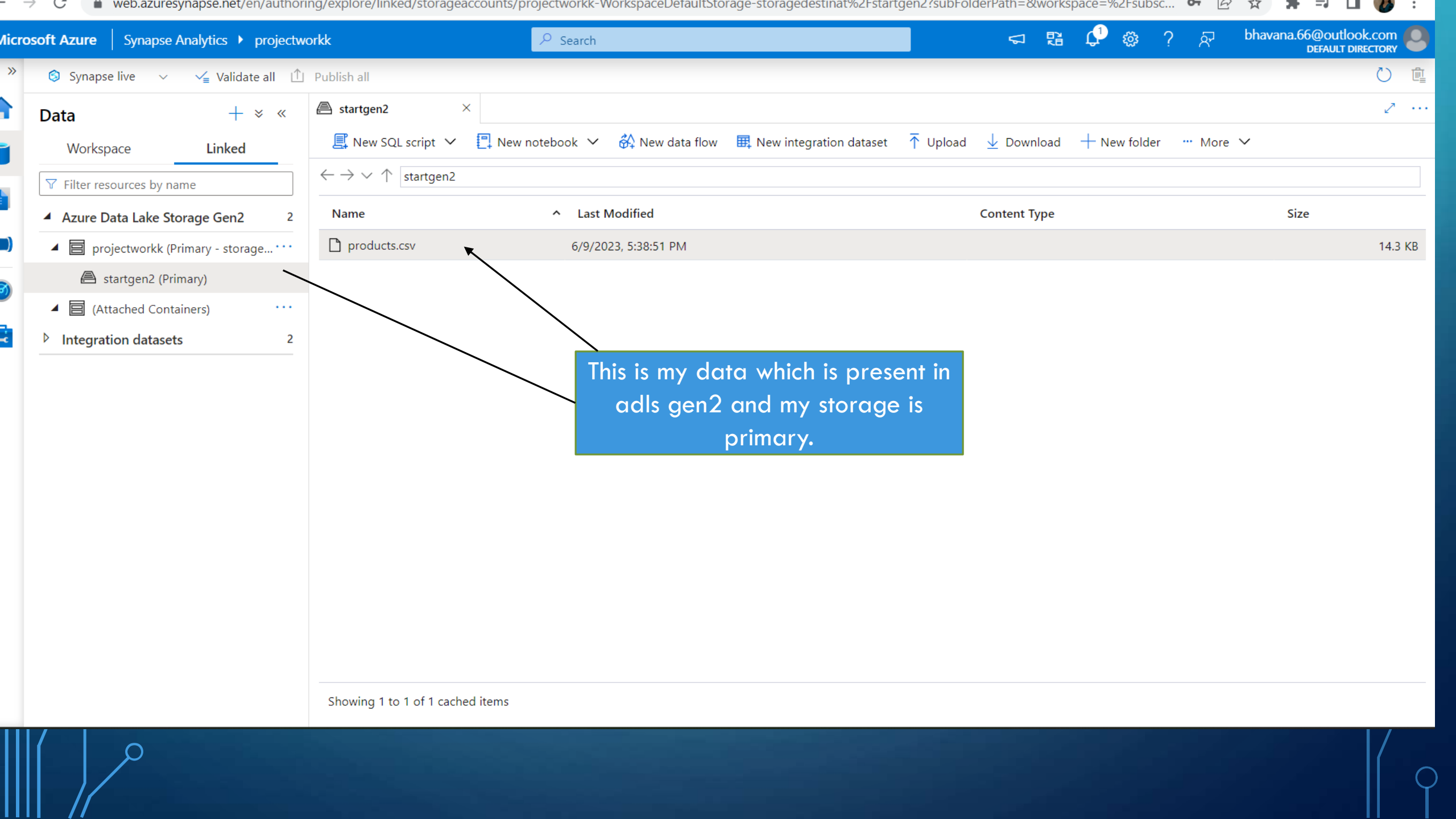[ Finish ]    [ Edit pipeline ]    [ Monitor ]

# PROBLEM SECOND

- **PROBLEM STATEMENT 2:** YOU NEED TO FIND THE TOP 100 ROWS FROM NEW SQL SCRIPT FROM YOUR DATA IN THE WORKSPACE, RUN THE CODE AND CHECK THE RESULT DATASETS. THEN, UPDATE THE QUERY BY SELECTING THE CATEGORY AND COUNT AS PRODUCT NUMBERS. FINALLY, MAKE NECESSARY CHANGES TO THE CHART VIEW.

