

STOCK PRICE PREDICTION

Loading and pre-processing the provided stock price dataset typically involves the following steps:

Data Collection: We need to obtain historical stock price data. We provide from Kaggle website the link for the dataset on which tasks are performed is given below:

<https://www.kaggle.com/datasets/prasoonkottarathil/microsoft-lifetime-stocks-dataset>

Data Format: First we should ensure that the data is in a format that you can work with. Typically, this data will be in CSV (Comma-Separated Values) format or some other structured format.

Import Libraries: We will need to import relevant libraries for data manipulation and analysis. Common libraries for this task include *pandas*, *numpy*, and *matplotlib* for visualization. We may also need libraries to fetch data from the web.



The screenshot shows a Google Colab notebook interface. The browser address bar at the top displays the URL: `colab.research.google.com/drive/1f90FZiXoxlPMdbMOgBe7-oBgHUTA-26A?authuser=0#scrollTo=z1jpB9XBT41L`. The notebook is titled "phase3ab.ipynb". Below the title bar, there is a menu with options: "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help". The main area of the notebook shows a code cell with the following Python code:

```
[1] import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

The code cell is marked with a green checkmark and the number "1" in the left margin, indicating it has been executed successfully. On the left side of the notebook, there is a file explorer showing a folder named "sample_data" and a file named "MSFT.csv". On the right side, there is a "Release notes" panel with a close button (X) and a message: "Please follow our [blog](#) to see more information about tricks, and featured notebooks such as [Analyzing](#)". The date "2023/09/22" is displayed at the bottom right of the panel.

Load the Data: Use pandas to read the dataset into a DataFrame. Here we load a CSV file:

```
' > sample_data
  MSFT.csv
]
```

```
✓ [2] # dataset is saved as CSV in the working directory
0s df = pd.read_csv('/content/MSFT.csv')
```

Data Cleaning: Check for missing values, duplicate rows, and outliers in the dataset. We need to handle missing data through imputation, drop duplicates, and handle outliers based on our analysis requirements.

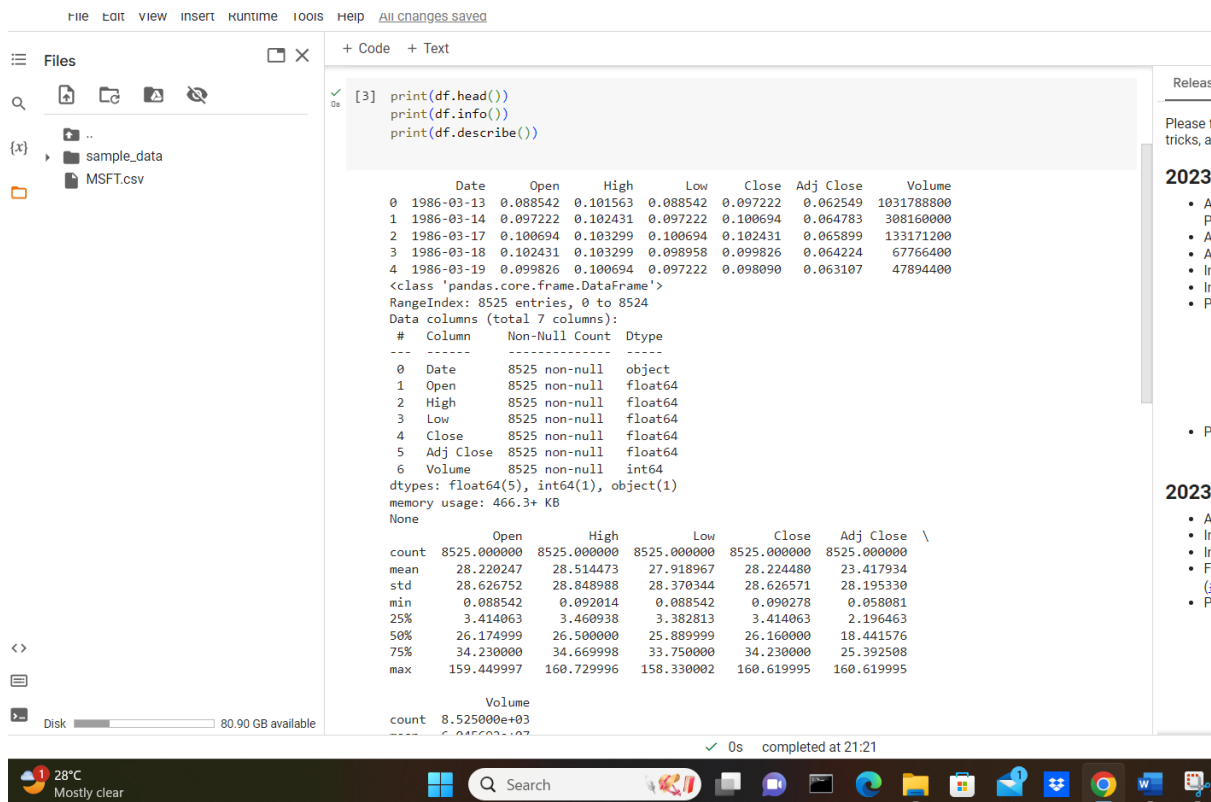
```
{x} > ..
  sample_data
  MSFT.csv
```

