

Exploratory Data Analysis (EDA) and data visualization

Exploratory Data Analysis (EDA) and data visualization are crucial steps in understanding and gaining insights from your data. Databricks provides a powerful platform for performing EDA and visualizing data using its integrated tools and libraries such as Apache Spark, Matplotlib, Seaborn, and Plotly. Here's a basic guide on how to perform EDA and visualize data in Databricks:

1.Importing Data:

Start by importing your dataset into Databricks. You can do this by uploading your data file directly to Databricks or by connecting Databricks to your data source (like Azure Blob Storage, AWS S3, or any other data lake).

2.Data Exploration:

Once your data is imported, begin exploring it to understand its structure, contents, and quality. You can use various Spark DataFrame operations and functions to perform basic exploration tasks such as checking the schema, summarizing statistics, and sampling the data.

```
python
```

```
# Read data into a Spark DataFrame
```

```
df = spark.read.csv("path_to_your_data.csv", header=True,  
inferSchema=True)
```

```
# Display the schema
```

```
df.printSchema()
```

```
# Summary statistics
```

```
df.describe().show()
```

Sample data

df.sample(False, 0.1).show()

1.Data Visualization:

After getting a basic understanding of the data, visualize it to gain deeper insights. Databricks supports multiple visualization libraries like Matplotlib, Seaborn, and Plotly. You can create various types of plots such as histograms, scatter plots, box plots, etc., to visualize the distribution, relationships, and trends in your data.

python

Import required libraries

import matplotlib.pyplot as plt

import seaborn as sns

Example: Histogram

sns.histplot(df.select("column_name").toPandas(), bins=20)

plt.title("Histogram of Column")

plt.xlabel("Values")

plt.ylabel("Frequency")

plt.show()

Example: Scatter Plot

sns.scatterplot(x="column1", y="column2", data=df.toPandas())

plt.title("Scatter Plot")

plt.xlabel("Column 1")

plt.ylabel("Column 2")

`plt.show()`

- **Interactive Visualization:**

For interactive and more sophisticated visualizations, you can use libraries like Plotly. Plotly allows you to create interactive plots that can be customized and explored interactively within Databricks notebooks.

python

Import Plotly

import plotly.express as px

Example: Interactive Scatter Plot

**fig = px.scatter(df.toPandas(), x="column1", y="column2",
color="category_column", title="Interactive Scatter Plot")**

fig.show()

- **Dashboard Creation (Optional):**

If you want to create interactive dashboards with multiple visualizations, you can use Databricks' built-in dashboarding capabilities or integrate with external dashboarding tools like Tableau or Power BI.

- **Iterate and Refine:**

EDA is an iterative process. Continue exploring and visualizing your data, refine your analysis based on insights gained, and repeat until you have a good understanding of your data.

Activity:-

- Create cluster hexa-deb-1066 personalize policy and terminate after 30 min

The image consists of two screenshots from a Windows desktop environment.

The top screenshot shows the Microsoft Azure portal. The browser address bar displays the URL: `portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions%2F984f097c-963c-4eb6-a20d-839457ae9f08%2FresourceGroups%2Frg-azuser1066...`. The page title is "rg-azuser1066_mml.local-C3G5P_hexa-deb-1066 | Overview". The left sidebar shows the "Overview" tab selected. The main content area displays "Deployment is in progress" with the following details:

- Deployment name: rg-azuser1066_mml.local-C3G5P_hexa-deb-1066
- Subscription: Azure subscription 1
- Resource group: rg-azuser1066_mml.local-C3G5P
- Start time: 2/21/2024, 10:37:37 AM
- Correlation ID: 121061d4-3d2c-4dcd-84fa-7c8da0559e7f

Below this, a table shows the deployment details:

Resource	Type	Status	Operation details
hexa-deb-1066	Azure Databricks Service	Created	Operation details

On the right side of the deployment details, there are three sections:

- Microsoft Defender for Cloud**: Secure your apps and infrastructure. [Go to Microsoft Defender for Cloud >](#)
- Free Microsoft tutorials**: Start learning today >
- Work with an expert**: Azure experts are service provider partners who can help manage your assets on Azure and be your first line of support. [Find an Azure expert >](#)

The bottom screenshot shows the Databricks cluster configuration page. The browser address bar displays the URL: `adb-1690845575385826.azure.databricks.net/?o=1690845575385826#setting/clusters/0221-051220-try67f6l/configuration`. The page title is "azuser1066_mml.local's Cluster". The left sidebar shows the "Compute" tab selected. The main content area displays the cluster configuration:

- Policy**: Personal Compute
- Access mode**: Single user access
- Single user**: azuser1066_mml.local
- Performance**: Databricks Runtime Version: 14.3 LTS ML (includes Apache Spark 3.5.0, Scala 2.12). ☐ Use Photon Acceleration
- Node type**: Standard_DS3_v2 (14 GB Memory, 4 Cores)
- Terminate after**: ☒ 30 minutes of inactivity
- Tags**: No custom tags. > Automatically added tags
- Advanced options**: (collapsed)

On the right side of the configuration page, there is a "Summary" box:

Summary

- 1 Driver, 14 GB Memory, 4 Cores
- Runtime: 14.3 x-cpu-m-ml-scala2.12
- Standard_DS3_v2, 0.75 DBU/h

Launch workspace for databricks:-

The screenshot shows the Microsoft Azure portal interface. The main heading is 'deb-hexa-1066' with the subtitle 'Azure Databricks Service'. The left sidebar contains navigation options like Overview, Activity log, Access control (IAM), Tags, Diagnose and solve problems, Settings, Virtual Network Peering, Encryption, Networking, Properties, Locks, Monitoring, Diagnostic settings, Automation, CLI / PS, Tasks (preview), and Export template. The main content area is titled 'Essentials' and displays the following information:

- Status: Active
- Resource group: rg-azuser1066_mml.local-C3GSP
- Location: East US
- Subscription: Azure subscription 1
- Subscription ID: 984f097c-963c-4eb6-a20d-839457ae9f08
- Tags: Add tags
- Managed Resource Group: databricks-rg-deb-hexa-1066-y34rudj0z6ms
- URL: https://adb-5691141690344570.10.azure.databricks.net
- Pricing Tier: Premium (+ Role-based access controls) (Click to change)

Below this information is a large red Databricks logo and a blue 'Launch Workspace' button. At the bottom, there are several tiles for 'Documentation', 'Getting Started', 'Import Data from File', 'Import Data from Azure Storage', 'Notebook', 'Admin Guide', and 'Link Azure ML workspace'.

