

Volume: 09 Issue: 07 | July - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

A Review on Statistical Machine Learning Models for Stock Market Forecasting

Gaurav Mourya¹, Prof. Preetish Kshirsagar²

Abstract: Machine learning is transforming stock market prediction by leveraging vast datasets, advanced algorithms, and computational power to provide more accurate forecasts. While challenges remain, continuous advancements in artificial intelligence and deep learning are improving predictive models, making them an essential tool for traders and investors. Stock market prediction extremely challenging due to the dependence of stock prices on several financial, socio-economic and political parameters etc. For real life applications utilizing stock market data, it is necessary to predict stock market data with low errors and high accuracy. This needs design of appropriate artificial intelligence (AI) and machine learning (ML) based techniques which can analyze large and complex data sets pertaining to stock markets and forecast future prices and trends in stock prices with relatively high accuracy. This paper presents a comprehensive review on the various techniques used in recent contemporary papers for stock market forecasting.

Keywords: Time Series Models, Stock Market Forecasting, Artificial Intelligence, Artificial Neural Networks, Forecasting accuracy.

I. INTRODUCTION

Stock prices are influenced by numerous factors, including company performance, macroeconomic indicators, geopolitical events, and investor sentiment. The interplay of these factors makes stock markets highly nonlinear and difficult to predict using conventional statistical models. Machine learning can process vast amounts of structured and unstructured data, including news articles, social media trends, and financial reports, to extract meaningful insights that drive stock prices. Stock market prediction is fundamentally a regression problem in which patterns in previous data and its associated variables need to be found. Stock prices can be mathematically modelled as a time series function as [1]:

Prices = function (time, variables)

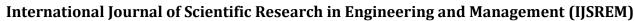
Mathematical modeling is a fundamental approach used to analyze and predict stock market trends. These models help traders, analysts, and investors make informed decisions based on historical data and market patterns. However, due to the dynamic and unpredictable nature of financial markets, developing accurate and reliable models is highly challenging. Several factors, including market volatility, data limitations, and external influences, create significant obstacles in constructing precise mathematical models for stock market prediction.

The stock prices depend both on time as well as associated variables and finding patterns in among the variables aid forecasting future stock prices which is often termed as stock market forecasting. Stock market prediction extremely challenging due to the dependence of stock prices on several financial, socio-economic and political parameters etc. For real life applications utilizing stock market data, it is necessary to predict stock market data with low errors and high accuracy. This needs design of appropriate artificial intelligence (AI) and machine learning (ML) based techniques which can analyze large and complex data sets pertaining to stock markets and forecast future prices and trends in stock prices with relatively high accuracy.

II. MACHINE LEARNING MODELS FOR STOCK MARKET FORECASTING

Machine Learning models are employed for data analysis where the data to be analyzed is extremely large and complex or both. Primarily, Artificial Intelligence and Machine Learning (AI and ML) have been extensively used for financial and business applications where large data has to be analyzed. One such major area is investment banking [2].

Supervised learning algorithms like regression models and neural networks can learn from past stock price movements and make predictions based on historical trends. Unsupervised learning methods, such as clustering, help identify patterns among stocks with similar behaviors. Additionally, reinforcement learning



IJSREM e-Journal

Volume: 09 Issue: 07 | July - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

can be used to optimize trading strategies by continuously learning from market conditions and adjusting decisions accordingly.

Additionally, Deep learning techniques, particularly recurrent neural networks (RNNs) and long short-term memory (LSTM) networks, have shown promise in time-series forecasting, making them valuable for predicting stock prices. Moreover, natural language processing (NLP) enables sentiment analysis, allowing traders to assess the impact of news, earnings reports, and social media discussions on stock movements.

Despite its potential, machine learning in stock market prediction faces several challenges. Market conditions change rapidly, and past trends may not always be indicative of future performance. Overfitting, where a model performs well on training data but poorly on realworld data, remains a major concern. Additionally, external factors like economic crises, political instability, and black swan events can disrupt even the most sophisticated models.

