

8th pipeline : Blob Storage → Azure Databricks → Azure SQL Database :

Services required :

Azure Blob Storage
Azure Databricks
Azure SQL Database
Azure key vault

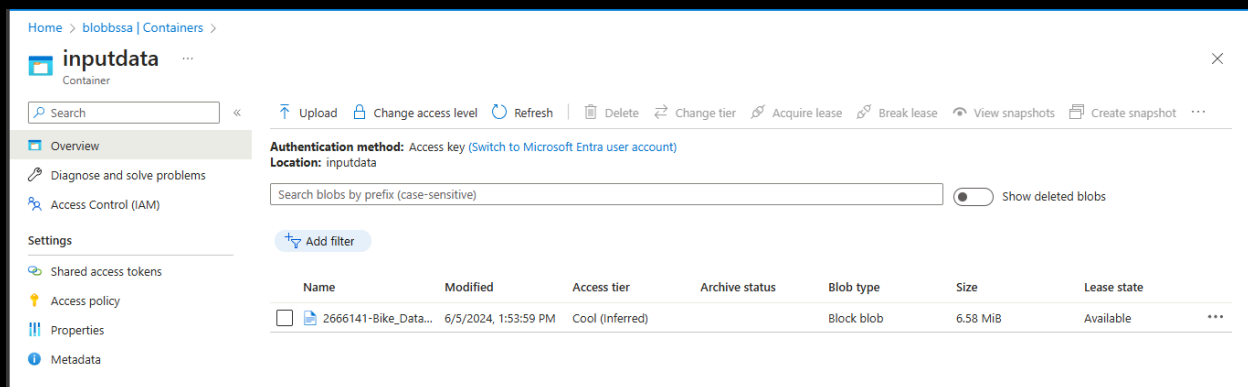
Creation of azure blob storage :

→ The creation of azure blob storage is shown in the document below.

<https://docs.google.com/document/d/1vZCMfM9ieALTQm6JdwL2JzxcMkwnyvkJHCrY/edit2kVlo?usp=sharing>

→ Create a container : inputData

→ Upload bike data in the inputData container



Creation of Azure SQL Database.

→ Creation of Azure SQL Database is shown in the document below.

https://docs.google.com/document/d/16iB1EsGKHc6-bcQTPSfqkK6BVf3n8_fpbat42uNOvXc/edit?usp=sharing

Creation of Azure Databricks

→ Creation of Azure SQL Databricks is shown in the document below.

