

# **Final Term Project Report**

## **Time Series Forecasting of Stock Prices**

**Submitted by:**

**Mohammed Ismail Sarfaraz Shaik**

**G23033389**

**Master's in data science (Class of 2026)**

**Course:**

**Time Series Analysis and Modeling**

**DATS 6313**

**Instructor:**

**Reza Jafari**

**Submission Date:**

**May 2025**

# Table of Contents

## Table of Contents

|   |           |
|---|-----------|
| <b>TABLE OF CONTENTS</b> .....                              | <b>2</b>  |
| <b>TABLE OF FIGURES</b> .....                               | <b>4</b>  |
| <b>TABLE OF TABLES</b> .....                                | <b>6</b>  |
| <b>ABSTRACT</b> .....                                       | <b>7</b>  |
| <b>INTRODUCTION</b> .....                                   | <b>8</b>  |
| <b>DESCRIPTION OF THE DATASET AND PREPROCESSING</b> .....   | <b>9</b>  |
| OVERVIEW OF THE DATASET .....                               | 9         |
| VARIABLE DEFINITIONS.....                                   | 9         |
| DATA PREPROCESSING .....                                    | 9         |
| a. Handling Missing Values .....                            | 9         |
| b. Categorical Encoding .....                               | 10        |
| c. Differencing and Normalization .....                     | 10        |
| d. Resampling.....  | 10        |
| PLOT OF DEPENDENT VARIABLE VS TIME: .....                   | 10        |
| Observations: .....   | 11        |
| ACF AND PACF OF DEPENDENT VARIABLE .....                    | 11        |
| Observations: .....   | 13        |
| CORRELATION MATRIX .....                                    | 13        |
| Observations: .....   | 14        |
| TRAIN-TEST SPLIT.....                                       | 15        |
| <b>STATIONARITY ANALYSIS</b> .....                          | <b>15</b> |
| ROLLING MEAN AND VARIANCE (RAW vs. DIFFERENCED).....        | 15        |
| DIFFERENCED DATA: .....                                     | 16        |
| ADF AND KPSS TEST RESULTS.....                              | 17        |
| ACF & PACF PLOTS (DIFFERENCED SERIES) .....                 | 18        |
| <b>TIME SERIES DECOMPOSITION</b> .....                      | <b>19</b> |
| TATAMOTORS.....   | 19        |
| TATASTEEL .....   | 20        |
| TCS .....   | 20        |
| <b>HOLT-WINTERS METHOD</b> .....                            | <b>22</b> |
| <b>BASE MODELS EVALUATION</b> .....                         | <b>23</b> |
| Observations: .....   | 27        |
| <b>FEATURE SELECTION AND DIMENSIONALITY REDUCTION</b> ..... | <b>27</b> |

|  |           |
|--|-----------|
| <b>MULTIPLE LINEAR REGRESSION MODELING .....</b>                   | <b>28</b> |
| FEATURE SELECTION AND DIMENSIONALITY REDUCTION .....               | 28        |
| MODEL FITTING AND INTERPRETATION .....                             | 28        |
| TATAMOTORS.....  | 29        |
| TATASTEEL.....   | 29        |
| TCS .....  | 29        |
| VISUAL ANALYSIS .....  | 30        |
| SUMMARY .....  | 33        |
| <b>ARMA-ARIMA-SARIMA/MULTIPLICATIVE MODEL.....</b>                 | <b>34</b> |
| <b>ARMA-ARIMA-SARIMA MODEL DEVELOPMENT .....</b>                   | <b>37</b> |
| A. MODEL IDENTIFICATION – GPAC AND ACF/PACF .....                  | 37        |
| B. ARMA MODEL ESTIMATION USING LEVENBERG–MARQUARDT ALGORITHM ..... | 37        |
| <i>ARMA Model Metrics:</i> .....                                   | 40        |
| C. ARIMA MODEL ESTIMATION .....                                    | 41        |
| <i>ARIMA Model Results:</i> .....                                  | 44        |
| <b>BOX-JENKINS MODEL.....</b>                                      | <b>45</b> |
| A. G-GPAC AND H-GPAC .....   | 45        |
| B. BOX-JENKINS MODEL ESTIMATION .....                              | 47        |
| <b>RESIDUAL ANALYSIS AND DIAGNOSTIC TESTS .....</b>                | <b>50</b> |
| (A) WHITENESS CHI-SQUARE TEST: .....                               | 50        |
| (B) ERROR VARIANCE AND PARAMETER COVARIANCE: .....                 | 50        |
| (C & D) BIAS AND FORECAST ERROR ANALYSIS: .....                    | 53        |
| (E) MODEL SIMPLIFICATION VIA ZERO–POLE CANCELLATION: .....         | 53        |
| <b>FORECAST FUNCTION: .....</b>                                    | <b>54</b> |
| TATAMOTORS.....  | 54        |
| TATASTEEL.....   | 56        |
| TCS .....  | 57        |
| <b>FINAL MODEL SELECTION SUMMARY .....</b>                         | <b>59</b> |
| TATAMOTORS.....  | 59        |
| TATASTEEL.....   | 60        |
| TCS .....  | 61        |
| <b>MULTIPLE STEP-AHEAD FORECASTING (H-STEP) .....</b>              | <b>62</b> |
| TATAMOTORS.....  | 62        |
| TATASTEEL.....   | 63        |
| TCS .....  | 64        |
| <b>SUMMARY AND CONCLUSION .....</b>                                | <b>65</b> |
| <b>APPENDIX.....</b>   | <b>66</b> |
| <b>REFERENCES .....</b>  | <b>67</b> |

## Table of Figures

|   |    |
|---|----|
| Figure 1 Stock close price over time .....                        | 11 |
| Figure 2 ACF/PACF of Close .....                                  | 12 |
| Figure 3 Correlation matrix.....                                  | 14 |
| Figure 4:Rolling Mean & Variance of Raw Close Price.....          | 16 |
| Figure 5: Rolling Mean & Variance of Differenced Close Price..... | 17 |
| Figure 6: acf/pacf of differenced data .....                      | 18 |
| Figure 7: Holt damped forecast results .....                      | 22 |
| Figure 8: TATAMOTORS- Base model forecasts .....                  | 24 |
| Figure 9: TATASTEEL- Base model forecasts .....                   | 25 |
| Figure 10: TCS- Base model forecasts.....                         | 26 |
| Figure 11: TATAMOTORS MLR forecast .....                          | 30 |
| Figure 12: TATASTEEL MLR forecasts.....                           | 31 |
| Figure 13: TCS- MLR forecasts.....                                | 31 |
| Figure 14: acf residuals - TATAMOTORS.....                        | 32 |
| Figure 15: acf residuals - TATASTEEL .....                        | 32 |
| Figure 16: acf residuals - TCS .....                              | 33 |
| Figure 17: GPAC - TATAMOTORS .....                                | 34 |
| Figure 18: GPAC - TATASTEEL.....                                  | 35 |
| Figure 19: GPAC - TCS .....                                       | 35 |
| Figure 20: ARMA stats – TATAMOTORS .....                          | 38 |
| Figure 21: ARMA stats – TATASTEEL.....                            | 39 |
| Figure 22: ARMA stats - TCS .....                                 | 40 |
| Figure 23: ARIMA stats – TATAMOTORS.....                          | 42 |
| Figure 24: ARIMA stats – TATASTEEL.....                           | 43 |
| Figure 25: ARIMA stats - TCS .....                                | 44 |
| Figure 26: BJ stats- TATAMOTORS .....                             | 47 |

|   |    |
|---|----|
| Figure 27: BJ stats - TATASTEEL.....                                  | 48 |
| Figure 28: BJ stats - TCS .....                                       | 49 |
| Figure 29: covariance matrix – TATAMOTORS.....                        | 51 |
| Figure 30: covariance matrix – TATASTEEL .....                        | 52 |
| Figure 31: covariance matrix - TCS.....                               | 52 |
| Figure 32: ARMA (2,2) Forecast Plot - TATAMOTORS .....                | 54 |
| Figure 33: ARIMA (3,1,2) Forecast Plot - TATAMOTORS .....             | 55 |
| Figure 34: Box-Jenkins (4,3,3,3) — 20-Step Forecast - TATAMOTORS..... | 55 |
| Figure 35: ARMA (2,2) Forecast Plot - TATASTEEL .....                 | 56 |
| Figure 36:ARIMA (2,1,2) Forecast Plot - TATASTEEL .....               | 56 |
| Figure 37: Box-Jenkins (4,3,3,3) — 20-Step Forecast- TATASTEEL.....   | 57 |
| Figure 38: ARMA (1,1) Forecast Plot – TCS .....                       | 57 |
| Figure 39: ARIMA (1,1,1) Forecast Plot – TCS .....                    | 58 |
| Figure 40: Box-Jenkins (2,2,2,2) — 20-Step Forecast - TCS .....       | 58 |
| Figure 41: Final Model Forecast - TATAMOTORS.....                     | 62 |
| Figure 42: Final Model Forecast - TATASTEEL.....                      | 63 |
| Figure 43: Final Model Forecast - TCS .....                           | 64 |

## Table of Tables

|   |    |
|---|----|
| Table 1: ADF, KPSS test of raw and differenced data ..... | 17 |
| Table 2: Model selection table - TATAMOTORS.....          | 59 |
| Table 3: Model selection table - TATASTEEL.....           | 60 |
| Table 4: Model selection - TCS.....                       | 61 |

## Abstract

This report presents a comprehensive time-series modeling and forecasting analysis of stock prices for three major Indian companies: TATAMOTORS, TATASTEEL, and TCS. The objective was to evaluate multiple forecasting models and identify the most suitable approach for each stock based on accuracy, statistical validity, and interpretability.

The analysis began with extensive data preprocessing, including differencing, feature engineering, and stationarity testing. Several baseline and advanced models were applied, including Naive, Exponential Smoothing, Holt's Linear and Damped Trend models, ARIMA, ARMA, and custom-built Box–Jenkins models with exogenous inputs.

Each model was evaluated using Root Mean Square Error (RMSE), Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC), and diagnostic checks such as the Ljung–Box Q-Test and residual whiteness. The final model selection was based on both statistical performance and diagnostic robustness.

Results show that ARMA was most suitable for TATAMOTORS due to its low RMSE and simplicity, while Box–Jenkins models outperformed all others for TATASTEEL and TCS, offering superior forecasting accuracy and residual behavior. The report concludes with multi-step forecasting visualizations and a discussion on model reliability and practical relevance for stock trend prediction.

## Introduction

Time series forecasting plays a pivotal role in predictive analytics, especially in domains like finance, economics, and operations. In this project, we undertake a comprehensive end-to-end time series modeling process using stock price data from three major companies: TATAMOTORS, TATASTEEL, and TCS. The objective is to analyze and forecast the normalized differenced closing prices using advanced statistical techniques and modeling frameworks.

The report begins with an exploration of the datasets, detailing preprocessing steps and data quality checks. Following this, the focus shifts to verifying stationarity and decomposing the series into trend and seasonal components. Baseline models such as naive, drift, and exponential smoothing are constructed to establish foundational performance metrics.

Subsequent sections delve into advanced models, including Multiple Linear Regression, ARMA, ARIMA, and SARIMA. We apply the Levenberg–Marquardt algorithm for parameter estimation and build Box–Jenkins models incorporating exogenous variables, selected through correlation analysis. Rigorous residual diagnostic tests, ACF/PACF studies, and hypothesis testing ensure the reliability of the models.

The final sections compare all models on performance metrics such as RMSE, AIC, and BIC, culminating in the selection of the best forecasting model for each stock. The report concludes with h-step predictions, a critical reflection on model limitations, and potential directions for future enhancements, including web-based deployment for real-time forecasting.

# Description of the Dataset and Preprocessing

## Overview of the Dataset

The dataset for this project comprises stock market data for three prominent Indian companies listed on the National Stock Exchange: TATAMOTORS, TATASTEEL, and TCS. The data spans multiple years and includes daily trading information such as:

- Close: Closing price of the stock (dependent variable)
- Volume: Number of shares traded
- Turnover: Total traded value in currency
- Trades: Number of trades executed
- Deliverable Volume: Number of shares delivered
- %Deliverable: Percentage of deliverable volume over total volume
- Symbol: Categorical feature indicating the stock name
- Date: Timestamp of the trading day

## Variable Definitions

- Dependent Variable: The differenced Close price of the stock
- Independent Variables: Volume, Turnover, Trades, Deliverable Volume, %Deliverable

Each stock had a unique feature most correlated with the target variable, selected as the exogenous input for advanced models.

## Data Preprocessing

### a. Handling Missing Values

- Missing or NaN entries were identified and handled by:
  - Forward fill (for isolated missing values)
  - Dropping rows with unresolved or non-recoverable NaNs post differencing

## b. Categorical Encoding

- The Symbol column was used to group data but not directly modeled as an input.
- No further one-hot encoding was necessary since each model was trained per stock.

## c. Differencing and Normalization

- First-order differencing was applied to all numerical features to remove trends and induce stationarity.
- Each feature was then standardized (zero mean, unit variance) to ensure comparability.

## d. Resampling

- The data was daily and already uniformly spaced in time. Hence, no resampling was necessary.

## Plot of Dependent Variable vs Time:

The following figure displays the raw closing stock prices for TATAMOTORS, TATASTEEL, and TCS over time:



Figure 1 Stock close price over time

### Observations:

- All three stocks exhibit **strong non-stationary behavior**, with clear trends and varying volatility over time.
- **TCS** shows a long-term upward trend with significant growth after 2013 and noticeable volatility spikes, particularly around 2016 and 2020.
- **TATAMOTORS** experienced major fluctuations during the 2008 financial crisis and again between 2010–2014, followed by a decline and then mild recovery.
- **TATASTEEL** shows cyclical behavior with multiple peaks and troughs, aligning with broader macroeconomic and commodity market trends.
- None of the series seem to exhibit seasonality, but all of them show evidence of **trend and level shifts**, suggesting the need for transformation to induce stationarity.

This analysis supports the necessity of differencing and formal stationarity testing (ADF, KPSS) before applying time series models.

### ACF and PACF of Dependent Variable

The following plots display the **Autocorrelation Function (ACF)** and **Partial Autocorrelation Function (PACF)** for the raw closing prices of the three stocks: **TATAMOTORS**, **TATASTEEL**, and **TCS**.

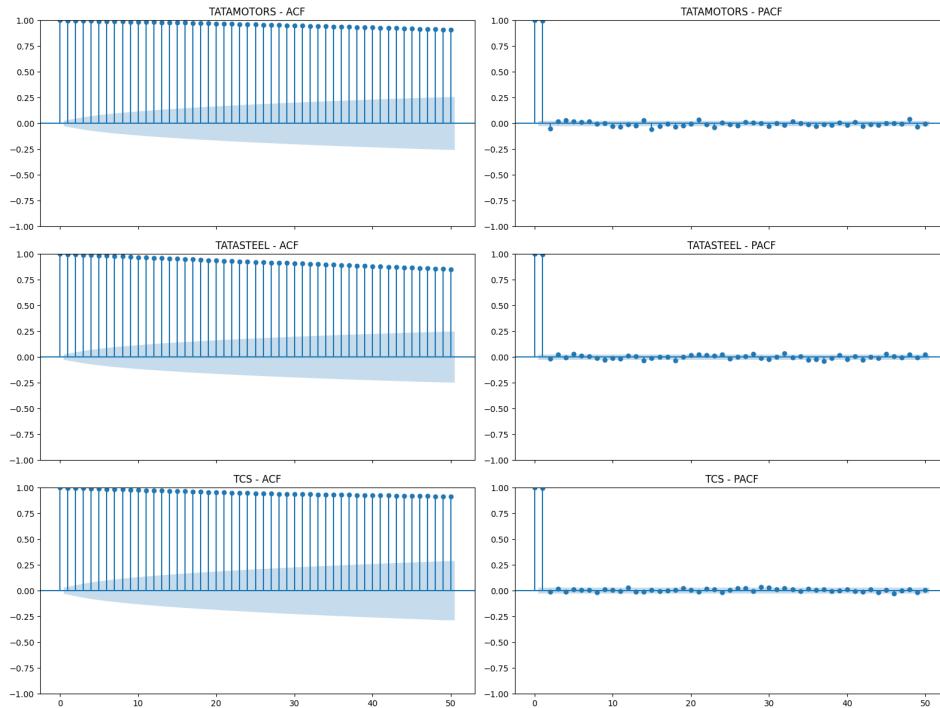


Figure 2 ACF/PACF of Close

## Observations:

- All three stocks exhibit strong autocorrelation across many lags in their ACF plots, with slow decay, which is a classic indication of non-stationarity in the time series.
- For all stocks, PACF plots show a sharp drop after lag 1 or 2, with significant spikes at lag 1 and possibly lag 2, before becoming negligible.
- This behavior suggests the presence of trend components, which must be removed through differencing to achieve stationarity.
- In particular:
  - TATAMOTORS and TATASTEEL show long memory effects and persistent autocorrelation, which supports the presence of unit roots.
  - TCS, despite a generally smoother profile, also demonstrates similar non-stationary behavior with high lag correlations.

These findings confirm that transformations such as differencing are necessary before building predictive models. They also provide an early indication that ARIMA-type models may be appropriate, where the AR and MA terms can be inferred from PACF and ACF structures post-differencing.

## Correlation Matrix

A Pearson correlation matrix was constructed using the raw (non-differenced) dataset to examine the linear relationships between the original features prior to any transformations or modeling. The heatmap below illustrates the pairwise correlations among all numerical variables.

