# TCS Stock Market Forecasting

## 1. Introduction

The dataset used for this analysis consists of minute-wise stock data for Tata Consultancy Services (TCS). This data is crucial for understanding the minute-to-minute fluctuations in stock prices, which can be pivotal for traders and investors. The primary objectives of this analysis are to identify trends within the data, visualize key metrics, and forecast future stock prices.

## 2. Objective:

- Explore and understand the features of the TCS stock dataset, including open, high, low, close prices, and volume.
- Perform data preprocessing, such as handling missing values and outliers to ensure data integrity.
- Identify key factors influencing stock prices using statistical analysis techniques.
- Develop predictive models to forecast stock prices accurately based on historical data.

## 3. Scope:

1. Data Exploration: Understanding the dataset and its features (open, high, low, close, volume).
2. Data Preprocessing: Cleaning the dataset by handling missing values, detecting, and removing outliers.
3. Feature Analysis: Identifying significant features that impact stock prices.
4. Data Visualization: Creating visual representations to analyze the relationships between features and stock prices.
5. Model Building: Developing and evaluating predictive models for stock price forecasting.
6. Reporting: Documenting findings and providing recommendations to stakeholders.

## 4. Methodology:

- Data Collection:
  - The dataset will be sourced from a local CSV file containing minute-wise stock data for TCS.
- Data Preprocessing:
  - Handle missing data through visualization and imputation techniques.
  - Detect and remove outliers using clustering and statistical methods to ensure clean data for analysis.
- Exploratory Data Analysis (EDA):
  - Utilize descriptive statistics and visualizations such as histograms, heatmaps, and scatter plots to explore the dataset and understand key trends.
- Modeling:
  - Implement ARIMA for time series forecasting to predict future stock prices based on historical data trends.

## i. Data Preprocessing

- Loading the Data:
  The dataset comprises several columns, including timestamp, open, high, low, close, and volume. The timestamp column is essential for tracking the chronological order of price movements.

- Handling Missing Values:

    In this analysis, missing values and duplicates are removed to ensure the integrity of the dataset, providing a clean slate for subsequent analyses.

- Feature Extraction"

    From the timestamp, various date and time components are extracted, including year, month, day, hour, minute, and second. This helps in analyzing trends over specific periods.
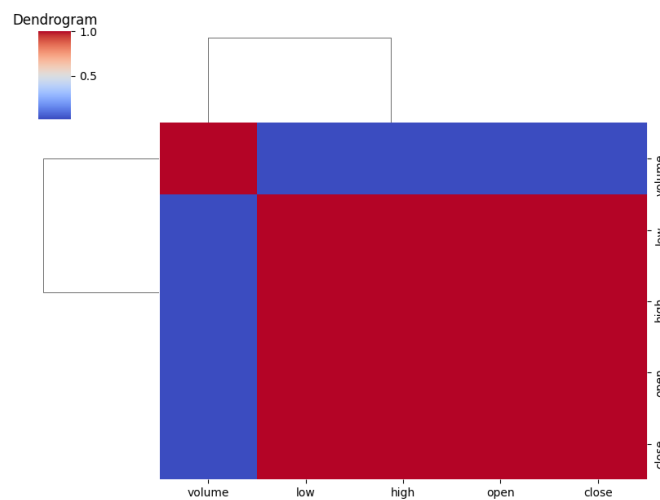
## ii. Exploratory Data Analysis (EDA)

- Summary Statistics

    A detailed description of summary statistics for numerical columns such as open, close, high, low, and volume was generated to understand the central tendency and dispersion of the data.
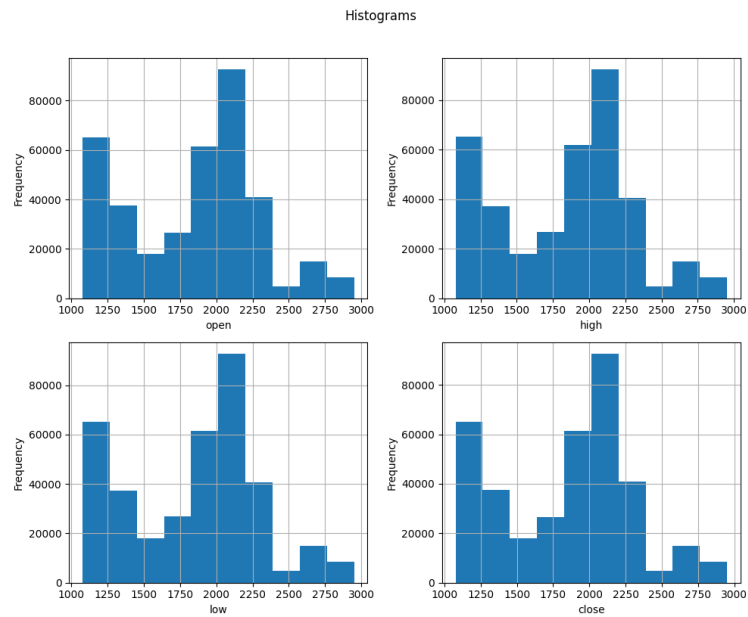
- Correlation Matrix and Heatmap