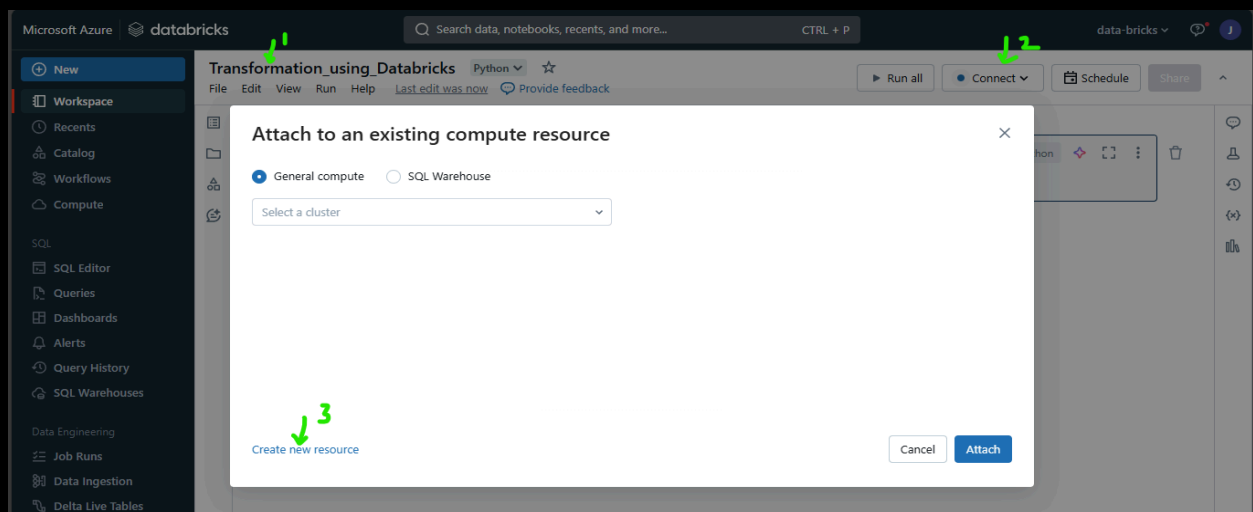
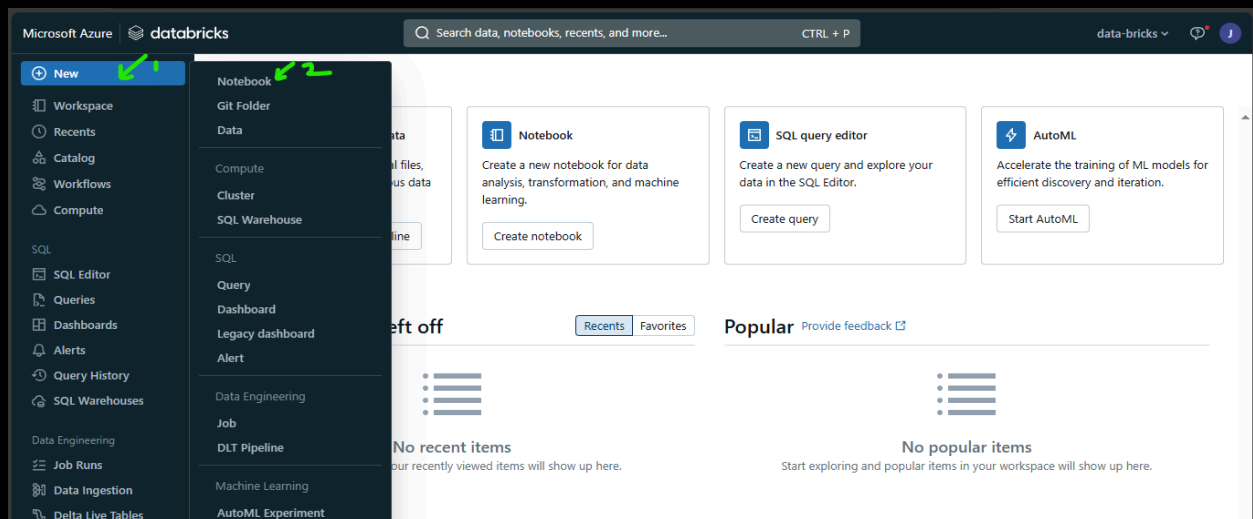
<https://docs.google.com/document/d/1s7Vqs4gcZbGMPK7vEWGFBXIMC5AoaOvC-E57qILwRw/edit?usp=sharing>

