

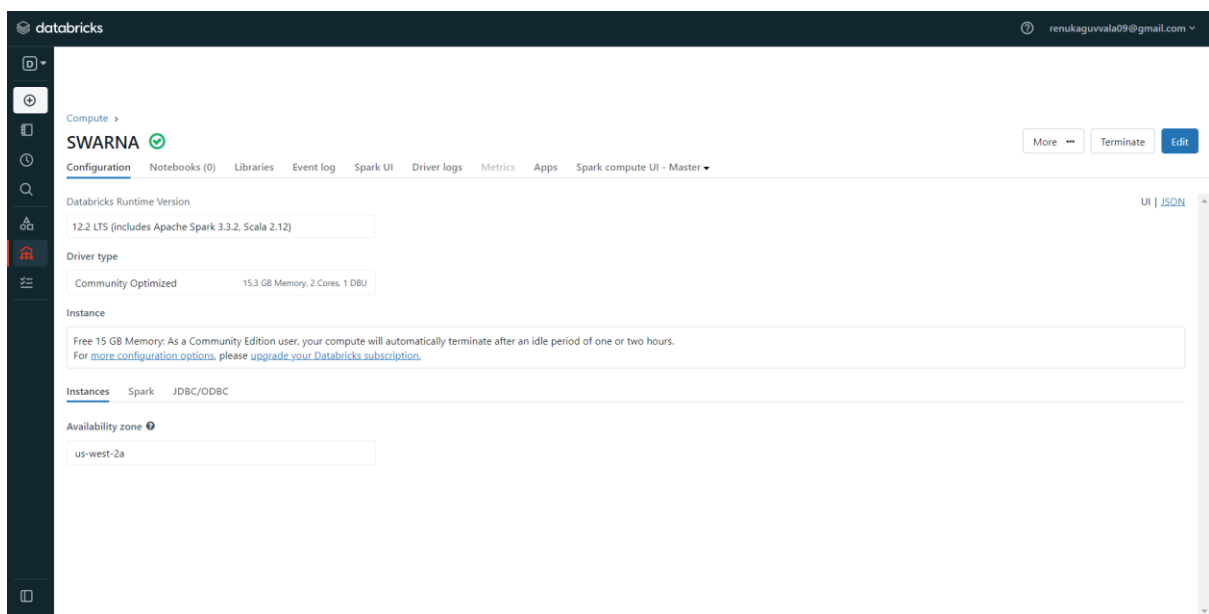
AZURE DATARICKS CODING ASSESSMENT

NAME-SWARNA LATHA GUVVALA

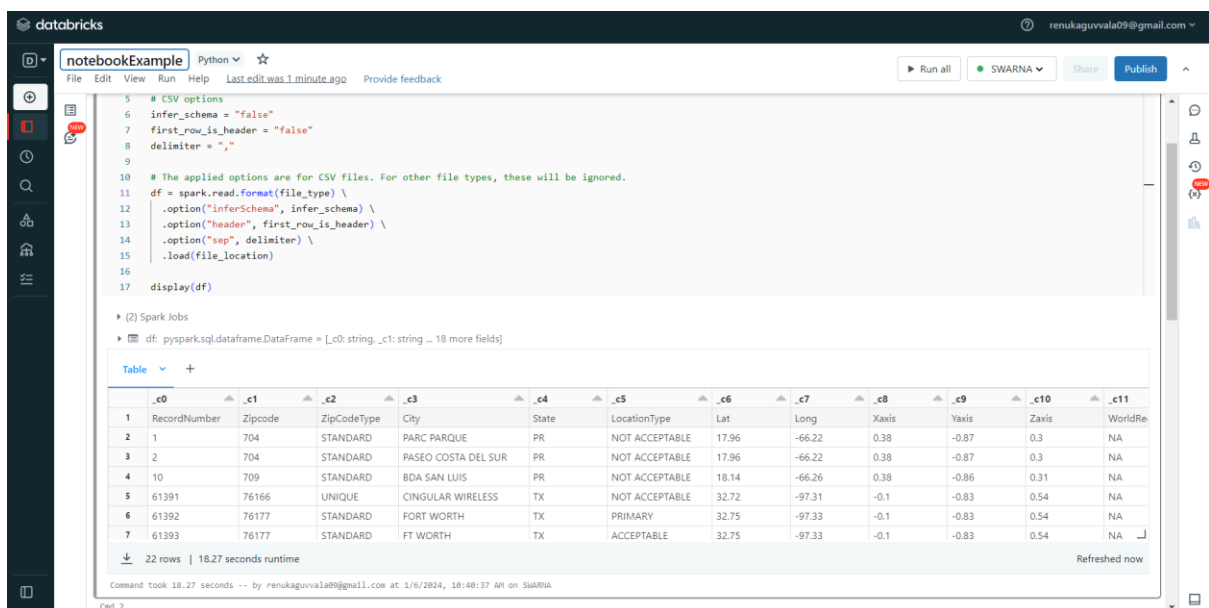
DATE-06/01/2024

1.create a cluster &attach the notebook to the cluster and run all commands in the notebook &creates a dataframe from a databricks dataset&create a visualizations in databricks notebooks &rename, duplicate, or remove a visualization or data profile?

Create a Cluster:



Attach the notebook to the cluster and run all commands in the notebook:



Creates a dataframe from a databricks dataset:

The screenshot shows a Databricks notebook interface. The top bar includes the Databricks logo, a user profile icon, and a dropdown menu. The notebook is titled "dataframe" and is in Python mode. The code in the notebook is as follows:

```
1 data = [[295, "South Bend", "Indiana", "IN", 101190, 112.9]]
2 columns = ["rank", "city", "state", "code", "population", "price"]
3
4 df1 = spark.createDataFrame(data, schema="rank LONG, city STRING, state STRING, code STRING, population LONG, price DOUBLE")
5 display(df1)
```

The output of the code is a table with 7 columns: rank, city, state, code, population, and price. The table contains one row of data: rank 295, city South Bend, state Indiana, code IN, population 101190, and price 112.9. The table is displayed with a "Table" button and a "1 row" indicator. The runtime is 5.04 seconds.

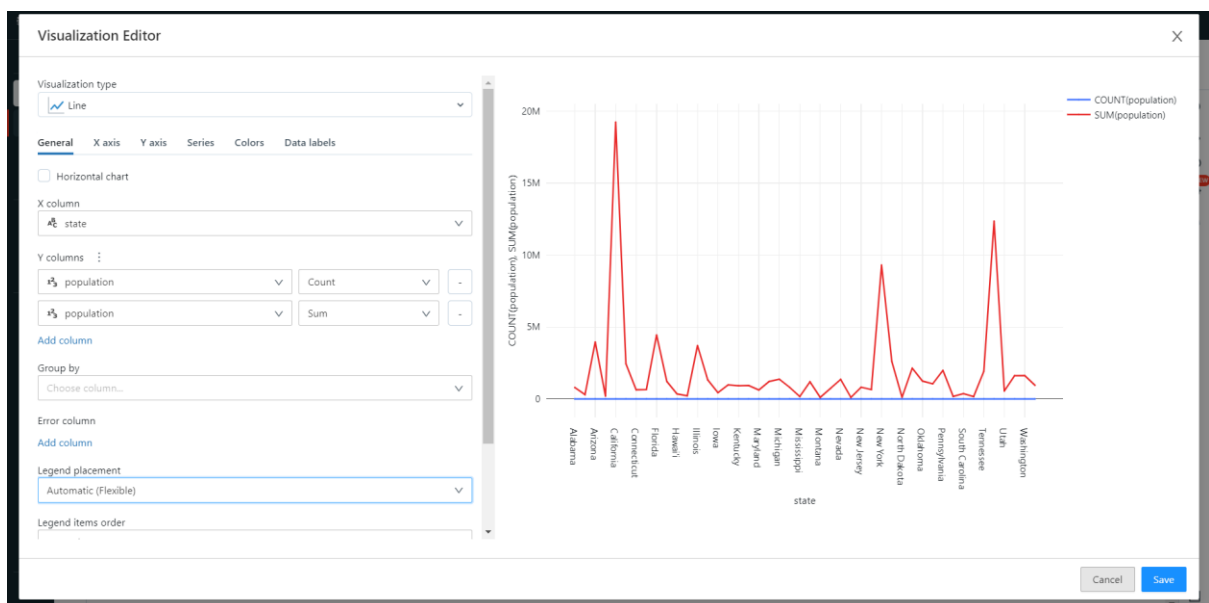
Below the table, the command took 5.84 seconds to execute by reneaguvala09@gmail.com at 1/6/2024, 18:44:49 AM on SWARNA.