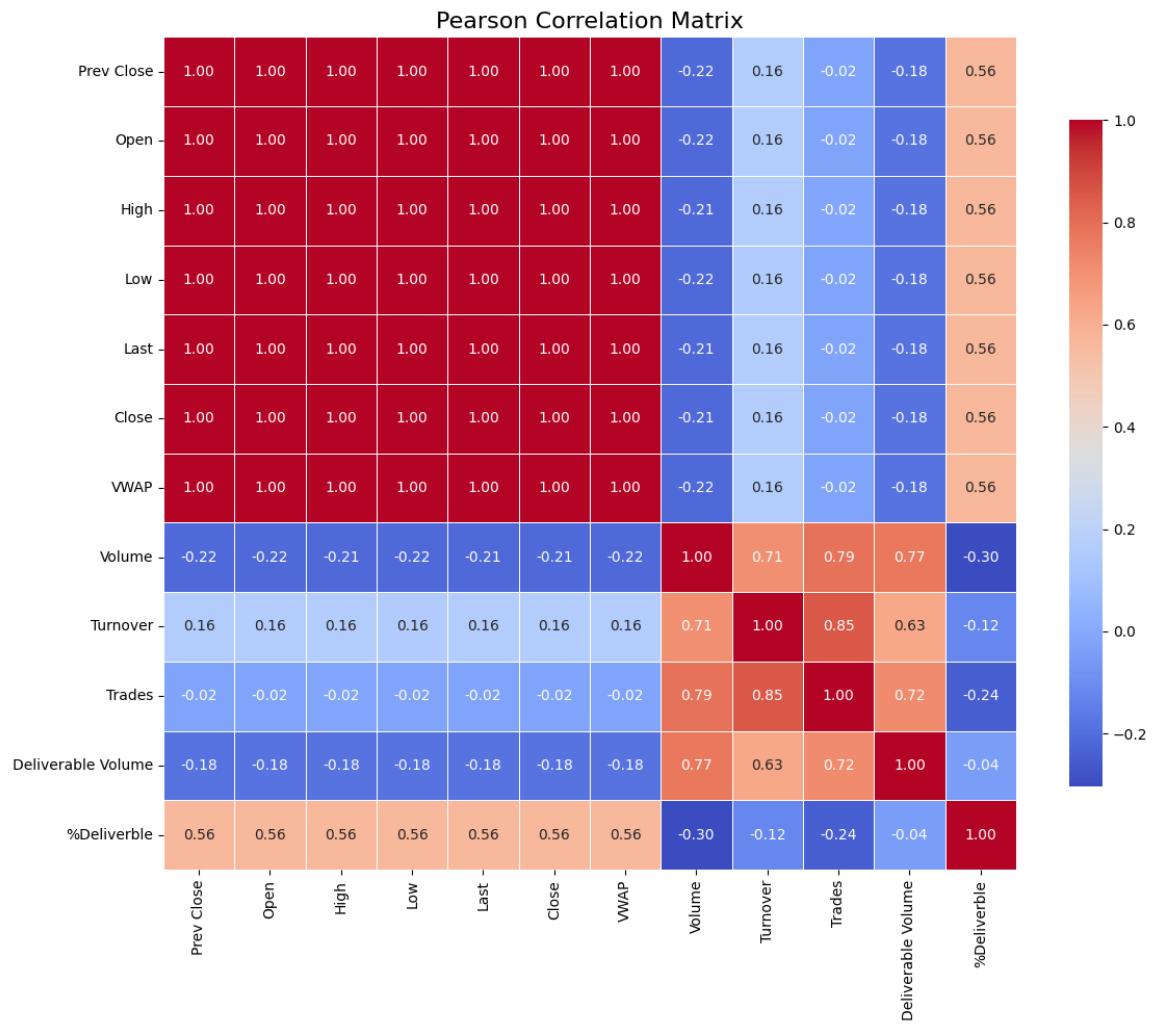


Figure 3 Correlation matrix

## Observations:

- A **perfect or near-perfect correlation** was observed among Open, High, Low, Last, Close, and Prev Close. This is expected as these are all price-based variables that move together in financial time series.
- Turnover, Volume, and Trades also showed **strong positive correlations** with each other, indicating that they capture similar aspects of trading activity.
- %Deliverable showed **moderate correlation** with price features (~0.56 with Close), suggesting it could hold independent predictive power in some contexts.

- Given the presence of **high multicollinearity**, especially among the price-based and volume-based variables, **feature elimination** was performed to retain only the most representative features for modeling:
  - Highly correlated features such as High, Low, Last, Prev Close, and VWAP were excluded to avoid redundancy.
  - Final selected features included: Close, Volume, Turnover, Trades, Deliverable Volume, and %Deliverable.

This step ensured cleaner input data for downstream time series modeling while preserving essential information for forecasting stock behavior.

## Train-Test Split

The final dataset for each stock was split into:

- Training set: 80% of the time series (chronologically earlier dates)
- Test set: 20% of the remaining data (latest dates)

This split ensured no look-ahead bias and preserved temporal ordering critical for time series validation.

## Stationarity Analysis

To verify stationarity, both statistical tests and visual diagnostics were applied to the raw and differenced stock closing prices of **TATAMOTORS**, **TATASTEEL**, and **TCS**.

## Rolling Mean and Variance (Raw vs. Differenced)

### Raw Close Data:

As shown in the rolling mean and variance plots of the raw close prices, the data exhibits non-stationary behavior — both mean and variance are non-constant over time, particularly around major structural shifts or volatility spikes.

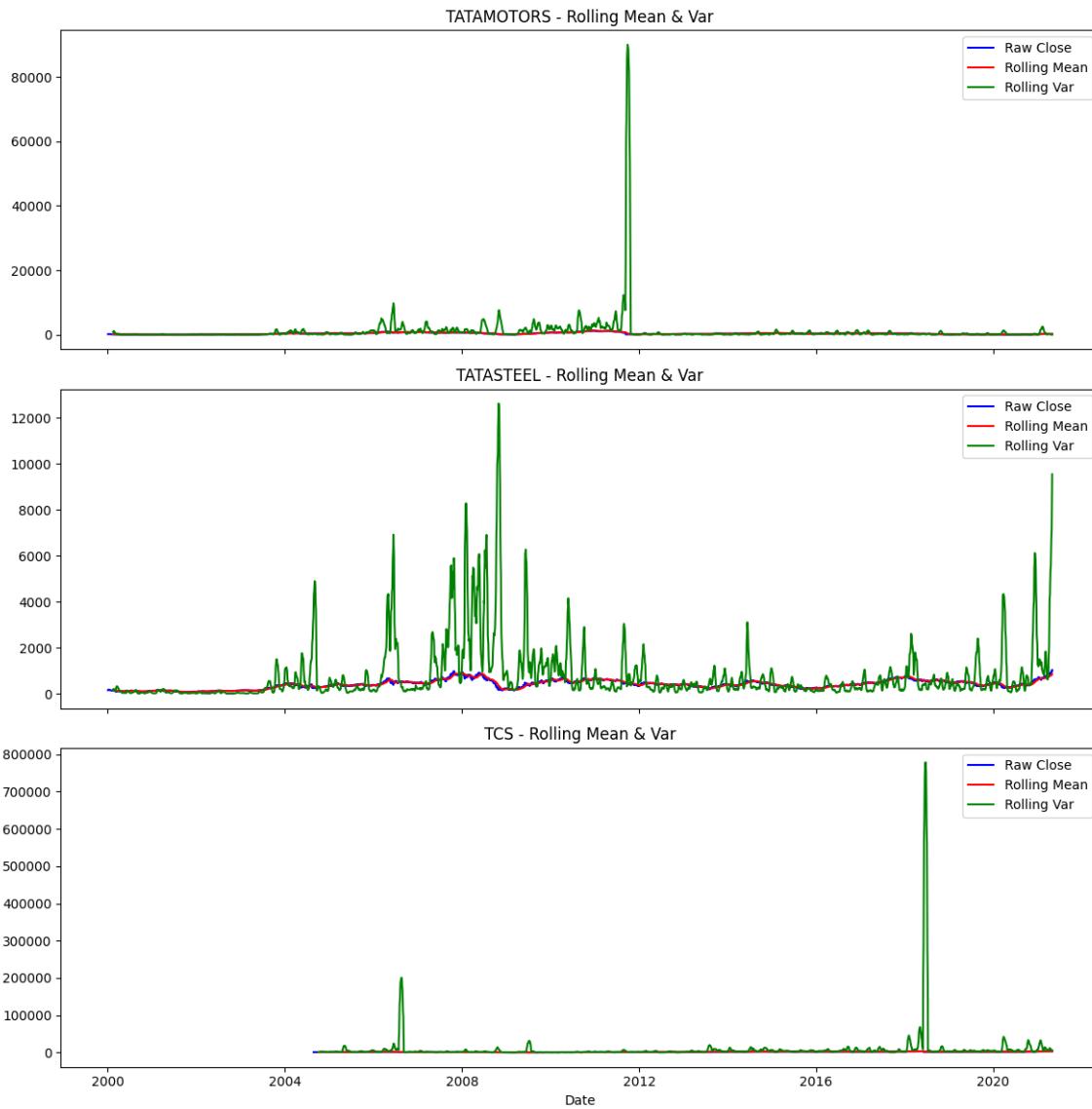


Figure 4: Rolling Mean & Variance of Raw Close Price

## Differenced Data:

Post first-order differencing, the mean stabilizes considerably, but the variance still appears to fluctuate, especially during market shocks. While this suggests that the variance is not perfectly constant, it is substantially more stable than in the raw series. Minor heteroskedasticity may persist, but the series is sufficiently stationary for linear modeling techniques like ARIMA and Box-Jenkins.

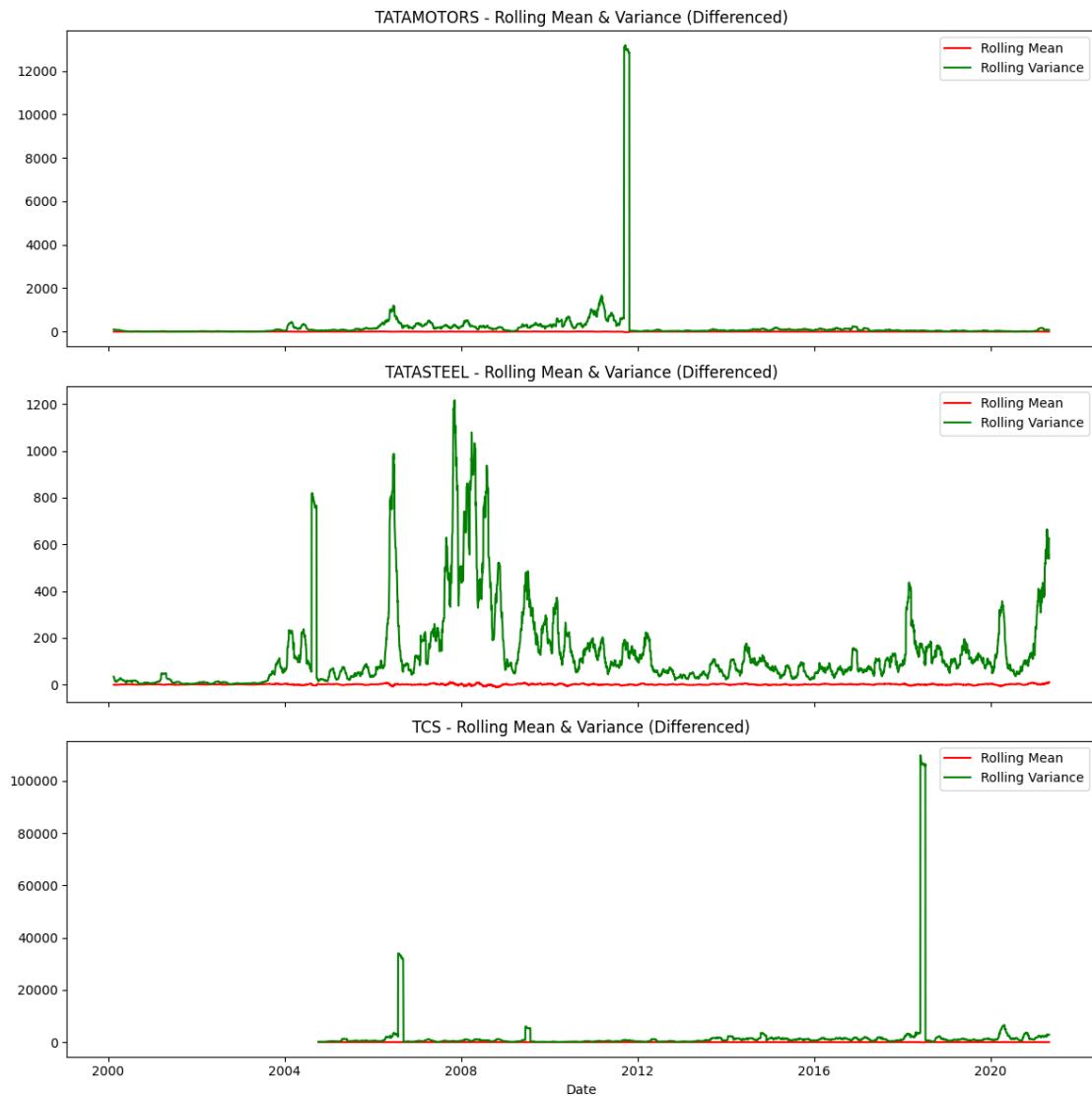


Figure 5: Rolling Mean & Variance of Differenced Close Price

## ADF and KPSS Test Results

Table 1: ADF, KPSS test of raw and differenced data

| STOCK      | ADF (raw)    | ADF (diff)   | KPSS (raw)   | KPSS (diff)  |
|------------|--------------|--------------|--------------|--------------|
| TATAMOTORS | <b>0.170</b> | <b>0.000</b> | <b>0.010</b> | <b>0.100</b> |
| TATASTEEL  | <b>0.525</b> | <b>0.000</b> | <b>0.010</b> | <b>0.100</b> |
| TCS        | <b>0.459</b> | <b>0.000</b> | <b>0.010</b> | <b>0.100</b> |

- ADF Test (Augmented Dickey-Fuller): High p-values on raw data indicate non-stationarity. After differencing, all p-values drop to 0.000, confirming strong evidence of stationarity.
- KPSS Test: Confirms the ADF results. Raw data shows significant trend stationarity violations. Post differencing, the null hypothesis of stationarity is not rejected ( $p=0.1$ ), indicating adequate stationarity.

## ACF & PACF Plots (Differenced Series)

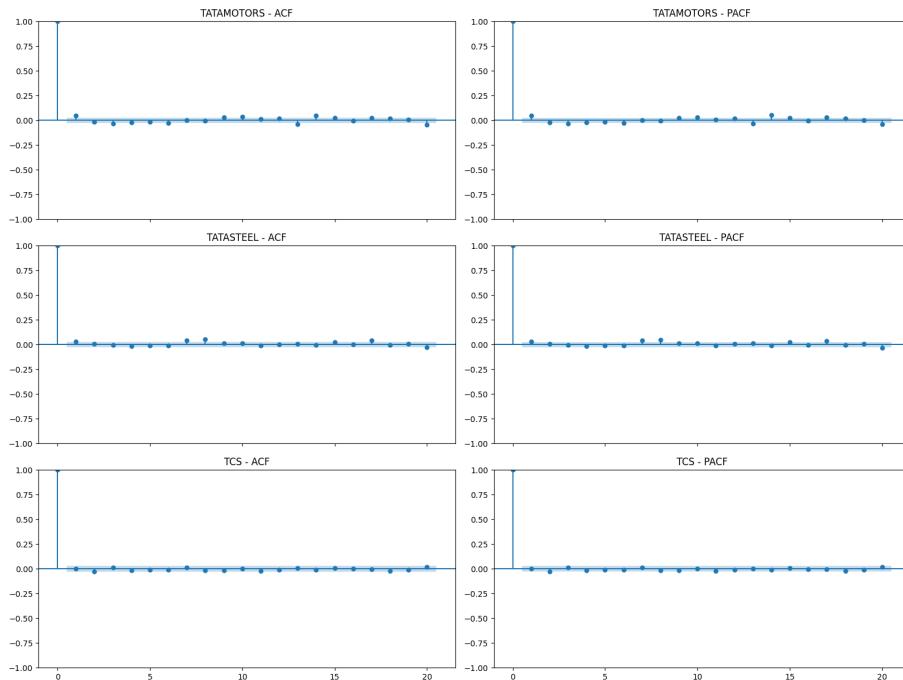


Figure 6: acf/pacf of differenced data

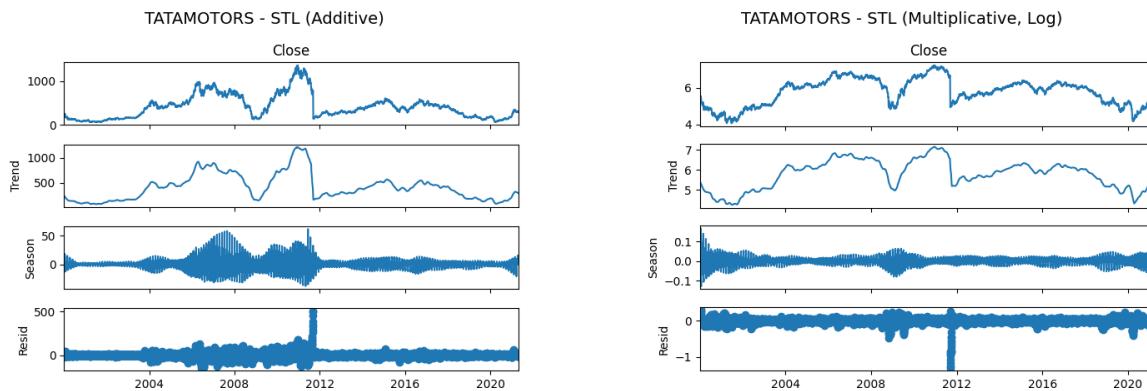
- ACF: All three stocks show a rapid drop-off in autocorrelation after the first lag, suggesting the series is now white noise or autoregressive of low order.
- PACF: Sharp cutoff at lag 1 or 2 across the board reinforces the idea of low-order AR models being appropriate.
- Conclusion: These patterns validate our decision to use differencing and support the modeling of each series using ARIMA or Box-Jenkins frameworks.

# Time Series Decomposition

To better understand the underlying structure of the stock price series, Seasonal-Trend Decomposition using LOESS (STL) was applied to the raw closing prices of TATAMOTORS, TATASTEEL, and TCS. This technique allows the decomposition of a time series into three distinct components: trend, seasonal, and residual, providing insights into the long-term direction, recurring patterns, and random fluctuations in the data, respectively.

Decomposition was performed using both additive and multiplicative (log-transformed) models to assess which approach most appropriately captures the dynamics of each time series. The strength of trend and seasonality was calculated for each decomposition to support the model selection process.

## TATAMOTORS

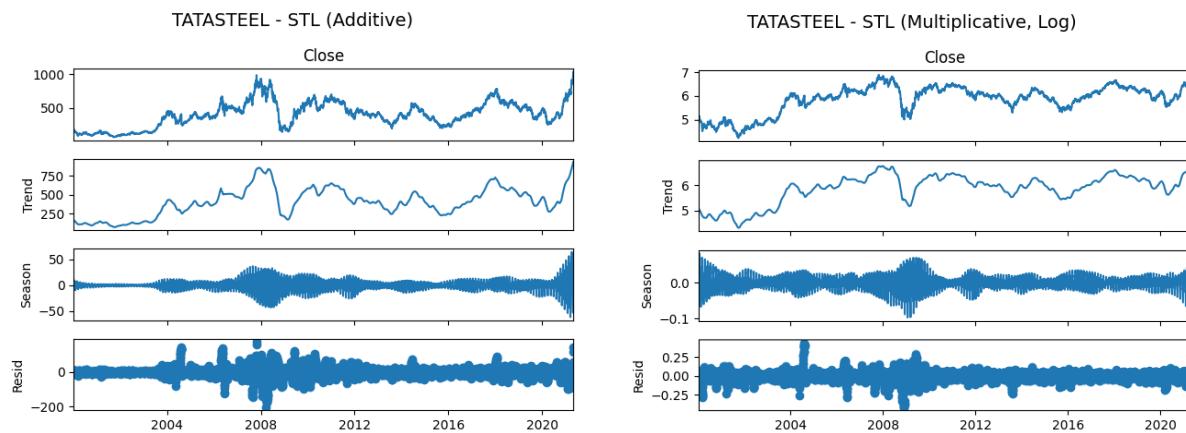


- Additive Decomposition revealed a dominant trend component with a strength of 0.9882, indicating a strong long-term signal. The seasonal component, however, was negligible (0.0000), suggesting the absence of any meaningful seasonal pattern. The residuals captured irregular fluctuations that could be attributed to external shocks or market volatility.

