

# **Stock Price Prediction Using NIFTY50**

A Project Report Submitted in Partial Fulfillment

of

The Bachelor of Technology Degree

In

Information Technology

Submitted By

Alavya Singh

(CSJMA20001390136)

Sejal Yadav

(CSJMA20001390180)

Suyash Pandey

(CSJMA20001390190)



**University Institute of Engineering & Technology**

**Chhatrapati Shahu Ji Maharaj University, Kanpur**

**April 2024**

# **Stock Price Prediction Using NIFTY50**

A Project Report Submitted in Partial Fulfillment

of

The Bachelor of Technology Degree

In

Information Technology

Submitted By

Alavya Singh

(CSJMA20001390136)

Sejal Yadav

(CSJMA20001390180)

Suyash Pandey

(CSJMA20001390190)



**University Institute of Engineering & Technology**

**Chhatrapati Shahu Ji Maharaj University, Kanpur**

**April 2024**

**Er. Prateek Srivastava**

Project Guide

**Dr Alok Kumar**

H.O.D

# CERTIFICATE



## University Institute of Engineering and Technology

I hereby certify that the work which is being presented in the B.Tech. Project Report entitled **"Stock Price Prediction Using NIFTY50"**, in partial fulfillment of the requirements for the Award of the Bachelor of Technology in Information Technology and submitted to the Department of Information Technology of the **University Institute of Engineering & Technology**, Kanpur UP is an authentic record of my work carried out during a period from July 2023 to April 2024 under the supervision of **Er. Prateek Srivastava** IT Department.

**We have not submitted the matter presented in this thesis for the award of any other degree elsewhere.**

Signature of Candidate  
Alavya Singh

Sejal Yadav

Suyash Pandey

**This is to certify that the above statement made by the candidate is correct to the best of my Knowledge.**

Signature of Project Guide  
**Er. Prateek Srivastava**

Signature of H.O.D  
**Dr Alok Kumar**

## Acknowledgment

First and foremost, we would like to thank our guides, Dr Alok Kumar and Er Prateek Srivastava for having suggested the topic of our project and for their constant support and guidance, without which we would not have been able to attempt this project.

Thanks to Er Prateek Srivastava Sir for being a constant help throughout the work for our current research. Furthermore, we acknowledge the contributions of researchers, scholars, and practitioners whose pioneering work has laid the foundation for this study. The amount of knowledge and insights shared through academic literature, research papers, and online resources have greatly enriched our understanding of stock market prediction methodologies.

Lastly, we wish to express our heartfelt thanks to our friends and family for their unwavering support, understanding, and encouragement throughout this academic journey. Their steadfast belief in our abilities & unwavering encouragement has served as a constant source of motivation.

Thank You.

Alavya Singh  
CSJMA20001390136

Sejal Yadav  
CSJMA20001390180

Suyash Pandey  
CSJMA20001390190

## Table Of Contents

Sno.	Topic	Page No
1	Introduction and Background	
2	Problem Statement	
3	Objective	
4	Literature Review	
5	Methodology	
6	Experimental Results	
7	Conclusion	
8	References	

## Introduction and Background

The stock market is viewed as an unpredictable, volatile, and competitive market. The prediction of stock prices has been a challenging task for many years. In fact, many analysts are highly interested in the research area of stock price prediction. Various forecasting methods can be categorized into linear and non-linear algorithms. In this paper, we offer an overview of the use of deep learning networks for the Indian National Stock Exchange time series analysis and prediction.

The networks used are Long Short-Term Memory Network (LSTM) and Convolutional Neural Networks (CNN) to predict future trends of NIFTY 50 stock prices. Comparative analysis is done using different evaluation metrics. Their analysis led us to identify the impact of feature selection process & hyper-parameter optimization on prediction quality and metrics used in the prediction of stock market performance and prices.

Additionally, we have employed a variety of traditional forecasting algorithms, including k-Nearest Neighbors (KNN), Holt Linear, Support Vector Regression (SVR), Ordinary Least Squares (OLS), Random Forest, and Autoregressive Integrated Moving Average (ARIMA) models. These algorithms complement the deep learning networks and provide a comprehensive analysis of stock price prediction techniques.

