

ISE 537 Project Report

Price Prediction: Comparison between ARIMA Model and LSTM

Siyu Wu

Abstract

In this study, I delved into the predictive capabilities of two distinct methodologies: the Long Short-Term Memory (LSTM) and AutoRegressive Integrated Moving Average (ARIMA) time series models. Employing a dataset centered around the stock prices of Tesla and Rivian, two key players in the electric vehicle market, I aimed to forecast the subsequent day's stock prices within the time span from January 1, 2020, to December 1, 2023. The finding suggest that while the ARIMA model may excel in capturing certain patterns during training, the LSTM model exhibits improved adaptability to diverse patterns, as evidenced by its superior performance on the test set.

Data. The data was gathered using Yahoo's free APIs, which can be found at <https://pypi.org/project/yfinance/>. The dataset covers the period from January 2020 to December 2023 and includes two columns: date and the closing prices.

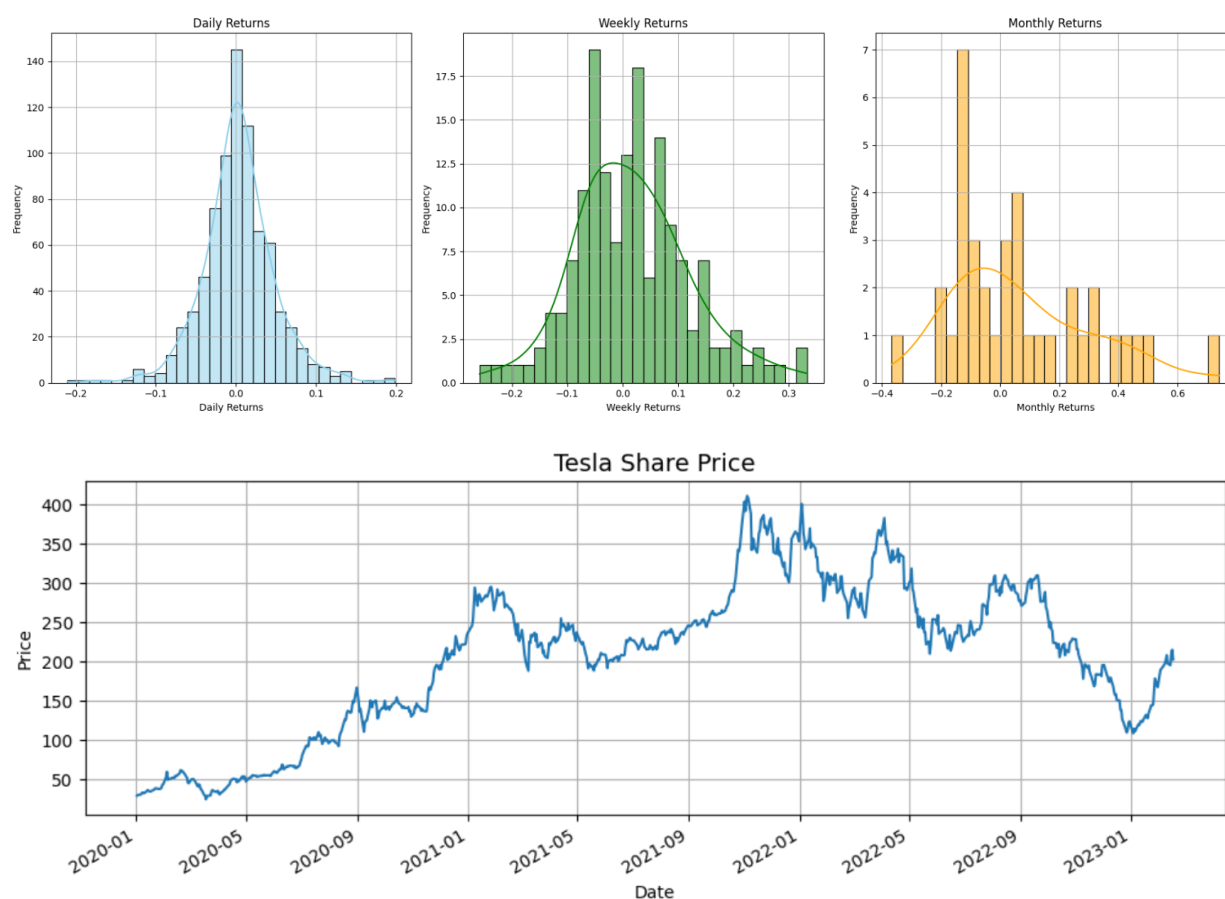


Figure 1: Distribution of Tesla Stock Returns and Price Plot

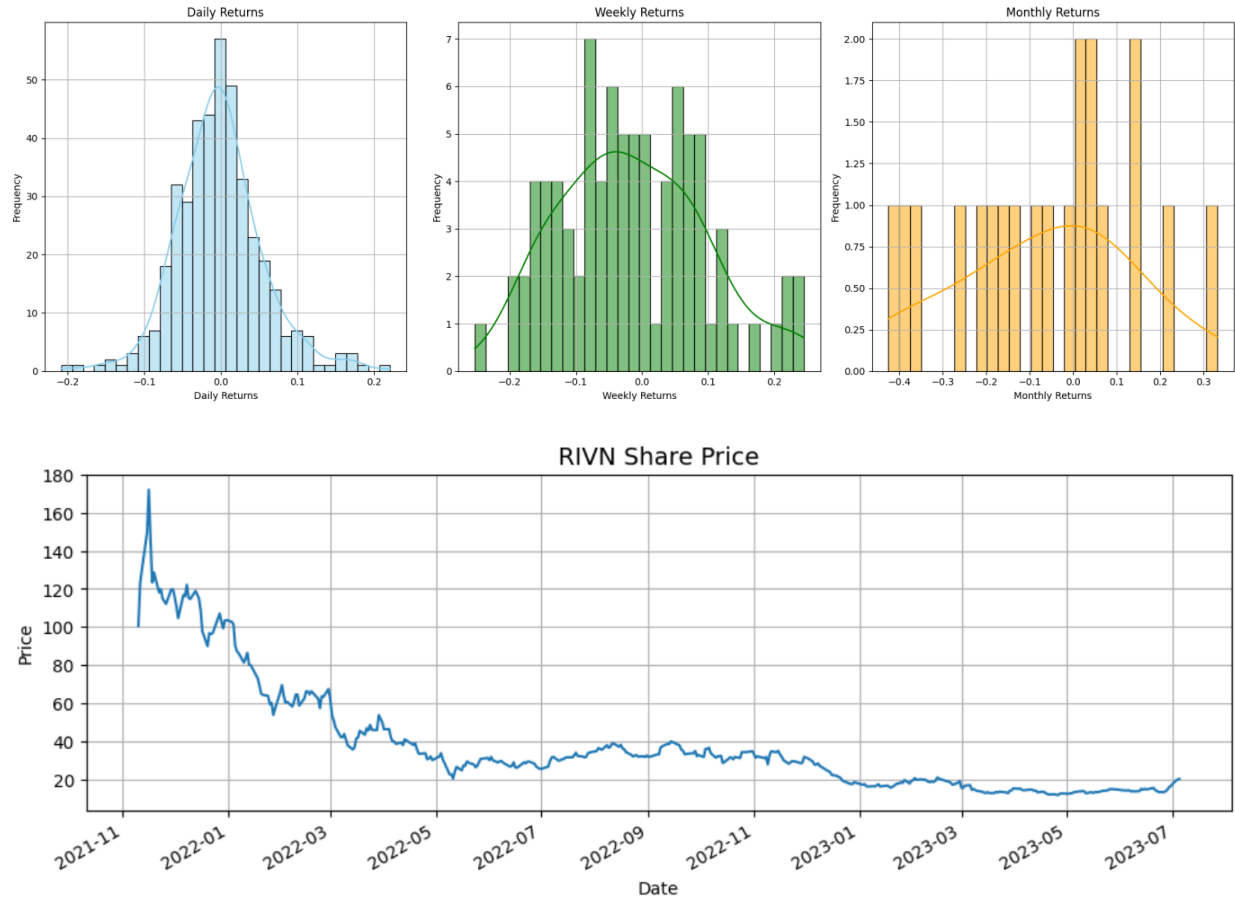


Figure 2: Distribution of Rivan Stock Returns and Price Plot

The selection of this timeframe is driven by the consideration that stock prices at the early stages can exhibit extreme highs or lows, influencing the modeling process. Additionally, a shorter period, covering the past three years, was chosen to focus on recent trends and avoid potential distortions arising from historical market dynamics. This approach enables a more relevant and accurate examination of stock price behavior in a contemporary context.

