# Case Study Documentation
# Retail Sales Analysis

**NAME: HARSHA M R**

# Objective

Objective of the Case Study:

- To build an end-to-end data pipeline for retail sales data.

- To ingest CSV data from Azure Blob Storage into Databricks.

- To clean, transform, and summarize the dataset using PySpark.

- To store the transformed data in Snowflake as raw and summary tables.

- To visualize insights using Power BI.

# Dataset Description

Source: Azure Blob Storage

File: Retail_Sales.csv

# FLOW:

```
Azure Blob Storage (Retail_Sales.csv)
      │
      ▼
  Databricks (PySpark)
     ├── Data Cleaning
     ├── Date Transformation
     ├── Type Casting (Sales, Profit, Quantity)
     └── Summary Table (Category-wise Sales)
         │
         ▼
  Snowflake (Data Warehouse)
     ├── RETAIL_SALES_RAW
     └── RETAIL_SALES_SUMMARY
         │
```

```
▼
Power BI
├── Sales Trend Analysis
├── Category-wise Sales
└── Top Products Analysis
```

## Environment Setup

Tools & Platforms Used:

- Databricks: Apache Spark environment for data processing

- Snowflake: Cloud Data Warehouse for storing raw and summary tables

- Azure Blob Storage: Source for CSV dataset

- Power BI Desktop: Visualization tool

- Python (PySpark): For data transformation

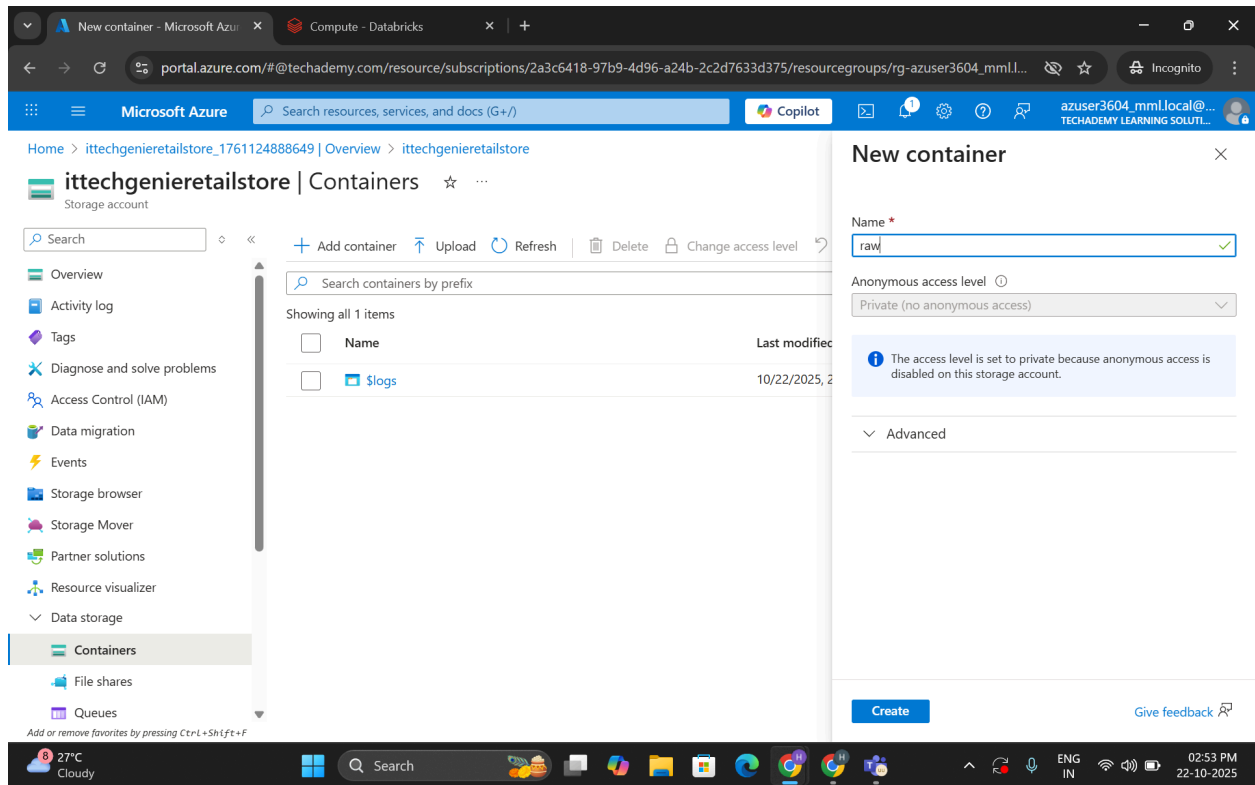## Step 1: Azure Blob Storage Setup

1. **Create Storage Account** in Azure Portal:

   ○ Name: ittechgenieretailstore

   ○ Type: Standard, StorageV2

   ○ Region: Central India

## 2. **Create Container**:

- Container Name: `raw`

- Public access level: Private



## 3. Upload CSV:

- File: Retail_Sales.csv

- Path: raw/Retail_Sales.csv

## Step 2: Load CSV into Databricks

```python
from snowflake.snowpark import Session

from snowflake.snowpark.functions import col, trim, upper, to_date, sum as
_sum



connection_parameters = {

    "account": "EBDQDYC-NB21316",

    "user": "HARSHA",

    "password": "HarshaRadhakrishnan02",

    "role": "ACCOUNTADMIN",

    "warehouse": "CASE_WH",

    "database": "ITTECHGENIE_DB",

    "schema": "PUBLIC"

}

session = Session.builder.configs(connection_parameters).create()



columns = [

    "OrderID", "OrderDate", "MonthOfSale", "CustomerID", "CustomerName",

    "Country", "Region", "City", "Category", "Subcategory",

    "Quantity", "Discount", "Sales", "Profit"

]
```

```python
df = session.read.option("skip_header",
1).csv("@retail_stage/Retail_Sales.csv")



df = df.to_df(*columns)



df_clean = (

    df.with_column("OrderDate", to_date(trim(col("OrderDate")), "YYYY-MM-DD"))

    .with_column("Category", upper(col("Category")))

    .with_column("Quantity", col("Quantity").cast("INTEGER"))

    .with_column("Discount", col("Discount").cast("FLOAT"))

    .with_column("Sales", col("Sales").cast("FLOAT"))

    .with_column("Profit", col("Profit").cast("FLOAT"))

)



df_summary =
df_clean.group_by("Category").agg(_sum(col("Sales")).alias("Category_Sales"))



df_clean.write.save_as_table("RETAIL_SALES_RAW", mode="overwrite")

df_summary.write.save_as_table("RETAIL_SALES_SUMMARY", mode="overwrite")
```

## Explanation:

- Removed spaces from OrderDate.

- Converted OrderDate to DATE type (yyyy-MM-dd).

- Uppercased Category.

- Casted numeric columns to correct types.

- Created summary table by category.

## Step 4: Snowflake Setup

- Create Warehouse: CASE_WH
- Create Database: ITTECHGENIE_DB
- Create Schema: PUBLIC
- Create Stage

```
ITTECHGENIE_DB.PUBLIC ∨        Settings ∨                                                    ⬈  Ope

 1    CREATE OR REPLACE DATABASE ITTECHGENIE_DB;
 2    USE DATABASE ITTECHGENIE_DB;
 3
 4    CREATE OR REPLACE SCHEMA PUBLIC;
 5
 6    CREATE OR REPLACE WAREHOUSE CASE_WH
 7      WITH WAREHOUSE_SIZE = 'XSMALL'
 8      AUTO_SUSPEND = 60
 9      AUTO_RESUME = TRUE
10      INITIALLY_SUSPENDED = TRUE
11      COMMENT = 'Warehouse for Databricks integration';
12
```

```
CREATE OR REPLACE FILE FORMAT csv_file_format
TYPE = CSV
FIELD_DELIMITER = ','
SKIP_HEADER = 1
EMPTY_FIELD_AS_NULL = TRUE
TRIM_SPACE = TRUE;

CREATE OR REPLACE STAGE retail_stage
URL='azure://ittechgenieretailstore.blob.core.windows.net/raw'
CREDENTIALS=(AZURE_SAS_TOKEN='?sv=2024-11-04&ss=bfqt&srt=sco&sp=rwdlacupiytfx&se=2025-10-
22T17:29:56Z&st=2025-10-
22T09:14:56Z&spr=https&sig=QUglqNAlSNQm%2FGNspwCv8yBNmQDLBYtMWauDolaUO84%3D')
FILE_FORMAT = csv_file_format;

LIST @retail_stage;
```

# Step 6: Verify Data in Snowflake

## Step 7: Power BI Connection

- Open Power BI Desktop → Get Data → Snowflake.
- Enter Server, Warehouse, Database, Schema.
- Load RETAIL_SALES_RAW and RETAIL_SALES_SUMMARY.

