

# ***Comparative Analysis of Forecasting Chevron's Crude Oil Stock Performance with Machine Learning Techniques***

**Muyang Chen<sup>1,a,\*</sup>**

<sup>1</sup>*Bachelor of Commerce, Monash University, Blackburn Road, Clayton, Victoria, 3800, Australia  
a. mche0237@student.monash.edu*

*\*corresponding author*

**Abstract:** The objective of this study is to predict the Chevron's Corporation stock market performance by conducting a comparative analysis of contemporary and conventional machine learning approaches, with a particular focus on the CNN-LSTM and ARIMA models. Given the unpredictable characteristics of the crude oil industry, forecasting stock prices with precision has emerged as a pivotal dilemma for both investors and analysts. This research utilizes ARIMA, which is representative of conventional time series forecasting methods, and CNN-LSTM, which embodies the latest advancements in deep learning techniques, to address the intricacies associated with predicting stock prices in the energy sector. Through a comprehensive data preparation process and the application of sophisticated modeling techniques, this study aims to rigorously assess the predictive capabilities of both models in forecasting Chevron's stock prices. Traditional statistical analysis often relies on the ARIMA model as a benchmark, while the CNN-LSTM model seeks to identify the complex, non-linear patterns prevalent in financial market time series data. This research conducts a comparative evaluation of the two models, focusing on their accuracy, strengths, and limitations. The findings carry important implications for the realm of financial forecasting, shedding light on how modern deep learning techniques stack up against traditional approaches in predicting stock market movements. Beyond contributing to scholarly debates on financial prediction, this study also provides actionable insights for financial analysts.

**Keywords:** Machine learning, CNN-LSTM, ARIMA, stock forecasting, time series

## **1. Introduction**

In the global energy sector, Chevron Corporation stands out for its significant market presence and vital contribution to the global oil supply. As a multinational entity with profound influence in the energy domain, Chevron's stock performance not only reflects the financial health of the company but also serves as a crucial indicator of global energy market trends. However, in 2023, Chevron experienced a notable decline in financial performance and the profits fell by 40% compared to the previous year [1]. In such a volatile market environment, accurately forecasting Chevron's stock price has become an urgent need.

This study employs two distinct models to forecast the stock performance of Chevron Corporation: the Auto-Regressive Integrated Moving Average (ARIMA) model and the Convolutional Neural Networks-Long Short-Term Memory (CNN-LSTM) hybrid model [2]. George and Gwilym firstly introduced the ARIMA model in the early 1970s, this approach to time series analysis has since

become one of the most widely used forecasting methods and its accuracy has been verified by predicting the stock prices, electricity prices and epidemic datasets [3]. The ARIMA model, as a classic tool for time series analysis, has long been utilized to forecast stock performances and other financial indicators. In the process of predicting Chevron's stock price, the traditional ARIMA model encounters limitations, particularly when dealing with high seasonal adjustments or when diagnostic tests fail to confirm the stationarity of the time series after such adjustments. This challenge highlights a critical drawback: the classical ARIMA model's reliance on static parameters restricts its effectiveness in forecasting Chevron's stock price fluctuations, which may exhibit significant seasonal variability. Additionally, the ARIMA approach's efficacy is further constrained by its need for a substantial volume of historical data to accurately select the optimal model for analysis [4]. Meanwhile, The CNN-LSTM Model utilizes the unique ability of Convolutional Neural Networks (CNN) to emphasize the most important features within its input, making it a crucial tool in feature engineering. LSTM networks are highly effective at analyzing time series data due to their ability to process data across temporal sequences [5]. This research combines the distinctive features of CNN and LSTM to create a model designed to forecast stock prices, with a focus on capturing the intricacies of Chevron's stock behavior [6]. This study aims to compare the effectiveness of these two models in predicting Chevron's stock prices, exploring which method proves more efficient [7].