    A correlation analysis was conducted to explore relationships between numerical features. The heatmap visualizes these correlations, helping to identify which variables are positively or negatively related.
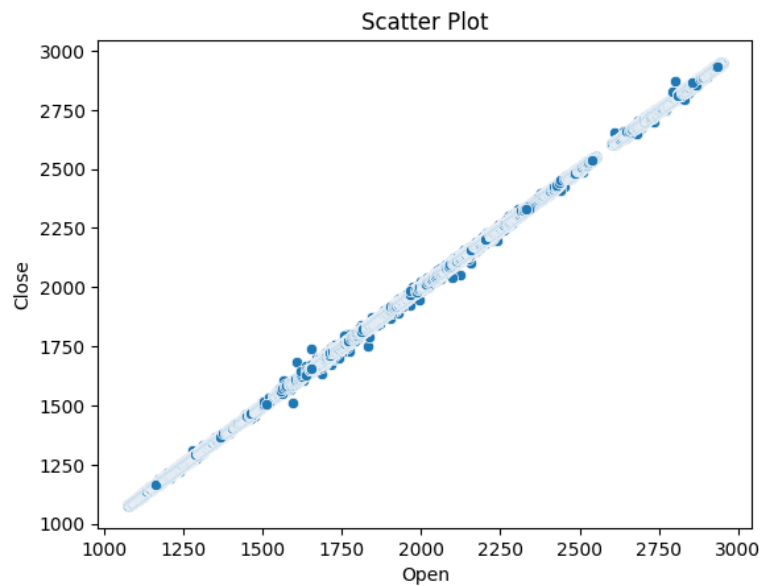


- Histograms

    Histograms were plotted for the numerical columns to visualize the distribution of values for open, close, high, low, and volume.
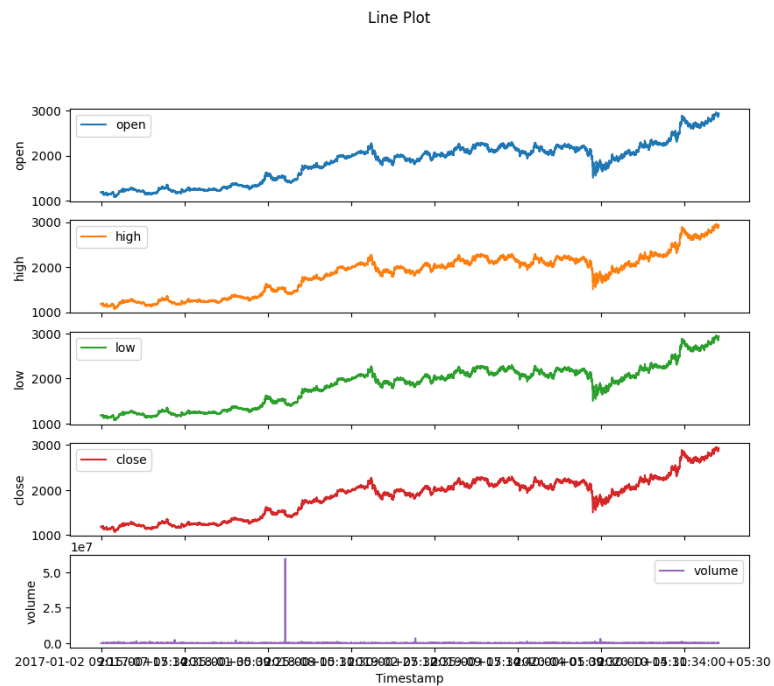
Histograms



- Scatter Plot

Scatter plots were created to demonstrate the relationship between open and close prices, allowing for a visual inspection of how closely related these two metrics are.



