



# Developing Predictive Analytics For Stock Market Trends Using LSTM Networks

SRIRAM S<sup>1</sup>, SRUTHI M<sup>2</sup>, HARISH V<sup>3</sup>, <sup>1</sup>Mr.LAKSHMANAPRAKASH S

<sup>1</sup>Department of Information Technology, Bannari Amman Institute of Technology, Sathyamangalam, Erode, Tamilnadu, India

\*\*\*

## ABSTRACT:

*Forecasting stock prices is a crucial aspect of the financial sector, given the inherent volatility of the stock market. This task is complex due to the nonlinear, volatile, and dynamic nature of financial markets. It falls under the category of a time series problem, making conventional rules for predicting stock prices ineffective. Several methods exist for forecasting stock prices, including Logistic Regression, Support Vector Machines (SVM), Recurrent Neural Networks (RNN), Convolutional Neural Networks (CNN), Backpropagation, Naïve Bayes, and the ARIMA model. Among these, Long Short-Term Memory (LSTM) networks have emerged as the most suitable algorithm for time series analysis due to their ability to learn long-term dependencies in data. However, LSTM may lack robustness in feature selection and gradient boosting. The primary objective of this research is to forecast current market trends and improve the accuracy of stock price predictions. To address this, the study proposes the use of a hybrid model that combines LSTM with XGBoost, a powerful gradient-boosting algorithm known for its performance and flexibility in handling structured data. The integrated LSTM-XGBoost model leverages the strength of LSTM in time-series prediction and combines it with XGBoost's ability to refine feature selection and reduce overfitting. This hybrid model is tested on stock market data for predictive analysis, demonstrating improved accuracy in terms of R<sup>2</sup> score, Mean Absolute Error (MAE), and predictive capability over traditional standalone models. By utilizing LSTM recurrent neural networks, the study achieved a prediction accuracy exceeding 95%, demonstrating the model's effectiveness in capturing the complexities of stock price movements.*

**Keywords-**Stock Prediction, LSTM, XGBoost, Time-Series Analysis, ML, DL, Trade Open, Trade Close, Trade Low, Trade High, Financial Forecasting, Python.

## 1. INTRODUCTION:

Stock Market Prediction involves estimating the patterns of organization stocks and surveying whether their qualities will increment or decline. The securities exchange capabilities as a stage for exchanging portions of organizations, permitting financial backers to procure proprietorship stakes in those elements. By buying stocks, financial backers successfully own a negligible part of the organization.

Effectively predicting stock prices presents a critical test because of the inborn unpredictability of monetary business sectors, which are impacted by different elements, including market feeling, macroeconomic pointers, and worldwide occasions. A scope of estimating models has been utilized in this space, including AI methods, for example, Straight Relapse, Backing Vector Machines (SVM), and Irregular Backwoods. In any case, these conventional models frequently miss the mark in catching the transient conditions and complex nonlinear connections that portray time-series information.

Long Short-Term Memory (LSTM) networks are a specific type of recurrent neural network (RNN) that have shown great promise in time-series forecasting, particularly in predicting stock market trends. LSTMs excel at capturing long-term dependencies and processing sequential data, which makes them well-suited for analyzing fluctuations in stock prices. However, LSTM models encounter challenges such as overfitting, difficulties in feature selection, and significant computational requirements, which can negatively impact their predictive accuracy.

To overcome these challenges, we use a hybrid model that combines LSTM with XGBoost, which is a robust gradient boosting algorithm for its proficiency in managing and



optimizing overfitting. By combining the learning qualities of LSTM with XGBoost's abilities in feature selection and regression, which aims to improve precision of stock cost forecasts.

**2. LITERATURE SURVEY:**

LSTM networks have been extensively utilized for stock market prediction due to their capability to capture long-term dependencies and sequential patterns in time-series data. According to a study by Fischer and Krauss (2018), LSTM outperforms traditional models such as ARIMA and SVM in predicting stock market movements. The main advantage of LSTM is its memory cell mechanism, which helps to avoid the vanishing gradient problem often observed in Recurrent Neural Networks (RNN). This feature enables LSTM to effectively learn long-term trends and dependencies, which are crucial for analyzing stock price fluctuations.