III. LITERATURE REVIEW

The contemporary work in the domain and the noteworthy contributions is cited in this section.

Subakkar et al. [3] proposed proposed that the autoregressive integrated moving average (ARIMA) is a type of moving average model that is used to for building a stock price forecasting model. In this paper, ARIMA models have been used to predict the stock price of daily trading stock prices published in the Bombay Stock Exchange (BSE) and National Stock Exchange (NSE). ARIMA model is analyzed for time series prediction, and the results obtained using this model show strong accuracy for short term and daily stock prediction, and this engages with other model for predicting stock price.

Soun et al. [4] proposed SLOT (Self-supervised Learning of Tweets for Capturing Multi-level Price Trends), an accurate method for stock movement prediction. SLOT has two main ideas to address the limitations of previous tweet-based models. First, SLOT learns embedding vectors of stocks and tweets in the same semantic space through self-supervised learning. The embeddings allow us to use all available tweets to improve the prediction for even unpopular stocks, addressing the sparsity problem. Second, SLOT learns multi-level relationships between stocks from tweets, rather than using them as direct evidence for prediction, making it robust to the unreliability of tweets. Extensive experiments on real world datasets show that SLOT

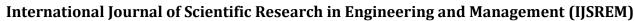
provides the state-of-the-art accuracy of stock movement prediction.

Sen et al. [5] proposed volatility models based on generalized autoregressive conditional heteroscedasticity (GARCH) framework for modeling the volatility of ten stocks listed in the national stock exchange (NSE) of India. The stocks are selected from the auto sector and the banking sector of the Indian economy, and they have a significant impact on the sectoral index of their respective sectors in the NSE. The historical stock price records from Jan 1, 2010, to Apr 30, 2021, are scraped from the Yahoo Finance website using the DataReader API of the Pandas module in the Python programming language. The GARCH modules are built and fine-tuned on the training data and then tested on the out-of-sample data to evaluate the performance of the models. The analysis of the results shows that asymmetric GARCH models yield more accurate forecasts on the future volatility of stocks.

S. Kim et al. [6] developed a technique termed as effective transfer entropy (ETE) to be used in conjugation with existing ML algorithms such as LR, MLP, LSTM etc. The ETE metric served as an exogenous feature which helped the training performance of the standard training models based on the entropy measure of the data set which is a stochastic variable of the training data set, while the data set used for the study was the US stick market dataset.

B. Bouktif et al. [7] designed an N-grams based approach utilizing the semantic analysis of data related to stock movement for the prediction problem. The sentiment polarity was utilized to predict the impact of the users of different social media platforms on the stock prices. The polarities used were positive, negative and neutral which served as tokenized impacts on the feature values of the dataset.

X.Li et al. in [8] devised a deep learning model employing sentiment analysis results to predict the stock market behaviour. The individual and cumulative impact of the sentiment features were used for designing the sentiment vector for the forecasting model. Textual normalization and opinion mining techniques were incorporated as features to gauge the sentiments of the common public regarding reputations of the firms since previous prices alone may not always render the moving trends in the market



IJSREM

Volume: 09 Issue: 07 | July - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

Gaurang Bansal et al. [9] proposed a de-centralized forecasting model incorporating block chain which acts as a distributed ledger for stock market behaviours. The use of block chain was done to relate the variables or features for training the system model to find trends or visible patterns in the data blocks. The performance evaluation of the system was done in terms of the accuracy of prediction.

Jithin Eapen et al. [10] proposed a pipeline approach of CNNs along with a bi-directional LSTM model. The authors were able to gain significant performance improvement in the prediction accuracy of the system using the pipeline CNN model as compared to the baseline regression models for the same S & P dataset. The bidirectional LSTM model was also tested for prediction accuracy of the system for the same data base. It was shown that the pipelined CNN based approach outperformed the conventional techniques.

Min Wen et al. [11] proposed a stacked CNN based approach for the analysis of noisy time series data for stock market behavioural patterns. The stacked CNN structure was able to extract different levels of features for the different layers of the data set. The proposed technique was shown to perform better than existing techniques for temporal stock market behavioural data patterns.