# AZURE DATA ANALYTICS WITH SQL POOL ENVIRONMENT

- CLOUD-BASED DATA WAREHOUSING SOLUTION FOR STORING AND MANAGING LARGE AMOUNTS OF DATA.

- HIGHLY SCALABLE AND SECURE ENVIRONMENT FOR ANALYZING LARGE DATASETS.

- PROVIDES ADVANCED SECURITY FEATURES AND SEAMLESS INTEGRATION WITH OTHER AZURE SERVICES.

# TECH STACK USED IN 2<sup>ND</sup> PROBLEM

- NEW SQL SCRIPT: A SCRIPT WRITTEN IN SQL (STRUCTURED QUERY LANGUAGE) USED TO QUERY DATA FROM A DATABASE.

- WORKSPACE: A DATA ANALYTICS WORKSPACE IN AZURE SYNAPSE ANALYTICS WHERE DATA CAN BE INGESTED, TRANSFORMED, AND ANALYZED.

- CHART VIEW: A GRAPHICAL REPRESENTATION OF DATA THAT CAN BE CREATED IN AZURE SYNAPSE ANALYTICS.

Microsoft Azure | Synapse Analytics ▸ projectworkk

Search

bhavana.66@outlook.com
DEFAULT DIRECTORY

Synapse live ∨ | ✓ Validate all | ↑ Publish all

**Data** + ⌄ «

Workspace | Linked

Filter resources by name

◢ Azure Data Lake Storage Gen2   2

  ◣ 🗎 projectworkk (Primary - storage... • • •

    🗄 startgen2 (Primary)

 ◣ 🗎 (Attached Containers) • • •

▷ Integration datasets   2

🗄 startgen2 ✕

🖳 New SQL script ∨   📄 New notebook ∨   🔀 New data flow   ▦ New integration dataset   ↑ Upload   ↓ Download   + New folder   ⋯ More ∨

← → ⌄ ↑ | startgen2

| Name | | Last Modified | Content Type | Size |
|------|---|---------------|--------------|------|
| 📄 products.csv | | 6/9/2023, 5:38:51 PM | | 14.3 KB |

This is my data which is present in adls gen2 and my storage is primary.

Showing 1 to 1 of 1 cached items

Microsoft Azure | Synapse Analytics ▶ projectworkk

Search

bhavana.66@outlook.com
DEFAULT DIRECTORY

Synapse live ⌄    ✓ Validate all    ⬆ Publish all 1

**Data**    + ⌄ «

startgen2    📄 SQL script 1 ●

Workspace    **Linked**

▶ Run    ↺ Undo ⌄    🖉 Publish    Query plan    **Connect to**  ✓ Built-in ⌄    **Use database** master ⌄

▽ Filter resources by name

◢ Azure Data Lake Storage Gen2    2

◢ 📄 projectworkk (Primary - storage... ⋯
  ◣ 📄 startgen2 (Primary)

◣ 📄 (Attached Containers)    ⋯

▷ Integration datasets    2

```
1   -- This is auto-generated code
2   SELECT
3       TOP 100 *
4   FROM
5       OPENROWSET(
6           BULK 'https://storagedestinat.dfs.core.windows.net/startgen2/products.csv',
7           FORMAT = 'CSV',
8           PARSER_VERSION = '2.0'
9       ) AS [result]
10
```

This is built-in sql pool and top 100 row.

**Results**    Messages

View    [ Table    Chart ]    ⤷ Export results ⌄

🔍 Search

| C1 | C2 | C3 | C4 |
|---|---|---|---|
| ProductID | ProductName | Category | ListPrice |
| 771 | Mountain-100 Silver, 38 | Mountain Bikes | 3399.9900 |
| 772 | Mountain-100 Silver, 42 | Mountain Bikes | 3399.9900 |
| 773 | Mountain-100 Silver, 44 | Mountain Bikes | 3399.9900 |
| 774 | Mountain-100 Silver, 48 | Mountain Bikes | 3399.9900 |