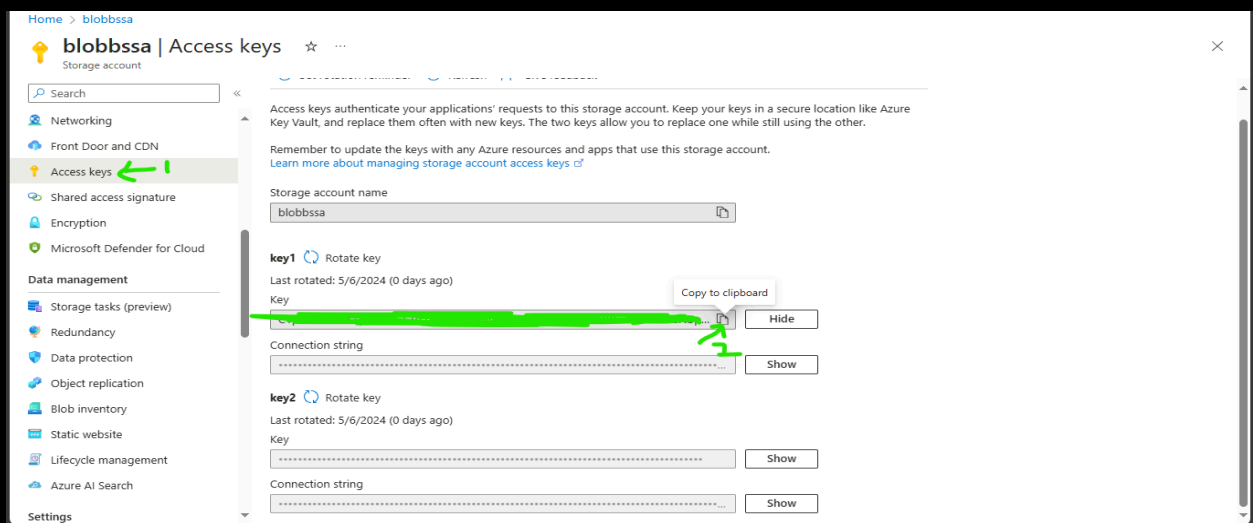
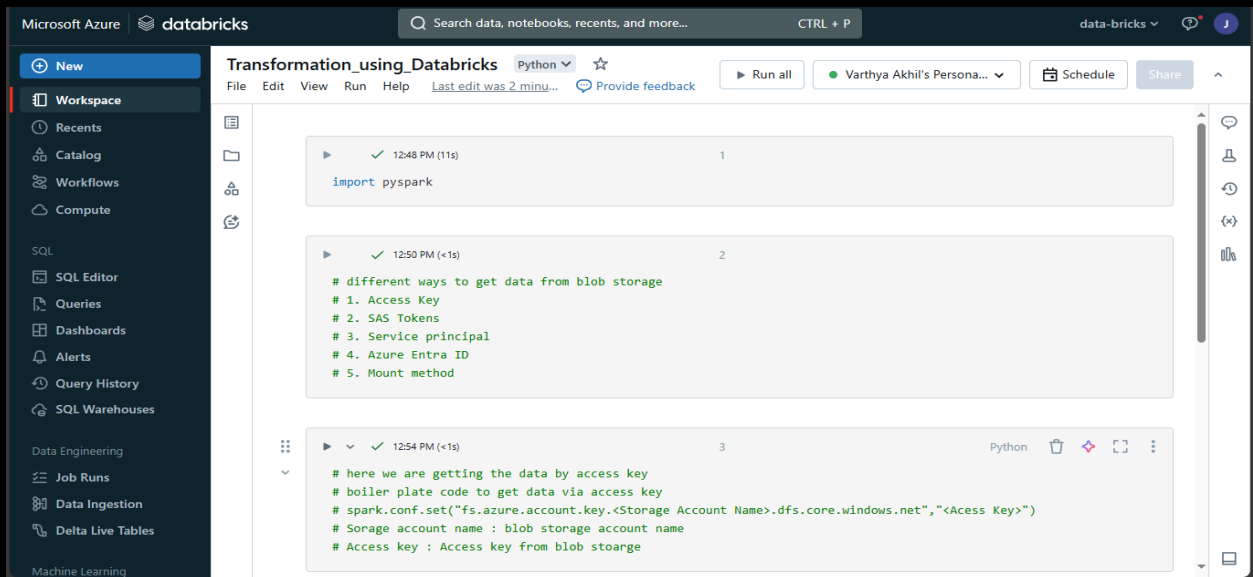
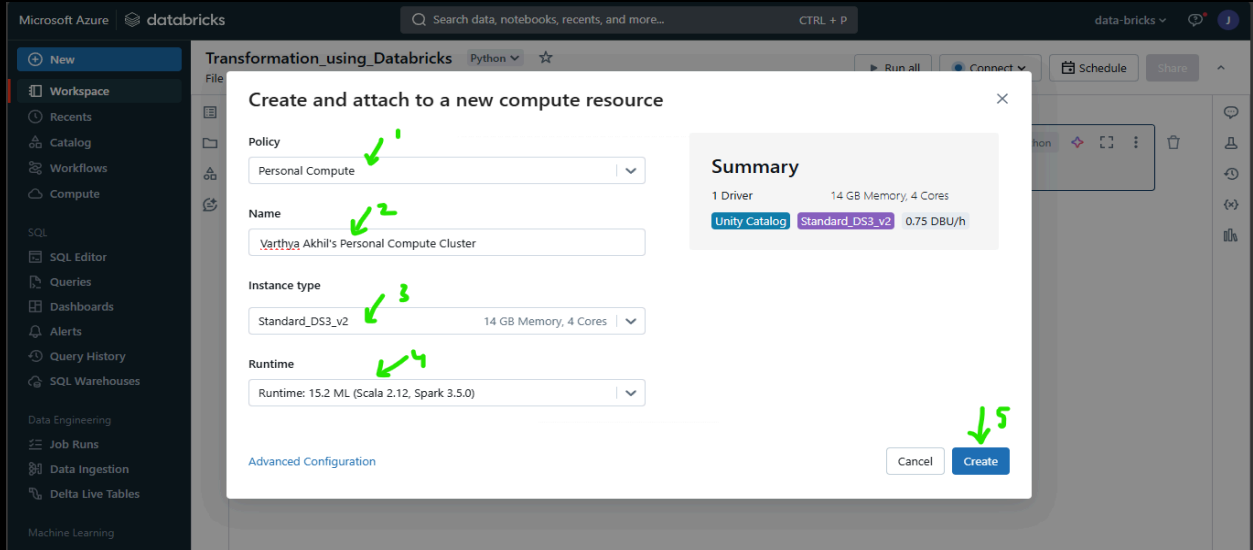
Creation of azure key vault :