Stationarity Check. To check if TESLA stock data is suitable for the ARIMA time series model, we need to ensure that the data is stationary. We used the Augmented Dickey–Fuller (ADF) test on different versions of the data—original, log-transformed, and daily returns. The ADF test showed strong evidence, with a 99% confidence level, supporting the idea that daily returns are stationary. This conclusion was backed up when looking at the ADF plot, which showed a clear characteristic of stationary time series—rapid decrease in the AutoCorrelation Function (ACF) to zero with increasing lag. These results confirm that daily returns are stationary, making it a good fit for applying the ARIMA time series model to TESLA stock data.

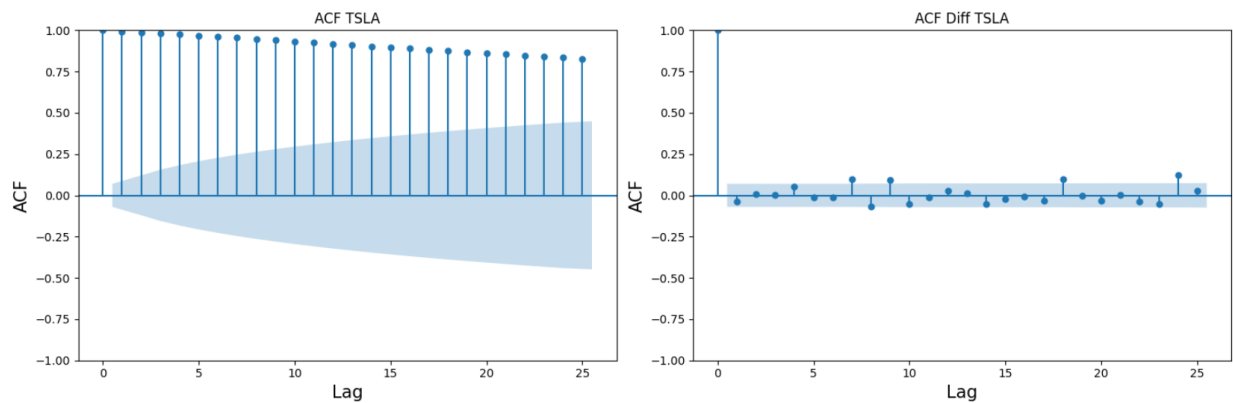


Figure 3: ACF for Close Price and Daily Return

Likewise, I replicated the procedure for RIVN stock. The test results indicate that the closing price is already stationary, affirming its suitability for constructing the ARIMA model.

Model Choice. I developed a search function for ARIMA models, systematically testing different combinations of parameters (p, d, q) to identify the optimal model based on lower AIC and BIC values. Upon analysis, the chosen ARIMA model for TSLA stock price is (0, 0, 0), indicating no autoregressive, differencing, or moving average components. For RIVAN stock price, the selected ARIMA model is (1, 0, 0), denoting one autoregressive component and no differencing or moving average components. These parameter configurations represent the best-fitting models according to the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC). A lower AIC or BIC indicates a better-performing model. AIC penalizes complex models less severely, making it sensitive to smaller improvements in fit. On the other hand, BIC imposes a stronger penalty for model complexity, favoring simpler models.

The LSTM model is composed of two main layers. The first layer is an LSTM layer with 4 units, designed to capture temporal dependencies in sequential data. The second layer is a Dense layer with a single unit responsible for generating the model's output. The model is compiled using the Adam optimizer and employs mean squared error (MSE) as the loss function. Throughout training, which spans 100 epochs with a batch size of 256, the model adjusts its parameters to minimize the discrepancy between predicted and actual values, thereby enhancing its ability to make accurate predictions on sequential data. I choose this model because it performs good in the previous assignment.