Y Guo et al. [12] proposed a modified version of the support vector regression (SVR) model with weight updating mechanism based on evolutionary algorithms such as the PSO. The inclusion of PSO helped in finding the local and global best feature values while optimizing the objective function simultaneously. It was shown that the proposed approach could outperform the existing regression or backpropagation models.

MS Raimundo et al. [13] proposed a technique that was the amalgamation of the wavelet transform and the support vector regression. The technique used the DWT as a data processing tools and retained the approximate co-efficient values of the multi-tree level DWT analysis of the raw data thereby enabling more noise immunity for the SVR algorithm. The DWT-SVR hybrid was

shown to perform better in terms of performance accuracy w.r.t. to SVR alone.

Y Baek et al. [14] designed a deep neural network named ModAugNet. The deep neural network was again an amalgamation of two LSTM layers. The first layer avoided the chances of overfitting while the second LSTM block was used purely for prediction. The approach was novel in the sense that a similar network with different Hyperparameters were used for the optimization and prediction purposes.

S.Selving et al. [15] utilized different data fitting algorithms for stock movement estimates. The data fitting approaches utilized were both linear and nonlinear in approach such as ARIMA, GARCH etc. The exogenous input feature vector was the closing price of the day which served as a separate feature value. The performance for the same data set with and without the closing price as an exogenous input was tested on the system performance.

Z. Zhao et al. [16] developed an approach which utilized different time-weighted feature vectors to train an LSTM neural net. The essence of the proposed approach was the fact that recent time or temporal sample has different weighted values compared to the generalized weighted values of a normal feature vector. The performance of the system was evaluated in terms of the accuracy achieved. The designed system achieved an accuracy score of 83.91% while feeding the system with refined feature values.

Nelson et al. [17] proposed an LSTM based model for stock market prediction along with technical analysis indicators. The fundamental approach of the system was to find the correlation among different variables for stock market movement. The price indicators of a particular company in a specific stock market were linked to the stock market in other stock markets listed globally. The effect of closing prices of one stock in a particular stock exchange was linked to the opening price of the same stock in some other stock exchange. Thus along with the historical data, the correlation among other variables was also evaluated.



International Journal of Scientific Research in Engineering and Management (IJSREM)

Volume: 09 Issue: 07 | July - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

Billah et al. in [18] designed a back-propagation based neural network training algorithm with data structuring. The Levenberg Marquardt (LM) weight updating rule was used to forecast closing stock prices of stocks for the Dhaka Stock Exchange. It was shown that the LM algorithm needed lesser memory consumption as well as iterations compared to conventional neural networks and ANFIS systems. The performance evaluation metric was the accuracy of the system.

Sadei et al in [19] proposed a fuzzy time series predictor based on the concept of fuzzy expert systems. The fuzzy set creation based on temporal data was dine followed by the design of membership functions. Finally, the fuzzy relationships were computed and the defuzzification block was used to predict future trends in the stock prices. Different membership functions were used for the purpose of designing the fuzzy sets.

.Lincy et al. in [20] proposed a model based on multiple fuzzy inference systems and applied it to the NASDAQ stock exchange data. The proposed system was pitted against the conventional ANFIS systems and it was shown that the proposed system outperformed the conventional techniques based on expert systems.



International Journal of Scientific Research in Engineering and Management (IJSREM)

Volume: 09 Issue: 07 | July - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