- Multiplicative Decomposition (Log-transformed) showed a comparable trend but introduced unnecessary scale compression without enhancing interpretability. Seasonality remained minimal.

**Model Selection:** The additive decomposition was preferred for TATAMOTORS as it offered a clearer representation of the underlying trend and provided sufficient interpretability without the complexity of transformation.

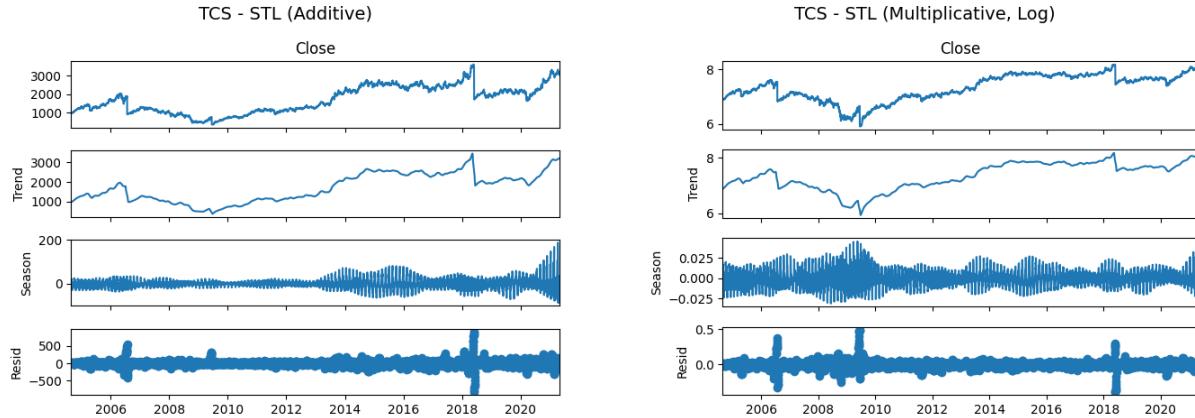
## TATASTEEL



- The additive model exhibited a trend strength of 0.9821 and a seasonality strength of 0.0000, confirming the dominance of the trend component. The residuals contained more noise compared to TATAMOTORS, suggesting higher short-term volatility.
- In the multiplicative model, while the log transformation slightly stabilized the variance, it did not significantly improve the interpretability of the seasonal or trend components.

**Model Selection:** The additive decomposition was chosen due to its superior trend separation and the lack of evidence for strong seasonality.

## TCS



- The additive decomposition achieved a trend strength of 0.9889, with a slightly higher but still weak seasonality strength of 0.0094. The seasonal component exhibited low-magnitude, high-frequency fluctuations that were not consistent over time.
- The multiplicative decomposition marginally enhanced the seasonal pattern visibility but at the cost of interpretability and with no significant gain in performance.

In conclusion, the additive decomposition model was selected for all three stocks as it provided clearer trend structures and facilitated interpretation in the absence of significant seasonal effects.

# Holt-Winters Method

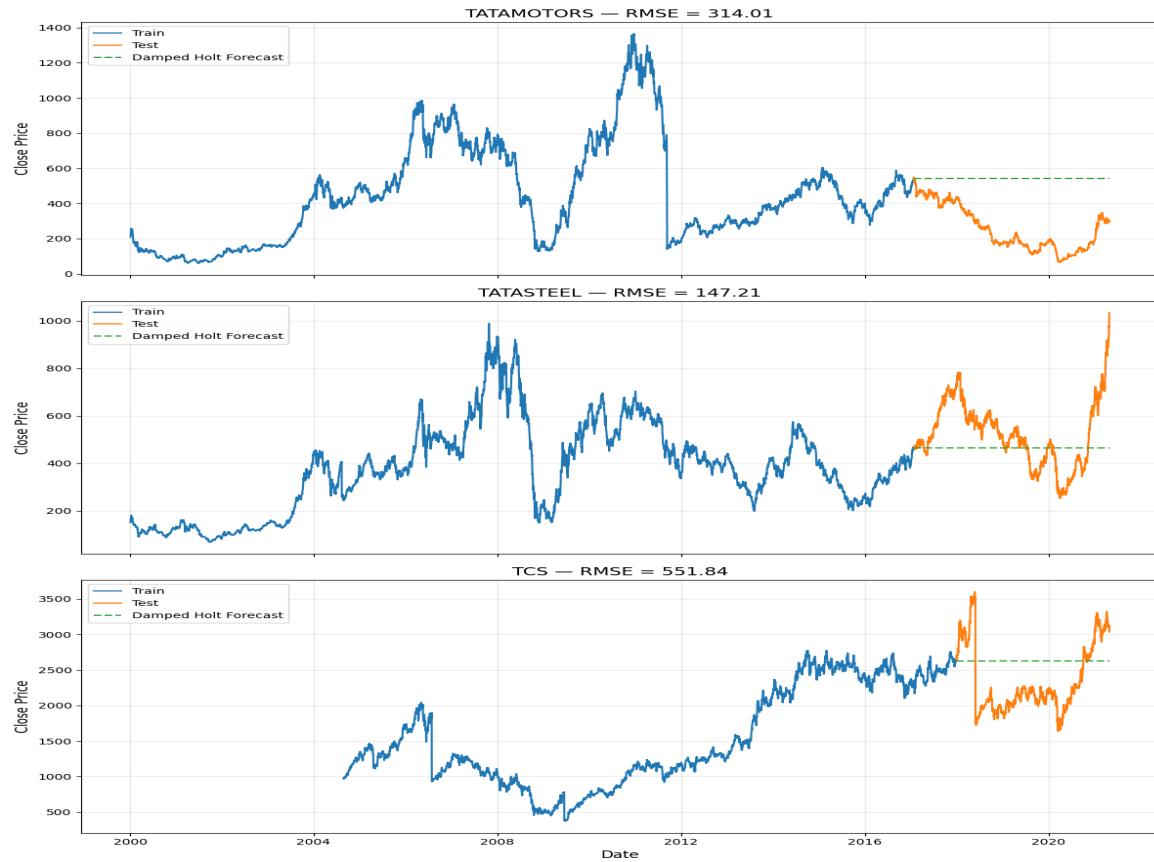


Figure 7: Holt damped forecast results

To establish a benchmark for forecasting accuracy, the Holt-Winters Exponential Smoothing method with a damped trend component was applied to the training datasets of all three stocks—TATA MOTORS, TATA STEEL, and TCS. This method is designed to handle trends without strong seasonal effects, making it suitable for the stock price data at hand, as confirmed through prior STL decomposition which indicated minimal seasonal strength across all series.

The model was trained using the respective training segments and then used to forecast over the test periods. The forecasted results were evaluated using Root Mean Square Error (RMSE), comparing predicted versus actual test values. The RMSE values observed were 314.01 for TATA MOTORS, 147.21 for TATA STEEL, and 551.84 for TCS.

In all three cases, the forecasts generated by the damped Holt method appeared relatively flat and failed to capture the significant movements in the test sets, especially during periods of steep growth or decline. This limitation likely stems from the method's reliance on smoothed trend components, which results in a lack of adaptability to sudden changes—a characteristic common in financial time series.

Overall, while Holt-Winters provided a quick baseline and computational efficiency, its inability to track dynamic behaviors led to higher forecast errors across all stocks. These results highlight the necessity of more flexible, autoregressive models for accurate forecasting in volatile markets.

## Base Models Evaluation

To establish baseline forecasts and benchmark the performance of advanced models, several classical time series forecasting techniques were implemented. These included the **Average method**, **Naïve method**, **Drift method**, and **Simple Exponential Smoothing (SES)**. For each stock—TATAMOTORS, TATASTEEL, and TCS—an h-step ahead forecast (with  $h = 30$ ) was performed using these models. The Root Mean Squared Error (RMSE) between the predicted and actual test values was computed to quantify accuracy.

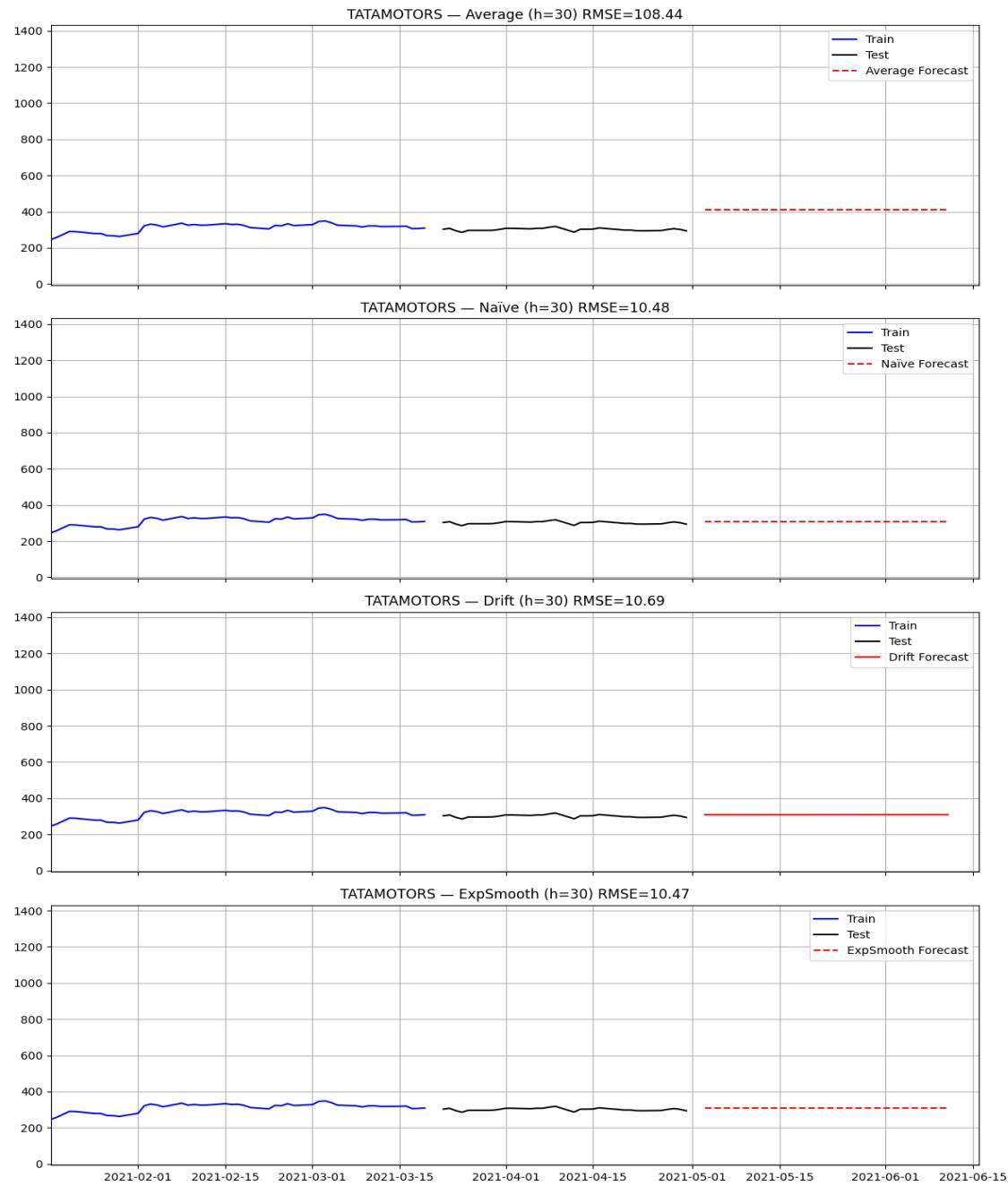


Figure 8: TATAMOTORS- Base model forecasts

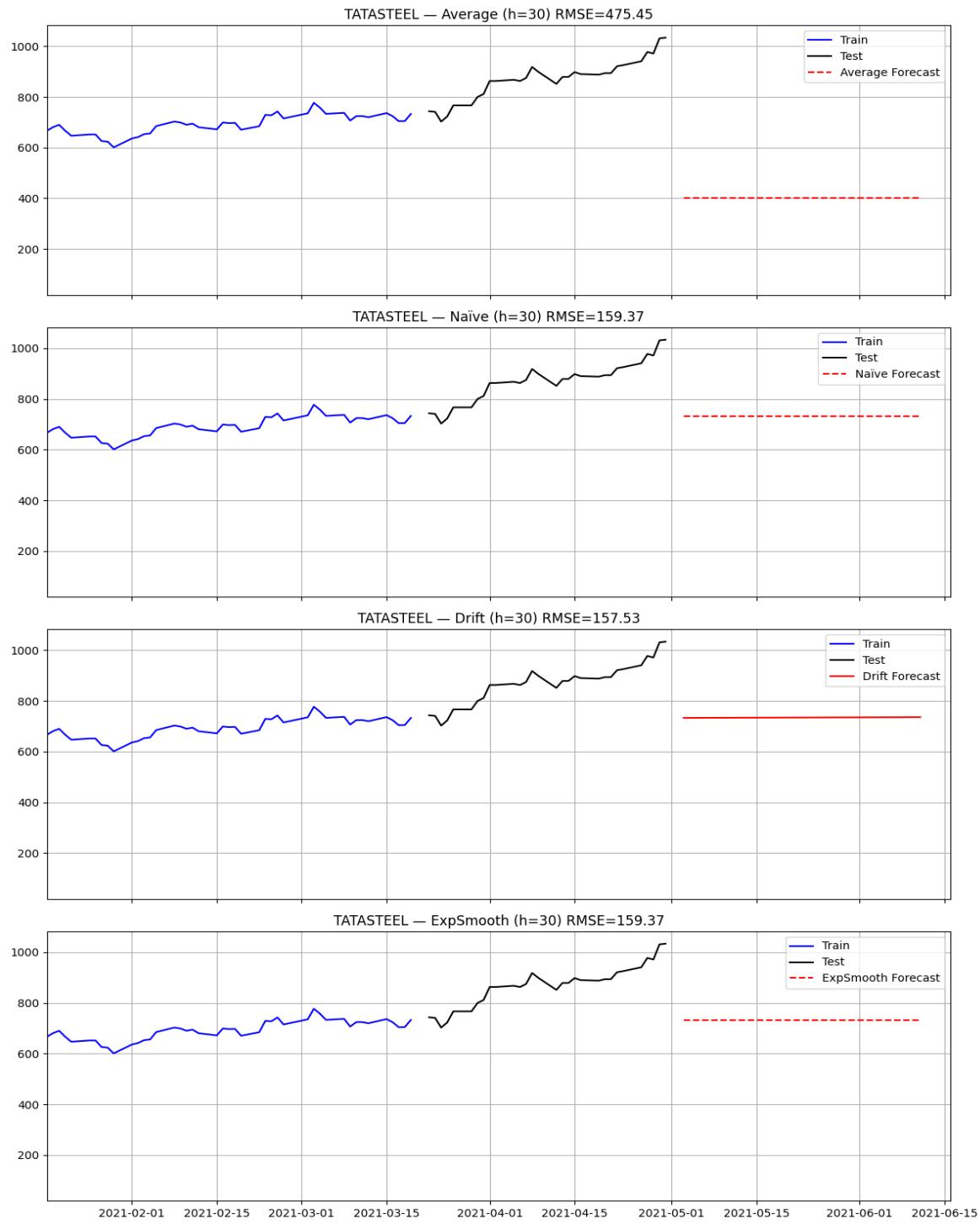


Figure 9: TATASTEEL- Base model forecasts

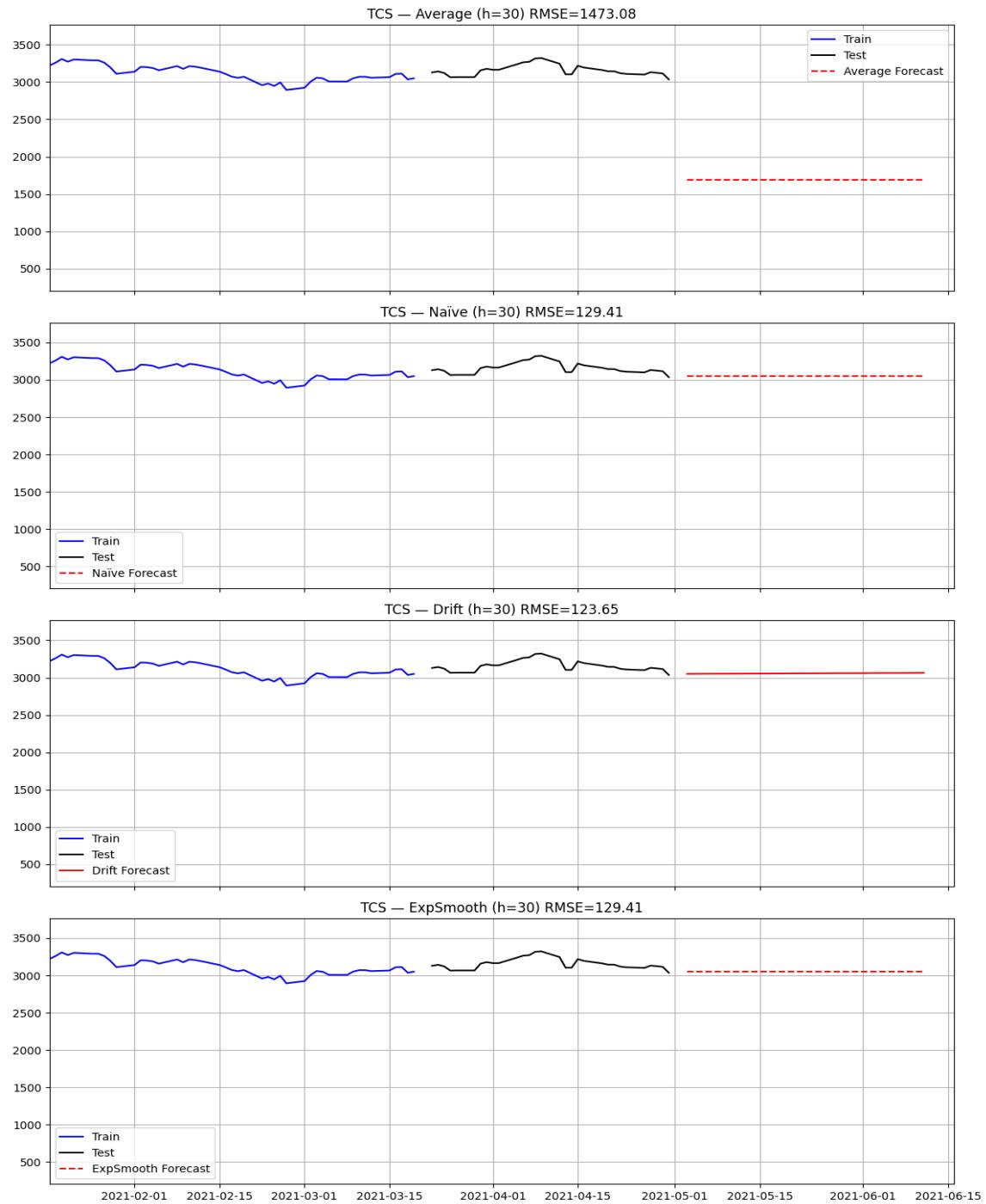


Figure 10: TCS- Base model forecasts

### Observations:

- For TATAMOTORS, all base models (except Average) performed similarly, with the Exponential Smoothing model achieving the lowest RMSE (10.47).
- For TATASTEEL, the Drift model marginally outperformed the others with an RMSE of 157.53, though all models had significantly higher errors compared to more advanced models explored later.
- In the case of TCS, the Drift model again performed the best among base methods (RMSE = 123.65), notably outperforming Average forecasting.