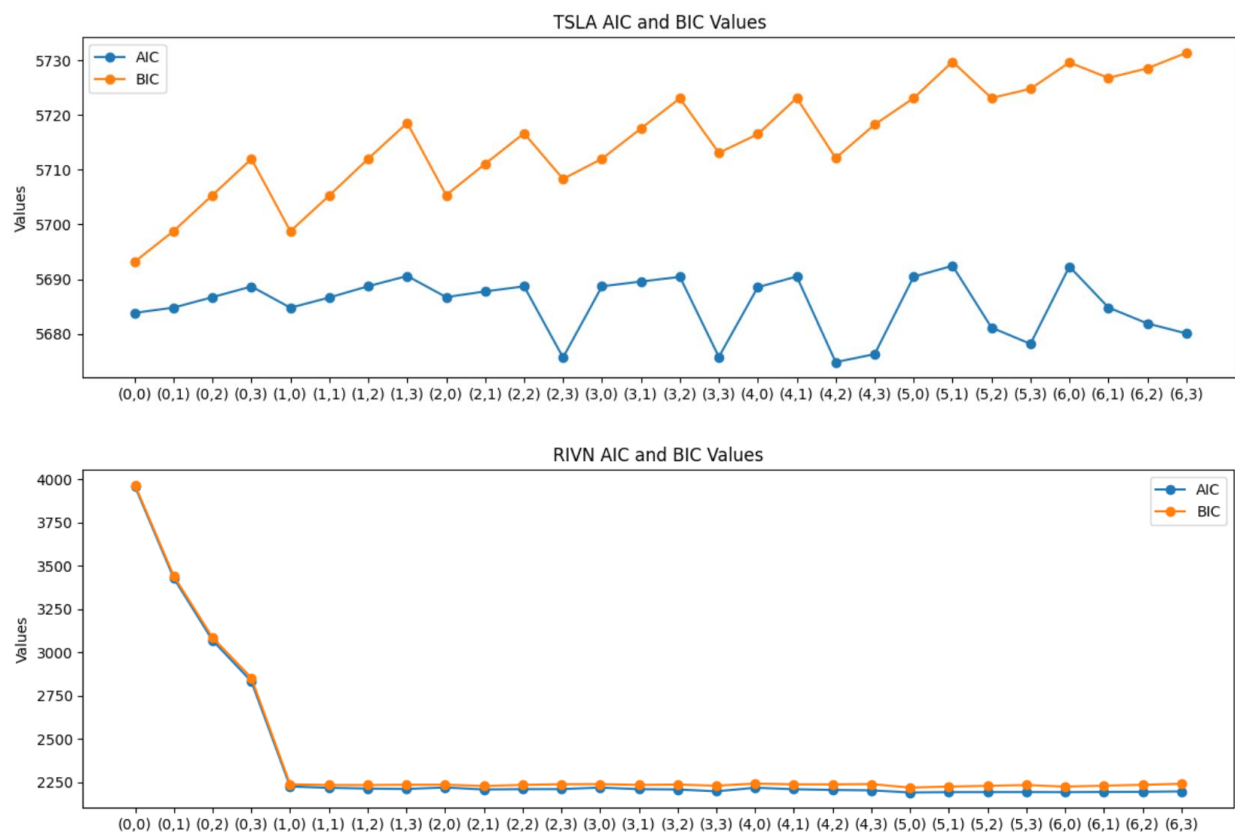


Figure 4: TSLA/RIVN AIC and BIC Values

Prediction. In the prediction phase, both models were evaluated using a rolling window spanning 10 days to forecast the subsequent day's data. The Root Mean Squared Error (RMSE) was employed as the metric for assessing predictive performance. For the ARIMA model, the training RMSE values for TSLA and RIVN were 17.8 and 2.86, respectively, while the corresponding test RMSE values were 17.2 and 20.95. In contrast, the LSTM model exhibited training RMSE values of 41.16 for TSLA, and 11.37, and test RMSE values of 17.8 and 6.72, respectively. These results provide a quantitative measure of the predictive accuracy of each model, with lower RMSE values indicating better performance.