The boosting techniques used are XG Boosting, Gradient Boosting, and Random Forest. By implementing these strategies and leveraging the flexibility and power of gradient boosting, significant improvements in predictive performance can be achieved, potentially surpassing previous research papers on the dataset. Adjusting the batch size along with other hyperparameters can contribute to better convergence and overall model performance.

The performance of the models was quantified using MSE metric. The study pioneers the fusion of time series, econometric, statistical, and learning-based techniques, expertly applied to predict stock prices across major sectors with exceptional precision. It also delves into the realm of classification models, leveraging the same multivariate data to predict whether stock prices will rise or fall the following day.

In summary, our research navigates the intersection of empirical insights and computational finesse, shedding light on the intricate path of stock price prediction. In a world where informed decisions hinge on such forecasts, our study underscores the potential of merging diverse modeling approaches to decode the intricate trajectory of stock price movements.

## **Problem Statement**

The research aims to develop a robust predictive model for forecasting stock prices of companies enlisted in the NIFTY50 index. This endeavor faces the challenge of effectively leveraging a significant historical dataset spanning, from Jan 2017 to Dec 2023. The primary objective is to generate accurate and reliable predictions for leading organizations within selected sectors based on this extensive historical data.

Moreover, the incorporation of NIFTY50 index data supplements the analysis by providing a comprehensive view of market sentiment. Utilizing Long Short-Term Memory Network (LSTM) and Convolutional Neural Networks (CNN), the study endeavors to predict future trends in NIFTY50 stock prices. A critical aspect of the research involves conducting a comparative analysis using various evaluation metrics. This analysis aims to discern the impact of feature selection processes and hyper-parameter optimization on prediction quality. By exploring different metrics employed in predicting stock market performance and prices, the study seeks to offer valuable insights into the effectiveness of these methodologies.

The research extends to the inclusion of traditional forecasting algorithms such as k-Nearest Neighbors (KNN), Holt Linear, Support Vector Regression (SVR), Ordinary Least Squares (OLS), Random Forest, Autoregressive Integrated Moving Average (ARIMA), XG Boost, Gradient Boosting, and other pertinent models. These conventional models are integrated to provide a comprehensive analysis of stock price prediction techniques, complementing the deep learning networks utilized in the study.

Furthermore, stock price data for NIFTY50, a benchmark index, is also incorporated to capture an encompassing market sentiment.

## Objectives

The primary objectives of this research are:

**Model Development:** Construct an array of predictive models spanning the realms of time series analysis, econometrics, statistics, machine learning, and deep learning. This collection of models will serve as tools for forecasting stock prices of the ten prominent organizations in each of the specified sectors, as well as for predicting NIFTY50's overall sentiment.

**Robust Prediction:** Harness the power of the diverse range of models to ensure robust and accurate predictions of stock prices. By exploring a wide spectrum of modeling techniques, we aim to capture complex patterns, trends, and relationships inherent in the historical data.

**Comparison with NIFTY50:** Compare the predictions of the Organizations with the overall market sentiment reflected by NIFTY50's stock price predictions. This comparative analysis will provide insights into how the organizations' stock prices align with the broader market trend.

A critical aspect of the research involves conducting a comparative analysis using various evaluation metrics. This analysis aims to discern the impact of feature selection processes and hyper-parameter optimization on prediction quality. By exploring different metrics employed in predicting stock market performance and prices, the study seeks to offer valuable insights into the effectiveness of these methodologies. Furthermore, the research extends to the inclusion of traditional forecasting algorithms such as k-Nearest Neighbors (KNN), Holt Linear, Support Vector Regression (SVR), Ordinary Least Squares (OLS), Random Forest, Autoregressive Integrated Moving Average (ARIMA), XG Boost, Gradient Boosting, and other pertinent models. These conventional models are integrated to provide a comprehensive analysis of stock price prediction techniques, complementing the deep learning networks utilized in the study.



# Literature Review

A comparative study of state-of-the-art of Machine learning & deep learning model for stock market prediction.