The study of predicting stock performance begins by looking back at its beginnings, which involved using classic statistical and econometric models [8]. The ARIMA model, moving averages, and exponential smoothing are fundamental tools used in early attempts to understand market dynamics. The ARIMA model comprises the autoregressive (AR) model, moving average (MA) model, and seasonal autoregressive integrated moving average (SARIMA) model [9]. The ARIMA model is well-known for its ability to capture different market circumstances using integrated, autoregressive, and moving average components, making it a fundamental tool in time series analysis. Moving averages are a simple but powerful tool for detecting trends by reducing short-term variations, offering a greater insight into the fundamental trajectory of stock prices [10]. Exponential smoothing improves the method by giving decreasing weights to past observations, making projections more responsive to recent changes. These models have established the foundation for comprehending market dynamics and providing insights into the temporal patterns influencing stock performance. Their lasting importance highlights the complex connection between statistical methods and the changing field of financial analysis, paving the way for the later introduction of machine learning approaches in predicting stock performance [11].

Forecasting approaches have evolved significantly in line with advancements in computing technology. This evolution signifies a shift from conventional methods to advanced analytical approaches, substantially transforming the field of financial prediction. With the growth in computational powers, there was a demand for methodologies that could utilize this increased power, resulting in the incorporation of machine learning and deep learning into the tools available to financial analysts [12]. Machine learning, rooted in regression analysis, offers a flexible structure for forecasting stock performance by allowing models to analyze past data and detect trends that traditional statistical methods may miss. Neural networks, a component of deep learning, have advanced by analyzing data through layers of interconnected nodes, imitating the human brain's capacity to identify intricate patterns and connections in extensive datasets. This architecture excels at capturing the complexities of market behavior and providing predictions that consider various aspects affecting stock prices [13].

Reinforcement learning is a key aspect of technology advancement that brings a dynamic method to predicting outcomes. Reinforcement learning algorithms optimize decision-making processes in real-time trading settings by learning from market interactions and modifying strategies depending

on rewards or penalties. This strategy highlights the move towards adaptive, self-enhancing models that can navigate the volatile financial markets more accurately.

The progress in machine learning and deep learning has greatly improved the tools available to financial analysts. These methodologies provide enhanced accuracy and the ability to analyze intricate market trends, marking a significant advancement in the pursuit of dependable stock performance prediction. This advancement demonstrates the improvement in computing technology and marks a new phase in financial analysis, where the level of understanding and ability to anticipate outcomes are constantly being redefined.

## **2. Methods**

### **2.1. Data Source**

This study examined Chevron Corporation's stock data spanning from December 30, 2019, to December 29, 2023. The dataset contains trading records spanning three years, with six main metrics: Date, Open, High, Low, Adjusted Close, and Volume. These indicators together offer a thorough summary of the stock market's performance. To ensure the effective training and testing of our models, the dataset was split in a 7:3 ratio: approximately the first 70% of the data (spanning about the first two years and nine months) was designated as the training set, used for model training and fine-tuning. The latter 30% of the data (covering the final approximately one year and three months) served as the test set, intended for assessing the models' forecasting accuracy. Special consideration was focused on exceptional cases in the dataset during the data preprocessing stage, especially when the volume showed unusual values or zeros. Linear interpolation was used to patch gaps in order to maintain data continuity and completeness. To improve the generalization power and forecast accuracy of our models, the data was normalized to make it more appropriate for deep learning model inputs. Our research seeks to investigate and compare the effectiveness of ARIMA and CNN-LSTM models in predicting stock prices through a detailed data processing and analysis methodology. This paper aims to offer more accurate and dependable methodological assistance for financial research and predicting stock market trends.