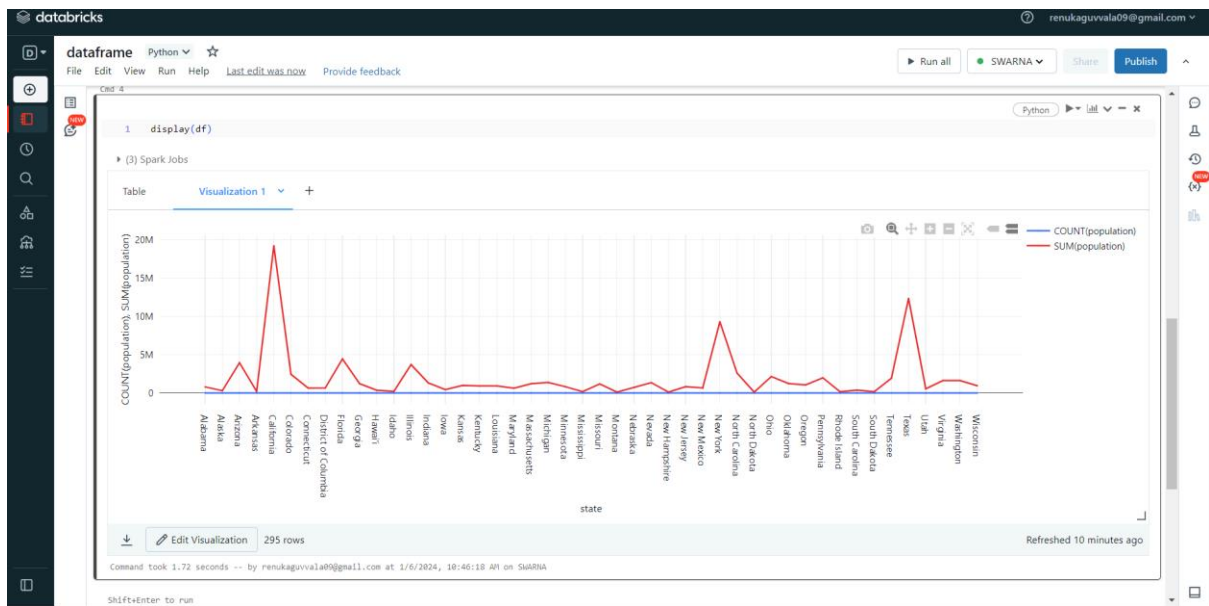
The notebook also shows a second code cell with the following code:

```
1 df2 = (spark.read
2       .format("csv")
3       .option("header", "true")
4       .option("inferSchema", "true")
5       .load("/databricks-datasets/samples/population-vs-price/data_geo.csv")
6 )
```