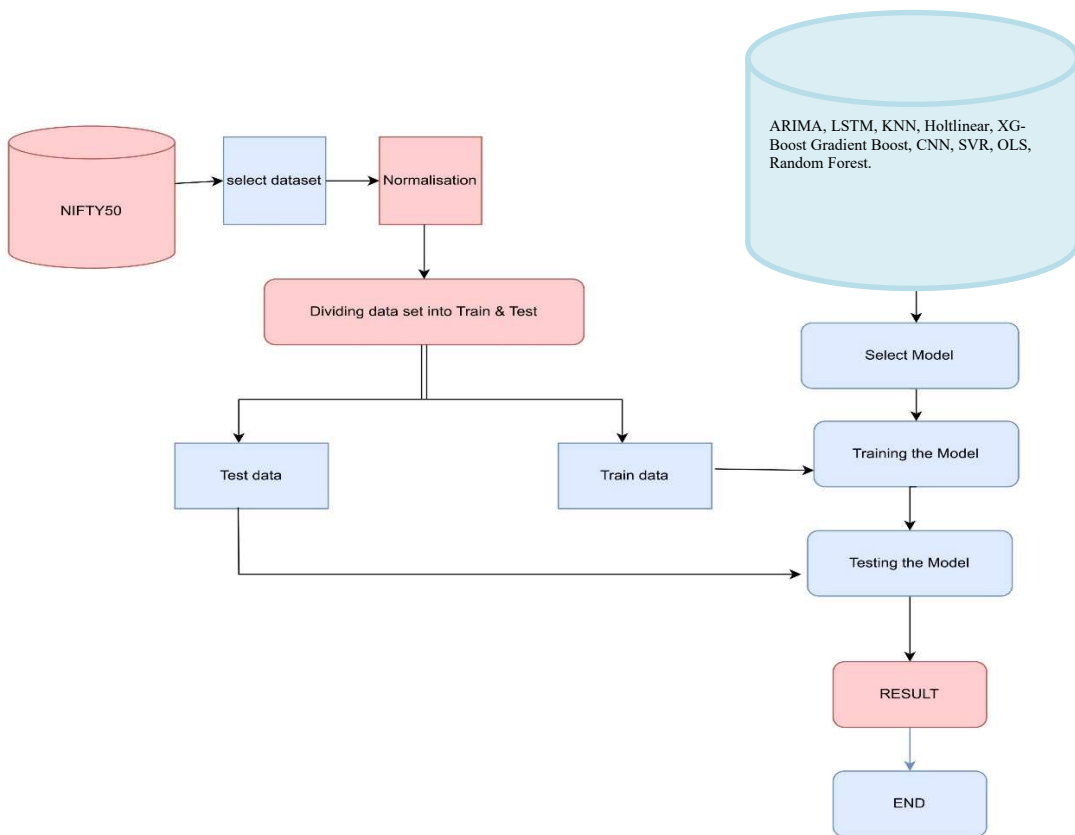
References	Dataset	Prediction Method /Technique Used LSTM CNN KNN ARIMA OLS HoltLinear XG Gradient RandomForest										Referred Dataset Results (RMSE)	Our Results (RMSE)
<a href="#">1</a>	Nifty 50	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	ARIMA:0.11 LSTM:0.012 CNN:0.0120 KNN:0.01 OLS:0.0108 Gradboost:0.009147 Holt linear:0.2069 XGBoost:0.0193	ARIMA:0.005718 LSTM:0.011461 KNN:0.0120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553
<a href="#">2</a>	Nifty 50	✓	✓		✓		✓				✓	ARIMA:0.11 LSTM:0.012 CNN: - Holt linear:0.04	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553

<a href="#">3</a>	Nifty 50	✓	✓	✓	✓	ARIMA:0.11 LSTM:0.037279 KNN:0.7120994 Holtlinear:0.0566	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553
<a href="#">4</a>	Nifty 50	✓	✓	✓		LSTM:0.012 CNN:0.001 ARIMA:0.121	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553
<a href="#">5</a>	Nifty 50 Index	✓	✓			LSTM:0.012 KNN:0.87002	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553