Write code for display data :-

The screenshot shows the Databricks workspace interface. The main heading is 'deb-hexa-1066' with the subtitle 'Azure Databricks Service'. The left sidebar contains navigation options like New, Workspace, Recent, 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, and Features. The main content area is titled 'Untitled Notebook 2024-02-21 11:09:49' and displays the following code:

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

Below the code is a table showing the results of the query. The table has 12 columns: instant, dteday, season, yr, mnth, holiday, weekday, workingday, weather, temp, atemp, and hu. The table contains 7 rows of data.

instant	dteday	season	yr	mnth	holiday	weekday	workingday	weather	temp	atemp	hu
1	2011-01-01	1	0	1	0	6	0	2	0.344167	0.369025	0.8
2	2011-01-02	1	0	1	0	0	0	2	0.363478	0.353739	0.6
3	2011-01-03	1	0	1	0	1	1	1	0.196364	0.189405	0.4
4	2011-01-04	1	0	1	0	2	1	1	0.2	0.212122	0.5
5	2011-01-05	1	0	1	0	3	1	1	0.226957	0.22927	0.4
6	2011-01-06	1	0	1	0	4	1	1	0.204348	0.233209	0.5
7	2011-01-07	1	0	1	0	5	1	2	0.196522	0.208829	0.4

Below the table is a status bar indicating '731 rows | 12.71 seconds runtime' and 'Refreshed now'.

Click on the + icons and get option of visualization:-

The screenshot shows the Databricks notebook interface. The notebook is titled "Untitled Notebook 2024-02-21 11:09:49". The code in the notebook is:

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

The output shows a table with 731 rows and 12.71 seconds runtime. The table has columns: instant, season, yr, mnth, holiday, weekday, workingday, weathersit, temp, atemp, and humidity. A '+' icon is visible next to the table header, indicating the option to visualize the data.

Select scatter vizulation type and x column y column group by and placement get a scatterd graph :-

The screenshot shows the Databricks Visualization Editor. The visualization type is set to "Scatter". The X column is "atemp", the Y column is "casual", and the Group by is "workingday". The legend placement is set to "Automatic (Flexible)". The legend items order is set to "Normal". The scatter plot shows a positive correlation between atemp and casual, with data points colored by workingday (0: blue, 1: orange).

Overview of 3 level namespace and creating Unity Catalog objects:-

The Unity Catalog within Databricks offers a comprehensive solution for managing data and AI assets across Azure Databricks workspaces. It provides centralized access control, auditing, lineage, and data discovery capabilities. Here's an explanation of the key concepts and processes involved in using Unity Catalog, including the three-level namespace and creating Unity Catalog objects:

1.Three-Level Namespace:

Unity Catalog organizes data assets using a three-level namespace structure:

- **Metastore:** This is the top-level container for metadata. It provides a three-level namespace (catalog.schema.table) to organize data.
- **Catalog:** The first layer of the namespace hierarchy, used to categorize data assets.
- **Schema:** Also known as databases, schemas organize tables and views within a catalog.

2.Creating Unity Catalog Objects:

- **Metastore:** Administrators create a metastore for each region, attaching it to Databricks workspaces in the same region. It registers metadata about data and AI assets and manages access permissions.
- **Catalog:** Users can create catalogs to organize data assets. Each catalog contains schemas, and users with appropriate permissions can manage catalogs.

- **Schema:** Schemas, or databases, organize tables and views. They provide logical separation within catalogs. Users can create schemas and manage their contents.

3.Unity Catalog Objects:

- **Tables:** Tables reside within schemas and contain rows of data. They can be managed or external. Managed tables are managed by Unity Catalog, while external tables manage their own lifecycle and file layout.
- **Views:** Views are read-only objects created from tables and views in the metastore. They provide a logical representation of data without storing it physically.
- **Volumes:** Volumes are in Public Preview and provide access to non-tabular data stored in cloud object storage. They contain directories and files and offer governance for non-tabular data.
- **Models:** Although not strictly data assets, registered models can also be managed within Unity Catalog and reside at the lowest level in the object hierarchy.

4.Creating Objects:

- **Managed Storage:** Managed tables and volumes can be stored at the metastore, catalog, or schema levels. Managed storage locations are declared during object creation and determine where data and metadata files are stored.
- **Storage Credentials and External Locations:** To manage access to underlying cloud storage for external tables, external volumes, and managed storage, Unity Catalog uses storage credentials and external locations.

5.Identity Management:

- Unity Catalog uses identities from the Azure Databricks account to resolve users, service principals, and groups for enforcing permissions.
- Users, service principals, and groups must be added to workspaces to access Unity Catalog data. This process is called identity federation, and it enables access to Unity Catalog across workspaces.

In summary, Unity Catalog provides a structured approach to managing data and AI assets within Databricks, offering centralized governance, access control, auditing, and data discovery capabilities. Through its three-level namespace and various object types, users can organize, secure, and manage their data effectively.