

RETAIL SALES ANALYSIS CASE STUDY (CS-3)

Batch 4
Group 4

TEAM MEMBERS

Smriti Ravindran

Ayushi Jain

Srishti Vermani

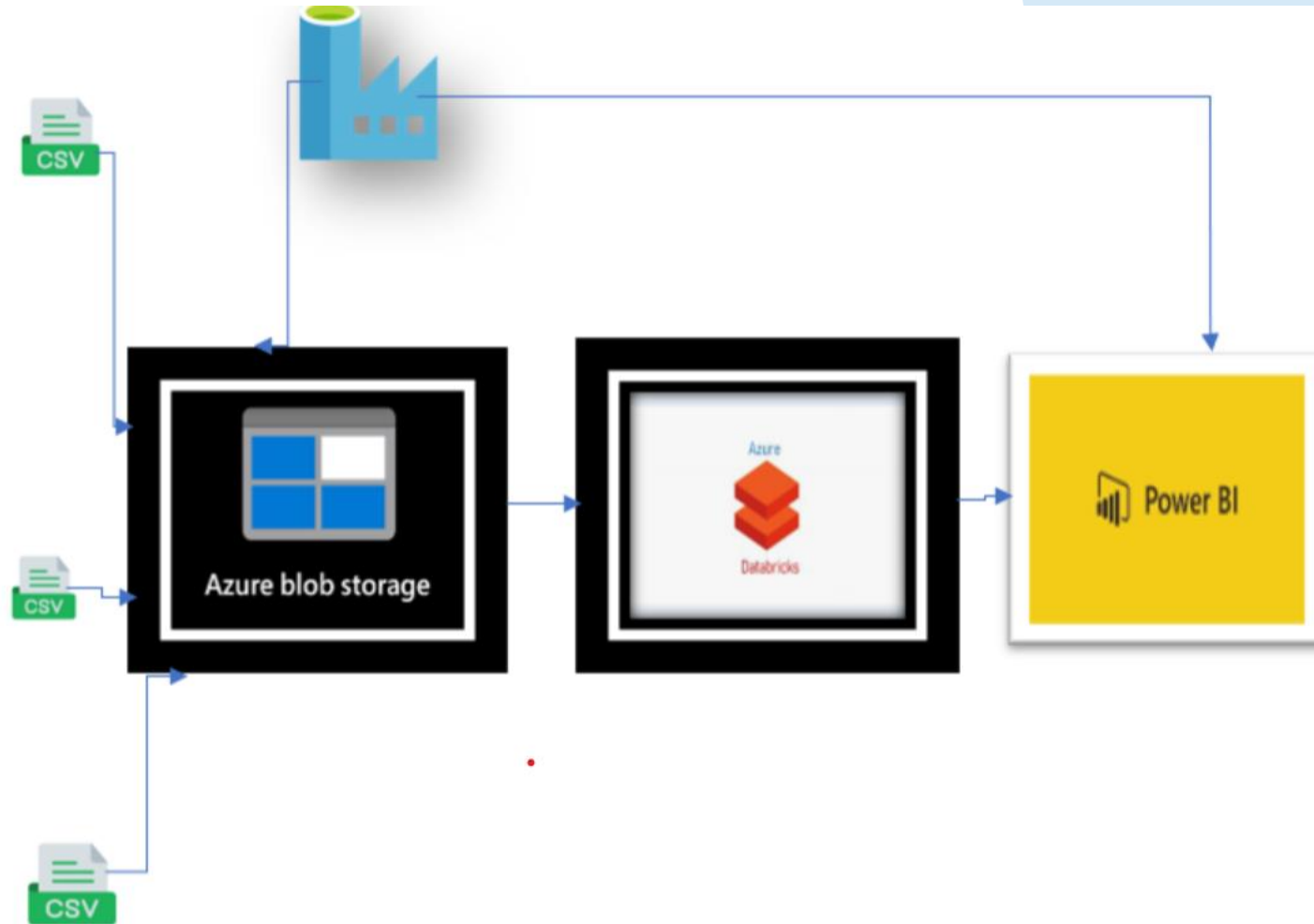
Shivam Jindal

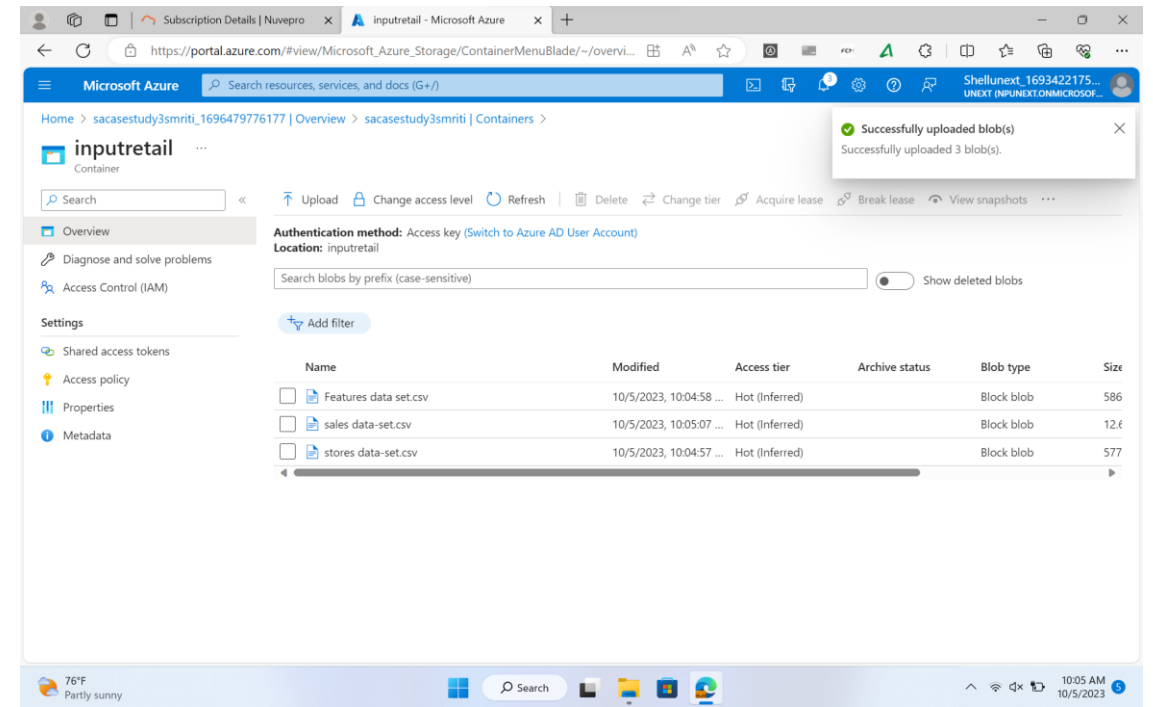
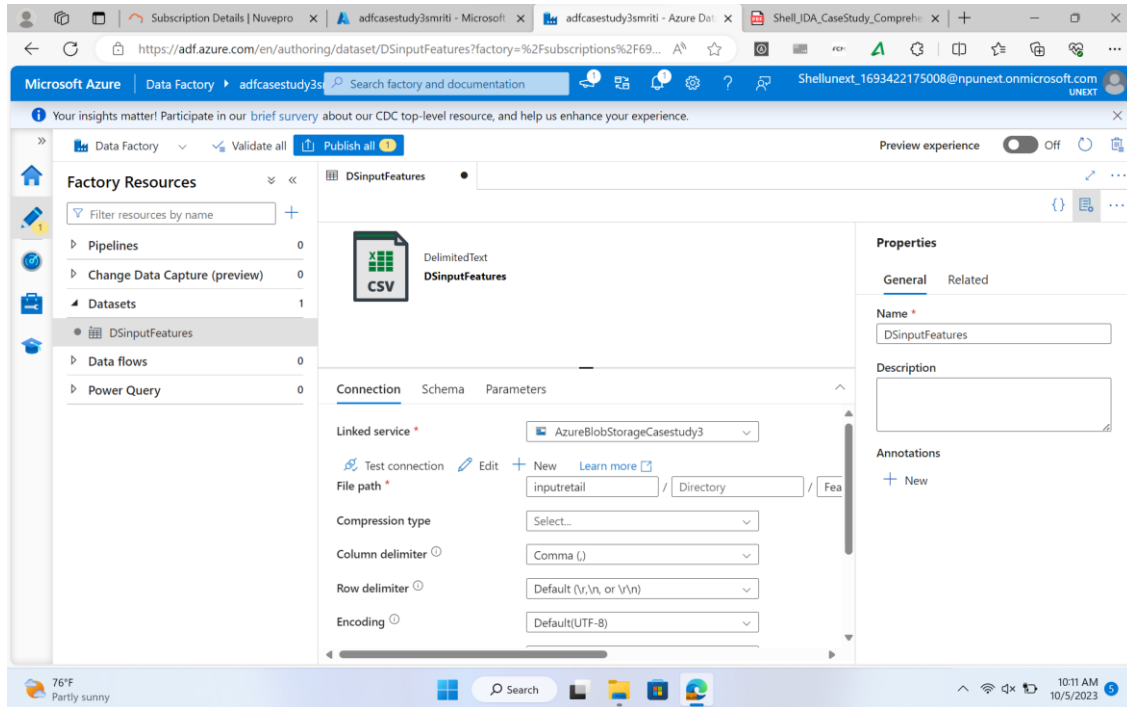
Shruti Katyal

In the realm of retail, making informed decisions based on historical data is crucial. The challenge arises when strategic choices are influenced by limited historical records, especially for events that occur annually, such as holidays and major promotions. This challenge is exacerbated by the fact that significant events, like holidays and markdowns, have a considerable impact on sales, making it necessary to anticipate the effects on specific departments and overall performance.