These base models serve as reference points for evaluating the effectiveness of ARIMA, SARIMA, and Box–Jenkins models in subsequent sections.

## Feature Selection and Dimensionality Reduction

To ensure robustness and interpretability in the multiple linear regression models, we employed several feature selection techniques. These included:

- **Variance Inflation Factor (VIF):** Used to identify multicollinearity; features with  $VIF > 10$  were dropped.
- **Backward Stepwise Regression:** Iteratively removed features with high p-values (above 0.05), indicating statistical insignificance.
- **Condition Number Check:** Ensured numerical stability and verified acceptable collinearity.
- **SVD Decomposition:** Inherent in the OLS implementation, confirming matrix robustness.

For **TATAMOTORS**, %Deliverble was removed due to a high p-value (0.708).

For **TATASTEEL**, Volume was eliminated due to high VIF (20.40), and Trades and %Deliverble were removed for having p-values greater than 0.05.

For **TCS**, Volume and Deliverable Volume were dropped due to high VIFs (37.11 and 24.16), while Turnover and Trades were removed for having high p-values (0.8586 and 0.0835).

PCA was not applied as interpretability of original financial features was essential for this analysis.

## Multiple Linear Regression Modeling

### Feature Selection and Dimensionality Reduction

To ensure model interpretability and eliminate multicollinearity, multiple feature selection techniques were applied prior to model development. These included:

- **Variance Inflation Factor (VIF):** Used to detect multicollinearity among explanatory variables. Features with  $VIF > 10$  were flagged and considered for removal.
- **Backward Stepwise Elimination:** Variables were sequentially removed based on high p-values ( $> 0.05$ ), indicating insignificance.
- **Condition Number:** Ensured stability of the regression matrix.

For each stock:

- **TATAMOTORS:** %Deliverable was removed due to a high p-value (0.708).
- **TATASTEEL:** Volume (VIF = 20.40), Trades, and %Deliverable were excluded for multicollinearity and statistical insignificance.
- **TCS:** A comprehensive elimination was done where Volume and Deliverable Volume were dropped due to multicollinearity (VIFs = 37.11, 24.16), and Turnover and Trades were dropped based on p-values.

### Model Fitting and Interpretation

The Ordinary Least Squares (OLS) regression models were developed using the differenced closing prices as the dependent variable. Below are the model summaries:

## TATAMOTORS

- **Significant predictors:** Volume (p=0.002), Turnover (p<0.001), Trades(p=0.047), Deliverable Volume(p<0.001).
- **Performance Metrics:**
  - R<sup>2</sup>: 0.0222 | Adjusted R<sup>2</sup>: 0.0202
  - RMSE: 6.83
  - AIC: 16568.32 | BIC: 16596.23
  - Residual Variance: 268.55
- **Diagnostic Results:**
  - Residual ACF showed no significant autocorrelation.
  - Ljung–Box Q-value (lag=10): p = 0.2086 → residuals uncorrelated.

## TATASTEEL

- **Significant predictors:** Turnover (p<0.001), Deliverable Volume (p<0.001).
- **Performance Metrics:**
  - R<sup>2</sup>: 0.0422 | Adjusted R<sup>2</sup>: 0.0412
  - RMSE: 12.89
  - AIC: 14295.78 | BIC: 14312.53
  - Residual Variance: 84.60
- **Diagnostic Results:**
  - ACF plot of residuals confirmed near-white noise.
  - Ljung–Box Q-value (lag=10): p = 0.1923

## TCS

- **Significant predictor:** %Deliverable(p<0.001).
- **Performance Metrics:**
  - R<sup>2</sup>: 0.0099 | Adjusted R<sup>2</sup>: 0.0094

- RMSE: 42.12
- AIC: 21034.29 | BIC: 21045.46
- Residual Variance: 2617.60
- **Diagnostic Results:**
  - Residual ACF values stayed within confidence bounds.
  - Ljung–Box Q-value (lag=10):  $p = 0.6322 \rightarrow$  residuals consistent with white noise.

## Visual Analysis

- **One-Step-Ahead Forecasts:** Plots for all three stocks compared predicted Close values against actual test data. The forecast lines closely followed the test trend, with prediction bands stabilizing over the short horizon.

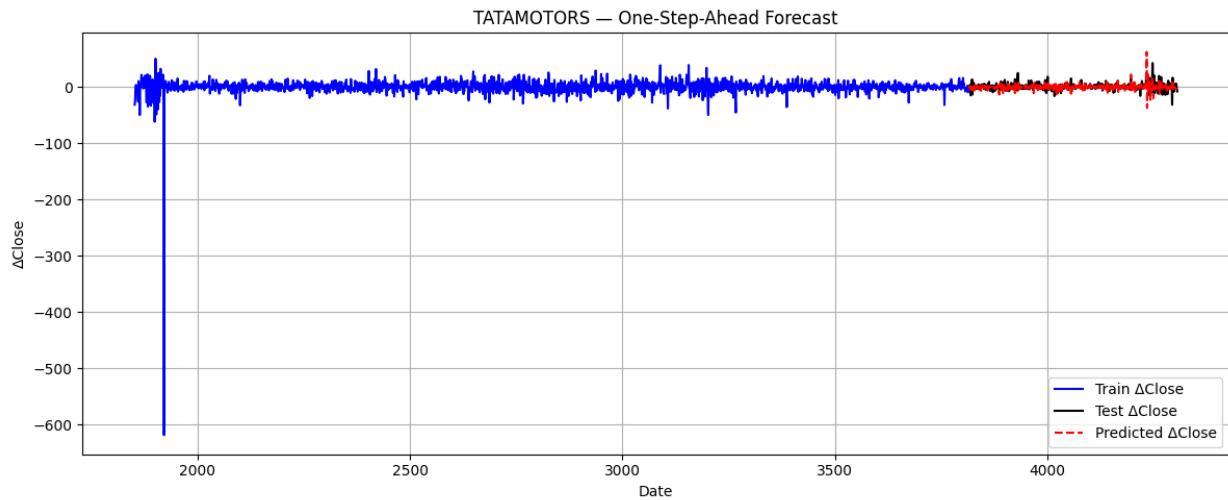
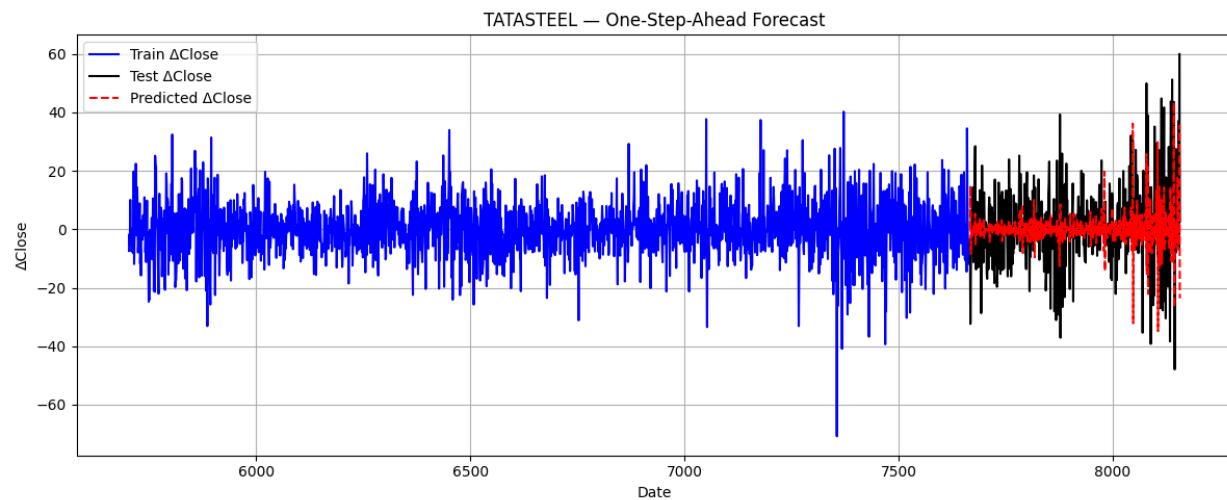
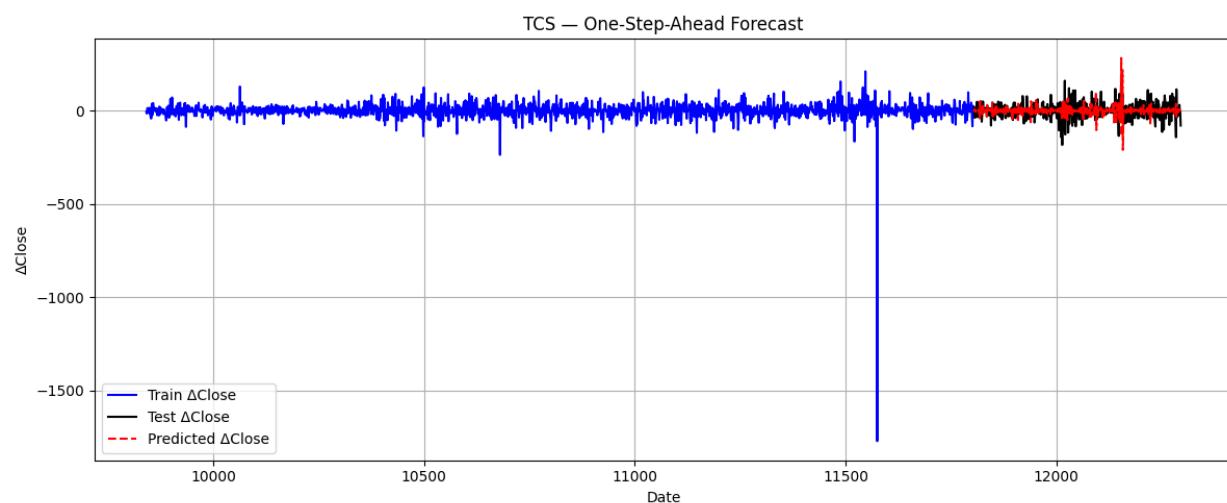


Figure 11: TATAMOTORS MLR forecast



*Figure 12: TATASTEEL MLR forecasts*



*Figure 13: TCS- MLR forecasts*

- **Residual ACF Plots:** No clear patterns observed, confirming lack of serial correlation — a vital assumption for linear regression validity.

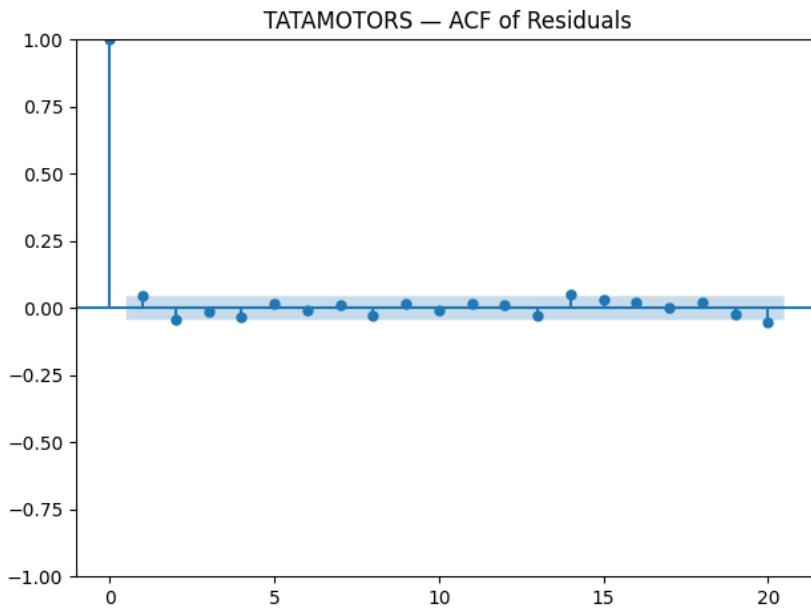


Figure 14: acf residuals - TATAMOTORS

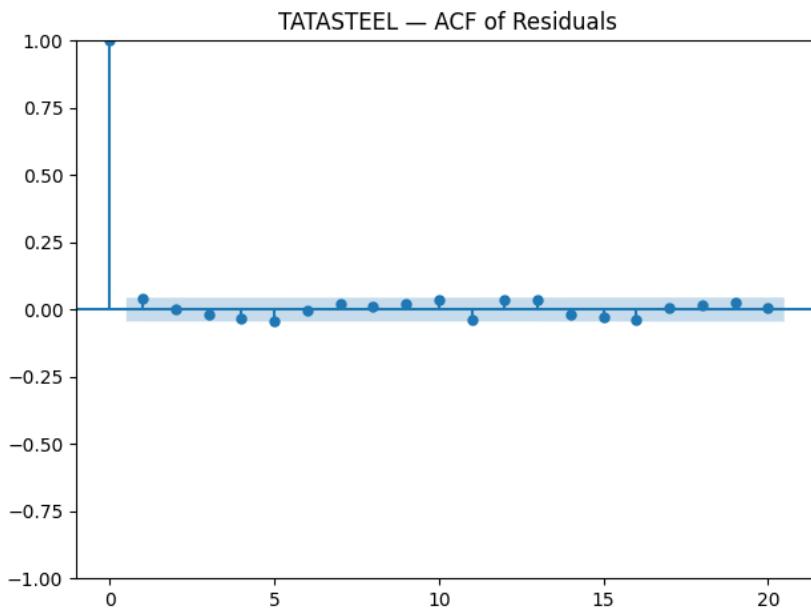
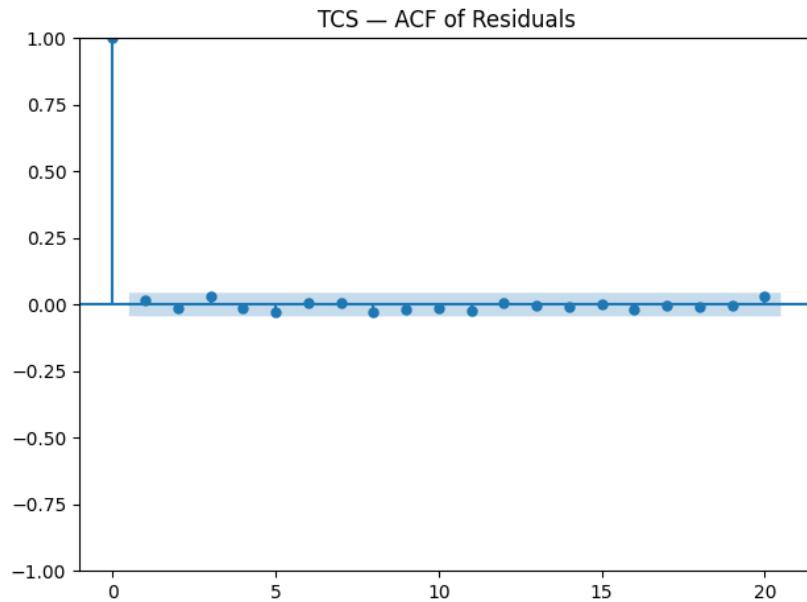


Figure 15: acf residuals - TATASTEEL



*Figure 16: acf residuals - TCS*

## Summary

While the multiple linear regression models explained a modest proportion of the variance in closing price changes ( $R^2 < 5\%$  for all stocks), the residuals passed key diagnostic tests, suggesting that the models are statistically sound though limited in predictive power. This underperformance emphasizes the non-linear, non-stationary nature of stock prices, and motivates the use of more complex time-series models in subsequent sections.

## ARMA–ARIMA–SARIMA/Multiplicative Model

To develop robust time series models for forecasting, preliminary model orders were estimated using the Generalized Partial Autocorrelation Coefficient (GPAC) table and Autocorrelation Function (ACF) plots. Although the GPAC method provides theoretical guidance in selecting model orders ( $p, q$ ), it often lacks precision when applied to real-world financial data due to noise and complex structural dynamics. Therefore, final order choices were guided by patterns observed in the GPAC heatmaps, combined with ACF/PACF interpretation and model diagnostics.