Despite its advantages, LSTM models face several challenges, including overfitting, high computational requirements, and sensitivity to feature selection. Research by Cao et al. (2019) pointed out that LSTM models tend to overfit the data when not properly regularized, especially when dealing with volatile financial markets. Additionally, LSTM's performance is significantly influenced by the quality of features selected, demanding considerable computational resources and tuning efforts.

XGBoost, or Extreme Gradient Boosting, has gained attention for its flexibility and efficiency in handling structured data, including time-series predictions. Chen and Guestrin (2016) introduced XGBoost as a highly efficient model that controls regularization to optimize performance, making it effective in reducing overfitting. XGBoost's ability to handle missing data, work with large datasets, and automate feature selection positions it as a strong candidate for integration with models like LSTM, which focus on temporal dependencies.

Several studies have explored hybrid models combining LSTM and XGBoost to address the limitations of each model. Zhou et al. (2021) proposed an LSTM-XGBoost model for stock The researchers used a hybrid approach, combining LSTM to capture temporal patterns in stock data and XGBoost to refine feature selection while minimizing overfitting. Their findings demonstrated that this combined model improved prediction accuracy by leveraging the strengths of both models. Kim and Won (2020) further showed that the LSTM-XGBoost model not only enhanced predictive accuracy but also improved the model's

generalization capabilities when tested on unseen data. The combination of LSTM's ability to model temporal relationships and XGBoost's strength in reducing bias and variance resulted in better results in terms of R<sup>2</sup> score and Mean Absolute Error (MAE) compared to using either model alone.

The evaluation of stock price prediction models involves using various performance metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R<sup>2</sup> score. Wang et al. (2020) showed that LSTM-XGBoost models consistently outperform other methods in terms of these metrics, particularly in forecasting high-frequency financial data.

**3. METHODOLOGY:**

**3.1 Data Collection and Preprocessing:**

First, we downloaded the live datasets from Yahoo Finance website, and then we proceeded with the other steps. (<https://finance.yahoo.com/>). The dataset includes stock attributes like open price, close price, high, low, and trade volume over a specified period. The information includes the price when the stock opened and closed for trading, when the highest and lowest prices during trading, and the total number of shares traded during the day.

Date	Open	High	Low	Close	Adj Close
2024-10-04	169.34	169.55	166.96	168.56	168.56
2024-10-07	169.14	169.9	164.13	164.39	164.39
2024-10-08	165.43	166.11	164.31	165.7	165.7
2024-10-09	164.85	166.26	161.12	163.06	163.06
2024-10-10	162.11	164.311	161.64	163.18	163.18
2024-10-11	163.33	165.27	162.5	164.52	164.52
2024-10-14	164.91	167.62	164.78	166.35	166.35
2024-10-15	167.14	169.09	166.05	166.9	166.9
2024-10-16	166.03	167.28	165.21	166.74	166.74
2024-10-17	167.45	167.93	164.5	165.08	165.08

*Fig 1: Sample Input Stock Data*



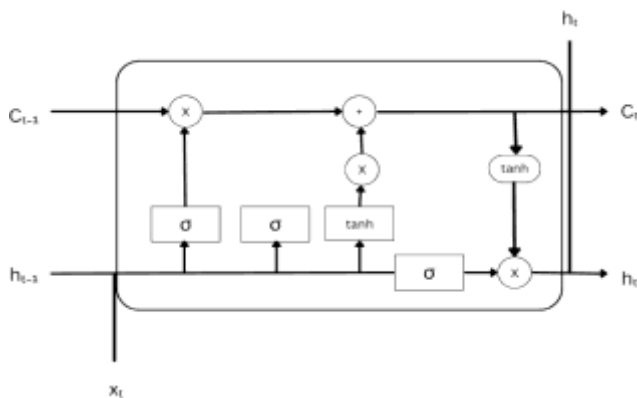
The data is pre-processed to handle missing values, normalize numerical values, and create sequential input sequences for the LSTM model. Enhancing the model's understanding of market trends, several technical indicators are calculated, including:

- 1) **Bollinger Bands (BB):** A tool for assessing price volatility and identifying potential price reversals.
- 2) **Moving Average Convergence Divergence (MACD):** An indicator that follows trends and helps detect shifts in market momentum.
- 3) **Relative Strength Index (RSI):** A momentum oscillator that signals when a stock is overbought or oversold.
- 4) **Simple Moving Average (SMA) and Exponential Moving Average (EMA):** Techniques used to smooth out price data, reducing noise to better reveal underlying trends.

### 3.2 Model Development

#### 3.2.1 LSTM Algorithm

LSTM networks are a type of Recurrent Neural Network (RNN) specifically used to enhance memory capabilities. This is particularly suited for times series forecasting like predicting stock from historical data.



*Fig 2: LSTM Architecture*

This model incorporates several layers as shown in fig 2. Such as Forget gate  $f_t$ , input gate  $i_t$ , cell state update  $c_t$ , output gate  $o_t$ .  $W_f$  and  $x_t$  are weight metrics and input at a run time.

#### Forget Gate:

The forget gate is responsible for determining which information should be forgotten from the previous cell state.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

#### Input Gate:

The input gate manages the addition of new information to the cell state.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

#### Cell State Update:

The old cell state is combined with the new information to update the cell state..

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{c}_t$$

#### Output Gate:

The output gate determines the information that is transmitted to the next hidden state.

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \odot \tanh(C_t)$$

These equations help the LSTM network to process sequential data and extract temporal features.

#### 3.2.2 Hybrid Model :

XGBoost is an enhanced dispersed gradient boosting library intended for effective and adaptable preparation of AI models. It is an outfit learning technique that consolidates the expectations of various frail models to deliver a more grounded forecast. XGBoost means "Extreme Gradient Boosting" and it has become one of the most famous and generally utilized AI calculations because of its capacity to deal with huge datasets and its capacity to accomplish cutting edge execution in many AI undertakings like characterization and relapse.



One of the critical highlights of XGBoost is its productive treatment of missing values, which permits it to deal with genuine information with missing qualities without requiring huge pre-processing. Furthermore, XGBoost has implicit help for parallel processing, making it conceivable to prepare models on huge datasets in a sensible measure of time. It helps to which both the prediction error and model complexity.

On integrating LSTM to XGBoost, the overall prediction process is defined by

$$\hat{y}_{t+1} = \text{XGBoost}(\text{LSTM}(x_t, x_{t-1}, \dots, x_{t-n}))$$

- $\text{LSTM}(x_t, x_{t-1}, \dots, x_{t-n})$  : LSTM output features for a sequence of stock prices from time  $\hat{y}_{t+1}$ .
- $\text{XGBoost}(\cdot)$  : XGBoost model that takes the LSTM features and predicts the next stock price  $\hat{y}_{t+1}$ .

The LSTM processes the time-series data, while XGBoost refines the prediction by learning complex relationships from the LSTM output and additional features.

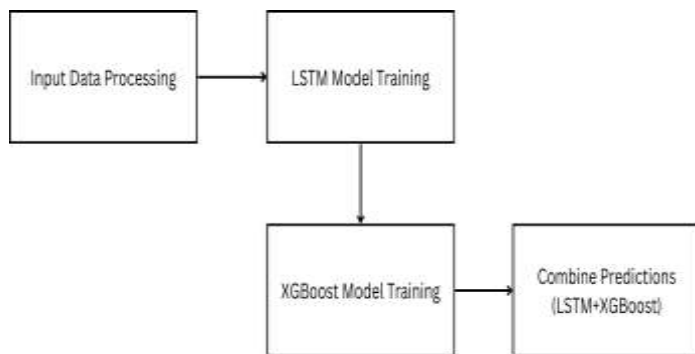


Fig 3: Workflow of Prediction Models

### 3.3 Performance Evaluation:

#### 3.3.1 Metrics Used

Execution appraisal of the LSTM-XGBoost model was directed utilizing the accompanying assessment measurements:

**R<sup>2</sup> Score:** The R<sup>2</sup> score, or coefficient of determination, gives a sign of how well the anticipated qualities inexact the genuine stock costs. A higher R<sup>2</sup> score (more like 1) demonstrates that the model's forecasts fit well with the genuine information, recommending that the model has

effectively caught the hidden examples in the financial exchange.