Figure 5: TSLA Model Performance

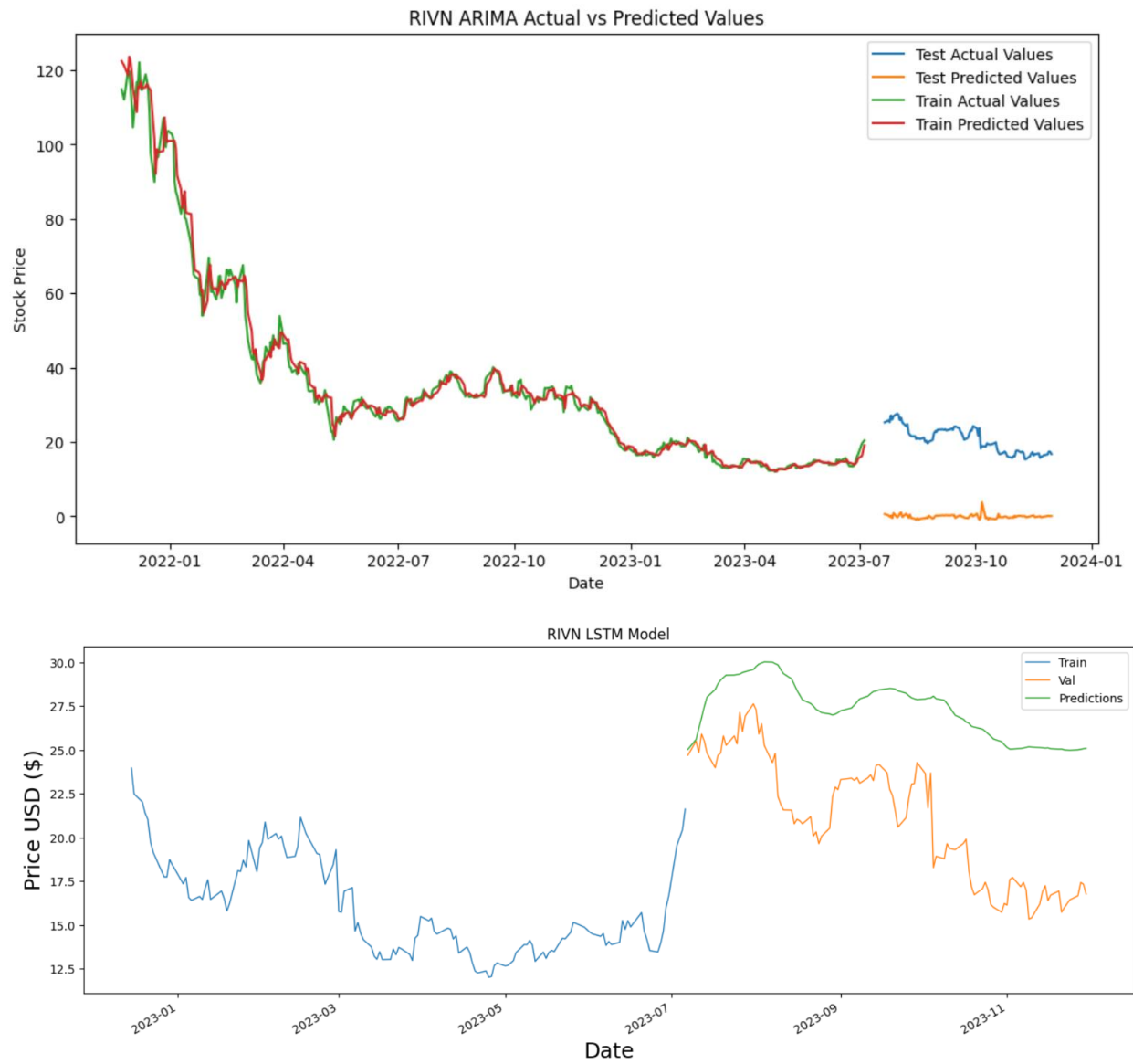


Figure 6: RIVN Model Performance

In the context of the LSTM model, a declining trend observed in the loss plot would suggest that the model is converging.