PROBLEM STATEMENT

DATA FLOW DIAGRAM





Stream 1 was created using Azure data factory. The above screenshots represent the same.

RAW DATA

Microsoft Azure

databricks

Search data, notebooks, recents, and more...

CTRL + P

ws05octshivam

shellunext_1693422177077@npune...

Data Science & Engi...

New

Workspace

Repos

Recents

Catalog

Compute

Workflows

Marketplace

0/3 Tasks Completed

Mounting_RetailSales_G4

Python

☆

File Edit View Run Help

Last edit was 3 hours ago

Provide feedback

▶ Run all

■ Terminated ▼

📅 Schedule

Share

^

Cmd 1

1 from pyspark.sql.functions import col,sum,isnan,when,count

Command took 0.02 seconds -- by shellunext_1693422177077@npunext.onmicrosoft.com at 10/5/2023, 1:39:50 PM on Shell102 Unext's Cluster

Cmd 2

1 strSource="wasbs://input@sa05octshivam.blob.core.windows.net/"

2 strMountPoint="/mnt/mp2"

3 strKey="fs.azure.account.key.sa05octshivam.blob.core.windows.net"

4 strValue="zNp5NWdcehwp55bwhu8KtcJa995RNonyglvLW9MJHvmG088g6tMP0lkx7Ec10AZu+9lcLrn86joW+ASTShqUdw=="

5

6 result=dbutils.fs.mount(

7 source=strSource,

8 mount_point=strMountPoint,

9 extra_configs={strKey:strValue}

10)

java.rmi.RemoteException: java.lang.IllegalArgumentException: requirement failed: Directory already mounted: /mnt/mp2; nested exception is:

Command took 1.18 seconds -- by shellunext_1693422177077@npunext.onmicrosoft.com at 10/5/2023, 12:32:37 PM on Shell102 Unext's Cluster

Cmd 3

1 %fs ls

Mounting_RetailSales_G4

Python

☆

File Edit View Run Help

Last edit was 3 hours ago

Provide feedback

▶ Run all

■ Terminated ▼

📅 Schedule

Share

^

1 %fs ls /mnt/mp2

Table +

New result table: OFF ▼

	path	name	size	modificationTime
1	dbfs:/mnt/mp2/Features data set.csv	Features data set.csv	600478	1696476990000
2	dbfs:/mnt/mp2/sales data-set.csv	sales data-set.csv	13264115	1696476992000
3	dbfs:/mnt/mp2/stores data-set.csv	stores data-set.csv	577	1696476989000

3 rows | 0.32 seconds runtime

Refreshed yesterday

Command took 0.32 seconds -- by shellunext_1693422177077@npunext.onmicrosoft.com at 10/5/2023, 1:32:57 PM on Shell102 Unext's Cluster

Cmd 6

1 df1=spark.read.csv("dbfs:/mnt/mp2/stores data-set.csv",header=True,inferSchema=True)

2 df2=spark.read.csv("dbfs:/mnt/mp2/sales data-set.csv",header=True,inferSchema=True)

3 df3=spark.read.csv("dbfs:/mnt/mp2/Features data set.csv",header=True,inferSchema=True)

▶ (6) Spark Jobs

▶ df1: pyspark.sql.dataframe.DataFrame = [Store: integer, Type: string ... 1 more field]

▶ df2: pyspark.sql.dataframe.DataFrame = [Store: integer, Dept: integer ... 3 more fields]

▶ df3: pyspark.sql.dataframe.DataFrame = [Store: integer, Date: date ... 10 more fields]

CURATED

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P ws05octshivam shellunext_1693422177077@npune...

Data Science & Engi...
New
Workspace
Repos
Recents
Catalog
Compute
Workflows
Marketplace
0/3 Tasks Completed
Menu options

Mounting_RetailSales_G4 Python ☆
File Edit View Run Help Last edit was now Provide feedback

Run all Terminated Schedule Share

Cmd 26

1 df_all.show()

(3) Spark Jobs

	Store	Date	IsHoliday	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	CPI
8.106	1	24924.5	A	151315							
8.106	1	46039.49	A	151315							
8.106	1	41595.55	A	151315							
8.106	1	19403.54	A	151315							
8.106	1	21827.9	A	151315							
8.106	1	21043.39	A	151315							
8.106	1	22136.64	A	151315							
8.106	1	51.45			2.732	0.0	0.0	0.0	0.0	0.0	211.01804

Command took 0.59 seconds -- by shellunext_1693422177077@npunext.onmicrosoft.com at 10/5/2023, 1:36:58 PM on Shell102 Unext's Cluster

Cmd 27

Stream 2 concentrated on further processing data over azure data bricks.