<a href="#">6</a>	Nifty50	✓ ✓ ✓	LSTM:0.0012 Random forest:0.0278 CNN:0.6547	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553
<a href="#">7</a>	Nifty 50 Index	✓ ✓ ✓	LSTM:0.0102 Random forest:0.014753 CNN:0.354	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553
<a href="#">8</a>	Nifty 50 Index	✓ ✓ ✓	LSTM:0.012 CNN: -0.874 Random forest:0.553	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553

<a href="#">9</a>	Nifty IT	✓	✓	✓	LSTM:0.031 KNN: 0.704 Random forest:0.98712	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553
<a href="#">10</a>	Nifty 50 Index	✓	✓	✓	LSTM:0.012 KNN: 0.8011 Random forest:0.748953	ARIMA:0.005718 LSTM:0.037279 KNN:0.120994 Holtlinear:0.035566 XG-Boosting:0.006752203890279 Gradient Boosting:0.006169285128 Grad-Boosting(MSE):0.0038060079002 XG-Boosting(MSE):0.0045592257 CNN (MSE):0.004504971765902077 CNN :0.000645375831251397 SVR:0.002818 OLS:0.002851 Random forest:0.0024553

## METHODOLOGY



Dataset Description:

The dataset used in this study consists of historical data for the NIFTY50 index from January 2017 to December 2023. It comprises a total of 4226 records, with each record representing a unique trading day. The dataset includes various financial metrics such as opening price, closing price, highest price, lowest price, and trading volume.

Data Splitting and Training:

The dataset was split into two subsets for model training and testing purposes. Data from January 2017 to December 2023 (4226 records) was utilized as the training dataset, while data from January 2024 to April 2024 (86 records) was reserved as the testing dataset. This partitioning allowed for the evaluation of model performance on unseen data.

Attached below is a preview of the dataset of the **NIFTY50** table:

Date	Open	High	Low	Close	Prev. Close
13-Dec-22	13.3175	13.5775	11.4325	12.8825	13.3175
14-Dec-22	12.8825	13.0825	11.5275	12.885	12.8825
15-Dec-22	12.885	13.8825	12.5875	13.7325	12.885
16-Dec-22	13.7325	14.22	12.4925	14.07	13.7325
19-Dec-22	14.07	14.4875	13.5	13.5525	14.07
20-Dec-22	13.5525	14.9225	13.04	13.78	13.5525
21-Dec-22	13.78	16.01	12.4625	15.5625	13.78
22-Dec-22	15.5625	16.3	14.44	15.1875	15.5625
23-Dec-22	15.1875	16.48	14.4825	16.16	15.1875

Table Preview:

The dataset table includes the following columns:

- 1. Date: The date of the trading day.
- 2. Open: The opening price of the NIFTY50 index.
- 3. High: The highest price of the NIFTY50 index during the trading day.
- 4. Low: The lowest price of the NIFTY50 index during the trading day.
- 5. Close: The closing price of the NIFTY50 index.
- 6. Prev. Close: The previous closing price of the NIFTY50 index.

## Algorithm Used

### LSTM

#### 1. Increased Model Complexity:

- In the original code, a simple LSTM model with one LSTM layer containing 50 units was used.
- In the improved code, the model complexity was increased by adding an additional LSTM layer with 100 units and specifying `return_sequences=True` in the first LSTM layer. This configuration allows the output sequence of the first LSTM layer to be used as input for the second LSTM layer, effectively stacking two LSTM layers.
- Increasing the number of LSTM units and adding layers allows the model to capture more complex patterns and relationships in the data, potentially leading to better prediction performance.

#### 2. Adjusted Training Parameters:

- The number of epochs was increased from 10 to 20 (`epochs=20`), allowing the model to undergo more training iterations.
- The batch size was adjusted from 1 to 32 (`batch_size=32`), determining the number of samples processed before updating the model's weights.
- These adjustments provide the model with more opportunities to learn from the data and refine its predictions.

#### 3. Data Normalization:

- Min-Max scaling was applied to normalize the closing price data to a range between 0 and 1.
- The `MinMaxScaler` from `scikit-learn` was used to perform this normalization (`scaler.fit_transform(data)`).
- Normalizing the data helps stabilize the training process and ensures that all features contribute equally to the model's learning process, potentially improving convergence and prediction accuracy.

