

Name: Pradip Bochare

Azure Databricks Coding Assessment Question 1



Exploratory data analysis (EDA) in Databricks & Visualizing data in Databricks

❖ Exploratory data analysis (EDA) in Databricks

- Exploratory data analysis (EDA) includes methods for exploring data sets to summarize their main characteristics and identify any problems with the data.
- Using statistical methods and visualizations, you can learn about a data set to determine its readiness for analysis and inform what techniques to apply for data preparation.
- EDA can also influence which algorithms you choose to apply for training ML models.
- Utilize PySpark for data loading, manipulation, and analysis within Databricks notebooks.
- Leverage PySpark functions such as `show()`, `describe()`, and `printSchema()` to explore the dataset.
- Perform data cleaning and preprocessing using PySpark DataFrame operations like `dropna()` and `withColumn()`.
- Use Databricks built-in visualizations or external libraries like Matplotlib and Seaborn for data visualization.
- Visualize data distributions, trends, and relationships through histograms, scatter plots, etc. Utilize Seaborn for more advanced visualizations and styling options.
- Databricks provides seamless integration with PySpark, enabling scalable and efficient data analysis workflows.

❖ Visualizing data in Databricks

- Here we are using data frame to perform visualization on it

The screenshot shows the Databricks Data Visualization interface. The top bar includes the Microsoft Azure logo, the Databricks logo, a search bar, and user information. The left sidebar contains navigation options like 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, Experiments, Features, Models, Serving, and Marketplace. The main area displays a Python code cell with the following code:

```
1 sparkDF = spark.read.csv("/databricks-datasets/bikeSharing/data-001/day.csv", header="true", inferSchema="true")
2 display(sparkDF)
3
```

Below the code, a table visualization is shown with the following columns: Instant, dteday, season, yr, mnth, holiday, weekday, workingday, weathersit, temp, atemp, and hum. The table contains 7 rows of data. The bottom status bar indicates the command took 15.98 seconds to run.

	Instant	dteday	season	yr	mnth	holiday	weekday	workingday	weathersit	temp	atemp	hum
1	1	2011-01-01	1	0	1	0	6	0	2	0.344167	0.369625	0.805833
2	2	2011-01-02	1	0	1	0	0	0	2	0.363478	0.353739	0.696087
3	3	2011-01-03	1	0	1	0	1	1	1	0.196364	0.189405	0.437273
4	4	2011-01-04	1	0	1	0	2	1	1	0.2	0.212122	0.590435
5	5	2011-01-05	1	0	1	0	3	1	1	0.226957	0.22927	0.436957
6	6	2011-01-06	1	0	1	0	4	1	1	0.204348	0.233209	0.518261
7	7	2011-01-07	1	0	1	0	5	1	2	0.196522	0.208839	0.498696

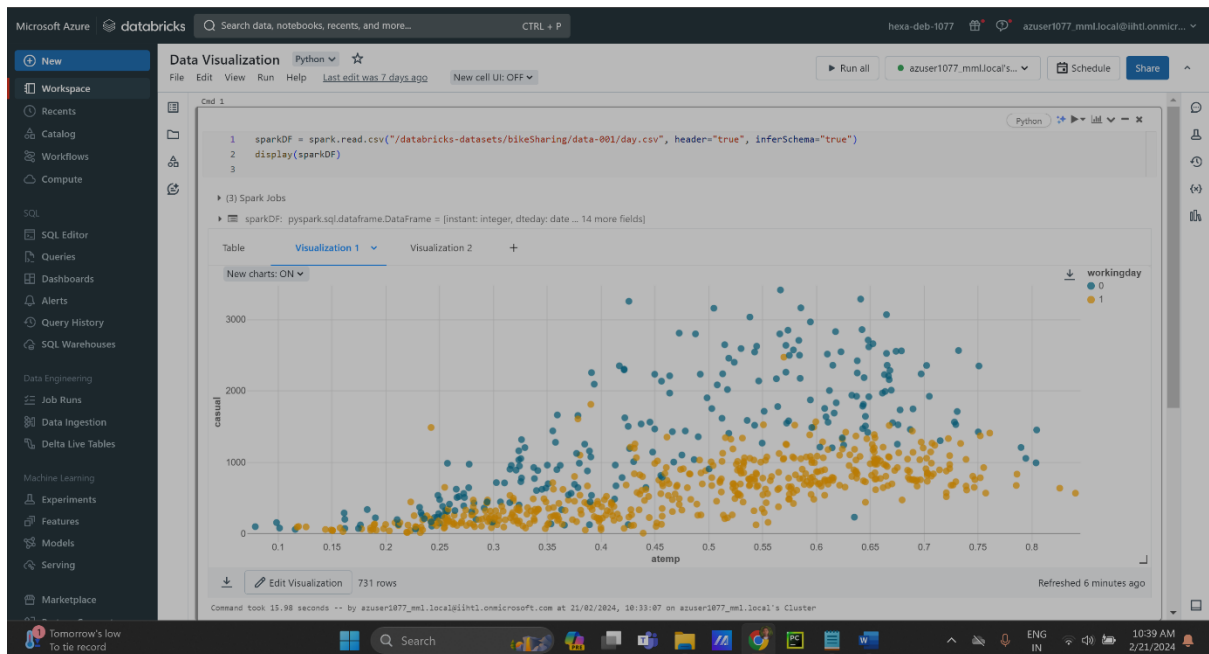
- Giving Parameters for Scatter plot according to columns in table