→ Creation of Azure Key Vault is shown in the document below.

<https://docs.google.com/document/d/1ScQ42B5c5ZuRnFjtLsdJpnV9067pLZkpcDf4WyKjaVw/edit?usp=sharing>

In Databricks workspace →





Microsoft Azure databricks

Search data, notebooks, recents, and more... CTRL + P

data-bricks

New

Workspace

Recents

Catalog

Workflows

Compute

SQL

SQL Editor

Queries

Dashboards

Alerts

Query History

SQL Warehouses

Data Engineering

Job Runs

Data Ingestion

Delta Live Tables

Machine Learning

Playground

Experiments

Transformation_using_Databricks Python

File Edit View Run Help Last edit was 1 minute ago Provide feedback

Run all Varthya Akhil's Persona... Schedule Share

01:03 PM (<1s) 4

spark.conf.set("fs.azure.account.key.blobssa.dfs.core.windows.net", "2666141-Bike_Data_xlsx")

1 Paste key

01:04 PM (17s) 5

display(dbutils.fs.ls("abfss://inputdata@blobssa.dfs.core.windows.net"))

(2) Spark Jobs

Table	path	name	size	modificationTime
1	abfss://inputdata@blobssa.dfs.core.windows.net/2666141-Bike_Data_xl...	2666141-Bike_Data_xl...	3340468	1717570990000

1 row | 16.60 seconds runtime Refreshed 50 minutes ago

data = spark.read.csv("abfss://inputdata@blobssa.dfs.core.windows.net/2666141-Bike_Data.xlsx") Python

Here, we exposed the access key. It is not a good practice to expose access key. Instead we will import the access key via azure key vault and scope.

Home > projects-key-vaults | Overview > projects-key-vaults

Key vault

projects-key-vaults | Secrets

Generate/Import Refresh Restore Backup View sample code Manage deleted secrets

Name	Type	Status	Expiration date
There are no secrets available.			