STAGING

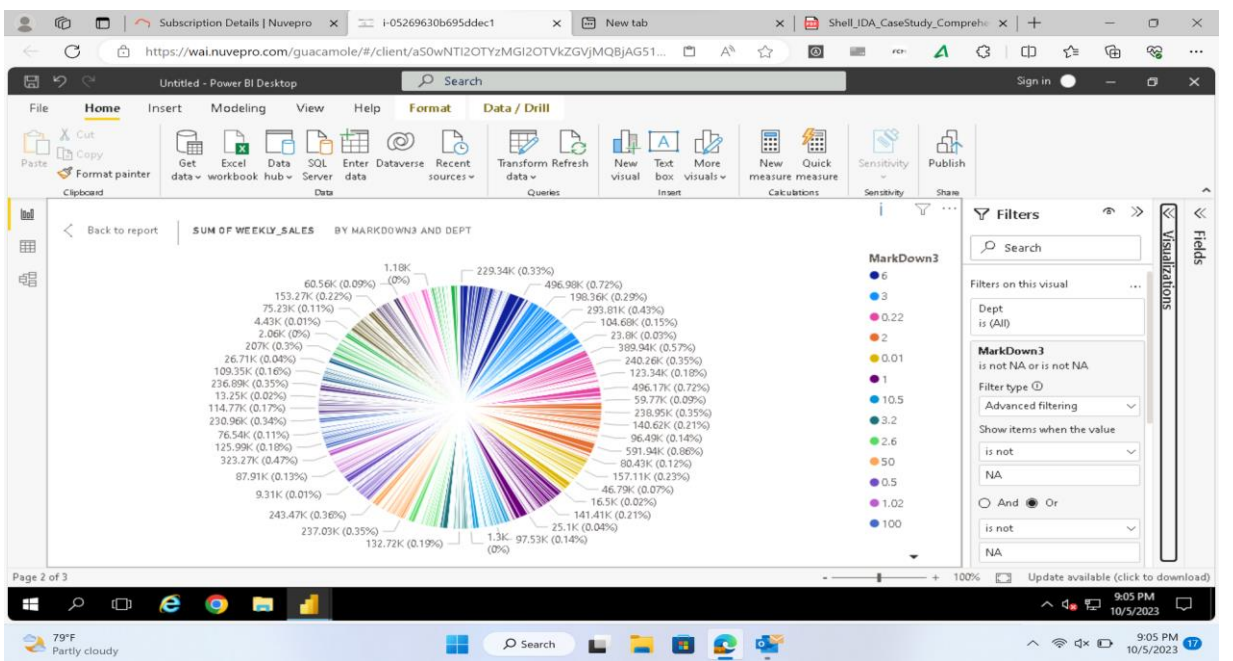
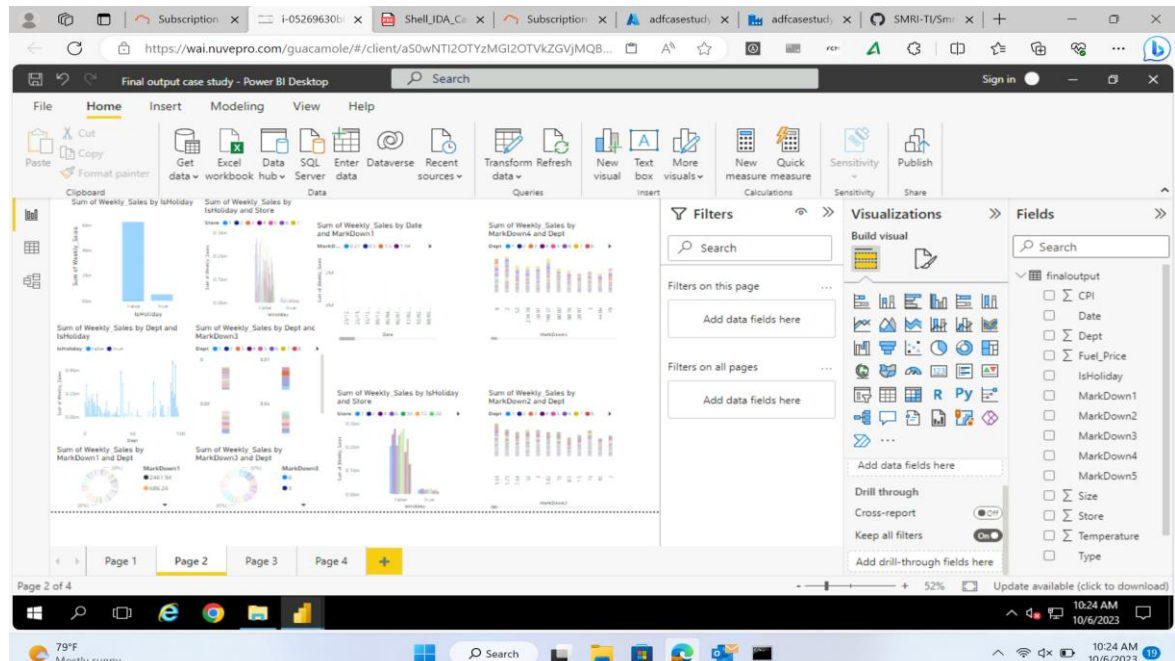
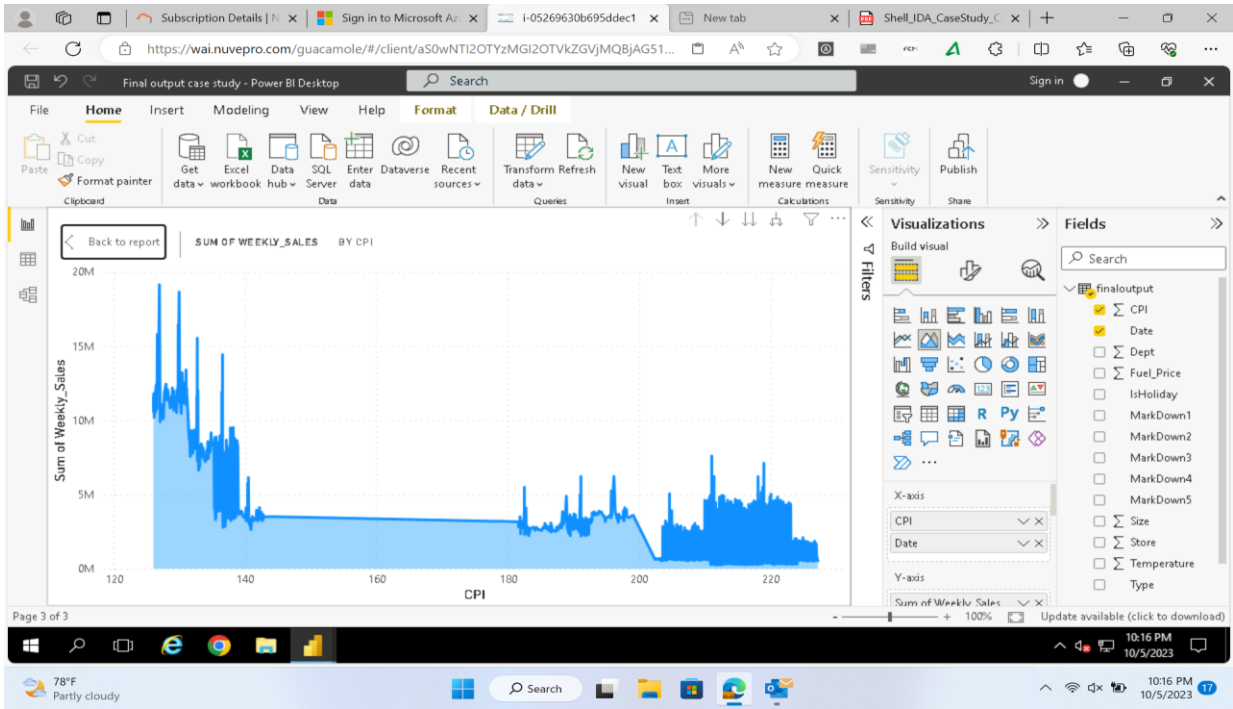
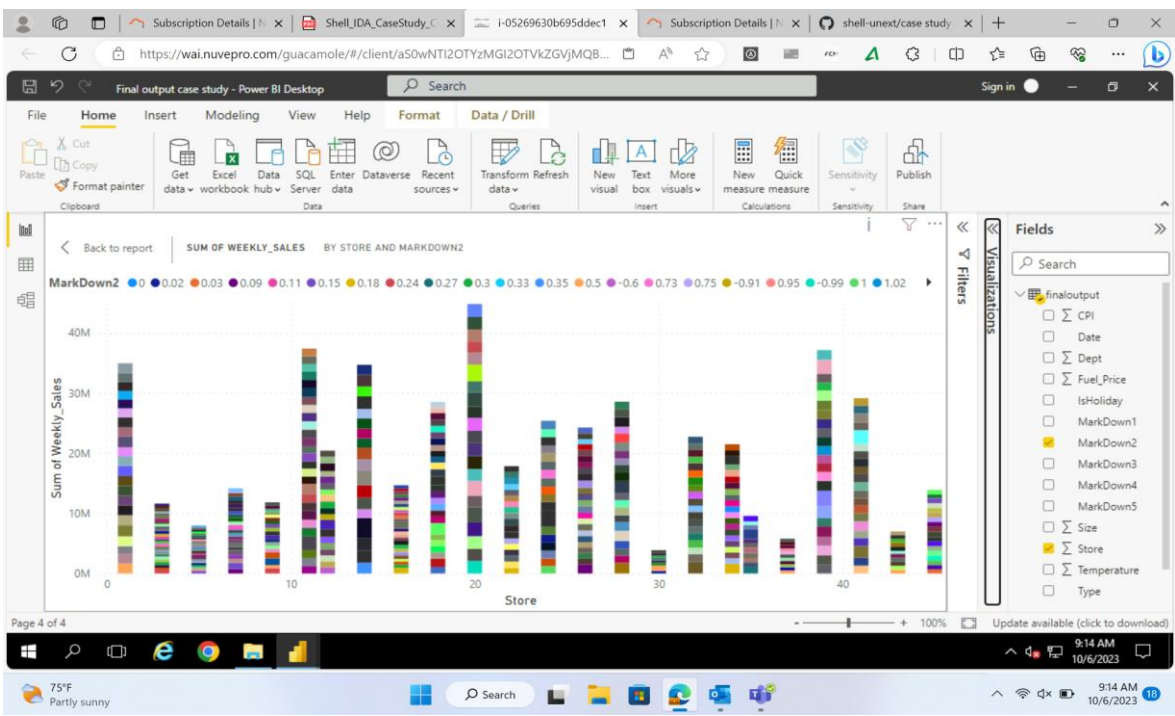
Microsoft Azure Data Factory portal showing the 'dataflow1' pipeline. The pipeline consists of two sources (source1 and source2) joined together, followed by a sink named 'salesfeatures'. The 'Inspect' tab is active, showing the schema with 14 columns. The 'Data preview' tab is also visible, showing a table with columns: Store, Date, Temperature, Fuel_P..., and MarkD...