### **2.2. Data Preparation and Preprocessing**

For a detailed view, the author showcased data between December 30, 2019, and January 6, 2020, focusing on key features like Volume, Open, High, Low, Close, and Adjusted Close. The absence of any lacking values in the dataset indicates that it is comprehensive and does not require any further augmentation or removal. This validation demonstrates that the data are suitable for CNN-LSTM/ARIMA model training and advanced analysis. For model training, the author divided the data into 70% training and 30% assessment sets. In order to exclude any potential absent or anomalous values, a comprehensive examination was conducted. By utilizing linear interpolation to bring in missing data points, the CNN-LSTM model was capable of analyzing trends and momentum. The ARIMA model prioritized autoregressive characteristics of time series and necessitated minimal feature engineering. The author normalized and standardized the data with MinMaxScaler and StandardScaler in order to enhance the model. Additionally, the author formatted and chronologically arranged the time series data in preparation for analysis.

## 2.3. Method Introduction

### 2.3.1. ARIMA Model

As a widely utilized statistical model in time series forecasting, ARIMA excels in capturing linear relationships and trends in historical data. It integrates autoregression, differencing to achieve stationarity, and a moving average component. The model's parameters (2, 1, 0) were systematically identified through a process of iterative testing and validation, ensuring the model's adequacy in capturing the underlying patterns of the stock price data.

### 2.3.2. CNN-LSTM Model

This advanced model combines the feature extraction capabilities of CNNs with the sequence prediction strengths of LSTMs. CNN layers first process the input data, identifying salient features without manual intervention. The extracted features are then fed into LSTM layers, which are adept at handling long-term dependencies and fluctuations in time series data. This model's architecture was optimized to balance complexity and performance, ensuring efficient processing of Chevron's stock data.

## 3. Results and Discussion

### 3.1. Stationary Test

The building of the ARIMA model begins by importing the target time series data, specifically the 'Close' prices, from a pandas DataFrame. To assess the appropriateness of the ARIMA model, the author performs an Augmented Dickey-Fuller (ADF) test to verify the presence of stationarity, which is a critical need for ARIMA modeling. When the ADF test shows that there is non-stationarity, the author does differencing on the series to make it stable. Afterward, the author retests for stationarity. After confirming stationarity, the author examines the Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF) plots to ascertain the order parameters (p, d, q) of the ARIMA model. These plots aid in determining the appropriate autoregressive (p) and moving average (q) terms, taking into account the applied differencing (d). Next, the author proceeds to build and apply the ARIMA model using the statsmodels package, while carefully choosing the suitable parameters. A detailed summary of the model is generated, which includes diagnostics and estimates of parameters, in order to evaluate the quality of fit and the statistical significance. When there are problems with the model converging, tweaks are made to the optimization process in order to enhance the accuracy of the fit.

Table 1: ARIMA (2,1,0) Model Parameter Table.

Item	Symbol	Coefficient	SE	z-value	p-value	95% CI
Constant	c	0.029	0.085	0.34	0.734	-0.138 ~ 0.196
AR	$\alpha_1$	-0.045	0.025	-1.825	0.068	-0.093 ~ 0.003
	$\alpha_2$	0.048	0.023	2.062	0.039	0.002 ~ 0.094

Table 2: ARIMA Ljung-Box Q Test Statistic Result

LB test	Value
Root mean square error RMSE	2.6501
Root mean square error RMSE	7.0230
Mean absolute error MAE	1.8891

Table 2: (continued).