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Access policies

Events

Objects

Keys

Secrets

Certificates

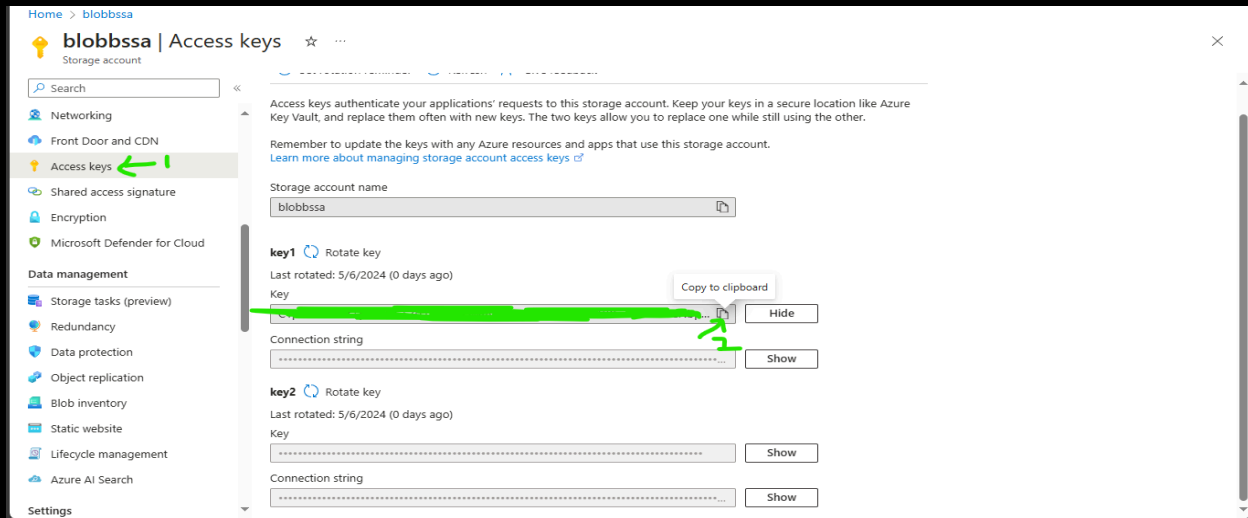
Settings

Access configuration

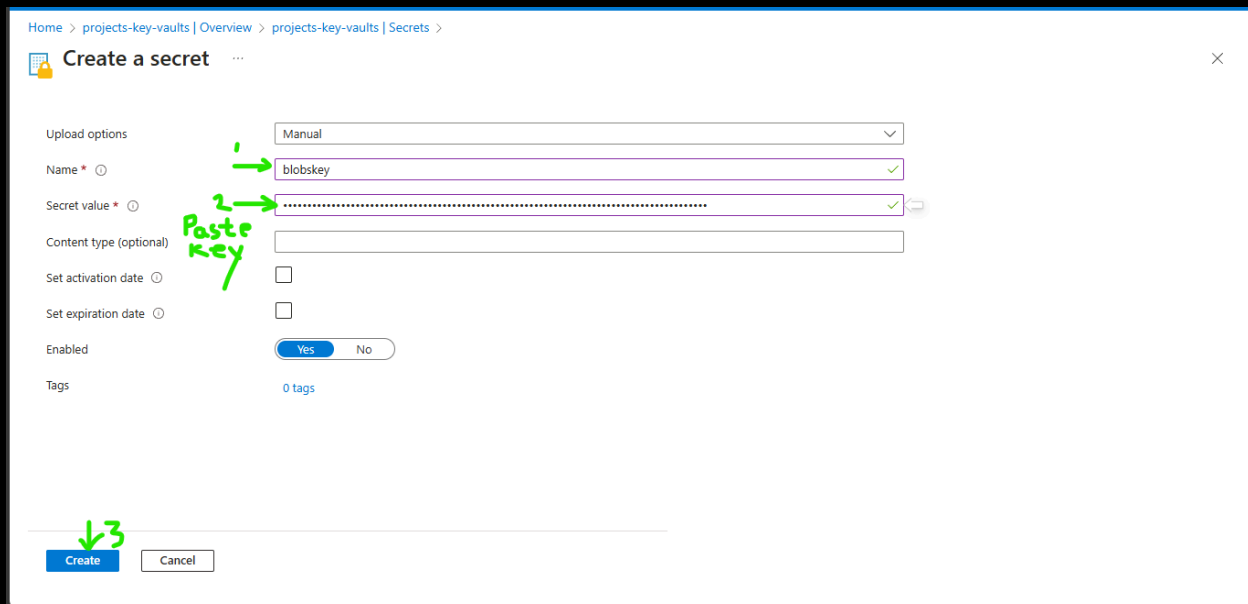
Networking

Microsoft Defender for Cloud

Properties



Copy the access key from blob storage and paste in the secret value.



Let's create a scope in databricks.

Using this scope we will import the key.

Copy the url, paste in the new tab and remove the characters after # and give secrets/createscope in the url.

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P data-bricks

New

- Workspace
- Recents
- Catalog
- Workflows
- Compute
- SQL
- SQL Editor
- Queries
- Dashboards
- Alerts
- Query History
- SQL Warehouses
- Data Engineering
- Job Runs
- Data Ingestion
- Delta Live Tables
- Machine Learning
- Playground
- Experiments

HomePage / Create Secret Scope

Create Secret Scope

Cancel Create

A store for secrets that is identified by a name and backed by a specific store type. [Learn more](#)

Scope Name

Manage Principal

Azure Key Vault

DNS Name

Resource ID

Home > projects-key-vaults | Overview > projects-key-vaults

projects-key-vaults | Properties

Key vault

Search Save Discard changes Refresh

Secrets

Certificates

Settings

- Access configuration
- Networking
- Microsoft Defender for Cloud
- Properties**
- Locks