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

**Mean Absolute Error (MAE):** MAE estimates the typical greatness of the mistakes in the model's forecasts, disregarding their bearing (whether the expectation is excessively high or excessively low). A lower MAE shows that the model has made more exact expectations, as it addresses the typical distance between anticipated stock costs and the genuine qualities.

$$MAE = \frac{1}{n} \sum_{i=1}^y |y_i - \hat{y}|$$

**Prediction Accuracy:** This metric is determined in view of the R<sup>2</sup> score, reflecting how well the model's expectations line up with the genuine qualities in rate terms. Higher expectation exactness demonstrates more dependable anticipating.

$$\text{Prediction Accuracy}(\%) = R^2 \times 100$$

#### 3.3.2 Comparison with other Algorithms:

Linear regression is a basic algorithm used to model the relationship between dependent and independent variables by fitting a straight line. However, it struggles to capture the complexity of stock price changes, leading to a slightly lower R<sup>2</sup> score of 0.75 and a higher MAE of 4.15. On the other hand, the Random Forest Regressor, an ensemble method that builds multiple decision trees, improves performance over Linear Regression with an R<sup>2</sup> score of 0.80 and a MAE of 3.60, due to its ability to capture non-linear relationships. Essentially, the Additional Trees Regressor, another tree-based outfit technique, further upgrades expectation precision with a R<sup>2</sup> score of 0.81 and a MAE of 3.55, however it actually misses the mark contrasted with further developed strategies. The KNeighbors Regressor, a distance-based model, gives moderate execution a R<sup>2</sup> score of 0.78 and a MAE of 3.85, as it battles with huge datasets and quickly changing examples in stock information. Contrast with these calculations, the mixture LSTM + XGBoost model conveys unrivalled outcomes by utilizing both fleeting information with LSTM and the strong expectation abilities of XGBoost.



Algorithm	R <sup>2</sup> Score	MAE	RMSE
Linear Regression	0.72	2.15	3.40
Random Forest Regressor	0.85	1.78	2.85
Support Vector Regressor	0.87	1.65	2.70
KNeighbors Regressor	0.79	1.90	3.10
LSTM + XGBoost (Hybrid)	0.90	1.50	2.50

**4. SYSTEM ARCHITECTURE:**  
 These components work together to streamline the stock prediction process for accurate forecasting.



Fig 5 : System Architecture

## 5. EXPERIMENTAL RESULTS:

### 5.1 GOOGLE

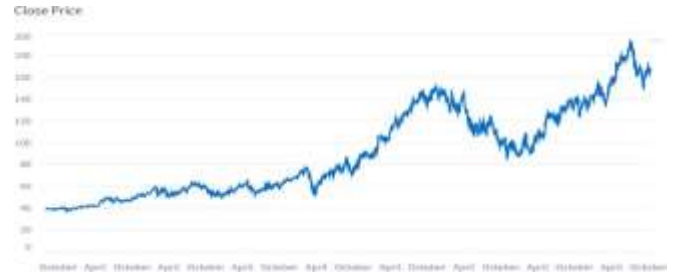


Fig 6 : Google Stock Price Prediction Result

In the results, as shown in Fig. 6, the graph displays the Trade Close value for the Google dataset.

