# Azure Databricks Coding Challenge
# Name : Mitushi Vishwakarma

## Question 1 : Exploratory data analysis (EDA) in Databricks & Visualizing data in Databricks

Exploratory Data Analysis (EDA) in Databricks involves utilizing its robust tools for data exploration and visualization. Databricks provides a unified analytics platform that integrates Apache Spark for big data processing with SQL, Python, R, and Scala. Here's how we can perform EDA and visualize data in Databricks:

1. **Data Import:** First, import the dataset into Databricks. We can upload data files directly or connect to various data sources like Azure Data Lake Storage, AWS S3, or a database.
2. **Data Exploration**:
   - Visualization: Use Databricks' built-in visualization tools or integrate with popular libraries like Matplotlib, Seaborn, or Plotly for creating plots and charts.
3. **Visualizing Data**:
   - Databricks Visualizations: Databricks provides interactive visualization capabilities. You can create charts directly from DataFrame results using Databricks notebooks. Supported chart types include line plots, bar charts, scatter plots, and more.
4. **Advanced Analysis:**
   - Machine Learning: Databricks supports building machine learning models using Spark MLlib or integrating with popular ML frameworks.
   - Time Series Analysis: For time series data, use Databricks' time-series functions and visualization tools to analyze trends, seasonality, and anomalies.
5. **Collaboration and Sharing**: Databricks allows collaborative work through shared notebooks. You can share notebooks with team members for collaborative analysis and decision-making.
6. **Exporting Results**: Export your analysis results and visualizations in various formats like CSV, Excel, or images for sharing with stakeholders or further analysis.
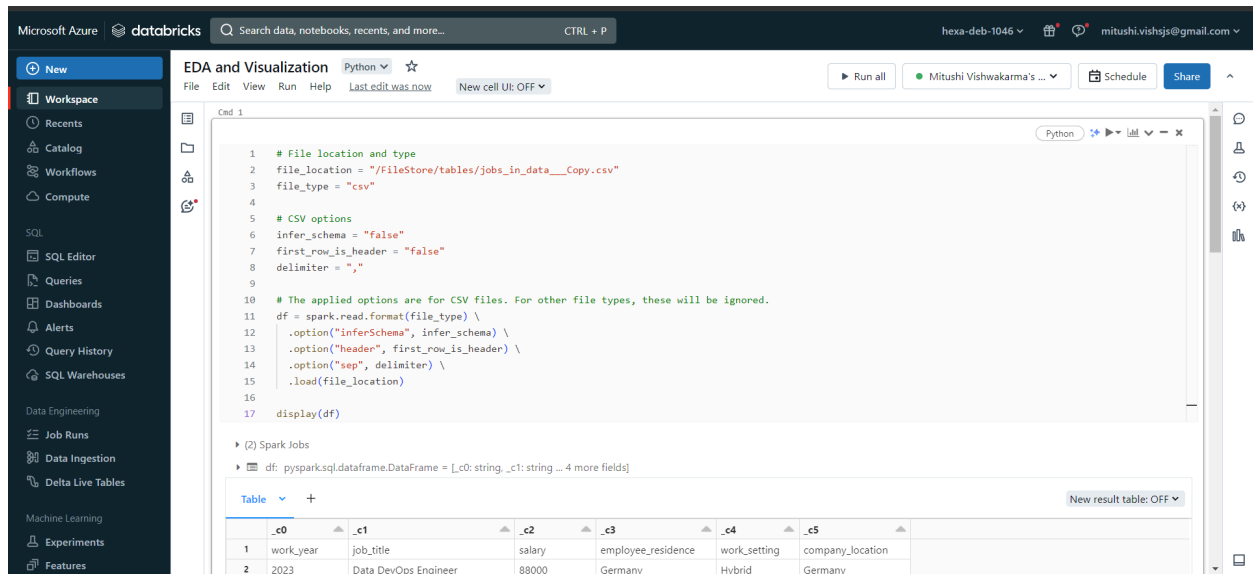
## Creating new visualizations:

In order to create visualizations, we need to have data.
Here I provided data and used the display() method.
- After creating a table
- Click on + symbol
- Click on visualization.
- Select the type of visualization, here I selected "Pie" and "Bar"

In the visualization editor, We can select which type of chart we want in the drop down. There are various options then to groupBy, change colors,etc.
You can also change the color of the chart visualization by clicking "edit" on the bottom left corner.
You can edit , delete , view and download visualization charts.

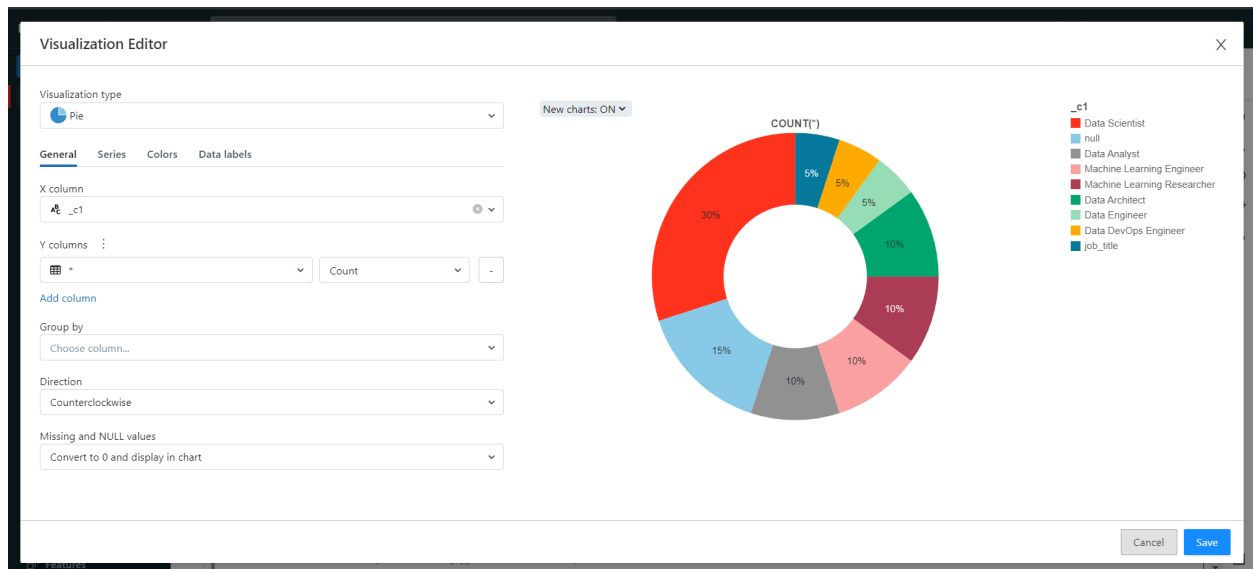## Uploaded CSV file into DBFS and displaying it :



## Creating Visualization 1 using "Pie" :

## Creating Visualization 2 using "Bar" :