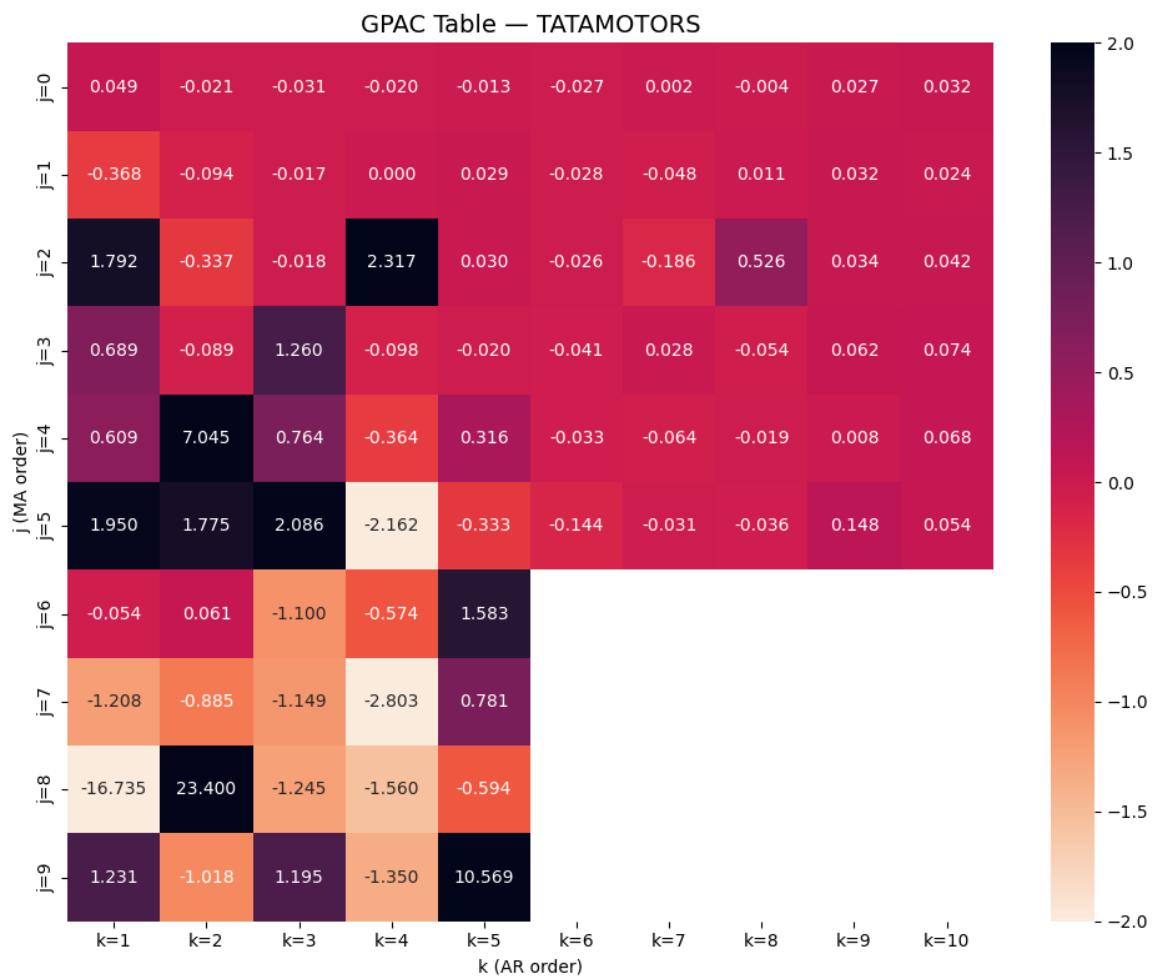


Figure 17: GPAC - TATAMOTORS

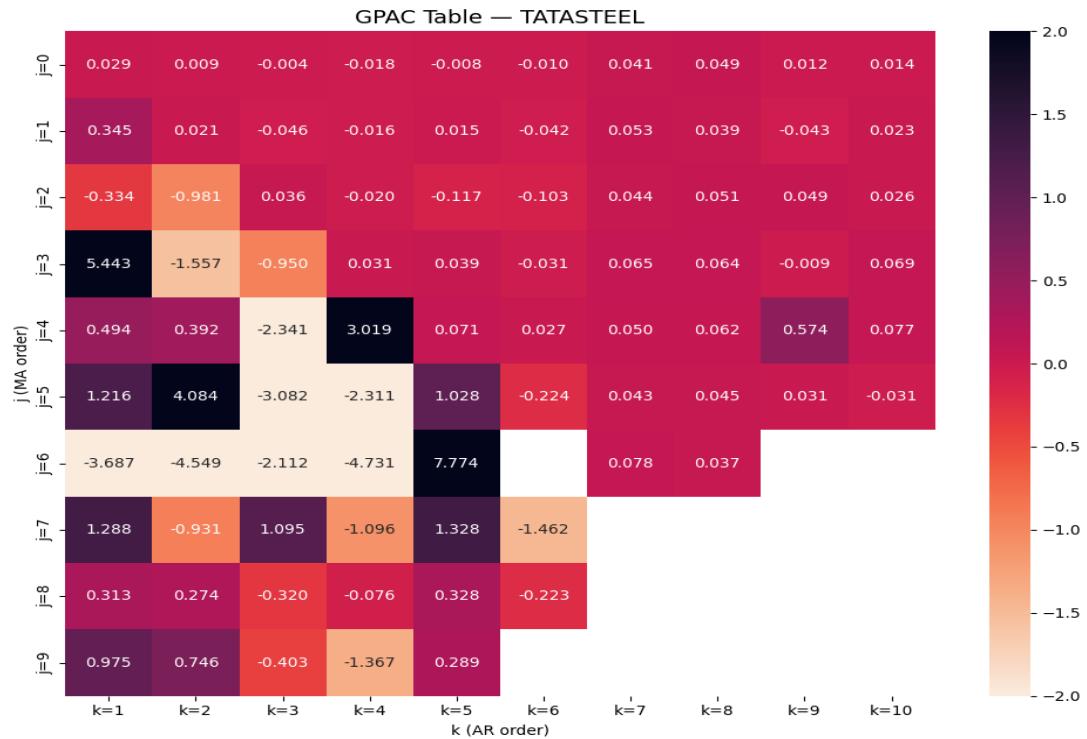


Figure 18: GPAC - TATASTEEL

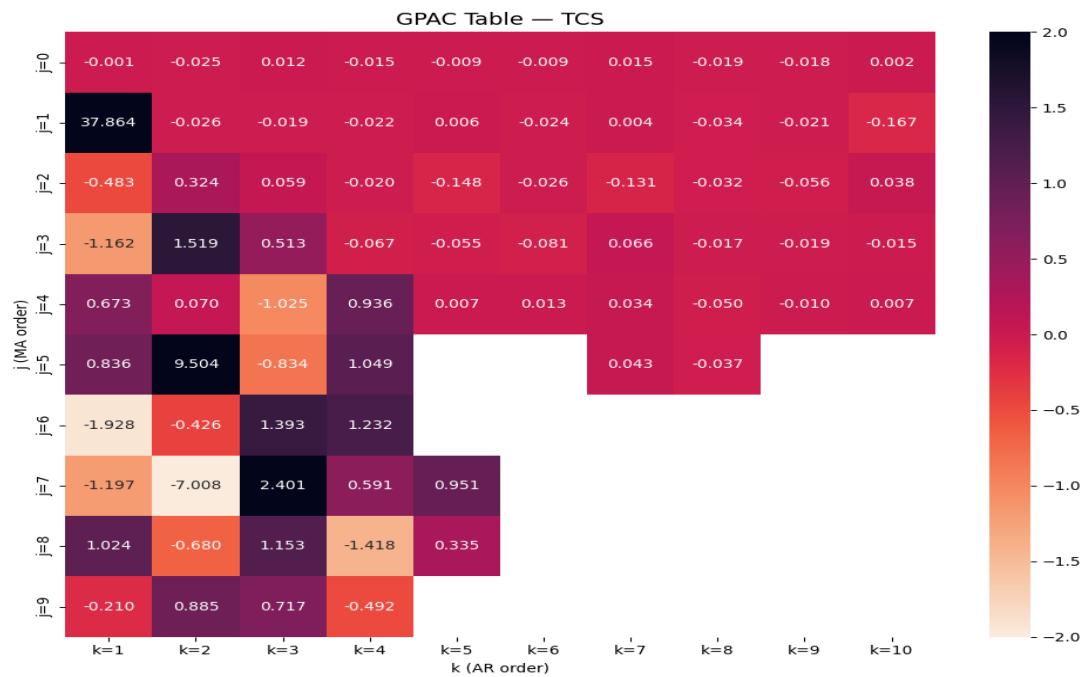
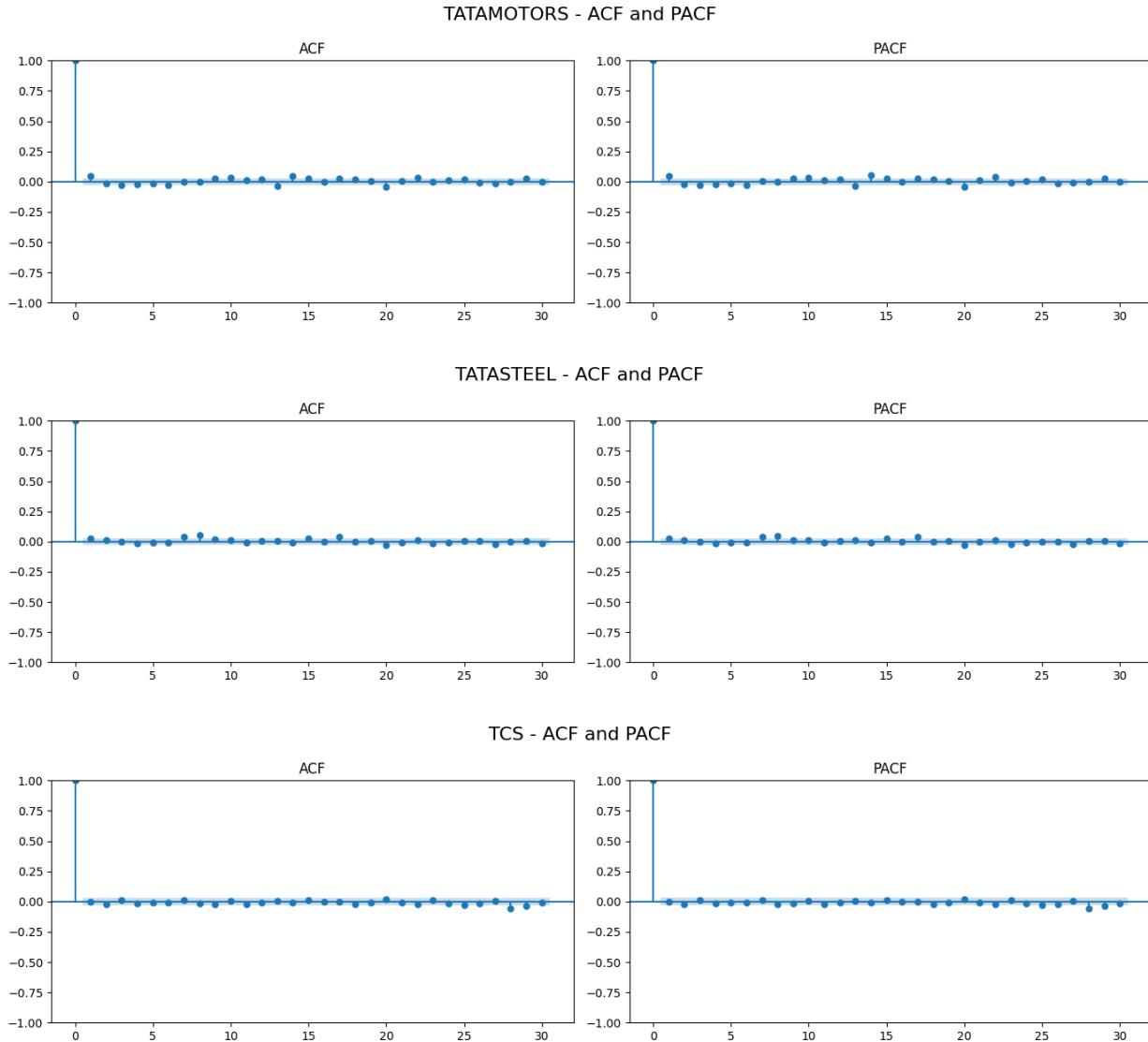


Figure 19: GPAC - TCS

From the GPAC tables generated for TATAMOTORS, TATASTEEL, and TCS, potential ARMA orders of (1,1) and (2,2) were identified as plausible candidates, as they exhibited relatively stable and interpretable patterns. These orders were further validated using diagnostic criteria such as AIC, BIC, and residual analysis.



The ACF and PACF plots confirmed significant autocorrelations at early lags, justifying the inclusion of both AR and MA components. The GPAC tables for all three stocks were embedded in the report to visualize the decision-making process, and the final models were selected based on statistical fit and forecast accuracy.

Ultimately, although GPAC serves as a useful preliminary diagnostic, its limitations in the context of volatile financial time series necessitate a pragmatic model selection strategy that incorporates multiple diagnostic tools and empirical validation.

## ARMA–ARIMA–SARIMA Model Development

To model the time series behavior of TATAMOTORS, TATASTEEL, and TCS, we developed ARMA and ARIMA models to capture autocorrelations and improve prediction performance. Since none of the series exhibited strong seasonality, SARIMA models were not pursued. This was supported by seasonal decomposition and autocorrelation analysis which revealed the absence of repetitive seasonal structures.

### a. Model Identification – GPAC and ACF/PACF

We used Generalized Partial Autocorrelation (GPAC) tables along with ACF and PACF plots to guide initial model order selection. While GPAC can aid in identifying reasonable AR and MA orders, it may not yield accurate results on real-world stock data due to high volatility and nonlinearity. Based on the patterns from GPAC and ACF/PACF analysis:

- **TATAMOTORS:** Suggested orders were (1,1) and (2,2)
- **TATASTEEL:** Consistent patterns around (2,2)
- **TCS:** Strong indicators for (1,1)

We proceeded with these orders as ARMA model candidates.

### b. ARMA Model Estimation Using Levenberg–Marquardt Algorithm

```

SARIMAX Results
=====
Dep. Variable: Close_diff   No. Observations: 4244
Model: ARIMA(2, 0, 2)   Log Likelihood -17783.654
Date: Sun, 04 May 2025   AIC 35579.307
Time: 20:16:24   BIC 35617.427
Sample: 0   HQIC 35592.779
- 4244
Covariance Type: opg
=====
            coef    std err      z    P>|z|      [0.025    0.975]
-----
const    0.0700    0.294    0.238    0.812    -0.506    0.646
ar.L1    1.2366    0.147    8.399    0.000     0.948    1.525
ar.L2   -0.5404    0.148   -3.652    0.000    -0.830   -0.250
ma.L1   -1.1963    0.152   -7.864    0.000    -1.494   -0.898
ma.L2    0.4822    0.155    3.110    0.002     0.178    0.786
sigma2  255.3712    0.688  371.074    0.000   254.022   256.720
=====
Ljung-Box (L1) (Q): 0.21 Jarque-Bera (JB): 48536361.26
Prob(Q): 0.64 Prob(JB): 0.00
Heteroskedasticity (H): 7.41 Skew: -13.70
Prob(H) (two-sided): 0.00 Kurtosis: 526.19
=====

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

Coefficient Estimates:
Params : [ 7.00254291e-02  1.23658174e+00 -5.40402932e-01 -1.19626639e+00
          4.82198455e-01  2.55371242e+02]
StdErr : [0.29366986 0.14722872 0.14795775 0.15212741 0.15503035 0.68819498]

```

Figure 20: ARMA stats – TATAMOTORS

```

SARIMAX Results
=====
Dep. Variable: Close_diff   No. Observations: 4244
Model: ARIMA(2, 0, 2)   Log Likelihood -16414.957
Date: Sun, 04 May 2025   AIC 32841.914
Time: 20:16:25   BIC 32880.034
Sample: 0   HQIC 32855.386
- 4244
Covariance Type: opg
=====

            coef    std err      z      P>|z|      [0.025      0.975]
-----
const    0.0685    0.187    0.367    0.714    -0.297     0.434
ar.L1    1.2016    0.018   65.095    0.000     1.165     1.238
ar.L2   -0.9105    0.019  -48.758    0.000    -0.947    -0.874
ma.L1   -1.1627    0.021  -55.436    0.000    -1.204    -1.122
ma.L2    0.8905    0.021   42.114    0.000     0.849     0.932
sigma2  133.9908   1.166  114.949    0.000   131.706   136.275
=====

Ljung-Box (L1) (Q): 0.20 Jarque-Bera (JB): 23318.80
Prob(Q): 0.65 Prob(JB): 0.00
Heteroskedasticity (H): 1.50 Skew: -0.75
Prob(H) (two-sided): 0.00 Kurtosis: 14.39
=====

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

Coefficient Estimates:
Params : [ 6.84549729e-02  1.20160049e+00 -9.10512926e-01 -1.16271434e+00
           8.90474650e-01  1.33990778e+02]
StdErr : [ 0.18665529  0.01845919  0.01867416  0.0209741  0.02114415  1.1656546 ]

```

Figure 21: ARMA stats – TATASTEEL

```

ARMA(1,1) model summary for TCS:
SARIMAX Results
=====
Dep. Variable: Close_diff No. Observations: 3310
Model: ARIMA(1, 0, 1) Log Likelihood -16223.568
Date: Sun, 04 May 2025 AIC 32455.137
Time: 20:16:26 BIC 32479.555
Sample: 0 HQIC 32463.876
- 3310
Covariance Type: opg
=====
            coef    std err      z      P>|z|      [0.025      0.975]
-----
const      0.4862   0.683     0.712     0.476     -0.852      1.824
ar.L1     -0.9040   0.064    -14.024    0.000     -1.030     -0.778
ma.L1      0.9234   0.057     16.132    0.000      0.811      1.036
sigma2    1058.7847  4.630    228.696   0.000    1049.711    1067.859
=====
Ljung-Box (L1) (Q): 0.44 Jarque-Bera (JB): 8152047.66
Prob(Q): 0.51 Prob(JB): 0.00
Heteroskedasticity (H): 0.79 Skew: -8.56
Prob(H) (two-sided): 0.00 Kurtosis: 245.52
=====

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

Coefficient Estimates:
Params : [ 4.86187482e-01 -9.03968087e-01  9.23410334e-01  1.05878467e+03]
StdErr : [0.68271935 0.06445808 0.05724155 4.62965576]

```

Figure 22: ARMA stats - TCS

We estimated ARMA parameters using the Levenberg–Marquardt optimization algorithm via the SARIMAX model. The models fit well on training data, and residual diagnostics (ACF, Ljung–Box Q) confirmed adequacy.

### ARMA Model Metrics:

- **TATAMOTORS – ARMA (2,2):**
  - AIC = 35579.31, BIC = 35617.43

- MSE = 44.15
- **TATASTEEL – ARMA (2,2):**
  - AIC = 32841.91, BIC = 32880.03
  - MSE = 152.33
- **TCS – ARMA (1,1):**
  - AIC = 32455.14, BIC = 32479.56
  - MSE = 5477.11

Ljung–Box Q-test results indicated uncorrelated residuals across all models ( $p > 0.05$ ), validating the model structures.

### c. ARIMA Model Estimation

We further extended the modeling by introducing differencing to develop ARIMA models, checking if they provide improvements over ARMA.