Mean absolute percentage error MAPE	0.0158
p-value	0.033

The table 1 and 2 illustrate the parameter estimates for an ARIMA(2,1,0) model that was applied to the stock data of Chevron Corporation. The text describes a constant term that is not statistically significant, as evidenced by a coefficient of 0.029 and a p-value of 0.734. This suggests that the constant term has a limited impact on the model's ability to make accurate predictions. It is worth noting that the AR terms, which represent the impact of previous data points on subsequent values, exhibit divergent levels of statistical significance. The coefficient of the first autoregressive (AR) component is -0.045, indicating a minimal statistical contribution. The p-value associated with this term is close to the standard significance threshold of 0.068. On the other hand, the second autoregressive (AR) term has statistical significance, as indicated by its coefficient of 0.048 and a p-value of 0.039, which above the conventional alpha threshold of 0.05. The precision of the coefficients is reflected by the standard errors associated with these estimations, where smaller numbers indicate more dependability. The relevance of the coefficients is further supported by the z-values, which measure the standard deviations from the mean estimate. Finally, the 95% confidence intervals represent a comprehensive range in where the actual parameter values are expected to lie, so providing valuable information regarding the dependability of the estimates. The statistical evidence emphasizes the significant contribution of the second autoregressive (AR) component to the model, emphasizing its importance in capturing the temporal dependencies of the stock.

### 3.2. Model Results

The Figure 1 illustrates a graphical representation that juxtaposes the observed stock prices with the projected values produced by a predictive model. The 'Fit' line, which is presumed to depict the model's predictions within the given sample, exhibits a strong correlation with the observed prices, suggesting a satisfactory alignment between the model and historical data. In general, the strong correlation observed between the 'Fit' and 'Actual' lines until the forecast initiation point suggests that the model has effectively captured the fundamental characteristics of the historical data. The forecast interval is a valuable tool for assessing the anticipated variability in projections, which is a critical factor in evaluating risks in financial decision-making.

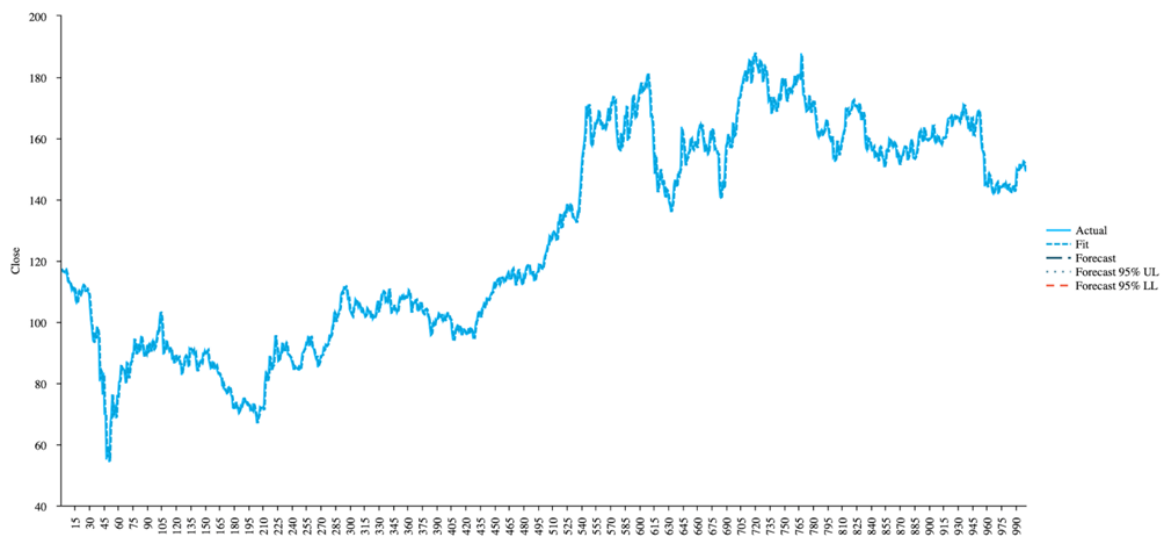


Figure 1: ARIMA Model fitting and prediction.