Table.1 Summary of Previous Work

S.No.	Authors	Findings	Research Gaps
		0	•
	~		
1.	Subbakar et al.	ARIMA based regression model for	No data pre-processing and
		forecasting Tesla stock prices.	optimization
2.	Soun et al.	Stock Movement Prediction with Self-	No feature optimization for the
		supervised Learning from Sparse Noisy	LSTM model resulting in
2	Sen et al.	GARCH model for modelling of	Data optimization or Deep Nets
		volatility movement of Indian stock.	not employed.
1	Kim et al.	Effective transfer entropy (ETE) used in	Lack of Data Optimization and
		conjugation with existing ML	Dimensional Reduction.
2	Bouktif et al.	Opinion Mining and Sentiment Analysis	No data filtering or optimization
		was used along with historical stock	approach used.
3	Li et al.	Textual normalization and opinion	No data filtering approach
		mining techniques were incorporated as	
4	Bansal et al.	Model proposed de-centralized	No data optimization
		forecasting model incorporating block	
5	Eapen et al.	A pipelined approach of CNNs along	Dimensional reduction and data
		with a bi-directional LSTM model.	optimization not used. Opinion
6	Wen at al.	Stacked CNN based approach for the	No estimation of overfitting for
		analysis of noisy time series data for	the CNN model.
7	Guo et al.	Adaptive Support vector regression	Opinion mining and filtering of
		(SVR) model with weight updating	data not employed. Performance
8	Raimundo et	Combination of discrete wavelet	SVR'r performance doesn't
	al.	transform (DWT) and the support vector	improve above a threshold.
9	Baek et al.	An amalgamation of two LSTM layers	No data optimization and
		performed. The first layer avoided the	sentiment analysis data used.
10	Selving et al.	The approach used made predictions	Only daily closing price chosen as
		based on the daily closing price. The	the time dependent feature.
11	Zhao et al.	Proposed model used time-weighted	No estimates of overfitting or
		feature vectors to train an LSTM neural	vanishing gradient.
	1	<u>L</u>	

Volume: 09 Issue: 07 | July - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

IV. PERFORMACNE METRICS

The parameters which can be used to evaluate the performance of the ANN design for time series models is given by [21]:

- 1) Mean Absolute Error (MAE)
- 2) Mean Absolute Percentage Error (MAPE)
- 3) Mean square error (MSE)

The above mentioned errors are mathematically expressed as:

$$MAE = \frac{1}{N} \sum_{t=1}^{N} |V_t - \widehat{V}_t|$$

$$MAE = \frac{1}{N} \sum_{t=1}^{N} |e_t|$$

$$MAPE = \frac{100}{N} \sum_{t=1}^{N} \frac{|V_t - \hat{V}_t|}{V_t}$$

$$MSE = \frac{1}{N} \sum_{t=1}^{N} e_t^2$$

Here,

N is the number of predicted samples V is the predicted value \hat{V}_t is the actual value e is the error value

CONCLUSION

This paper presents a comprehensive review and taxonomy on the use of machine learning based approaches for stock market prediction or forecasting. As the data to be analysed is extremely large and complex, hence it is mandatory to employ machine learning for regression analysis. The multiple machine learning and deep learning approaches used in contemporary work have been cited and the related research gaps have been identified. It is expected that this paper puts future research directions in better stead with an aim to enhance the forecasting accuracy.

References

- 1. Martin T. Hagan, Howard B. Demuth, Mark H. Beale, Orlando De Jesus, "Neural Network Design", 2nd edition, Cengage Publications.
- 2. Shalev-Shwartz, Shai, Ben-David, "Understanding Machine Learning: From Theory to Algorithms", Cambridge University Press.

- 3. A Subakkar, S Graceline Jasmine, L Jani Anbarasi, J Ganesh, CM Yuktha, "An Analysis on Tesla's Stock Price Forecasting Using ARIMA Model", Proceedings of the International Conference on Cognitive and Intelligent Computing, Springer, 2023, pp 83–89.
- 4. Y. Soun, J. Yoo, M. Cho, J. Jeon and U. Kang, "Accurate Stock Movement Prediction with Self-supervised Learning from Sparse Noisy Tweets," 2022 IEEE International Conference on Big Data (Big Data), Osaka, Japan, 2022, pp. 1691-1700.
- 5. J. Sen, S. Mehtab and A. Dutta, "Volatility Modeling of Stocks from Selected Sectors of the Indian Economy Using GARCH," IEEE Access, 2021, pp. 1-9.
- 6. S Kim, S Ku, W Chang, JW Song, "Predicting the Direction of US Stock Prices Using Effective Transfer Entropy and Machine Learning Techniques", IEEE Access 2022, Vol-8, pp. 111660 111682.
- 7. S Bouktif, A Fiaz, M Awad, Amir Mosavi, "Augmented Textual Features-Based Stock Market Prediction", IEEE Access 2021, Volume-8, PP: 40269 40282.
- 8. X Li, P Wu, W Wang, "Incorporating stock prices and news sentiments for stock market prediction: A case of Hong Kong", Information Processing & Management, Elsevier 2020.