```
✓ # Check for missing values
0s print(df.isnull().sum())

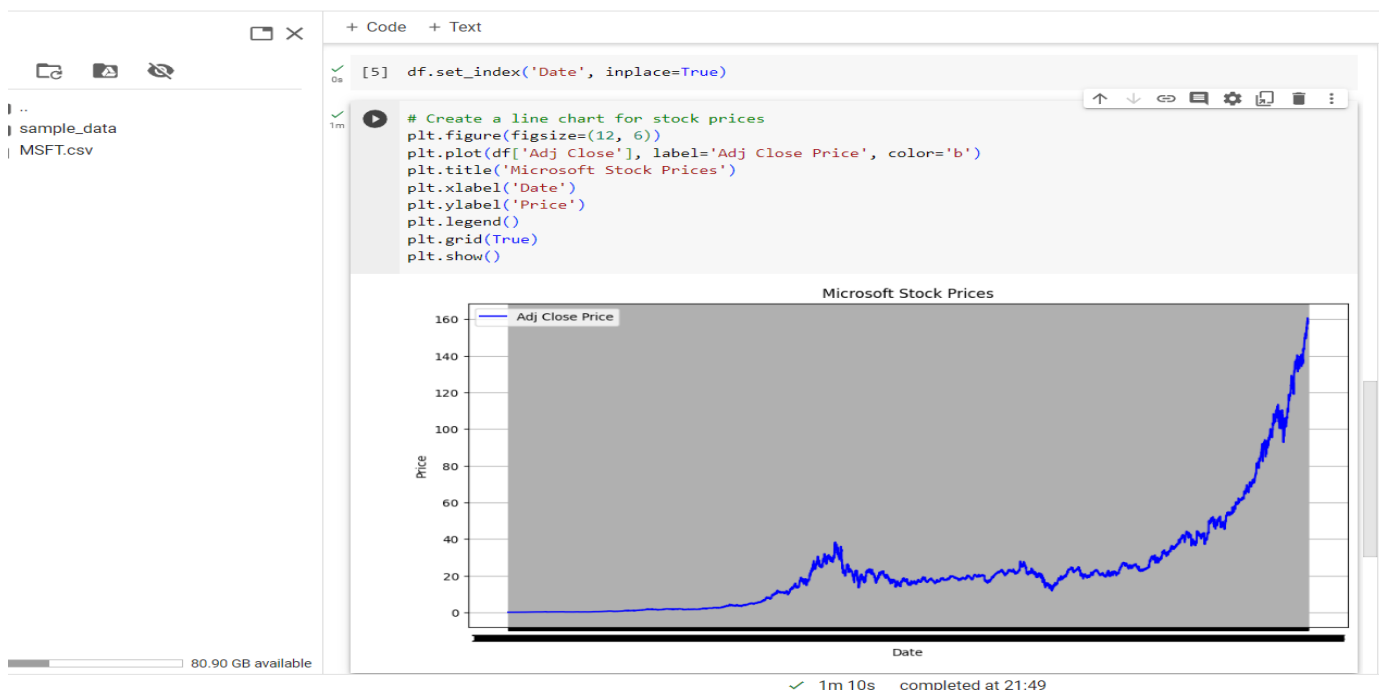
# Handle missing values (e.g., forward fill or interpolate)
df['Open'].fillna(method='ffill', inplace=True)
```

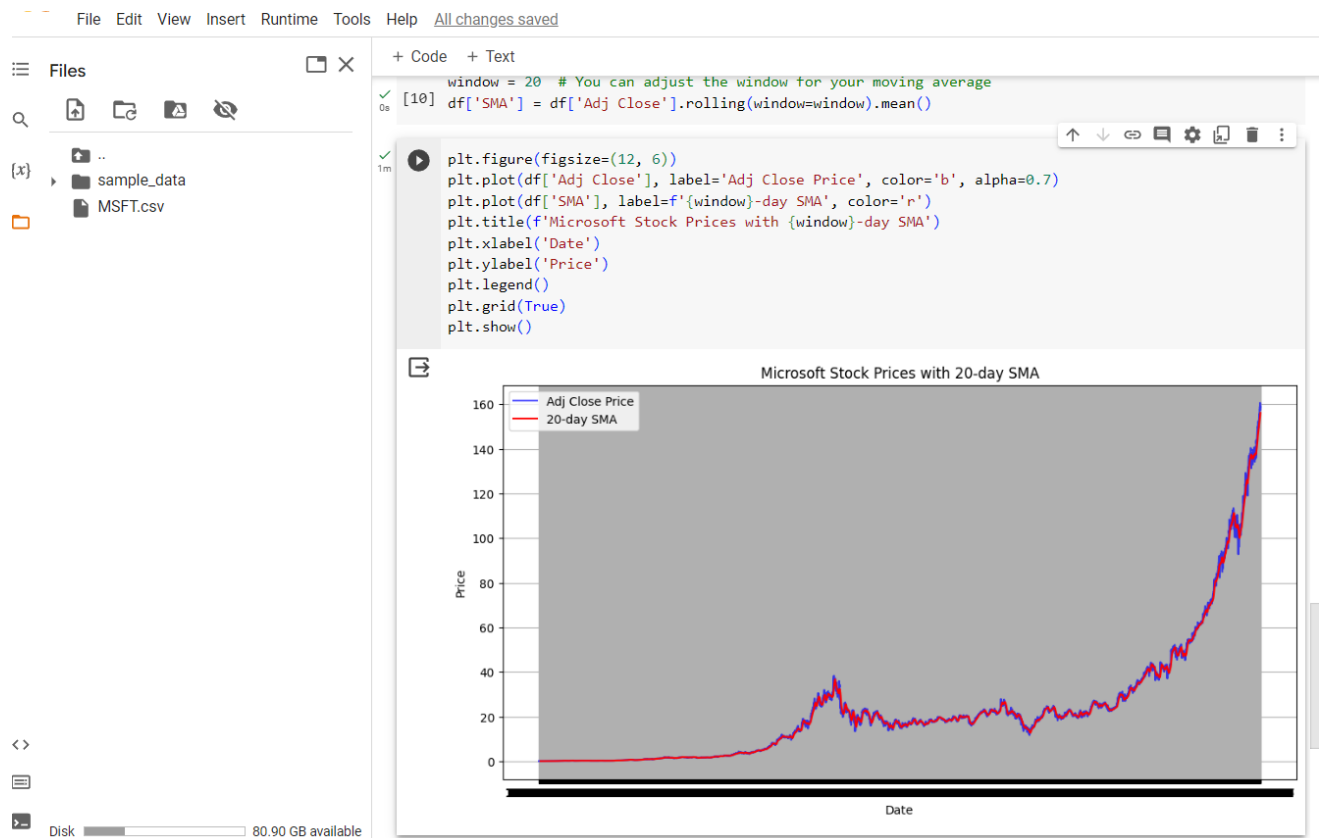
```
➡ Date      0
   Open      0
   High      0
   Low       0
   Close     0
   Adj Close  0
   Volume    0
   dtype: int64
```

Data Exploration: Explore the dataset to gain insights. We use methods like `describe()`, `info()`, and visualizations to understand the data better.



After loading and pre-processing stock price dataset, it's a good practice to create visualizations to gain insights and better understand the data. Common visualizations for financial time series data like stock prices include line charts, candlestick charts, volume plots, and moving averages. Here we can create these visualizations using Python.