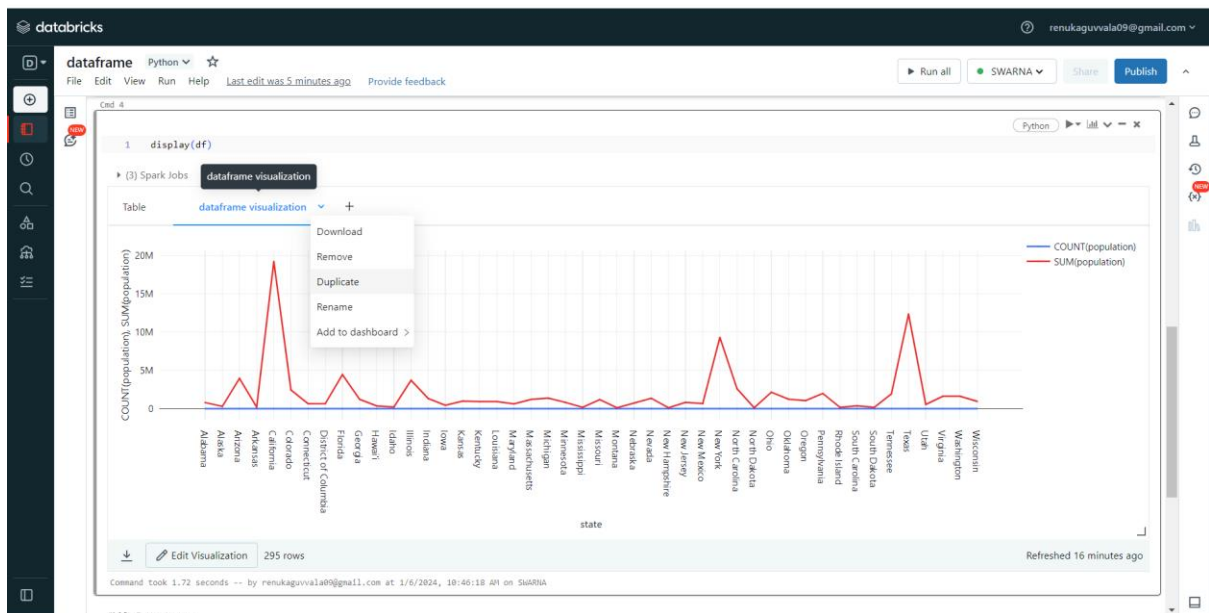
The output of the second code cell is a DataFrame with 2014 rows and 4 columns: rank (integer), city (string), state (string), and price (double).