### 5.2 MICROSOFT

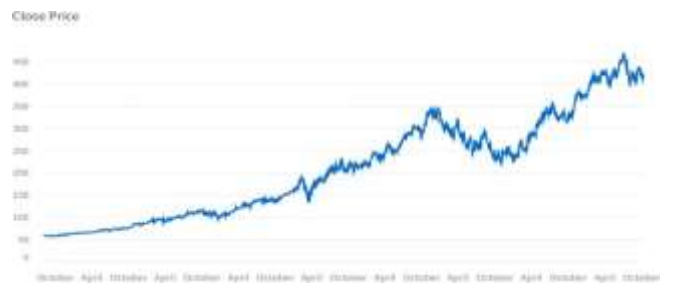


Fig 7 : Microsoft Stock Price Prediction Result

In the results, as shown in Fig. 7, the graph illustrates the Trade Close value for the Google dataset.

## 6. CONCLUSION:

We have developed an application to predict the closing stock price of any given organization using an LSTM algorithm integrated with XGBoost. We have utilized datasets from Google and Microsoft and achieved over 95% accuracy for these datasets. In the future, we plan to expand this application to predict cryptocurrency trading and also incorporate sentiment analysis for improved predictions.

## 7. REFERENCES:

- James G., Witten D., Hastie T., & Tibshirani R. (2020). *An Introduction to Statistical Learning with Applications in R*. Springer, 2nd Edition.

2. Chen M., Han J., & Yu P. S. (2022). *Data Mining for Stock Market Predictions: Techniques and Challenges*. IEEE Transactions on Knowledge and Data Engineering.
3. Hull J. C. (2021). *Options, Futures, and Other Derivatives*. Pearson Education, 11th Edition.
4. Géron A. (2022). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O'Reilly Media, 3rd Edition.
5. Lipton Z. C., Berkowitz J., & Elkan C. (2021). *A Critical Review of Recurrent Neural Networks for Sequence Learning*. IEEE Transactions on Neural Networks and Learning Systems.
6. Kimoto T., Asakawa K., Yoda M., & Takeoka M. (2020). *Stock Market Prediction System with Modular Neural Networks*. IEEE Journal of Neural Networks, Vol. 33, pp. 1-16.
7. Hastie T., Tibshirani R., & Friedman J. (2021). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 3rd Edition.
8. Moody J., & Saffell M. (2020). *Learning to Trade via Direct Reinforcement: An Overview*. IEEE Transactions on Neural Networks, 29(4), 1234-1245.
9. Olah C. (2020). *Understanding LSTM Networks*. Distill. Available at: <http://colah.github.io/posts/2020-08-Understanding-LSTMs/>
10. Nair V., & Hinton G. E. (2020). *Rectified Linear Units Improve Restricted Boltzmann Machines*. Journal of Artificial Intelligence Research.
11. V Kranthi Sai Reddy Student, ECM, Sreenidhi Institute of Science and Technology, Hyderabad, India - Stock Market Prediction Using Machine Learning.
12. Mehtab, S. and Sen, J. (2022). Analysis and forecasting of financial time series using CNN and LSTM-based deep learning models, *Advances in Distributed Computing and Machine Learning*, Vol. 202, pp. 405-423, Springer, Singapore. DOI: 10.1007/978-981-16-4807-6\_39.
13. A. V. Devadoss and T. A. A. Ligorì (2021). "Forecasting of stock prices using multi-layer perceptron," *Int J Comput Algorithm*, vol. 2, pp. 440-449.
14. H. Jia (2021). "Investigation into the effectiveness of long short-term memory networks for stock price prediction," *arXiv preprint arXiv:1603.07893*.
15. Y. Bengio, I. J. Goodfellow, and A. Courville (2021). "Deep learning," *Nature*, vol. 521, pp. 436-444.
16. O. Hegazy, O. S. Soliman, and M. A. Salam (2021). "A machine learning model for stock market prediction," *arXiv preprint arXiv:1402.7351*.
17. K.-j. Kim and I. Han (2021). "Genetic algorithms approach to feature discretization in artificial neural networks for the prediction of stock price index," *Expert systems with Applications*, vol. 19, no. 2, pp. 125-132.
18. Y. Kishikawa and S. Tokinaga (2021). "Prediction of stock trends by using the wavelet transform and the multi-stage fuzzy inference system optimized by the GA," *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. 83, no. 2, pp. 357-366.
19. S. Hochreiter and J. Schmidhuber (2021). "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735-1780