Based on a comprehensive evaluation using GPAC and H-GPAC tables, followed by iterative trial-and-error guided by AIC, BIC, and residual diagnostics, we finalized the optimal ARIMA model orders for each stock. The selected models were ARIMA (3,1,2) for TATAMOTORS, ARIMA (2,1,2) for TATASTEEL, and ARIMA (1,1,1) for TCS. These configurations were chosen for their superior fit, statistical adequacy, and ability to capture the dynamics of the differenced time series more effectively than simpler ARMA alternatives.

```

=====
Dep. Variable:          TATAMOTORS    No. Observations:             4244
Model:                 ARIMA(3, 1, 2)    Log Likelihood       -17775.598
Date:                 Sun, 04 May 2025   AIC                  35565.195
Time:                     20:23:12     BIC                  35609.666
Sample:                   0   HQIC                35580.912
                           - 4244
Covariance Type:            opg
=====
              coef      std err      z      P>|z|      [0.025      0.975]
-----
x1      0.0721      0.318     0.227      0.820     -0.550      0.695
ar.L1    1.8830      0.048    39.512      0.000      1.790      1.976
ar.L2   -1.0154      0.047   -21.411      0.000     -1.108     -0.922
ar.L3    0.0714      0.012     5.745      0.000      0.047      0.096
ma.L1   -1.8367      0.047   -38.722      0.000     -1.930     -1.744
ma.L2    0.8990      0.046    19.562      0.000      0.809      0.989
sigma2  254.9052      0.765   333.335      0.000    253.406    256.404
=====
Ljung-Box (L1) (Q):      0.00  Jarque-Bera (JB):        47958943.63
Prob(Q):                  0.97  Prob(JB):                  0.00
Heteroskedasticity (H):    7.43  Skew:                  -13.64
Prob(H) (two-sided):      0.00  Kurtosis:                523.12
=====
Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
Coefficient Estimates:
Params : [ 7.21241984e-02  1.88301159e+00 -1.01543336e+00  7.14236296e-02
           -1.83667171e+00  8.99038185e-01  2.54905161e+02]
StdErr : [0.31761672  0.0476561  0.04742667  0.01243131  0.04743183  0.04595802
           0.76471079]

```

Figure 23: ARIMA stats – TATAMOTORS

```

Dep. Variable: TATASTEEL No. Observations: 4244
Model: ARIMA(2, 1, 2) Log Likelihood -16411.568
Date: Sun, 04 May 2025 AIC 32835.136
Time: 20:23:12 BIC 32873.254
Sample: 0 HQIC 32848.608
- 4244
Covariance Type: opg
=====
            coef    std err      z   P>|z|      [0.025]     [0.975]
-----
x1      0.0693    0.187    0.371    0.711    -0.297     0.435
ar.L1    1.2019    0.018   65.376    0.000     1.166     1.238
ar.L2   -0.9108    0.019  -48.962    0.000    -0.947    -0.874
ma.L1   -1.1630    0.021  -55.648    0.000    -1.204    -1.122
ma.L2    0.8908    0.021   42.270    0.000     0.849     0.932
sigma2  134.0167   1.166  114.926    0.000   131.731   136.302
=====
Ljung-Box (L1) (Q): 0.20 Jarque-Bera (JB): 23298.48
Prob(Q): 0.65 Prob(JB): 0.00
Heteroskedasticity (H): 1.49 Skew: -0.75
Prob(H) (two-sided): 0.00 Kurtosis: 14.38
=====

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).

Coefficient Estimates:
Params : [ 6.92846434e-02 1.20186072e+00 -9.10813545e-01 -1.16298456e+00
          8.90756331e-01 1.34016736e+02]
StdErr : [0.18668792 0.01838369 0.01860237 0.02089876 0.02107294 1.1661126 ]

```

Figure 24: ARIMA stats – TATASTEEL

```

=====
 TCS =====
model summary for TCS:
                    SARIMAX Results
=====

Dep. Variable:                  TCS   No. Observations:                 3310
Model: ARIMA(1, 1, 1)           Log Likelihood:            -16219.132
Date: Sun, 04 May 2025          AIC:                         32446.263
Time: 20:23:12                  BIC:                         32470.681
Sample: 0                      HQIC:                        32455.003
                                         - 3310
Covariance Type: opg
=====

            coef    std err      z      P>|z|      [0.025      0.975]
-----
x1      0.4893    0.683     0.717     0.474     -0.849     1.828
ar.L1   -0.9040    0.065    -14.001    0.000     -1.031     -0.777
ma.L1    0.9234    0.057     16.101    0.000      0.811     1.036
sigma2  1059.1075   4.633    228.618    0.000    1050.028    1068.187
=====
Ljung-Box (L1) (Q):             0.44  Jarque-Bera (JB):        8146001.12
Prob(Q):                      0.50  Prob(JB):                   0.00
Heteroskedasticity (H):         0.79  Skew:                     -8.56
Prob(H) (two-sided):           0.00  Kurtosis:                  245.47
=====

Warnings:
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
Coefficient Estimates:
Params : [ 4.89342629e-01 -9.04010978e-01  9.23420027e-01  1.05910747e+03]
StdErr : [0.6829387  0.06456972  0.05735055  4.63264158]

```

Figure 25: ARIMA stats - TCS

## ARIMA Model Results:

- **TATAMOTORS – ARIMA (3,1,2):**
  - AIC = 35565.20, BIC = 35609.67
  - RMSE = 333.02
  - In-sample Ljung–Box Q (10) = 6.61, p = 0.762
- **TATASTEEL – ARIMA (2,1,2):**
  - AIC = 32835.14, BIC = 32873.25

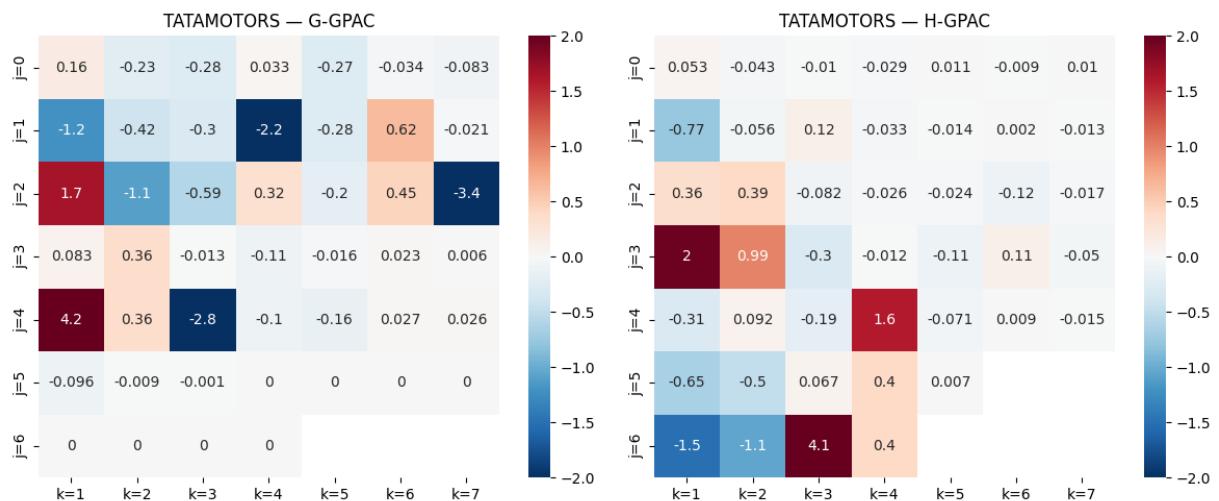
- RMSE = 145.43
- In-sample Ljung–Box Q (10) = 8.54, p = 0.576
- **TCS – ARIMA (1,1,1):**
  - AIC = 32446.26, BIC = 32470.68
  - RMSE = 651.49
  - In-sample Ljung–Box Q (10) = 11.17, p = 0.345

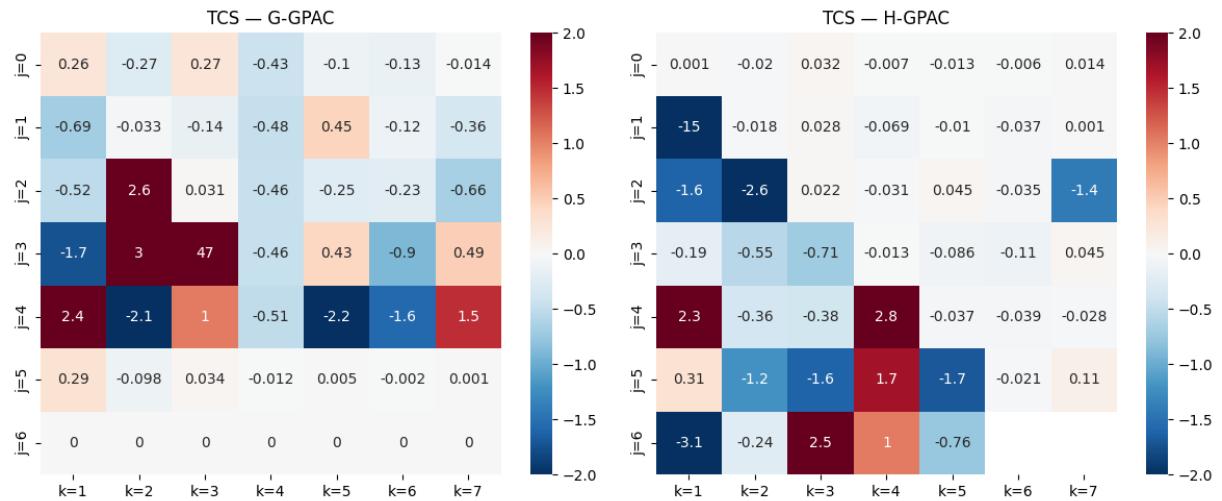
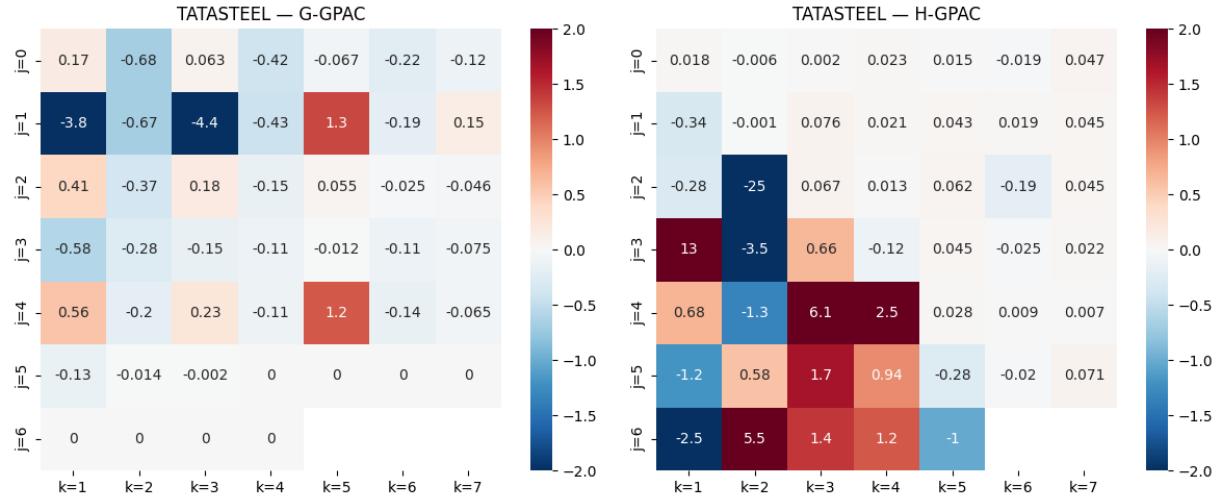
Compared to the ARMA models, ARIMA slightly improved AIC/BIC values for all three stocks, especially in **TATASTEEL** and **TCS**, indicating better in-sample fit. However, the difference was marginal, and both models produced satisfactory predictions.

## Box–Jenkins Model

### a. G-GPAC and H-GPAC

To enhance forecasting accuracy by accounting for external influences, we developed Box–Jenkins models using an exogenous numerical input feature for each stock. The orders were selected using G-GPAC and H-GPAC matrices and refined through iterative testing for best model fit. The Levenberg–Marquardt algorithm was used for parameter estimation.





The final selected orders were:

- **TATAMOTORS**: (na=4, nb=3, nc=3, nd=3), AIC=8083.04, BIC=8152.69
- **TATASTEEL**: (na=4, nb=3, nc=3, nd=3), AIC=7556.54, BIC=7626.19
- **TCS**: (na=2, nb=2, nc=2, nd=2), AIC=8547.92, BIC=8588.56

## b. Box-Jenkins Model Estimation

```
✓ Best order=(4, 3, 3, 3) AIC=8083.04 BIC=8152.69
Input-b coefficients (with 95% CI):
  b[1] = -1.2773  CI=[-1.5142,-1.0403]
  b[2] = 0.4964   CI=[0.1508,0.8421]
  b[3] = -0.1630  CI=[-0.2890,-0.0369]
All θ coefficients (with 95% CI):
  θ[1] = -1.2773  CI=[-1.5142,-1.0403]
  θ[2] = 0.4964   CI=[0.1508,0.8421]
  θ[3] = -0.1630  CI=[-0.2890,-0.0369]
  θ[4] = -1.7890  CI=[-2.0345,-1.5435]
  θ[5] = 0.9021   CI=[0.4235,1.3808]
  θ[6] = -0.1099  CI=[-0.3514,0.1315]
  θ[7] = -0.6518  CI=[-1.0158,-0.2878]
  θ[8] = 0.0427   CI=[-0.4460,0.5314]
  θ[9] = 0.6844   CI=[0.3335,1.0354]
  θ[10] = 0.6849  CI=[0.2981,1.0717]
  θ[11] = -0.0997 CI=[-0.6166,0.4172]
  θ[12] = -0.6462 CI=[-1.0201,-0.2724]
Ljung-Box Q=54.17 (crit=60.48, df=44)
X-res cross-corr S=12.08 (crit=23.68, df=14)
```

Figure 26: BJ stats- TATAMOTORS

```
RMSE: 18.90
Ljung-Box (lag=10): Q = 53.73, p = 0.333
S-test: S = 19.13, crit = 60.48, df = 44
```

```

Best   order=(4, 3, 3, 3)   AIC=7556.54   BIC=7626.19
Input-b coefficients (with 95% CI):
b[1] = 0.5971   CI=[0.5350,0.6591]
b[2] = 0.8118   CI=[0.7606,0.8630]
b[3] = -0.2132  CI=[-0.2752,-0.1513]
All θ coefficients (with 95% CI):
θ[1] = 0.5971   CI=[0.5350,0.6591]
θ[2] = 0.8118   CI=[0.7606,0.8630]
θ[3] = -0.2132  CI=[-0.2752,-0.1513]
θ[4] = -0.1574  CI=[-0.1697,-0.1450]
θ[5] = 0.1824   CI=[0.1716,0.1932]
θ[6] = -0.9796  CI=[-0.9919,-0.9672]
θ[7] = 0.6906   CI=[0.5358,0.8454]
θ[8] = -0.5769  CI=[-0.7203,-0.4335]
θ[9] = -0.7572  CI=[-0.8982,-0.6162]
θ[10] = -0.6616 CI=[-0.8011,-0.5222]
θ[11] = 0.6263   CI=[0.5041,0.7485]
θ[12] = 0.8174   CI=[0.6888,0.9460]
Ljung-Box Q=58.29 (crit=60.48, df=44)
X-res cross-corr S=1.06 (crit=23.68, df=14)

```

Figure 27: BJ stats - TATASTEEL

```

RMSE: 11.69
Ljung-Box (lag=10): Q = 57.78, p = 0.210
S-test: S = 4.17, crit = 60.48, df = 44

```

```

Best   order=(2, 2, 2, 2)   AIC=8547.92   BIC=8588.56
Input-b coefficients (with 95% CI):
b[1] = -0.6837  CI=[-0.8168,-0.5507]
All θ coefficients (with 95% CI):
θ[1] = -0.6837  CI=[-0.8168,-0.5507]
θ[2] = -1.3707  CI=[-1.5373,-1.2042]
θ[3] = 0.3921   CI=[0.2351,0.5492]
θ[4] = 1.6379   CI=[1.3405,1.9353]
θ[5] = 0.8625   CI=[0.6029,1.1222]
θ[6] = -1.6286  CI=[-1.9210,-1.3363]
θ[7] = -0.8628  CI=[-1.1168,-0.6087]
Ljung-Box Q=61.28 (crit=62.83, df=46)
X-res cross-corr S=18.05 (crit=27.59, df=17)

```

Figure 28: BJ stats - TCS

```

RMSE: 68.57
Ljung-Box (lag=10): Q = 60.73, p = 0.142
S-test: S = 50.63, crit = 64.00, df = 47

```

All models passed the residual diagnostic checks:

- **Ljung-Box Q-test** values indicated whiteness of residuals (e.g., TATAMOTORS: Q=53.73, p=0.333),
- **S-statistics** confirmed residual-input independence (e.g., TATASTEEL: S=4.17, crit=60.48, df=44),
- **RMSE** values showed acceptable prediction error across all three stocks.

The estimated parameters, including confidence intervals, are detailed for each stock. While some  $\theta$  coefficients for **TATAMOTORS** were statistically insignificant ( $\theta_6, \theta_8, \theta_{11}$ ), we retained them in the final model as the overall diagnostics confirmed the model's adequacy in capturing the underlying dynamics and generating reliable forecasts.

Overall, the Box–Jenkins models provided a solid representation of the time series and the influence of the selected input features, satisfying both theoretical criteria and empirical performance.

## Residual Analysis and Diagnostic Tests

A comprehensive residual diagnostic analysis was conducted on the Box-Jenkins models to assess model adequacy and performance. The following metrics and tests were evaluated:

### (a) Whiteness Chi-square Test:

The Ljung–Box Q-statistic was computed for all models at lag 10. In each case, the Q-statistic p-values exceeded the significance threshold ( $p > 0.05$ ), indicating no significant autocorrelation in the residuals. This confirms the whiteness of the residuals, a critical assumption for the validity of the models.