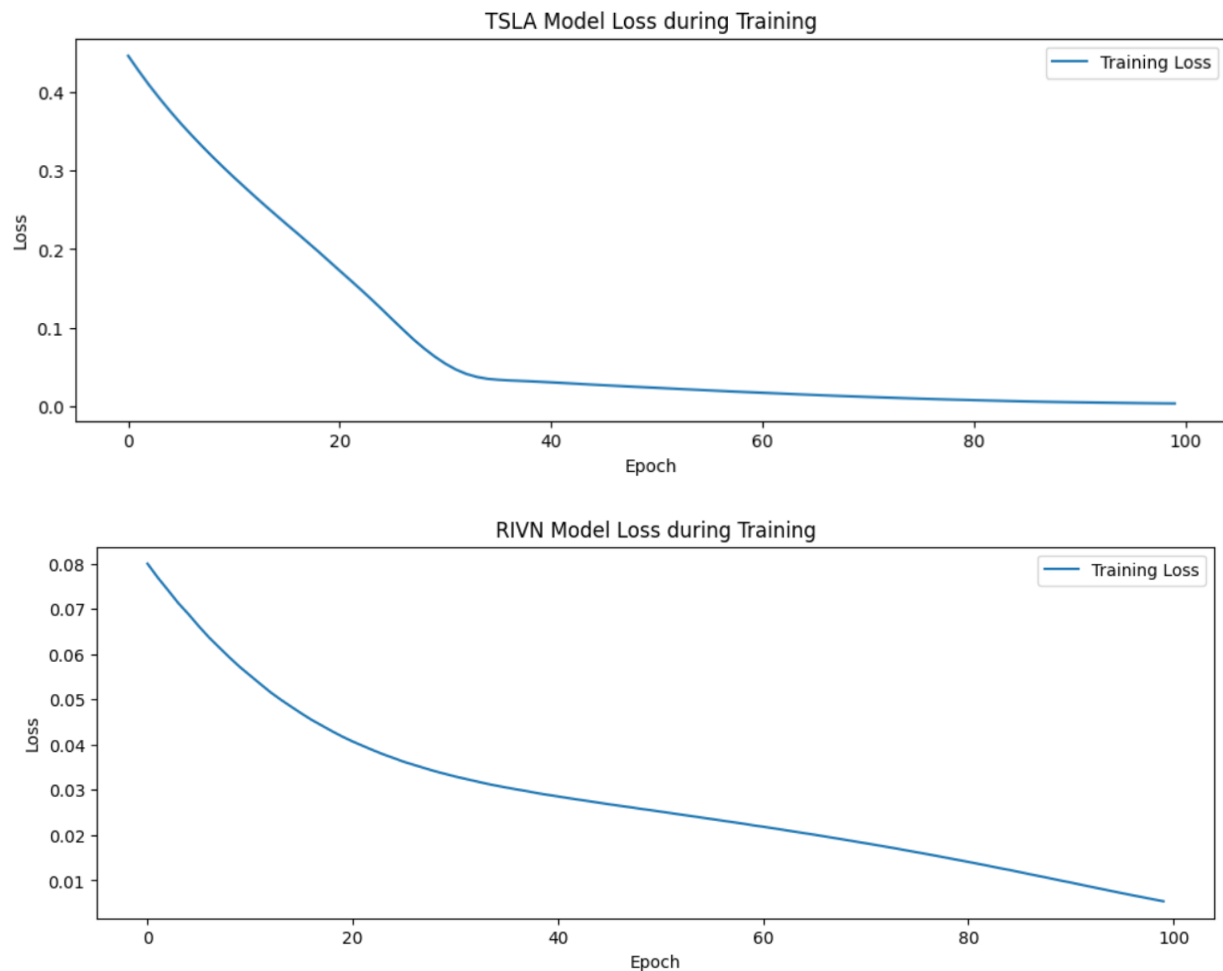


Figure 7: Model Convergence

Conclusion. The ARIMA model, as reflected in the RMSE values, demonstrated relatively lower training RMSE for both Tesla and RIVN compared to the LSTM model. However, its performance on the test set exhibited higher RMSE values. On the other hand, the LSTM model exhibited higher training RMSE values, indicating a larger error during the learning phase, but demonstrated better generalization with lower RMSE values on the test set for both Tesla and RIVN. Notably, RIVN's test RMSE for the LSTM model was substantially lower than that of the ARIMA model. In conclusion, the choice between ARIMA and LSTM depends on the specific characteristics of the data and the trade-off between training and generalization performance. While the ARIMA model may excel in capturing certain patterns during training, the LSTM model exhibits improved adaptability to diverse patterns, as evidenced by its superior performance on the test set.