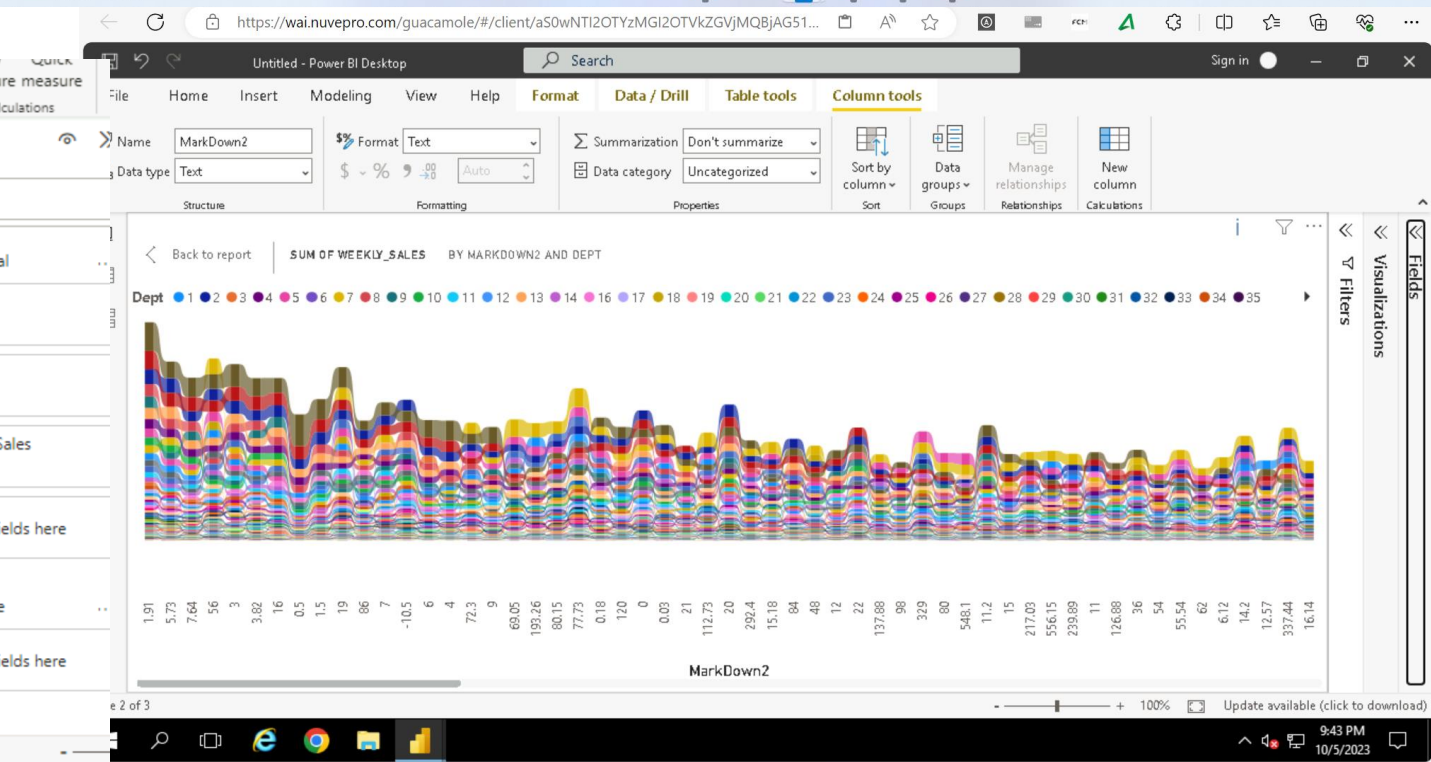
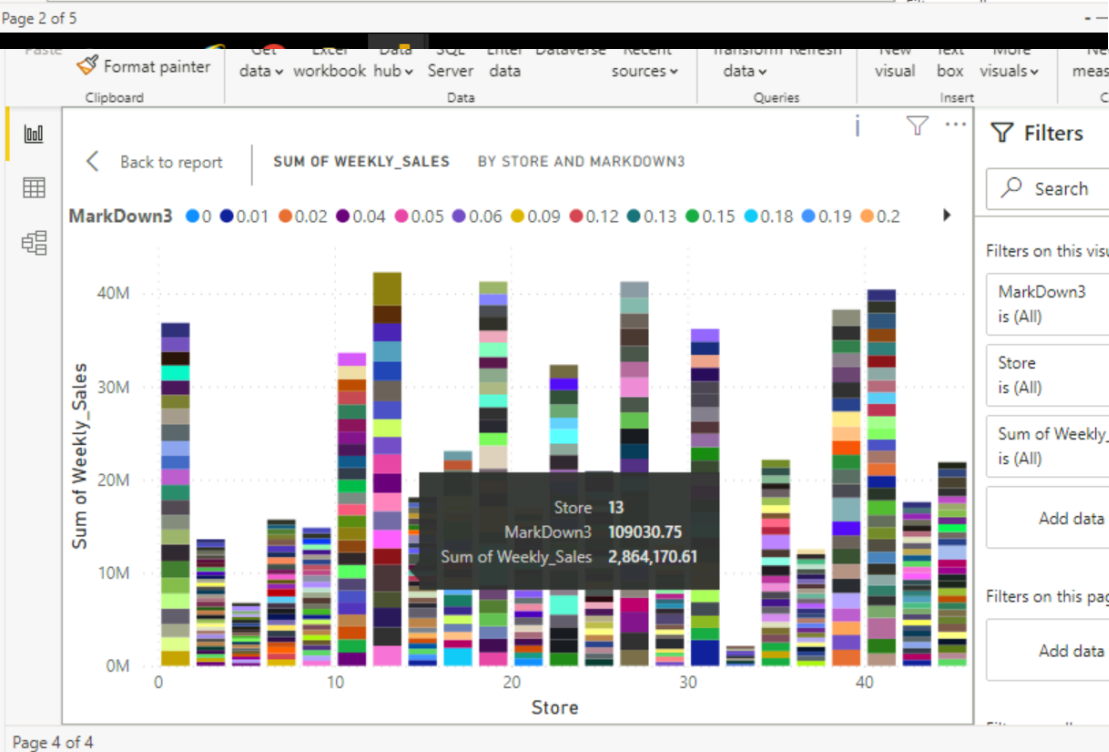
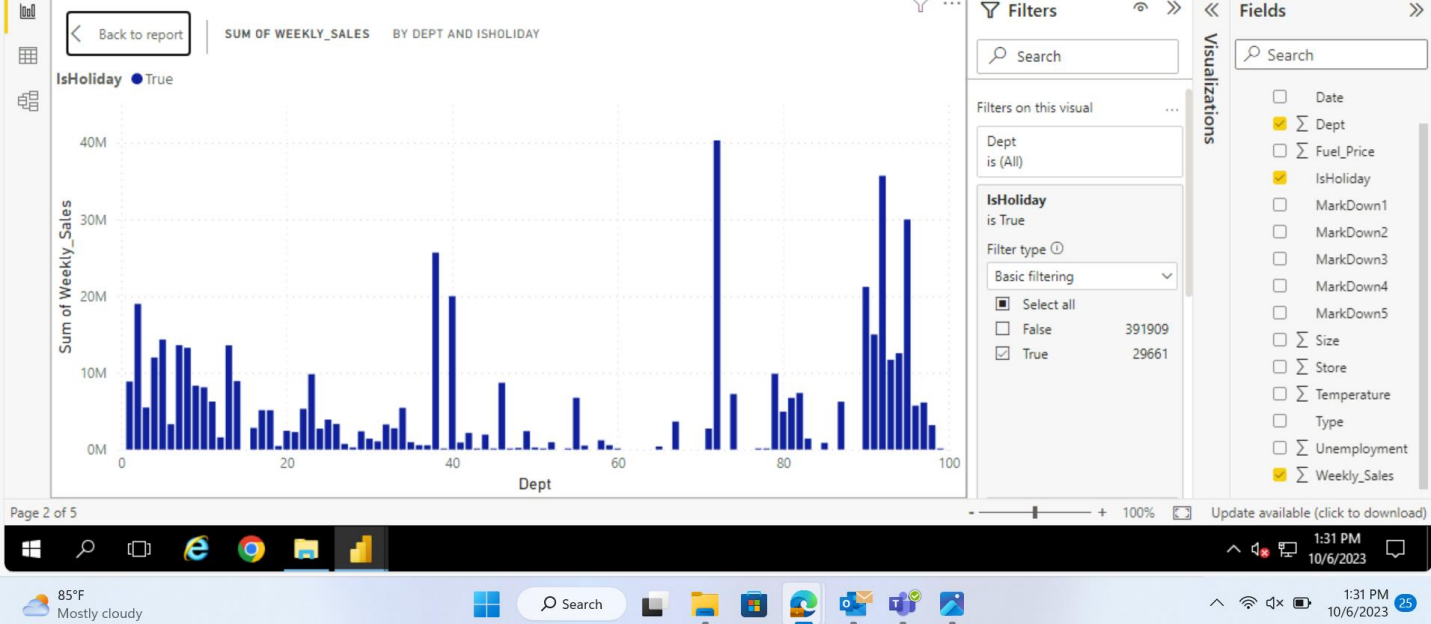