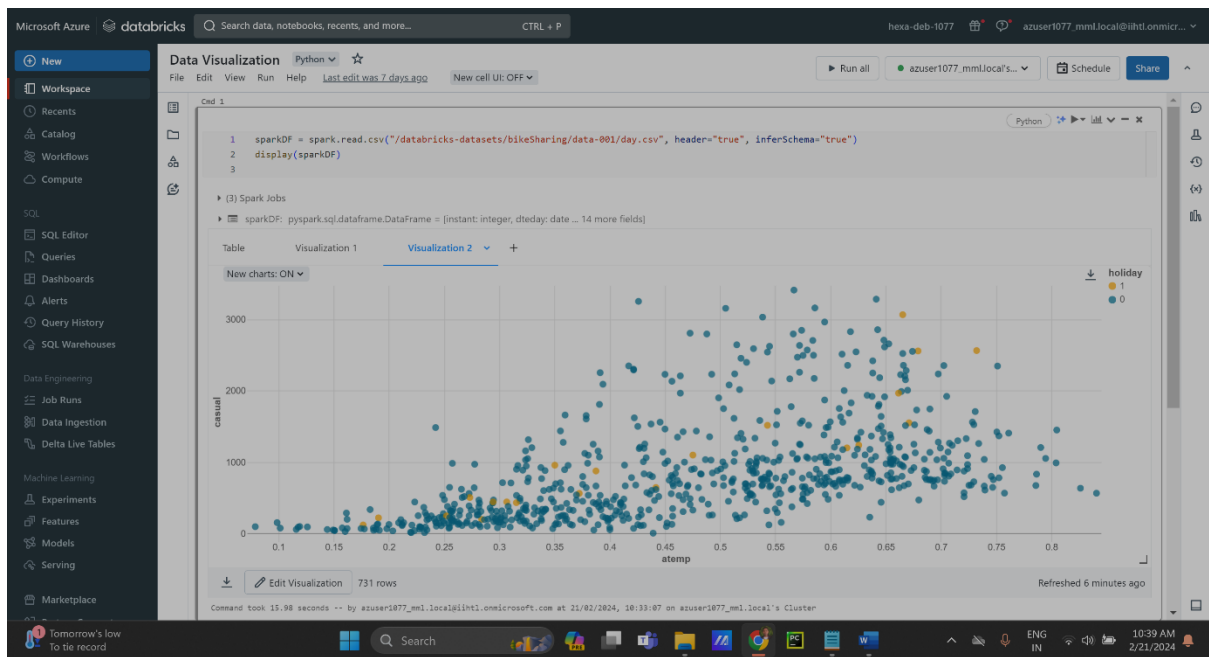


- Giving Parameters for Bubble plot according to columns in table



- Visualization In Azure Databricks





- After creating the visualization, we can also create data profile.

Data profile : Used to analyse the data based on some trends, you can also order and filter the data.