File    Home    Insert    Modeling

Share ▾

Get Data                                                    ✕

Search                        All

☐ All                ⬛ MySQL database
☐ File               ⬛ PostgreSQL database
☐ Database           ⬛ Sybase database
☐ Microsoft Fabric   ⬛ Teradata database
☐ Power Platform     ⬛ SAP HANA database
☐ Azure              ⬛ SAP Business Warehouse Application Server
☐ Online Services    ⬛ SAP Business Warehouse Message Server
☐ Other              ⬛ Amazon Redshift
                     ◆ Impala
                     ⬛ Google BigQuery
                     ⬛ Google BigQuery (Microsoft Entra ID)
                     ∨ Vertica
                     ❄ Snowflake
                     ⬛ Essbase
                     ▲ AtScale Models
                     ⬛ Power BI semantic models

Certified Connectors    Template Apps          Connect      Cancel

# Snowflake

Server

EBDQDYC-NB21316.snowflakecomputing.com

Warehouse

CASE_WH

▷ Advanced options

OK      Cancel

# Step 8: Create Power BI Visuals

**Screen 1 — Pie Chart**

Sum of PROFIT, Sum of QUANTITY and Sum of DISCOUNT by SUBCATEGORY and CATEGORY

SUBCATEGORY
- Laptops
- Networking
- Tables
- Storage
- Mobiles
- Appliances
- Furnishings
- Printers
- Paper
- Art
- Binders

Pie slice labels:
- 61.31K (33.94%)
- 31.27K (17.31%)
- 18.49K (10.24%)
- 16.09K (8.91%)
- 15.3K (8.47%)
- 14.74K (8.16%)
- 9.59K (5.31%)
- 8.74K (4.84%)
- 2.26K
- 2.26K (1.25%)

Data panel:

Data

RETAIL_SALES_RAW
- ☑ CATEGORY
- ☐ CITY
- ☐ COUNTRY
- ☐ CUSTOMERID
- ☐ CUSTOMERNA...
- ☑ Σ DISCOUNT
- ☐ MONTHOFSALE
- ☐ ⊞ ORDERDATE
- ☐ ORDERID
- ☑ Σ PROFIT
- ☑ Σ QUANTITY
- ☐ REGION
- ☐ Σ SALES
- ☑ SUBCATEGORY

RETAIL_SALES_SUMM...
- ☐ CATEGORY
- ☐ CATEGORY_SAL...

Tabs: BARCHART-CATEGORY+SALES | **PIECHART-SUBCATEGORY+PROFIT** | Page 1 | Page 2 | +

---

**Screen 2 — Line Chart**

Sum of SALES by Year and CATEGORY

CATEGORY ● FURNITURE ● OFFICE SUPPLIES ● TECHNOLOGY

Y-axis: Sum of SALES (0.0M, 0.1M, 0.2M, 0.3M, 0.4M, 0.5M, 0.6M, 0.7M)
X-axis: Year (2024, 2025)

Data

RETAIL_SALES_RAW
- ☑ CATEGORY
- ☐ CITY
- ☐ COUNTRY
- ☐ CUSTOMERID
- ☐ CUSTOMERNA...
- ☐ Σ DISCOUNT
- ☐ MONTHOFSALE
- ☑ ⊞ ORDERDATE
- ☐ ORDERID
- ☐ Σ PROFIT
- ☐ Σ QUANTITY
- ☐ REGION
- ☑ Σ SALES
- ☐ SUBCATEGORY

RETAIL_SALES_SUMM...
- ☐ CATEGORY
- ☐ CATEGORY_SAL...

Tabs: BARCHART-CATEGORY+SALES | PIECHART-SUBCATEGORY+PROFIT | **Page 1** | Page 2 | +

of 4          60%

Sum of QUANTITY by CUSTOMERID

**Data**

- RETAIL_SALES_RAW
  - CATEGORY
  - CITY
  - COUNTRY
  - ☑ CUSTOMERID
  - CUSTOMERNA...
  - Σ DISCOUNT
  - MONTHOFSALE
  - ORDERDATE
  - ORDERID
  - Σ PROFIT
  - ☑ Σ QUANTITY
  - REGION
  - Σ SALES
  - SUBCATEGORY
- RETAIL_SALES_SUMM...
  - CATEGORY
  - CATEGORY_SAL...

BARCHART-CATEGORY+SALES | PIECHART-SUBCATEGORY+PROFIT | Page 1 | Page 2 | +

of 4

60%