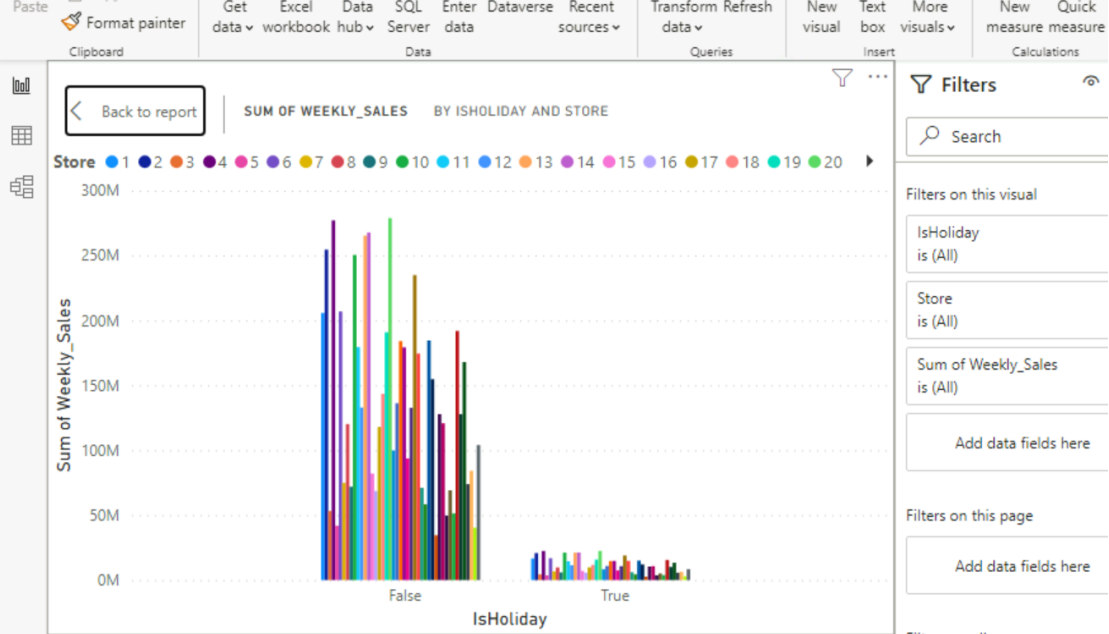
Store	Date	Temperature	Fuel_P...	MarkD...
1	05/02/...	42.31	2.572	NA
1	12/02/...	38.51	2.548	NA

