Name : Parth Nandedkar
Date : 21 Feb 2024
Topics : Azure Databricks
Batch : Data Engineering Batch-1

**Q1.Exploratory data analysis (EDA) in Databricks & Visualizing data in Databricks.**

Exploratory Data Analysis (EDA) is a crucial step in the data analysis process, and Databricks provides a powerful platform for performing EDA and visualising data. Here's a general guide on how you can perform EDA and visualise data using Databricks:

**Loading Data:**

Begin by loading your dataset into Databricks. You can do this from various sources such as Azure Data Lake Storage, Azure Blob Storage, Azure SQL Database, etc.
Databricks supports multiple file formats like CSV, Parquet, JSON, etc. You can use the appropriate reader to load your data into a DataFrame.

**Understanding Data:**

Once the data is loaded, you can use DataFrame operations to explore its structure and contents. Methods like display, show, describe, schema, etc., can be helpful.
Check for missing values, data types, summary statistics, unique values, etc., to get a better understanding of your data.

**Data Visualization:**

Databricks supports various visualisation libraries such as Matplotlib, Seaborn, Plotly, etc., which you can use directly in your notebooks.
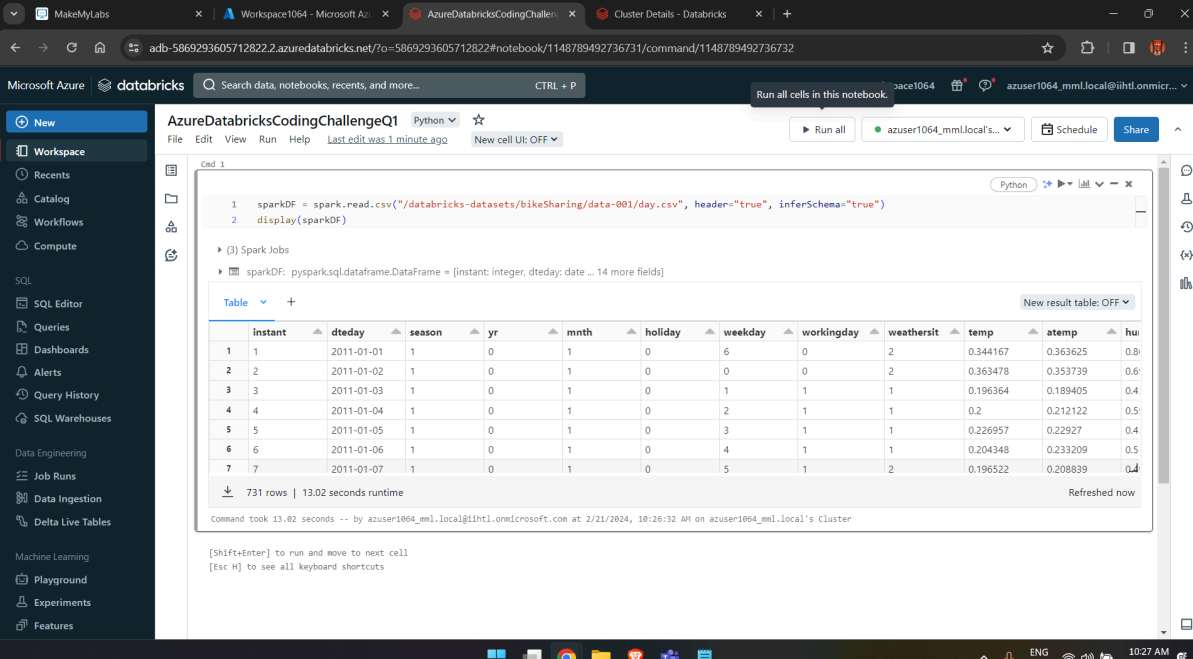You can create different types of plots like histograms, scatter plots, bar plots, line plots, etc., to visualise the distribution, relationships, and patterns in your data.

After creating a cluster and creating a new notebook.
By using PySpark using **spark.read.csv()**  we have stored/converted the csv
data which was already present in the databricks workspace into dataframe.

**sparkDF** is a data frame which stores all the values.

display() is the method to see whether the data is stored or not.