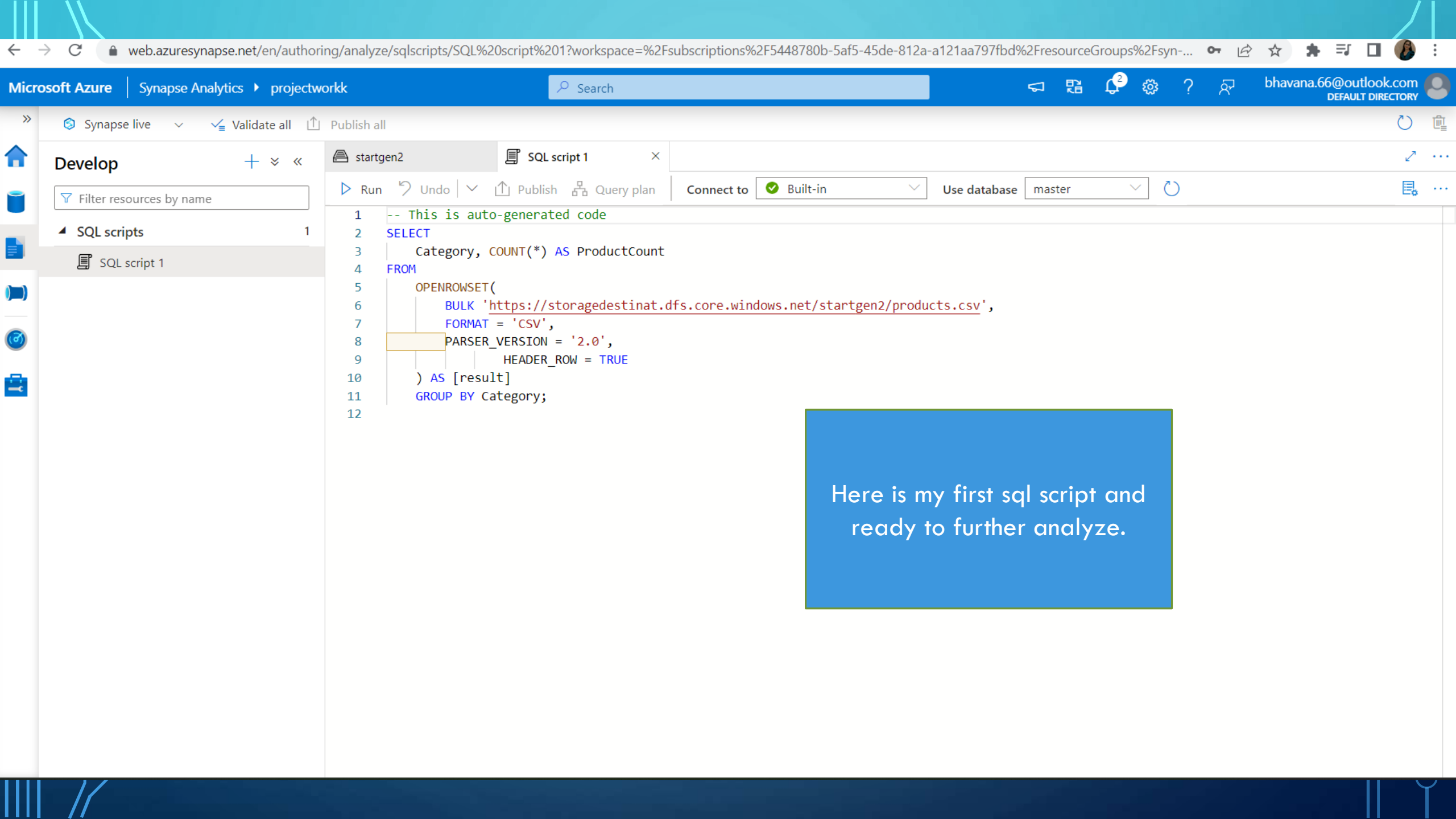✓ 00:00:15 Query executed successfully.