Microsoft Azure Data Factory portal showing the 'dataflow2' pipeline. The pipeline consists of two sources (source1 and source2) joined together, followed by a sink named 'sink1'. The 'Data preview' tab is active, showing a table with columns: Store, Date, Temperature, Fuel_P..., and MarkD...

Store	Date	Temperature	Fuel_P...	MarkD...
1	05/02/...	42.31	2.572	NA
1	12/02/...	38.51	2.548	NA

These snapshot represent the pipelines used for data processing over azure data factory.





GITHUB LINK

- Azure Data Factory (Stream 1)

[/Smriti-Case-Study-3 \(github.com\)](#)

- Notebook for Azure Databricks is with the name Mounting.ipnb (Stream 2)

[/RetailSales Unext FinalAssessment \(github.com\)](#)

CHALLENGES FACED

Data Source Complexity

- Dealing with complex and diverse data sources, including structured and unstructured data
- Ensuring data quality and consistency across these sources is critical.

Data Integration

Integrating data from various sources into Azure Data Factory was time-consuming and require in-depth knowledge of data connectors, APIs, and data integration patterns.

Data Volume and Performance

- Handling large volumes of data lead to performance issues.
- Optimizing data pipelines and queries to ensure efficient data processing and visualization is crucial.

LEARNINGS

Understanding
Data Sources

Data Ingestion

Data
Transformation

Join Strategies

Error handling

Business
Impact