As per the above screenshot we can say that the data is stored in the data frames and we can use this data for visualisation purposes.
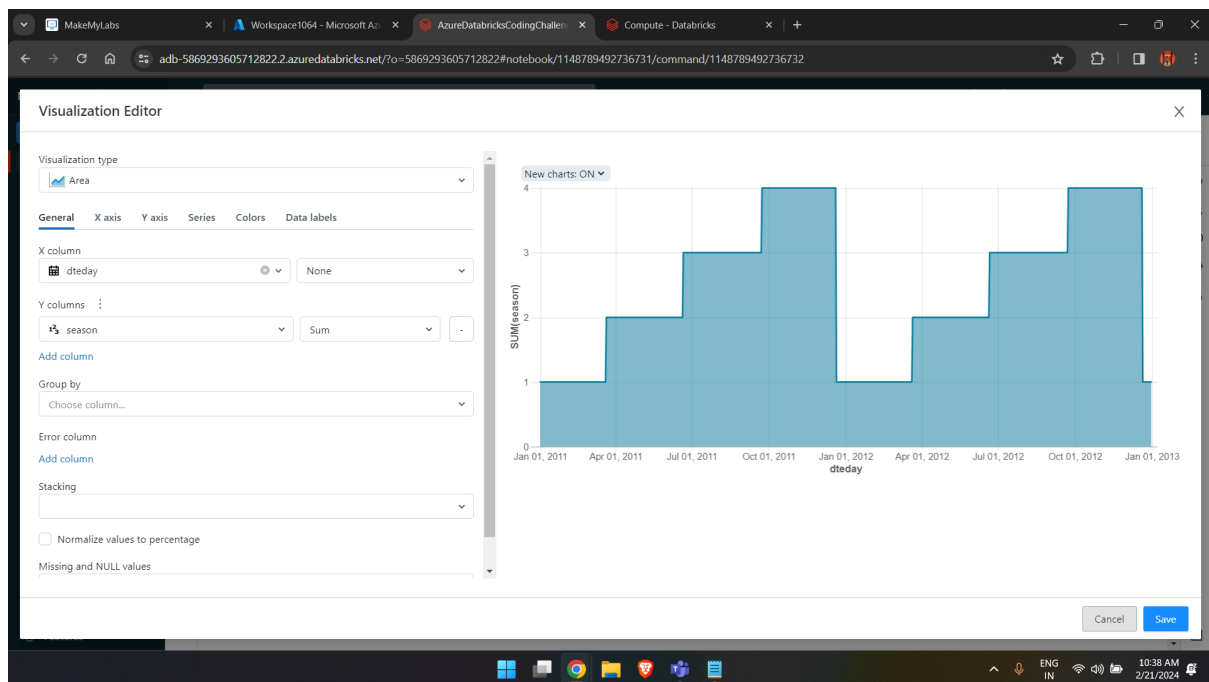


On table creation there is an option named as Visualization where there are options like below
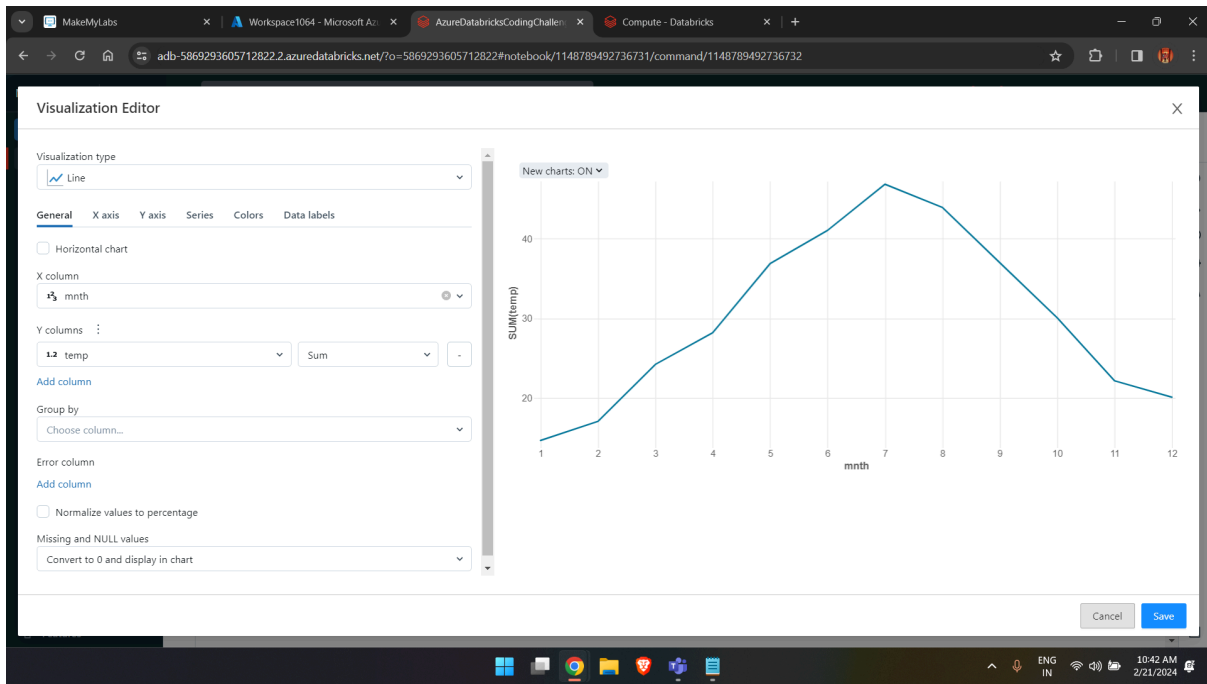
Now according to data we can use the different visualisations so as to analyse the data.