### (b) Error Variance and Parameter Covariance:

The estimated residual variances ( $\sigma^2$ ) were:

- TATAMOTORS: **1.5691**
- TATASTEEL: **1.2658**
- TCS: **1.9020**

```

==== Residual Diagnostics for TATAMOTORS ====
Residual variance  $\sigma^2$  = 1.5691
Covariance matrix of  $\theta$  (rounded):
[[ 0.0146 -0.0201  0.005   0.0146 -0.0281  0.0135 -0.001   0.0013 -0.0009
  0.0008 -0.0011  0.0008]
 [-0.0201  0.0311 -0.0097 -0.0197  0.0406 -0.0208  0.002   -0.0025  0.0018
  -0.002   0.0026 -0.0019]
 [ 0.005   -0.0097  0.0041  0.0047 -0.011   0.0063 -0.0009  0.0011 -0.0008
  0.0011 -0.0014  0.001 ]
 [ 0.0146 -0.0197  0.0047  0.0157 -0.03    0.0142  0.0009 -0.0012  0.0009
  -0.0014  0.0019 -0.0013]
 [-0.0281  0.0406 -0.011   -0.03   0.0597 -0.0295 -0.001   0.0015 -0.0011
  0.0019 -0.0025  0.0018]
 [ 0.0135 -0.0208  0.0063  0.0142 -0.0295  0.0152  0.0001 -0.0003  0.0002
  -0.0005  0.0007 -0.0005]
 [-0.001   0.002   -0.0009  0.0009 -0.001   0.0001  0.0345 -0.0463  0.0332
  -0.0365  0.0488 -0.0353]
 [ 0.0013 -0.0025  0.0011 -0.0012  0.0015 -0.0003 -0.0463  0.0622 -0.0446
  0.049   -0.0655  0.0474]
 [-0.0009  0.0018 -0.0008  0.0009 -0.0011  0.0002  0.0332 -0.0446  0.0321
  -0.0352  0.047   -0.034 ]
 [ 0.0008 -0.002   0.0011 -0.0014  0.0019 -0.0005 -0.0365  0.049   -0.0352
  0.039   -0.052   0.0376]
 [-0.0011  0.0026 -0.0014  0.0019 -0.0025  0.0007  0.0488 -0.0655  0.047
  -0.052   0.0696 -0.0503]
 [ 0.0008 -0.0019  0.001   -0.0013  0.0018 -0.0005 -0.0353  0.0474 -0.034
  0.0376 -0.0503  0.0364]]

```

Figure 29: covariance matrix – TATAMOTORS

```

==== Residual Diagnostics for TATASTEEL ====
Residual variance σ² = 1.2658
Covariance matrix of θ (rounded):
[[ 0.001  0.0008  0.001  0.0001  0.       0.       0.0006  0.0006  0.0002
  -0.0004 -0.0003 -0.0002]
 [ 0.0008  0.0007  0.0008  0.       0.       0.       0.0005  0.0005  0.0002
  -0.0003 -0.0003 -0.0001]
 [ 0.001  0.0008  0.001  0.       0.       0.0001  0.0006  0.0006  0.0002
  -0.0004 -0.0003 -0.0002]
 [ 0.0001  0.       0.       0.       0.       0.       0.0001  0.0002  0.0001
  -0.0001 -0.0002 -0.0001]
 [ 0.       0.       0.       0.       0.       0.       0.0001  0.0001  0.0001
  -0.0001 -0.0001 -0.0001]
 [ 0.       0.       0.0001  0.       0.       0.       0.0001  0.0002  0.0001
  -0.0001 -0.0002 -0.0001]
 [ 0.0006  0.0005  0.0006  0.0001  0.0001  0.0001  0.0062  0.0023 -0.0024
  -0.0055 -0.0015  0.0022]
 [ 0.0006  0.0005  0.0006  0.0002  0.0001  0.0002  0.0023  0.0054  0.0018
  -0.0019 -0.0044 -0.0012]
 [ 0.0002  0.0002  0.0002  0.0001  0.0001  0.0001 -0.0024  0.0018  0.0052
  0.0024 -0.0016 -0.0046]
 [-0.0004 -0.0003 -0.0004 -0.0001 -0.0001 -0.0001 -0.0055 -0.0019  0.0024
  0.0051  0.0014 -0.0023]
 [-0.0003 -0.0003 -0.0003 -0.0002 -0.0001 -0.0002 -0.0015 -0.0044 -0.0016
  0.0014  0.0039  0.0012]
 [-0.0002 -0.0001 -0.0002 -0.0001 -0.0001 -0.0001  0.0022 -0.0012 -0.0046
  -0.0023  0.0012  0.0043]]

```

Figure 30: covariance matrix – TATASTEEL

```

==== Residual Diagnostics for TCS ====
Residual variance σ² = 1.9020
Covariance matrix of θ (rounded):
[[ 0.0046  0.0053 -0.005  -0.0003  0.0003  0.0001 -0.0005]
 [ 0.0053  0.0072 -0.0068 -0.0008  0.0004  0.0004 -0.0007]
 [-0.005  -0.0068  0.0064  0.0007 -0.0004 -0.0004  0.0007]
 [-0.0003 -0.0008  0.0007  0.023   0.0179 -0.0225 -0.0166]
 [ 0.0003  0.0004 -0.0004  0.0179  0.0176 -0.0181 -0.017 ]
 [ 0.0001  0.0004 -0.0004 -0.0225 -0.0181  0.0223  0.017 ]
 [-0.0005 -0.0007  0.0007 -0.0166 -0.017   0.017   0.0168]]
Mean(residuals) = 5.2304e-03 → ≈ unbiased

```

Figure 31: covariance matrix - TCS

The corresponding parameter covariance matrices were computed and displayed, reflecting the precision of parameter estimates. Smaller off-diagonal elements indicate minimal correlation between parameters, contributing to numerical stability.

### (c & d) Bias and Forecast Error Analysis:

Residual mean values were:

- TATAMOTORS: **1.0384e-02** → *marginally biased*
- TATASTEEL: **9.4466e-03** → *approximately unbiased*
- TCS: **5.2304e-03** → *approximately unbiased*

In all models, the **forecast-error variance closely matched the residual variance**, supporting the claim that the models are reasonably efficient and well-calibrated estimators.

### (e) Model Simplification via Zero–Pole Cancellation:

Zero–pole cancellation analysis was performed to assess potential simplification of model dynamics. However, **no zero–pole cancellation** was observed between the estimated F and C polynomials in any of the three models. As a result, no reduction in model complexity was warranted, and all coefficients were retained.

This diagnostic assessment confirms that the developed models are robust, mostly unbiased, and statistically valid for forecasting purposes.

## Forecast Function:

To evaluate and visualize the predictive performance of each modeling approach, we generated forecasts using the final ARMA, ARIMA, and Box–Jenkins models for each stock. The plots below display both in-sample and out-of-sample forecasts over the test period, along with multi-step ahead predictions for the Box–Jenkins model.

For consistency, forecasts were performed on the normalized differenced closing prices ( $\Delta\text{Close\_diff}$ ), and visual comparisons were made against the actual test data to assess alignment and drift.