Volume 57, Issue 5, pp: 1-19.

- 9. Gaurang Bansal; Vikas Hasija; Vinay Chamola; Neeraj Kumar; Mohsen Guizani, "Smart Stock Exchange Market: A Secure Predictive Decentralized Model", 2019 IEEE Global Communications Conference (GLOBECOM), IEEE 2019 pp. 1-6.
- 10. Jithin Eapen; Doina Bein; Abhishek Verma, "Novel Deep Learning Model with CNN and Bi-Directional LSTM for Improved Stock Market Index Prediction", 2019 IEEE 9th Annual Computing and Communication Workshop and Conference (CCWC), IEEE 2019 pp. 0264-0270.
- 11. Min Wen; Ping Li; Lingfei Zhang; Yan Chen, "Stock Market Trend Prediction Using High-Order Information of Time Series", IEEE Access 2019, Volume 7, pp: 28299 28308.
- 12. Y Guo, S Han, C Shen, Y Li, X Yin, Y Bai, "An adaptive SVR for high-frequency stock price forecasting", Volume-6, IEEE Access 2018, pp: 11397 11404.
- 13. MS Raimundo, J Okamoto, "SVR-wavelet adaptive model for forecasting financial time series", 2018 International Conference on Information and Computer Technologies (ICICT), IEEE 2018, pp. 111-114.

International Journal of Scientific Research in Engineering and Management (IJSREM)



Volume: 09 Issue: 07 | July - 2025

SJIF Rating: 8.586

ISSN: 2582-3930

- 14. Y Baek, HY Kim, "ModAugNet: A new forecasting framework for stock market index value with
- Applications, Elsevier 2018, Volule-113, pp: 457-480.

 15. S Selvin, R Vinayakumar, E. A Gopalakrishnan; Vijay Krishna Menon; K. P. Soman, "Stock price prediction using LSTM, RNN and CNN-sliding window model", 2017 International Conference on Advances in Computing, Communications and Informatics (ICACCI), IEEE 2017, pp. 1643-1647.

an overfitting prevention LSTM module and a prediction LSTM module" Journal of Expert System and

- 16. Z Zhao, R Rao, S Tu, J Shi, "Time-weighted LSTM model with redefined labeling for stock trend prediction", 2017 IEEE 29th International Conference on Tools with Artificial Intelligence (ICTAI), pp. 1210-1217.
- 17. DMQ Nelson, ACM Pereira, Renato A. de Oliveira, "Stock market's price movement prediction with LSTM neural networks", 2017 International Joint Conference on Neural Networks (IJCNN), IEEE 2017, pp. 1419-1426
- 18. M Billah, S Waheed, A Hanifa, "Stock market prediction using an improved training algorithm of neural network", 2016 2nd International Conference on Electrical, Computer & Telecommunication Engineering (ICECTE), IEEE 2016, pp. 1-4,

doi: 10.1109/ICECTE.2016.7879611

- 19. HJ Sadaei, R Enayatifar, MH Lee, M Mahmud, "A hybrid model based on differential fuzzy logic relationships and imperialist competitive algorithm for stock market forecasting", Journal of Applied Soft Computing, Elsevier 2016, Volume 40, pp. 132-149.
- 20. GRM Lincy, CJ John, "A multiple fuzzy inference systems framework for daily stock trading with application to NASDAQ stock exchange", Journal of Expert Systems with Applications, Volume-44, Issue-C, ACM 2016.
- 21. M. Wen, P. Li, L. Zhang and Y. Chen, "Stock Market Trend Prediction Using High-Order Information of Time Series," in IEEE Access 2019, vol. 7, pp. 28299-28308.