## iii. Data Visualization

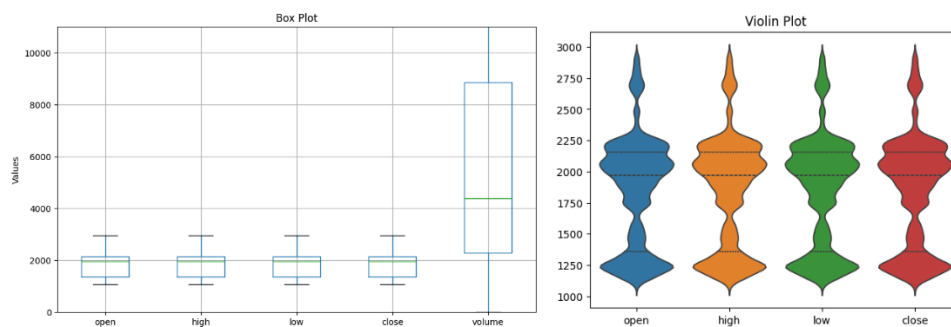- Line Plot:

Line plots were generated for the open, high, low, close, and volume columns over time, showcasing the price movements and trading volumes.

Line Plot

- **Box Plot and Violin Plot:**

    Box plots and violin plots were used to display the data distribution and identify outliers within the dataset.
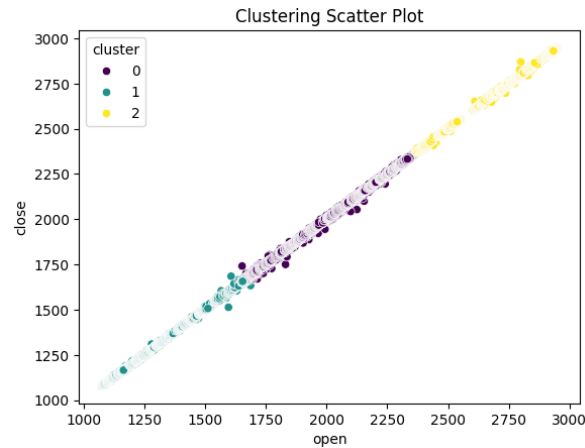


- **Candlestick Chart:**

    A candlestick chart was produced to visualize the stock's price movements over time, providing insights into market trends and volatility.

## iii a. Clustering

K-Means Clustering:

K-Means clustering was applied to identify patterns in stock prices, utilizing open-close and high-low values to classify different price movement behaviours.



## 5. Modelling

### Time Series Forecasting

- Stationarity Check (ADF Test):
  An Augmented Dickey-Fuller (ADF) test was performed to check the stationarity of the close price column, a crucial step for effective time series forecasting.



- ARIMA Modeling
  The data was differenced to achieve stationarity, followed by the selection of appropriate ARIMA model parameters (p, d, q). The model was then fitted to the data to forecast future stock prices.
- Forecasting
  Results of the forecasting were displayed, including a plot comparing historical values with forecasted values, allowing for an evaluation of the model's predictive power.

- Liner regression

    Linear regression is used to model the relationship between user behavior (e.g., past ratings or interactions) and the predicted rating or preference for a specific item. In the recommendation system, this method helps by creating a baseline prediction for how a user may rate a product or service based on their historical data. While straightforward, linear regression may struggle with non-linear relationships common in user behavior patterns.

- Decision tree

    Decision tree models help segment users based on their preferences and behavior. The tree structure breaks down the data into smaller subsets, allowing the recommendation system to identify patterns and make personalized suggestions. For example, users who prefer certain types of items can be categorized accordingly, and recommendations are tailored based on those segments. However, decision trees can sometimes overfit, particularly in cases with highly varied user preferences.

- Support Vector Machine (SVM):

    SVM, adapted for recommendation systems (SVR for regression tasks), works by finding the optimal boundaries between different user preference classes. This method excels in identifying complex patterns between users and items, such as implicit feedback or latent user-item interactions, making it a good choice for personalized recommendations. By maximizing the separation between preferences, SVM can deliver more accurate suggestions, though it may require tuning to handle large, sparse datasets.

- Ensemble:

    Ensemble learning combines the strengths of multiple models (e.g., linear regression, decision trees, SVM) to enhance the recommendation's accuracy. Techniques like bagging, boosting, or stacking leverage the diverse outputs of individual algorithms to improve predictions and reduce errors. In your recommendation system, ensemble methods help balance different types of user behaviors, leading to more reliable and tailored suggestions by mitigating the weaknesses of single-model approaches.

## 7. Conclusion

In summary,

- The project provided valuable insights into the factors influencing TCS stock prices.
- Predictive models developed can assist investors in making informed decisions.
- Future work could explore additional features and machine learning models for enhanced forecasting.