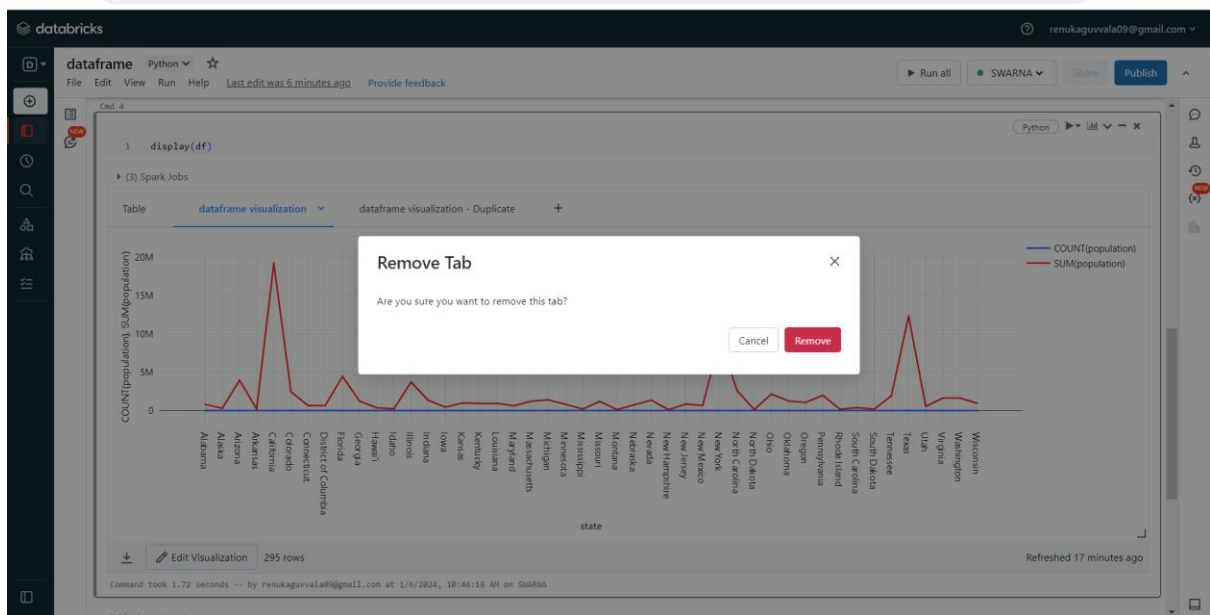
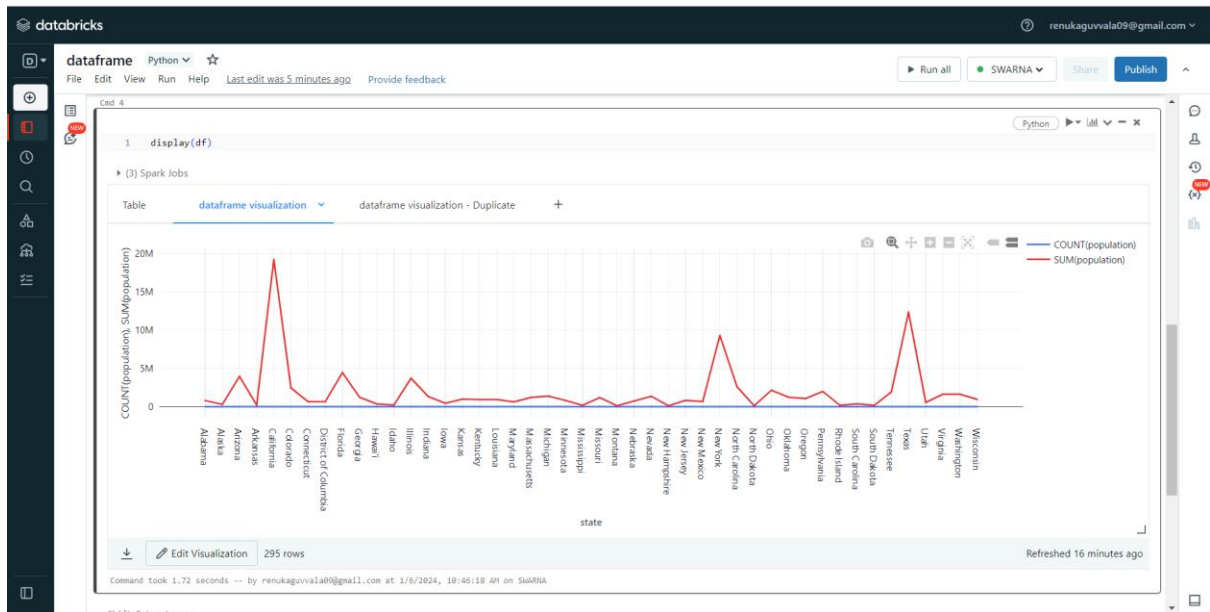
Create a visualizations in databricks notebooks:

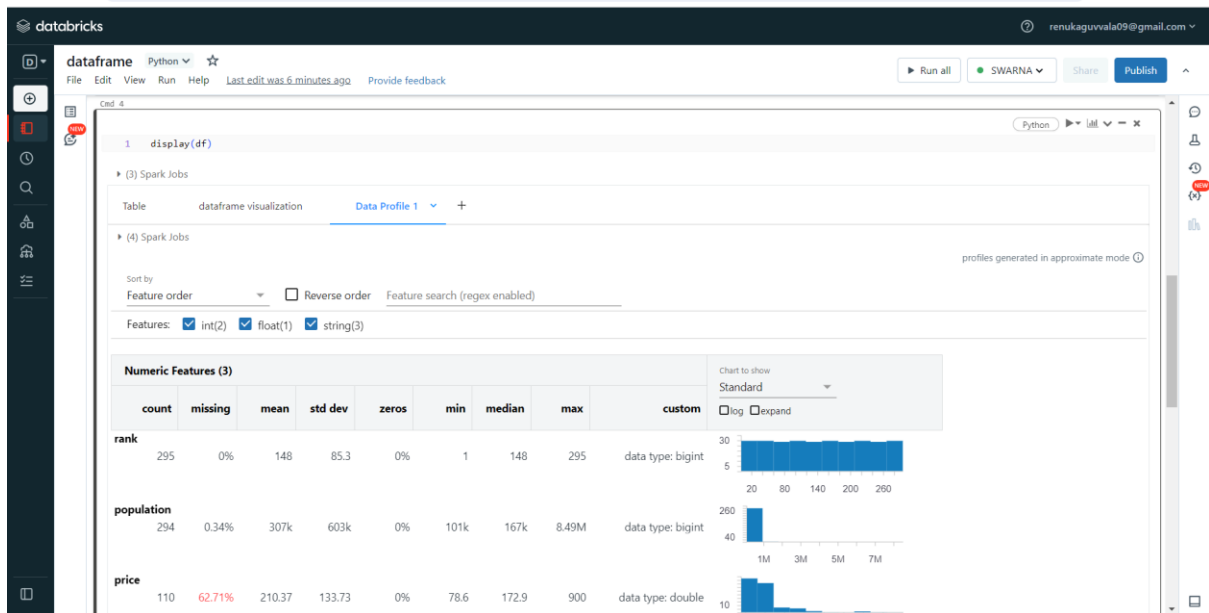




Rename, duplicate, or remove a visualization or data profile:







2.Explain the copy activity in Azure data factory.

In Azure Data Factory(ADF),the copy activity is a crucial component that enables the movement of data from source to destination. The primary purpose of the copy activity is to facilitate the extraction, transformation, and loading (ETL) of data.

Here is a breakdown of the Key aspects of the copy activity in Azure Data Factory:

1.Source and Destination:

- **Source Dataset:** This is the data store from which data is to be copied .It could be an on-premises database, a cloud-based storage solution (like Azure Blob Storage or Azure SQL Database)
- **Destination Dataset:** This is the target data store where the data will be copied .Similar to the source dataset, it can be on-premises or cloud-based data store.

2.Data Movement:

- The copy activity efficiently moves data from source to destination. It handles data transfer in a secure and performant manner, optimizing the process for large-scale data movement.

3.Supported Formats:

- The copy activity supports various date formats ,including CSV,JSON...etc. It can handle structured, semi-structured, and unstructured data.

4.Mapping and Transformation:

- The copy activity allows for simple data mappings between source and destination columns. It also supports limited transformations during the copying process.

5.Integration Runtime:

- Integration Runtimes are used to provide the necessary compute infrastructure to execute the copy activity .Azure data factory supports various types of integration runtimes, including Azure .These runtimes ensure that data movement occurs securely and efficiently.

6.Monitoring and Logging:

- Azure Data Factory provides monitoring capabilities for the copy activity. You can monitor the process, view logs and gain insights into the success or failure of the data movement operation.

7.Error Handling:

- The copy activity includes features for handling errors during the data movement process. This includes retry policies and configurable settings to manage how the system responds to failures.

8.Scheduling and Orchestration:

- The copy activity is often used within pipelines ,where it can be scheduled to run at specific intervals or triggered by external events. This enables you to orchestrate complex ETL workflows within Azure Data Factory.

In summary, the copy activity in Azure Data Factory is a versatile tool for efficiently moving and transforming data between various sources and destinations ,making it a key component for building data integration solutions in the Azure cloud.