**Properties**

General    Related (0)

Name *
SQL script 1

Description

Type
.sql script

Size
232 bytes

Results settings per query ⓘ
◉ First 5000 rows (default)
○ All rows

Microsoft Azure | Synapse Analytics ▶ projectworkk

Search

bhavana.66@outlook.com
DEFAULT DIRECTORY

Synapse live ∨ | ✓ Validate all | ⬆ Publish all 1

Data | + ≫ ≪

🗄 startgen2 | 📄 SQL script 1 ●

Workspace | Linked

▷ Run | ↩ Undo ∨ | ⬆ Publish | Query plan | Connect to | ✅ Built-in ∨ | Use database | master ∨

Filter resources by name

◢ Azure Data Lake Storage Gen2    2

  ◢ 📄 projectworkk (Primary - storage... ⋯

    📄 startgen2 (Primary)

  ◢ 📄 (Attached Containers) ⋯

▷ Integration datasets    2

```
1   -- This is auto-generated code
2   SELECT
3       TOP 100 *
4   FROM
5       OPENROWSET(
6           BULK 'https://storagedestinat.dfs.core.windows.net/startgen2/products.csv',
7           FORMAT = 'CSV',
8           PARSER_VERSION = '2.0',
9               HEADER_ROW = TRUE
10      ) AS [result]
11
```

In this place I put HEADER_ROW = TRUE To REMOVE C1,C2 IN THE header of the column

Results    Messages

View | Table | Chart | → Export results ∨

Search

| ProductID | ProductName | Category | ListPrice |
|-----------|-------------|----------|-----------|
| 771 | Mountain-100 Silver, 38 | Mountain Bikes | 3399.99 |
| 772 | Mountain-100 Silver, 42 | Mountain Bikes | 3399.99 |
| 773 | Mountain-100 Silver, 44 | Mountain Bikes | 3399.99 |
| 774 | Mountain-100 Silver, 48 | Mountain Bikes | 3399.99 |
| 775 | Mountain-100 Black, 38 | Mountain Bikes | 3374.99 |

✅ 00:00:02 Query executed successfully.

**Properties**

General    Related (0)

Name *
SQL script 1

Description

Type
.sql script

Size
232 bytes

Results settings per query ⓘ
◉ First 5000 rows (default)
○ All rows

Microsoft Azure | Synapse Analytics ▸ projectworkk

Search

bhavana.66@outlook.com
DEFAULT DIRECTORY

Synapse live ⌄    ✓ Validate all    ⬆ Publish all

**Develop**    + ⌄ «

🔍 Filter resources by name

◢ SQL scripts                                    1

    📄 SQL script 1

startgen2        📄 SQL script 1    ✕

▷ Run    ↶ Undo ⌄    ⬆ Publish    ⧉ Query plan    Connect to  ✅ Built-in ⌄    Use database  master ⌄

```
1   -- This is auto-generated code
2   SELECT
3       Category, COUNT(*) AS ProductCount
4   FROM
5       OPENROWSET(
6           BULK 'https://storagedestinat.dfs.core.windows.net/startgen2/products.csv',
7           FORMAT = 'CSV',
8           PARSER_VERSION = '2.0',
9               HEADER_ROW = TRUE
10      ) AS [result]
11      GROUP BY Category;
12
```

Here is my first sql script and ready to further analyze.

This page is a screenshot of the Azure Synapse Analytics web interface showing a SQL script and chart visualization.

Microsoft Azure | Synapse Analytics ▶ projectworkk

Search

bhavana.66@outlook.com
DEFAULT DIRECTORY

Synapse live    Validate all    Publish all

**Develop**

Filter resources by name

SQL scripts                                          1

SQL script 1

startgen2                    SQL script 1

▷ Run    Undo    Publish    Query plan    Connect to    Built-in    Use database    master

```
1   -- This is auto-generated code
2   SELECT
3       Category, COUNT(*) AS ProductCount
4   FROM
5       OPENROWSET(
6           BULK 'https://storagedestinat.dfs.core.windows.net/startgen2/products.csv',
7           FORMAT = 'CSV',
8           PARSER_VERSION = '2.0',
9               HEADER_ROW = TRUE
10      ) AS [result]
11      GROUP BY Category;
```

Here I have analyze my data and make visualization of it. As per the given in the task.

Results    Messages

View    Table    Chart    Save as image

Chart type
Column

Category column
Category

Legend (series) columns
ProductCount

Legend position:
bottom - center

Legend (series) label

ProductCount

⊘ 00:00:02 Query executed successfully.

# PROBLEM THIRD

- PROBLEM STATEMENT 3:

- THE TASK INVOLVES FINDING THE TOP 100 ROWS FROM NEW SQL SCRIPT FROM THE DATA IN THE WORKSPACE, RUNNING THE CODE, AND CHECKING THE RESULT DATASETS. THEN, THE QUERY WILL BE UPDATED BY SELECTING THE CATEGORY AND COUNT AS PRODUCT NUMBERS. FINALLY, NECESSARY CHANGES WILL BE MADE TO THE CHART VIEW.