Monitoring

- Alerts
- Metrics
- Diagnostic settings
- Logs
- Insights
- Workbooks

Name	projects-key-vaults
Sku (Pricing tier)	Standard
Location	eastus
Vault URI	https://projects-key-vaults.vault.azure.net/
Resource ID	/subscriptions/c92e5287-784e-4a20-9e6c-8549947f2ee6/resourceGroups/projects_rg/providers/Microsoft.KeyVault/vaults/project...
Subscription ID	c92e5287-784e-4a20-9e6c-8549947f2ee6
Subscription Name	Free Trial
Directory ID	2282df71-75c7-49bc-8440-7dc9452f264c
Directory Name	Default Directory
Soft-delete	Soft delete has been enabled on this key vault
Days to retain deleted vaults	<input type="text" value="90"/>
Purge protection	<input checked="" type="radio"/> Disable purge protection (allow key vault and objects to be purged during retention period) <input type="radio"/> Enable purge protection (enforce a mandatory retention period for deleted vaults and vault objects)

Handwritten notes: "COPY 2" with an arrow pointing to the Vault URI field.

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P data-bricks

New

- Workspace
- Recents
- Catalog
- Workflows
- Compute
- SQL
- SQL Editor
- Queries
- Dashboards
- Alerts
- Query History
- SQL Warehouses
- Data Engineering
- Job Runs
- Data Ingestion
- Delta Live Tables
- Machine Learning
- Playground
- Experiments

HomePage / Create Secret Scope

Create Secret Scope

Cancel Create

A store for secrets that is identified by a name and backed by a specific store type. [Learn more](#)

Scope Name

Manage Principal

Azure Key Vault

DNS Name

Resource ID

Handwritten notes: "1" with an arrow pointing to the Scope Name field, and "2 Paste" with an arrow pointing to the DNS Name field.

Home > projects-key-vaults | Overview > projects-key-vaults

projects-key-vaults | Properties

Key vault

Search

Save Discard changes Refresh

Secrets

Certificates

Settings

- Access configuration
- Networking
- Microsoft Defender for Cloud
- Properties**
- Locks

Monitoring

- Alerts
- Metrics
- Diagnostic settings
- Logs
- Insights
- Workbooks

Name	projects-key-vaults
Sku (Pricing tier)	Standard
Location	eastus
Vault URI	https://projects-key-vaults.vault.azure.net/
Resource ID	/subscriptions/c92e5287-784e-4a20-9e6c-8549947f2ee6/resourceGroups/projects_rg/providers/Microsoft.KeyVault/vaults/project...
Subscription ID	c92e5287-784e-4a20-9e6c-8549947f2ee6
Subscription Name	Free Trial
Directory ID	2282df71-75c7-49bc-8440-7dc9452f264c
Directory Name	Default Directory

Soft-delete

Days to retain deleted vaults

Purge protection

Soft delete has been enabled on this key vault

☒ Disable purge protection (allow key vault and objects to be purged during retention period)

☐ Enable purge protection (enforce a mandatory retention period for deleted vaults and vault objects)

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P data-bricks

New

- Workspace
- Recents
- Catalog
- Workflows
- Compute

SQL

- SQL Editor
- Queries
- Dashboards
- Alerts
- Query History
- SQL Warehouses

Data Engineering

- Job Runs
- Data Ingestion

Create Secret Scope

Cancel Create

A store for secrets that is identified by a name and backed by a specific store type. [Learn more](#)

Scope Name

Manage Principal

Azure Key Vault

DNS Name

Resource ID

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P data-bricks

New

- Workspace
- Recents
- Catalog
- Workflows
- Compute

SQL

- SQL Editor
- Queries
- Dashboards
- Alerts
- Query History
- SQL Warehouses

Data Engineering

- Job Runs
- Data Ingestion
- Delta Live Tables

Machine Learning

- Playground
- Experiments
- Features
- Models

Transformation_using_Databricks

Python

File Edit View Run Help Last edit was 2 minutes ago Provide feedback

Run all Varthya Akhil's Persona... Schedule Share

```
1 from pyspark.sql import *
  from pyspark.sql.functions import *

2 dbutils.secrets.listScopes()

3 spark.conf.set(
    "fs.azure.account.key.<storage-account>.dfs.core.windows.net",
    dbutils.secrets.get(scope="<scope>", key="<storage-account-access-key>"))

# scope = scopename (databricksscope)
# key = storage-account-access-key (blobkey (azure key vault -> secrets))

4 spark.conf.set("fs.azure.account.key.blobsssa.dfs.core.windows.net",dbutils.secrets.get(scope="databricksscope",key="blobkey"))
```