Figure 2 provides a visual representation of the Chevron stock prices in comparison to the predictions made by the CNN-LSTM model. Upon analyzing the graph, various insights are revealed: To begin with, it is noteworthy that both the observed (shown in red) and projected (represented in blue) prices demonstrate comparable patterns, indicating that the CNN-LSTM model adeptly reflects the broader fluctuations in prices. Nevertheless, there is a conspicuous delay between the projected and observed lines, which is characteristic of time series models. The projected prices exhibit a greater degree of smoothness, suggesting that the model may have accounted for some more noise. However, it is important to note that the projections do not fully reflect the dramatic variations in actual prices, as this is to be expected given the unpredictable nature of market movements. Substantial disparities between the lines indicate regions of less precise forecasts, necessitating additional examination. A strong correspondence between lines may signal overfitting on the training data, but a strong correspondence on unknown data indicates successful generalization. The presence of a divergence between lines towards the right end of the plot suggests the possibility of less precise forecasts, which could be attributed to the restricted availability of surrounding data for recent time steps. In general, the graph offers significant information regarding the performance of the model and identifies potential areas for enhancement.

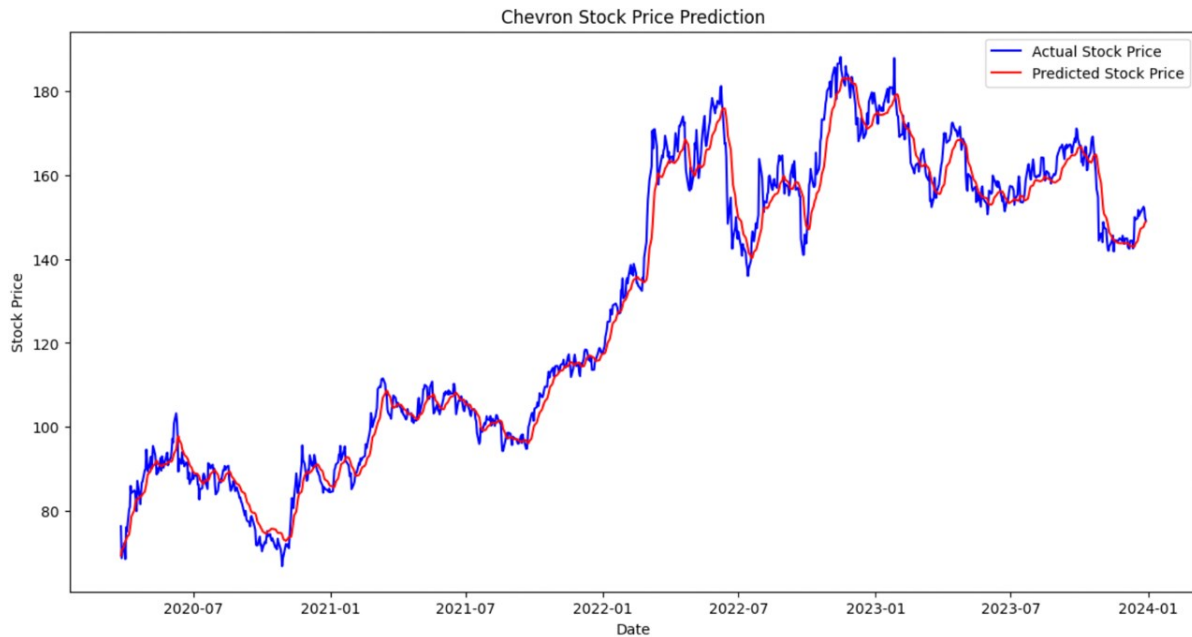


Figure 2: CNN-LSTM Model fitting and prediction.

#### 4. Conclusion

When evaluating the effectiveness of the Autoregressive Integrated Moving Average (ARIMA) and Convolutional Neural Network-Long Short-Term Memory (CNN-LSTM) models in forecasting Chevron stock prices, a thorough examination uncovers unique strengths and weaknesses associated with each model. The ARIMA model is widely praised for its simplicity and ease of comprehension. Time series data analysis is highly proficient at capturing linear trends and seasonal fluctuations, which are crucial elements in the field of financial forecasting. Nevertheless, the efficacy of the model decreases when faced with the complex and non-linear dynamics commonly observed in stock market fluctuations. This constraint can significantly affect the accuracy of the model, particularly when making projections over long durations. Moreover, the ARIMA model requires manual hyperparameter adjustment, a process that can be substantial in terms of time and effort, hence