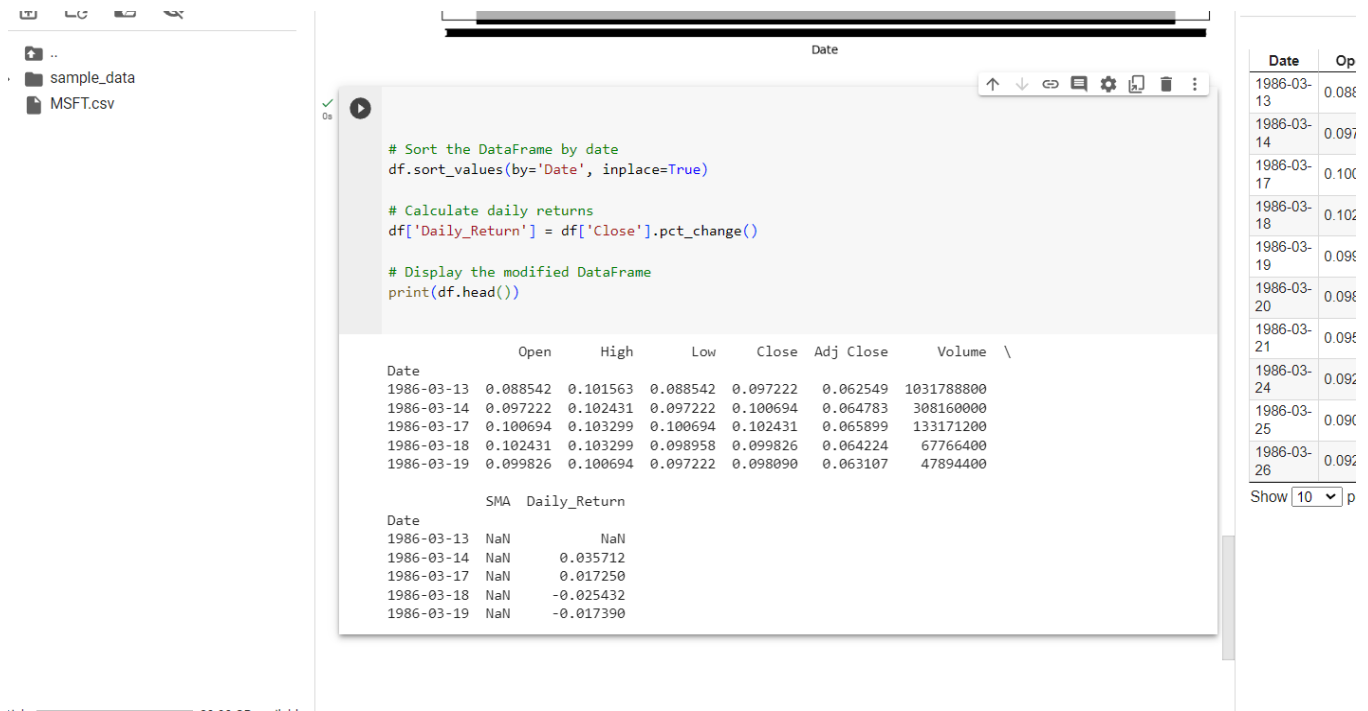
Data Pre-processing:

Date Conversion: Convert the date column to a datetime object if it's not already in that format.

Sorting: Sort the Data Frame by date if it's not already in chronological order.

Feature Engineering: Create additional features if needed, such as moving averages, daily returns, or technical indicators.

Normalization/Scaling: If you plan to use machine learning models, you might want to normalize or scale the data. Common techniques include Min-Max scaling or standardization.

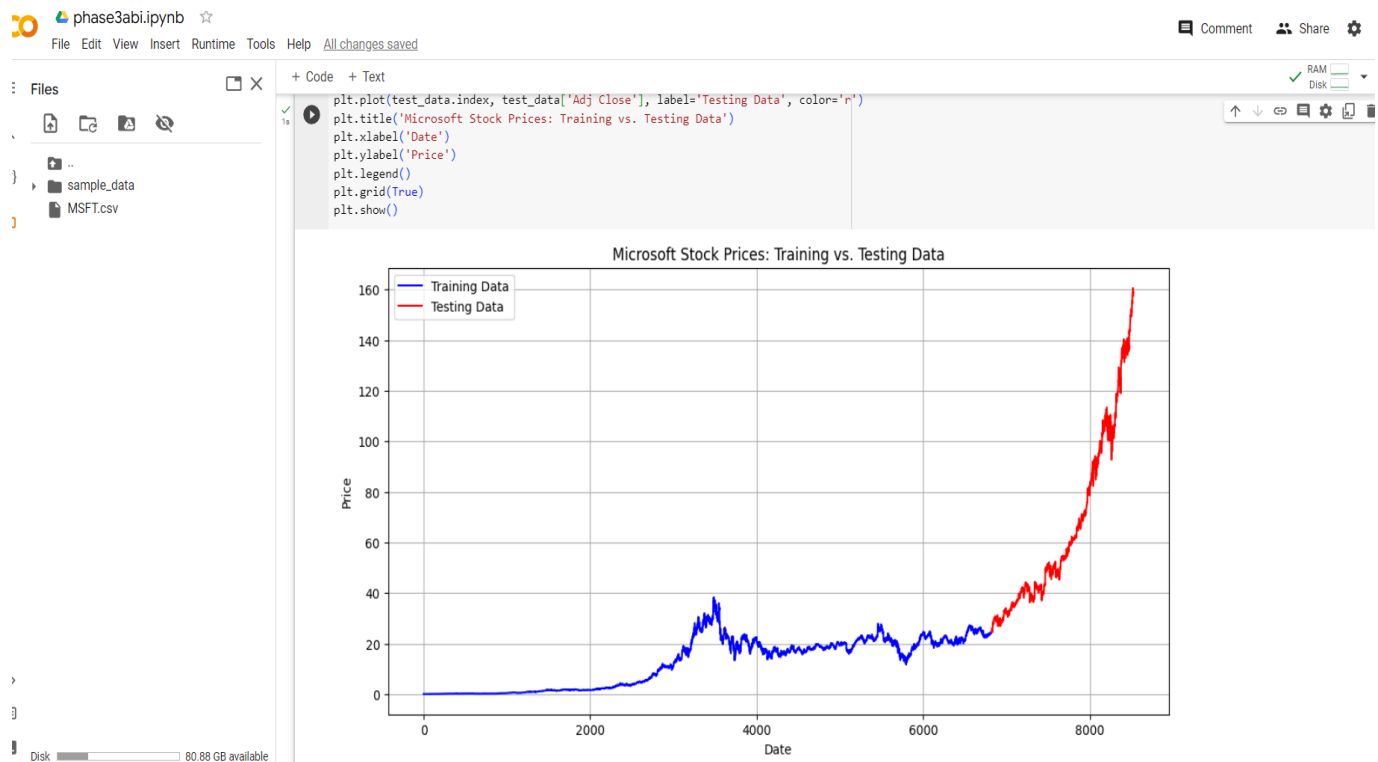


Data Splitting: If you plan to build a predictive model, split the dataset into training and testing sets. Make sure to maintain the chronological order.

```
[10]: from sklearn.model_selection import train_test_split

# Split the data into training and testing sets
train_size = 0.8 # 80% for training, 20% for testing
train_data, test_data = train_test_split(df, train_size=train_size, shuffle=False)
```

To visualize the training and testing datasets for stock price dataset, we create separate line charts for each dataset to compare them. This helps in understanding how our model's predictions (on the test dataset) compare to the actual stock price movements (from the training dataset).



We load the pre-processed data, containing the 'Adj Close' prices.

We have separate Data Frames for training and testing data, 'train_data' and 'test_data'.

We create line charts for both datasets, with different colours to distinguish between training and testing data.

The x-axis represents the date, and the y-axis represents the adjusted closing price.

This visualization helps compare the stock price movements in the training and testing datasets. It's useful for assessing how well our model generalizes to unseen data. In an ideal scenario, the testing data would closely follow the trends seen in the training data, indicating that our model has learned and generalizes well.

In addition to the basic preprocessing steps mentioned earlier, there are several other advanced preprocessing steps that can be performed

on a stock price dataset. These steps are often used to enhance the quality and relevance of the data for analysis and modelling.