Transformation_using_Databricks Python Run all Varthya Akhil's Persona... Schedule Share

File Edit View Run Help Last edit was 1 minute ago Provide feedback

10

display(dbutils.fs.ls("abfss://inputdata@blobbssa.dfs.core.windows.net"))

(2) Spark Jobs

Table	path	name	size	modificationTime
1	abfss://inputdata@blobbssa.dfs.core.windows.net/2666141-Bike_Data (1).c...	2666141-Bike_Data (1).c...	6895611	1717575839000

1 row | 2.07 seconds runtime Refreshed 4 minutes ago

11

data = spark.read.csv("abfss://inputdata@blobbssa.dfs.core.windows.net/2666141-Bike_Data (1).csv",header=True, inferSchema=True)

(2) Spark Jobs

data: pyspark.sql.dataframe.DataFrame = [Region: string, Country: string ... 9 more fields]

12

Activate Windows
Go to Settings to activate Windows.

Transformation_using_Databricks Python Run all Varthya Akhil's Persona... Schedule Share

File Edit View Run Help Last edit was 2 minutes ago Provide feedback

12

data.columns

```
[ 'Region',  
  'Country',  
  'Customer',  
  'Business Segment',  
  'Category',  
  'Model',  
  'Color',  
  'SalesDate',  
  'ListPrice',  
  'UnitPrice',  
  'OrderQty' ]
```

13

data.printSchema()

```
root  
 |-- Region: string (nullable = true)  
 |-- Country: string (nullable = true)  
 |-- Customer: string (nullable = true)  
 |-- Business Segment: string (nullable = true)  
 |-- Category: string (nullable = true)  
 |-- Model: string (nullable = true)  
 |-- Color: string (nullable = true)  
 |-- SalesDate: date (nullable = true)  
 |-- ListPrice: string (nullable = true)  
 |-- UnitPrice: string (nullable = true)  
 |-- OrderQty: integer (nullable = true)
```

Activate Windows
Go to Settings to activate Windows.

Transformation_using_Databricks Python Run all Varthya Akhil's Persona... Schedule Share

File Edit View Run Help Last edit was 2 minutes ago Provide feedback

14

data.select("Country","Customer","Business Segment","Category","Model","SalesDate","OrderQty").show(5,False)

(1) Spark Jobs

Country	Customer	Business Segment	Category	Model	SalesDate	OrderQty
United States	Advanced Bike Components	Components	Road Frames	LL Road Frame	2020-04-01	1
United States	Central Discount Store	Bikes	Mountain Bikes	Mountain-100	2020-04-01	1
United States	Leading Sales & Repair	Clothing	Jerseys	Long-Sleeve Logo Jersey	2020-04-01	6
United States	Paint Supply	Components	Mountain Frames	HL Mountain Frame	2020-04-01	2
United States	Scooters and Bikes Store	Bikes	Road Bikes	Road-450	2020-04-01	2

only showing top 5 rows

15

data = data.withColumn("Sales",data["ListPrice"]*data["OrderQty"])

data: pyspark.sql.dataframe.DataFrame = [Region: string, Country: string ... 10 more fields]

Activate Windows
Go to Settings to activate Windows.

New

Workspace

Recents

Catalog

Workflows

Compute

SQL

SQL Editor

Queries

Dashboards

Alerts

Query History

SQL Warehouses

Data Engineering

Job Runs

Data Ingestion

Delta Live Tables

Machine Learning

Playground

Transformation_using_Databricks

Python

File Edit View Run Help

Last edit was 3 minutes ago

Provide feedback

Run all

Varthya Akhil's Persona...

Schedule

Share

data.show(5,False)

(1) Spark Jobs

Region	Country	Customer	Business Segment	Category	Model	Color	SalesDate	ListPrice	UnitP
North America	United States	Advanced Bike Components	Components	Road Frames	LL Road Frame	Red	2020-04-01	337.22	183.9
4	1	337.22							
North America	United States	Central Discount Store	Bikes	Mountain Bikes	Mountain-100	Silver	2020-04-01	3,399.99	2,03
9.99	1	NULL							
North America	United States	Leading Sales & Repair	Clothing	Jerseys	Long-Sleeve Logo Jersey	Multi	2020-04-01	49.99	28.84
6		299.94							
North America	United States	Paint Supply	Components	Mountain Frames	HL Mountain Frame	Black	2020-04-01	1,349.60	714.7
12		NULL							
North America	United States	Scooters and Bikes Store	Bikes	Road Bikes	Road-450	Red	2020-04-01	1,457.99	874.7
9	2	NULL							