diminishing its overall effectiveness. On the other hand, the CNN-LSTM model exhibits enhanced proficiency in interpreting intricate nonlinear patterns and maintaining enduring dependencies that are typical of sequential data. This level of expertise makes it extremely well-suited for predicting stock prices, where such intricacies are typical. Despite the greater processing costs and the risk of overfitting associated with the CNN-LSTM model, its primary advantage comes in its ability to autonomously learn features from raw data. This particular capacity proves to be particularly advantageous when it comes to the management of multivariate and high-dimensional datasets. Nevertheless, the complex structure of the model presents interpretive obstacles, rendering it arduous to clarify the fundamental mechanisms that propel its projections. In conclusion, it is crucial to carefully consider the choice between ARIMA and CNN-LSTM models when predicting Chevron stock prices, taking into account the specific requirements of the forecasting task. The determination of the most successful methodology requires more empirical investigation and validation, considering the trade-offs between simplicity and analytical depth, as well as computational efficiency and forecast accuracy.

## References

- [1] Box, G.E., Jenkins, G.M., Reinsel, G.C. and Ljung, G.M. (2015) *Time series analysis: forecasting and control*. John Wiley & Sons.
- [2] Ariyo, A.A., Adewumi, A.O. and Ayo, C.K. (2014) *Stock price prediction using the ARIMA model*. In 2014 UKSim-AMSS 16th international conference on computer modelling and simulation, 106-112.
- [3] Contreras, J., Espinola, R., Nogales, F.J. and Conejo, A.J. (2003) *ARIMA models to predict next-day electricity prices*. *IEEE transactions on power systems*, 18(3), 1014-1020.
- [4] Benvenuto, D., Giovanetti, M., Vassallo, L., Angeletti, S. and Ciccozzi, M. (2020) *Application of the ARIMA model on the COVID-2019 epidemic dataset*. *Data in brief*, 29, 105340.
- [5] Fawaz, H.I., et al. (2019) *Deep learning for time series classification: a review*. *Data Min Knowl. Disc.*, 33(4), 917-963.
- [6] Livieris, I.E., Pintelas, E. and Pintelas, P. (2020) *A CNN-LSTM model for gold price time-series forecasting*. *Neural computing and applications*, 32, 17351-17360.
- [7] Yang, R., Liu, X., Yu, R., Hu, Z. and Duan, X. (2022) *Long short-term memory suggests a model for predicting shale gas production*. *Applied Energy*, 314, 118681.
- [8] Santur, Y. (2023) *A novel financial forecasting approach using deep learning framework*. *Computational Economics*.
- [9] Srijiranon, K., Lertratanakham, Y. and Tanantong, T. (2022) *A hybrid framework using PCA, EMD, and LSTM methods for stock market price prediction with sentiment analysis*. *Applied Sciences*, 12(21), 10823.
- [10] Yun, K.K., Yoon, S.W. and Won, D. (2023) *Interpretable stock price forecasting model using genetic algorithm-machine learning regressions and best feature subset selection*. *Expert Systems with Applications*, 212, 118450.
- [11] Fattah, J., et al. (2018) *Forecasting of demand using ARIMA model*. *International Journal of Engineering Business Management*, 10.
- [12] Li, Y., He, Y. and Zhang, M. (2020) *Prediction of Chinese energy structure based on convolutional neural network-long short-term memory (CNN-LSTM)*. *Energy Science & Engineering*, 8(8), 2680-2689.
- [13] Hauenstein, S., Wood, S.N. and Dormann, C.F. (2018) *Computing AIC for black-box models using generalized degrees of freedom: A comparison with cross-validation*. *Communications in Statistics-Simulation and Computation*, 47(5), 1382-1396.