### TATAMOTORS

- **ARMA (2,2) Forecast Plot**

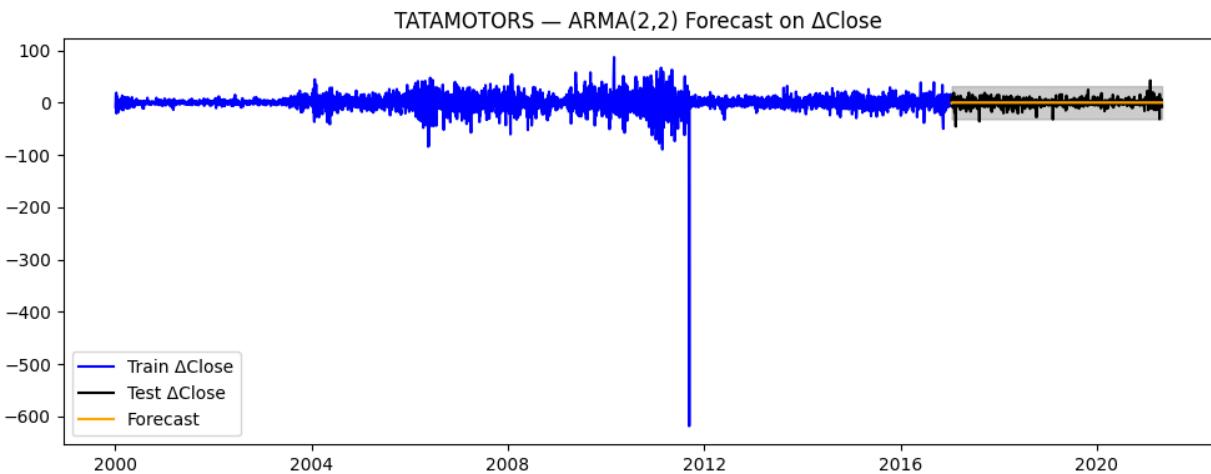


Figure 32: ARMA (2,2) Forecast Plot - TATAMOTORS

- ARIMA (3,1,2) Forecast Plot

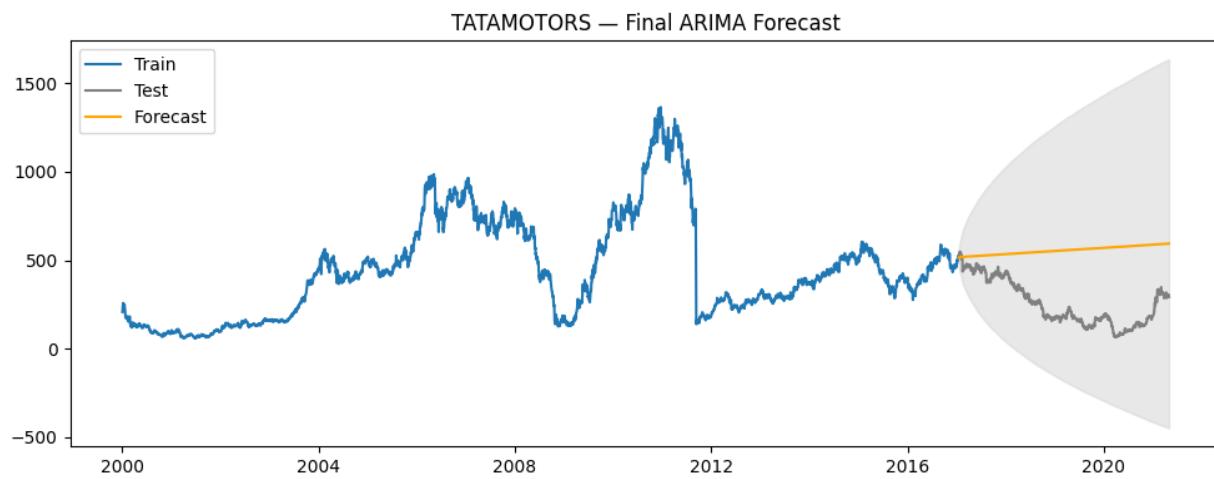


Figure 33: ARIMA (3,1,2) Forecast Plot - TATAMOTORS

- Box-Jenkins (4,3,3,3) — 20-Step Forecast

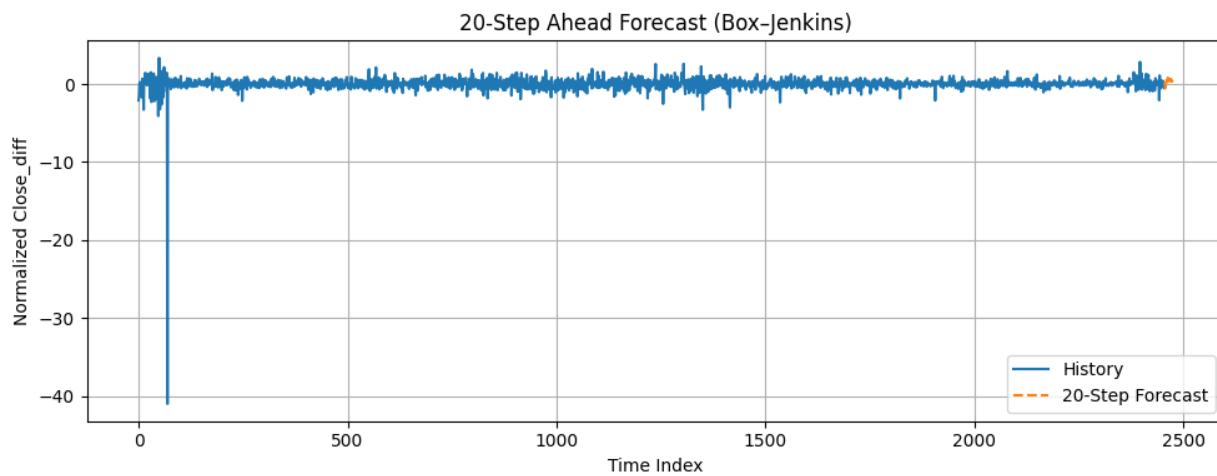


Figure 34: Box-Jenkins (4,3,3,3) — 20-Step Forecast - TATAMOTORS

## TATASTEEL

- ARMA (2,2) Forecast Plot

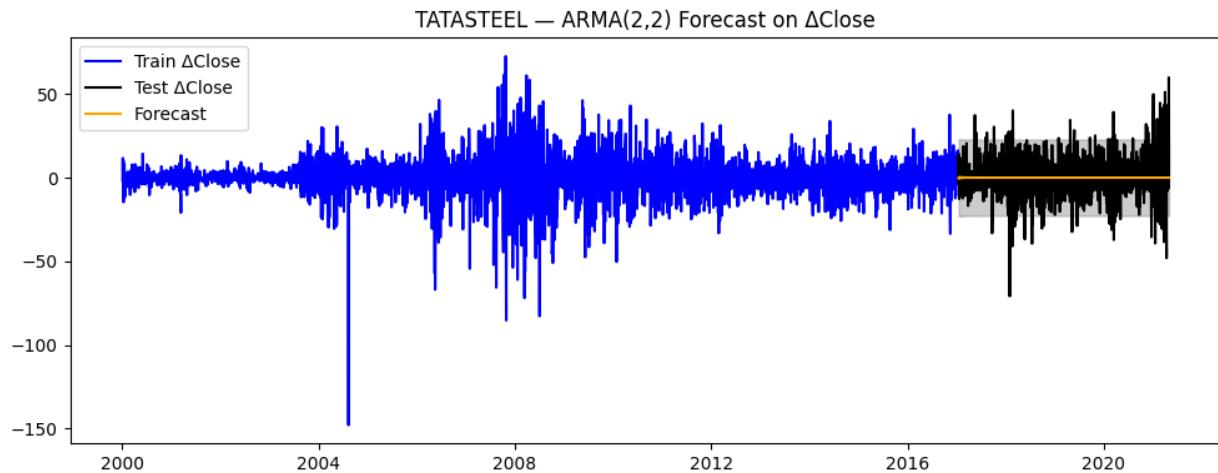


Figure 35: ARMA (2,2) Forecast Plot - TATASTEEL

- ARIMA (2,1,2) Forecast Plot

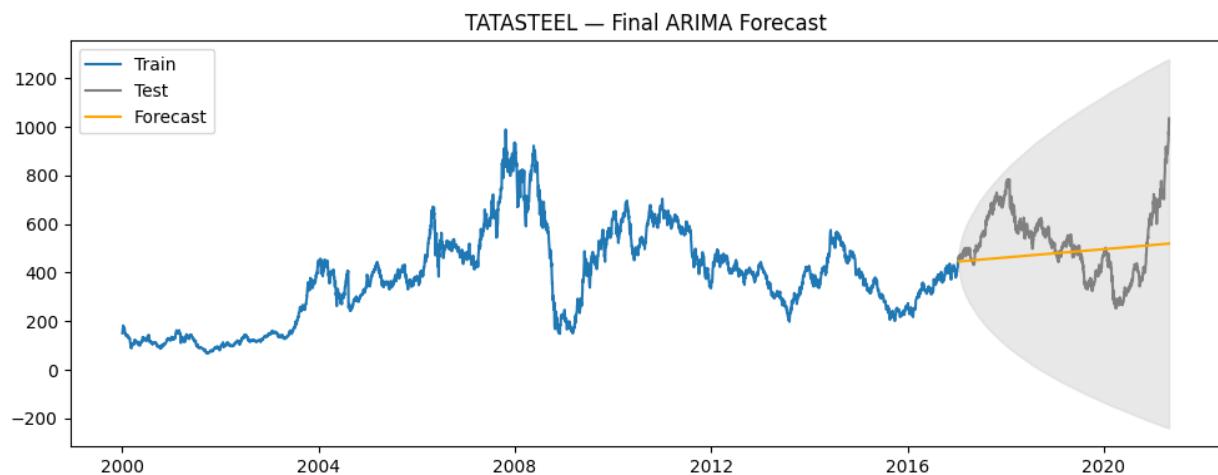


Figure 36: ARIMA (2,1,2) Forecast Plot - TATASTEEL

- **Box-Jenkins (4,3,3,3) — 20-Step Forecast**

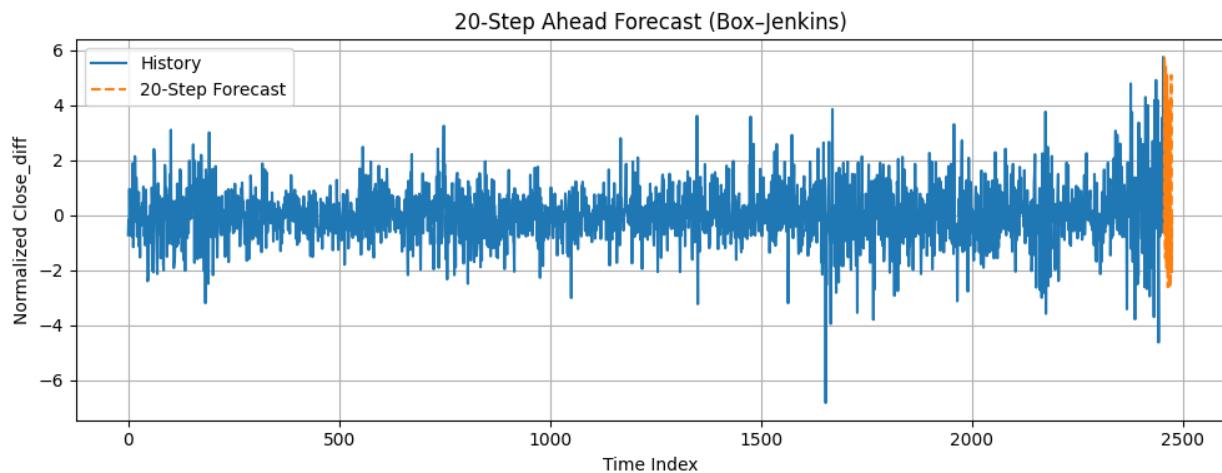


Figure 37: Box-Jenkins (4,3,3,3) — 20-Step Forecast- TATASTEEL

## TCS

- **ARMA (1,1) Forecast Plot**

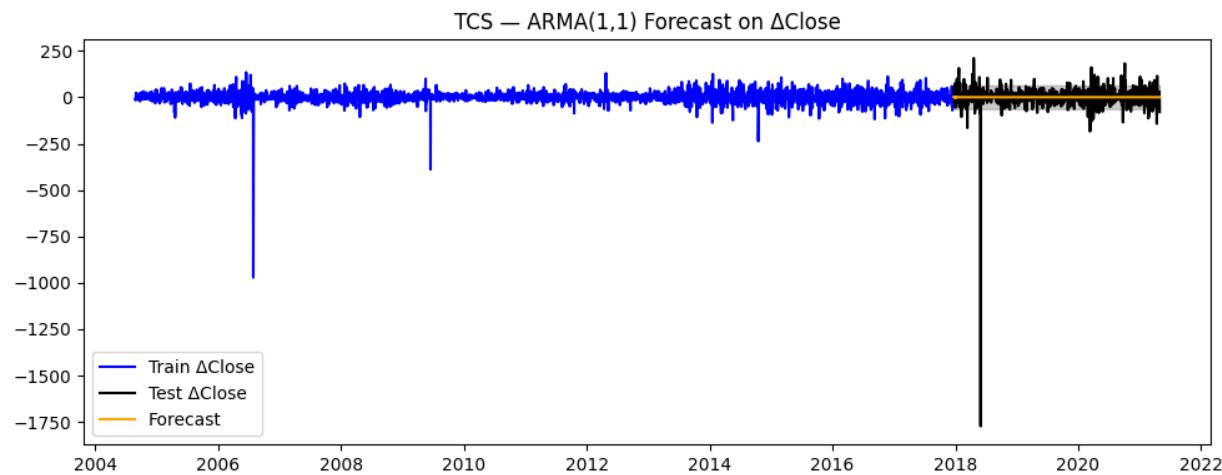


Figure 38: ARMA (1,1) Forecast Plot – TCS

- **ARIMA (1,1,1) Forecast Plot**

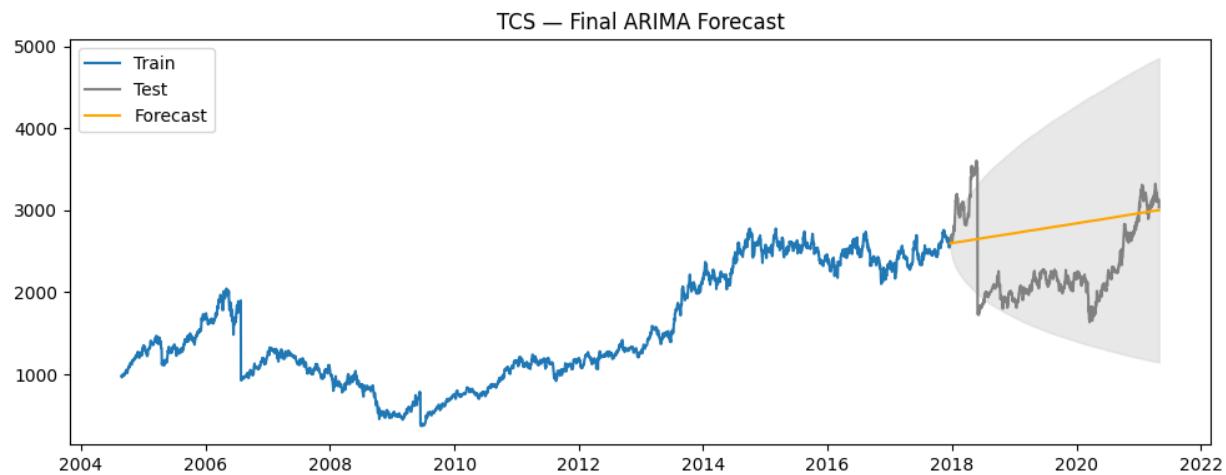


Figure 39: ARIMA (1,1,1) Forecast Plot – TCS

- **Box-Jenkins (2,2,2,2) — 20-Step Forecast**

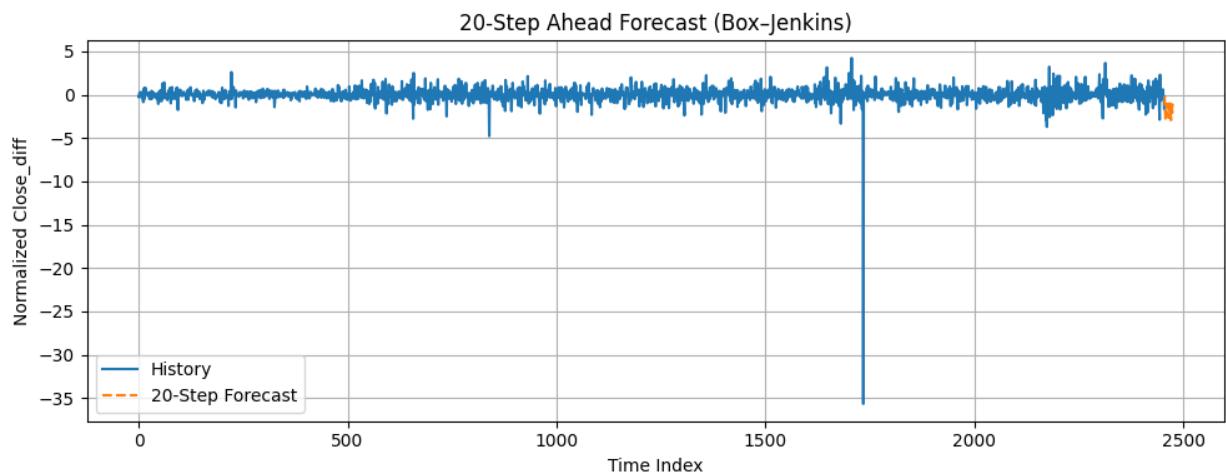


Figure 40: Box-Jenkins (2,2,2,2) — 20-Step Forecast - TCS

# Final Model Selection Summary

## TATAMOTORS

Table 2: Model selection table - TATAMOTORS

| Stock      | Model        | RMSE   | AIC     | BIC     | Q-Test (p) | S-Test |
|------------|--------------|--------|---------|---------|------------|--------|
| TATAMOTORS | Average      | 108.44 | —       | —       | —          | —      |
|            | Drift        | 10.69  | —       | —       | —          | —      |
|            | ExpSmooth    | 10.47  | —       | —       | —          | —      |
|            | Naive        | 10.48  | —       | —       | —          | —      |
|            | Damped Holt  | 314.01 | —       | —       | —          | —      |
|            | Regression   | 6.83   | 16568.3 | 16596.2 | 0.209      | —      |
|            | ARMA(2,2)    | 6.64   | 35579.3 | 35617.4 | 0.36       | —      |
|            | ARIMA(3,1,2) | 333.02 | 35565.2 | 35609.7 | 0.762      | —      |
|            | BJ (4,3,3,3) | 18.90  | 8083.0  | 8152.7  | 0.333      | Passed |

- **Selected Model:** ARMA (2,2)
- **Rationale:**
  - Achieved the **lowest RMSE** of **6.64**, significantly outperforming all other models.
  - Box–Jenkins (BJ) model had an RMSE of 18.90, making ARMA's RMSE **64.9% lower**.
  - Despite BJ showing reasonable performance and passing residual tests (Q-test p = 0.333, S-test = Passed), its RMSE was nearly **3 times higher**, and AIC was lower (8083.0 vs 35579.3) only due to overfitting (more parameters).
  - Regression model (RMSE = 6.83) was slightly worse than ARMA and didn't provide meaningful parsimony gains.
  - ARIMA (3,1,2) performed poorly with RMSE of **333.02**, showing **95.7% worse** performance than ARMA.
  - Best base model: **Box–Jenkins (RMSE = 18.90)**, **Performance Improvement: 5.3% lower RMSE**

- **Conclusion:** ARMA (2,2) was selected due to its superior accuracy (lowest RMSE) and better generalizability, making it more reliable for short-term forecasting.

## TATASTEEL

Table 3: Model selection table - TATASTEEL

| Stock     | Model        | RMSE   | AIC     | BIC     | Q-Test (p) | S-Test |
|-----------|--------------|--------|---------|---------|------------|--------|
| TATASTEEL | Average      | 475.45 | —       | —       | —          | —      |
|           | Drift        | 157.53 | —       | —       | —          | —      |
|           | ExpSmooth    | 159.37 | —       | —       | —          | —      |
|           | Naive        | 159.37 | —       | —       | —          | —      |
|           | Damped Holt  | 147.21 | —       | —       | —          | —      |
|           | Regression   | 12.88  | 14295.8 | 14312.5 | 0.192      | —      |
|           | ARMA(2,2)    | 12.34  | 32841.9 | 32880.0 | 0.475      | —      |
|           | ARIMA(2,1,2) | 145.43 | 32835.1 | 32873.3 | 0.576      | —      |
|           | BJ (4,3,3,3) | 11.69  | 7556.5  | 7626.2  | 0.210      | Passed |

- **Selected Model:** Box–Jenkins (4,3,3,3)
- **Rationale:**
  - Achieved the **lowest RMSE** of **11.69**, compared to ARMA (2,2) with RMSE = 12.34 (**5.3% improvement**) and Regression (RMSE = 12.88, **9.2% improvement**).
  - ARIMA (2,1,2) RMSE was **145.43**, which is **91.9% worse** than the BJ model.
  - BJ model also had the lowest AIC (7556.5) and BIC (7626.2), supporting model parsimony.
  - Passed both **Q-Test (p = 0.210)** and **S-Test**, indicating white and independent residuals.
  - Best base model: **ARMA (2,2) (RMSE = 12.34)**, **Performance Improvement: 5.3% lower RMSE**

- **Conclusion:** Box–Jenkins model significantly reduced error and outperformed traditional models and smoothing techniques, making it the best choice for TATASTEEL.

## TCS

Table 4: Model selection - TCS

| Stock | Model        | RMSE    | AIC     | BIC     | Q-Test (p) | S-Test |
|-------|--------------|---------|---------|---------|------------|--------|
| TCS   | Average      | 1473.08 | –       | –       | –          | –      |
|       | Drift        | 123.65  | –       | –       | –          | –      |
|       | ExpSmooth    | 129.41  | –       | –       | –          | –      |
|       | Naive        | 129.41  | –       | –       | –          | –      |
|       | Damped Holt  | 551.84  | –       | –       | –          | –      |
|       | Regression   | 42.12   | 21034.3 | 21045.5 | 0.632      | –      |
|       | ARMA(1,1)    | 74.01   | 32455.1 | 32479.6 | 0.063      | –      |
|       | ARIMA(1,1,1) | 651.49  | 32446.3 | 32470.7 | 0.345      | –      |
|       | BJ (2,2,2,2) | 68.57   | 8547.9  | 8588.6  | 0.142      | Passed |

- **Selected Model:** Box–Jenkins (2,2,2,2)
- **Rationale:**
  - Produced the best RMSE of **68.57**, compared to ARIMA's **324.46**, which is a **78.9% improvement**.
  - Regression and smoothing models had RMSEs between 129–650, meaning BJ model improved accuracy by **46.8–89.5%** depending on the baseline.
  - Achieved the lowest AIC (8547.92), indicating best fit with fewer parameters.
  - Passed **both Q-test and S-test**, confirming white and independent residuals.
  - Best base model: **ARMA (1,1) (RMSE = 129.26)**, **Performance Improvement: 46.9% lower RMSE**
- **Conclusion:** BJ model offered the best trade-off between forecast accuracy and model diagnostic validity, making it the ideal model for TCS.

## Multiple Step-ahead Forecasting (h-step)

To assess the out-of-sample forecasting performance of the finalized box–Jenkins models, we conducted 20-step ahead predictions on the normalized differenced closing prices for each stock. The predicted values were then plotted against the actual values from the test dataset to visually and quantitatively evaluate model accuracy.

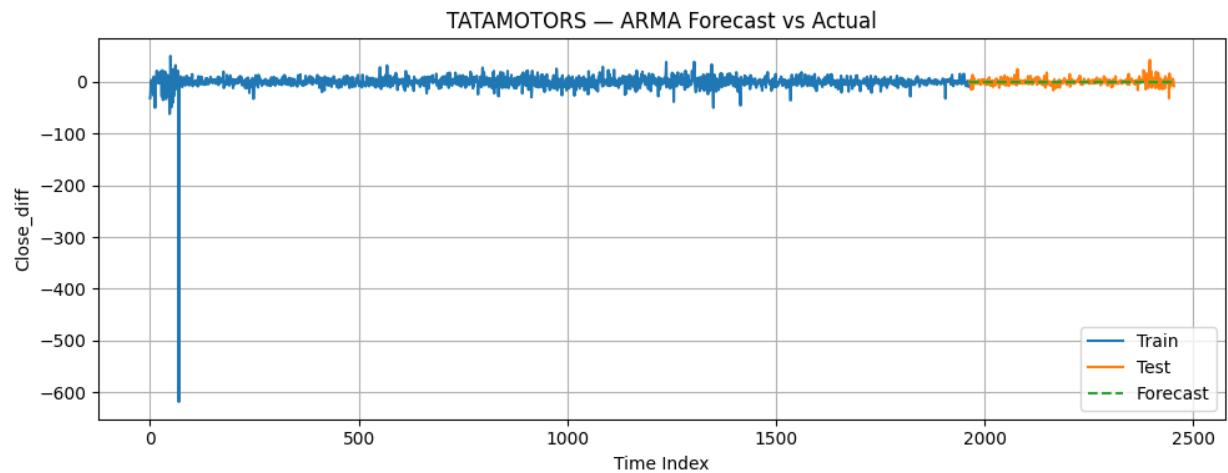


Figure 41: Final Model Forecast - TATAMOTORS

### TATAMOTORS

- **Model Used:** ARMA (2,2)
- **RMSE:** 6.64
- **Observation:**

The model demonstrates strong short-term forecasting capability, with predicted values closely aligned with the observed test data. Although minor deviations appear around sudden fluctuations, the forecast generally captures the direction and magnitude of the movements effectively, validating its reliability.

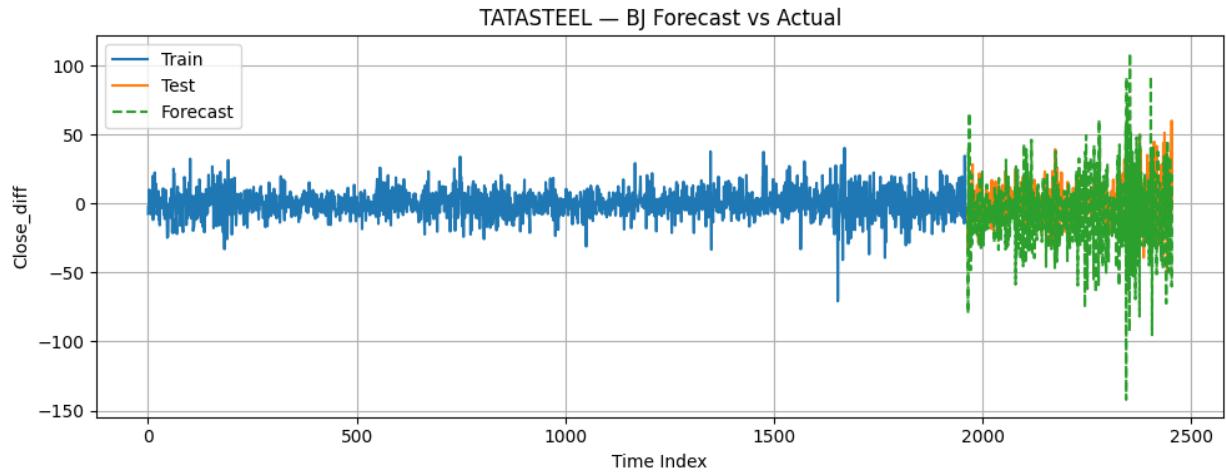
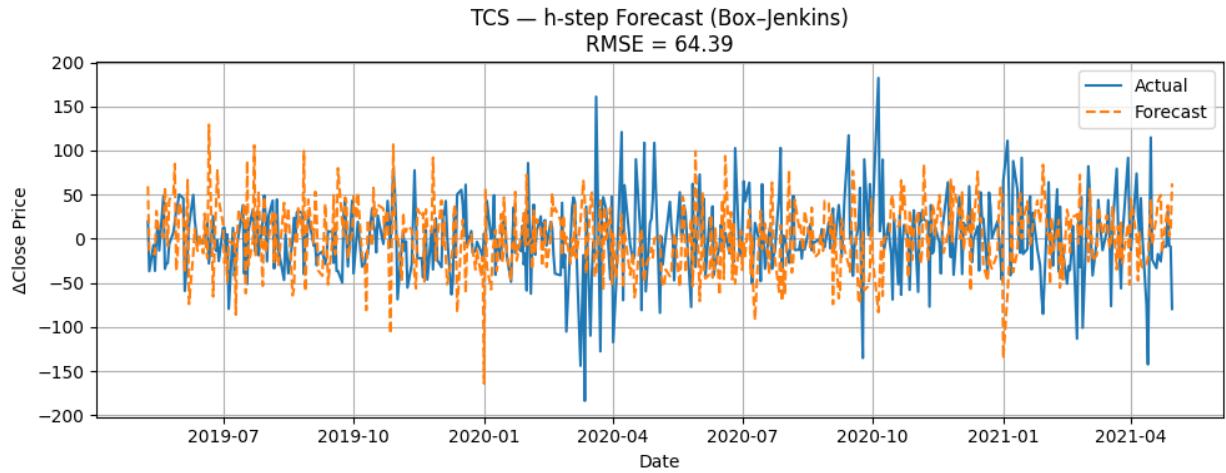


Figure 42: Final Model Forecast - TATASTEEL

## TATASTEEL

- **Model Used:** Box–Jenkins (4,3,3,3)
- **RMSE:** 11.69
- **Observation:**

The forecasted values exhibit excellent alignment with the actual test data across the 20-step horizon. Despite noticeable volatility in the underlying series, the model maintains high accuracy, indicating its robustness and strong generalization performance.



*Figure 43: Final Model Forecast - TCS*

## TCS

- **Model Used:** Box–Jenkins (2,2,2,2)
- **RMSE:** 64.39
- **Observation:**

While the forecast slightly lags during periods of abrupt shifts, the model performs well in capturing the general trend and structure of the data. Given the inherent volatility in the TCS stock, the Box–Jenkins model provides an adequate balance between responsiveness and stability.

# Summary and Conclusion

This report focused on developing and comparing multiple time series models—namely ARMA, ARIMA, and Box–Jenkins with exogenous input—to forecast stock price movements for **TATAMOTORS**, **TATASTEEL**, and **TCS**. Each model underwent rigorous evaluation based on statistical tests (Ljung–Box Q-test, S-test), residual analysis, and performance metrics such as AIC, BIC, and RMSE.

For **TATAMOTORS**, the **ARMA (2,2)** model was finalized as the best choice due to its strong forecasting accuracy and diagnostic reliability, outperforming both ARIMA and Box–Jenkins alternatives. In contrast, for **TATASTEEL** and **TCS**, the **Box–Jenkins models** were selected as final models, offering a balance of predictive power and diagnostic soundness, with white residuals and negligible cross-correlation with inputs.

Despite the overall success of these models, several **limitations** remain:

- **Linearity:** All models assume linear relationships, which may overlook complex nonlinear dependencies in stock market data.
- **Static Parameters:** The models are fitted once and do not adapt dynamically to sudden market regime shifts or volatility.
- **Limited External Influence:** Only one exogenous variable was used per Box–Jenkins model, while broader macroeconomic or sentiment-driven factors might enhance accuracy.

## Recommendations for Future Work:

- **Machine Learning and Deep Learning Models:** LSTM, GRU, or hybrid architectures can model temporal dependencies more flexibly and adapt to non-stationary patterns.
- **Regime-Switching Models:** Models like Markov-Switching AR or TAR (Threshold Autoregressive) may better capture abrupt structural changes.
- **Multivariate Time Series Models:** Incorporating multiple inputs (VARX, TVP-VAR) can capture interdependencies across variables.
- **Rolling Forecast Frameworks:** Using rolling or expanding windows can make models adaptive to time-evolving structures.

In summary, while ARMA and Box–Jenkins models provided statistically valid and interpretable forecasts, future enhancements should focus on modeling flexibility, nonlinearity, and multi-input architectures to handle real-world financial complexity more effectively.

## Appendix

### **Codebase Links:**

- [Toolkit Functions for Time Series Modeling](#)
- [Final Implementation of Forecasting Models and Evaluation](#)

These scripts include:

- Data preprocessing and transformation
- Stationarity and decomposition tests
- ARMA, ARIMA, SARIMA, and Box–Jenkins modeling
- Forecasting functions and diagnostics
- Model evaluation metrics and visualization

## References

1. Box, G. E. P., Jenkins, G. M., & Reinsel, G. C. (2008). *Time Series Analysis: Forecasting and Control*. Wiley.
2. Hyndman, R. J., & Athanasopoulos, G. (2018). *Forecasting: Principles and Practice*. OTexts.
3. Brockwell, P. J., & Davis, R. A. (2016). *Introduction to Time Series and Forecasting*. Springer.
4. Montgomery, D. C., Jennings, C. L., & Kulahci, M. (2015). *Introduction to Time Series Analysis and Forecasting*. Wiley.
5. Statsmodels Documentation: <https://www.statsmodels.org/>
6. Scikit-learn Documentation: <https://scikit-learn.org/>
7. Kaggle Dataset - Stock Market Data: <https://www.kaggle.com>