only showing top 5 rows

New

Workspace

Recents

Catalog

Workflows

Compute

SQL

SQL Editor

Queries

Dashboards

Alerts

Query History

SQL Warehouses

Data Engineering

Job Runs

Data Ingestion

Delta Live Tables

Machine Learning

Playground

Transformation_using_Databricks

Python

File Edit View Run Help

Last edit was 10 minutes ago

Provide feedback

Run all

Varthya Akhil's Persona...

Schedule

Share

only showing top 5 rows

data = data.withColumn("Cost",data["UnitPrice"]*data["OrderQty"])

data: pyspark.sql.dataframe.DataFrame = [Region: string, Country: string ... 11 more fields]

data = data.withColumn("Profit",data["Sales"]-data["Cost"])

data: pyspark.sql.dataframe.DataFrame = [Region: string, Country: string ... 12 more fields]

df = data.select("Country","Customer","Category","SalesDate","Sales","Cost","Profit")

df: pyspark.sql.dataframe.DataFrame = [Country: string, Customer: string ... 5 more fields]

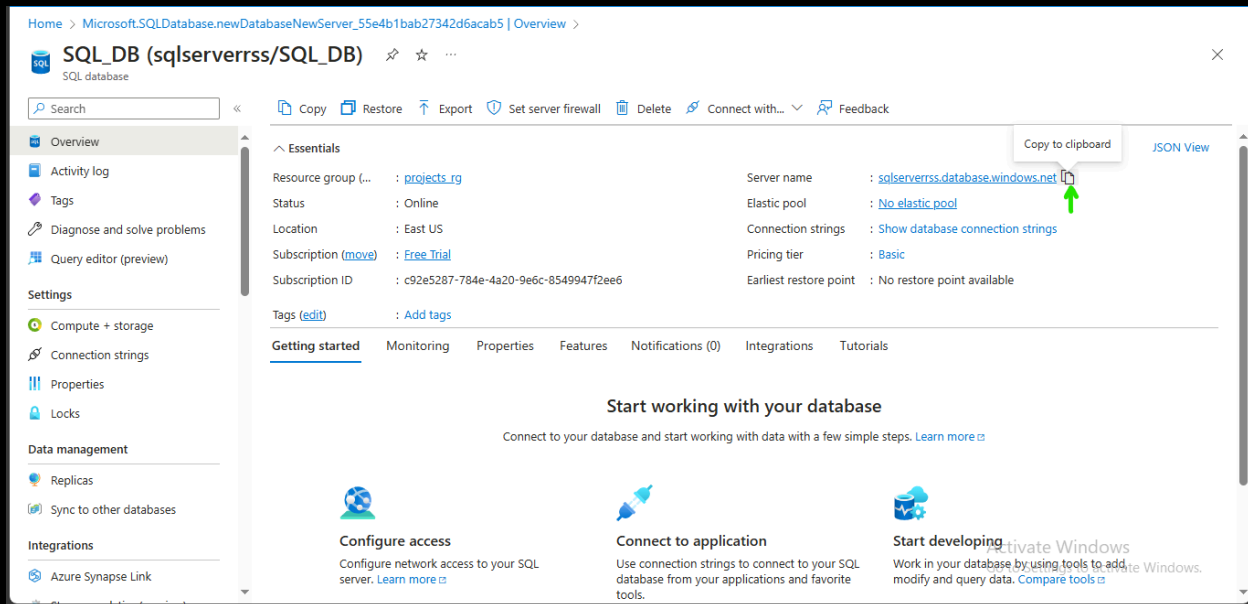
Loading dataframe into sql database.

Boilerplate code to connect sql database.

→ jdbc_url =

```
"jdbc:sqlserver://<your_server_name>.database.windows.net:1433;database=<your_database_name>"
```

Copy the servername from azure sql database overview



Paste the server name in the boilerplate code.



Boilerplate code to writeback to sql server database

```
→ df.write.jdbc(url=jdbc_url, table="your_table_name", mode="overwrite",  
properties=connection_properties)
```



Now data has been successfully loaded into an azure sql database.

Lets see data in azure sql database by querying .