This is the area wise graph of the every day present in the data and based on season.
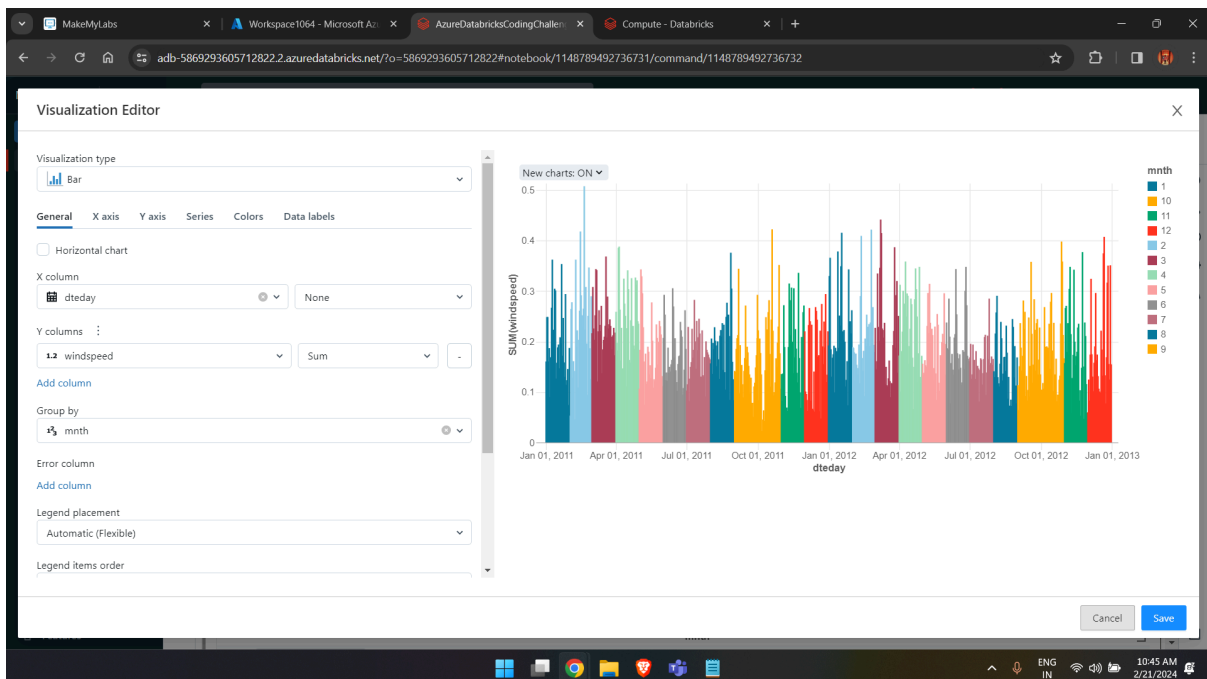


Above is the scattered graph of per day according to temperature.

Above is the line graph for the same.

We can analyse that during month 6-9 the temperature is at its peak and before that temp is rising and after that temperature is falling.



In above visualisation the wind speed can be analysed for every day according to month the colours are different as we have used group by months.