# Azure Data Analytics with Spark Pool Environment.

- Cloud-based big data processing solution using Apache Spark.

- Fully managed Spark environment for processing large datasets.

- Highly scalable and integrates seamlessly with other Azure services.

# TECH STACK USED IN 3RD PROBLEM

- SPARK POOL: A MANAGED APACHE SPARK SERVICE IN AZURE SYNAPSE ANALYTICS THAT ALLOWS YOU TO PROCESS BIG DATA WORKLOADS. IT CAN BE USED FOR DATA TRANSFORMATION, MACHINE LEARNING, AND DATA VISUALIZATION.

- APACHE SPARK: AN OPEN-SOURCE DISTRIBUTED COMPUTING SYSTEM USED FOR PROCESSING LARGE DATA SETS. IT IS DESIGNED TO BE FAST, EASY TO USE, AND SCALABLE.

- MANAGE HUB: A CENTRALIZED MANAGEMENT PORTAL IN AZURE SYNAPSE ANALYTICS WHERE YOU CAN MANAGE RESOURCES SUCH AS SQL POOLS, SPARK POOLS, AND DATA FLOWS.

Synapse live ⌄    ✓ Validate all    ⬆ Publish all

**Apache Spark pool**

Apache Spark pools can be tuned to run different kinds of Apache Spark workloads using specific configuration libraries, permissions, etc. Learn more ⬏

Analytics pools

- SQL pools
- **Apache Spark pools**
- Data Explorer pools (pre...

➕ New    ↻ Refresh

⧩ Filter by name

Showing 1-1 of 1 item

External connections

- Linked services
- Microsoft Purview

| Name | Node size family | Size |
|------|------------------|------|
| sparkpool | Memory Optimized | Small (4 vCores / 32 GB) - 3 to 5 nodes |

Integration

- Triggers
- Integration runtimes

Security

- Access control
- Credentials
- Managed private endpoi...

This is my spark pool and here also I can do the same thing as I do in sql pool environment.

Configurations + libraries

- Workspace packages
- Data flow libraries
- Apache Spark configurat...

Source control

Microsoft Azure | Synapse Analytics ▸ projectworkk

Search

bhavana.66@outlook.com
DEFAULT DIRECTORY

Synapse live ⌄    ✓ Validate all    📤 Publish all 3

**Data**    + ⌄ «

Workspace    **Linked**

🔻 Filter resources by name

🔺 Azure Data Lake Storage Gen2    2

🔻🔺 projectworkk (Primary - storage... ···
   📦 startgen2 (Primary)

🔻🔺 (Attached Containers)    ···

🔺 Integration datasets    2
   📋 DestinationDataset_ead
   📋 SourceDataset_ead

📦 startgen2    📄 SQL script 1    📖 Notebook 1 ●    📖 Notebook 2 ●    📖 **Notebook 3** ●

▷ Run all ⌄    ↶ Undo ⌄    📤 Publish    ☰ Outline    Attach to  sparkpool ⌄    Language  PySpark (Python) ⌄    ▦ Variables

● Ready

```
1  %%pyspark
2  df = spark.read.load('abfss://startgen2@storagedestinat.dfs.core.windows.net/products.csv', format='csv'
3  ## If header exists uncomment line below
4  ##, header=True
5  )
6  display(df.limit(10))
```

[10]    ⊗ <1 sec - Failed to create session

···

InvalidHttpRequestToLivy: Your Spark job requested 48 vcores. However, the workspace has a 0 core limit. Try reducing the numbers of vcores requested or increasing your vcore quota. HTTP status code: 400. Trace ID: 1b3ba1d1-854d-44bb-b877-f75e9852b35f.

+ Code    + Markdown

Getting error even I changed the core multiple times higher or lower also. But giving me same message.

# FINAL LOCATION OF THE DATA PRESENT IN ADLS_GEN2 STORAGE ACCOUNT

THE DATA IS FINALLY COME TO MY ADLS-GEN2

STORAGE ACCOUNT .

As you can see my transformed data stored in the destination path.

This is the preview of the stored data.

# Thank you