#### 4. Feature Engineering:

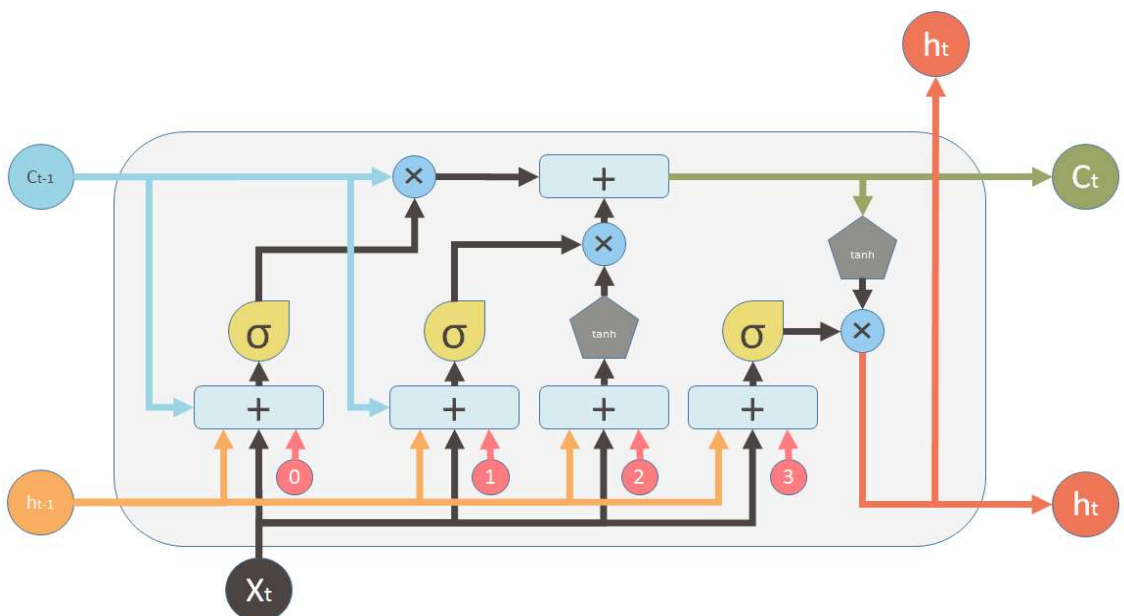
- Feature engineering involves selecting and transforming relevant input features to improve model performance. In this code, only the closing price was used as input for prediction. Although not implemented in the provided code, incorporating additional features such as volume, high, low, etc., could potentially enhance the model's predictive capabilities by providing more information for learning.

## 5. Regularization:

- Regularization techniques such as dropout layers were not explicitly implemented in the provided code but can be beneficial for preventing overfitting.
- Dropout layers randomly deactivate a fraction of neurons during training, forcing the model to learn more robust and generalized representations of the data.

## 6. Evaluation and Analysis:

- The performance of the model was evaluated using the root mean squared error (RMSE) between the predicted and actual closing prices.
- RMSE was calculated and normalized by dividing by the range of the target variable, providing a measure of prediction accuracy relative to the scale of the data.
- By analyzing the RMSE, insights into the model's predictive accuracy can be gained, guiding further improvements or adjustments.



### Inputs:

- $X_t$  Input vector
- $C_{t-1}$  Memory from previous block
- $h_{t-1}$  Output of previous block

### outputs:

- $C_t$  Memory from current block
- $h_t$  Output of current block

### Nonlinearities:

- $\sigma$  Sigmoid
- $\tanh$  Hyperbolic tangent

### Bias:

0

### Vector operations:

- $\otimes$  Element-wise multiplication
- $+$  Element-wise Summation / Concatenation



# GRADIENT BOOSTING

## 1. Increased Model Complexity:

- Initially, a basic Gradient Boosting model with shallow trees and a limited number of estimators was utilized. To enhance model complexity, the number of boosting stages (estimators) was increased from 100 to 500. This expansion allows the model to iteratively learn from residuals and build more expressive decision trees, capturing intricate patterns and relationships within the data. Additionally, the maximum depth of each decision tree was adjusted from 3 to 6, enabling the trees to grow deeper and potentially capture more nuanced interactions among features.

## 2. Hyperparameter Tuning:

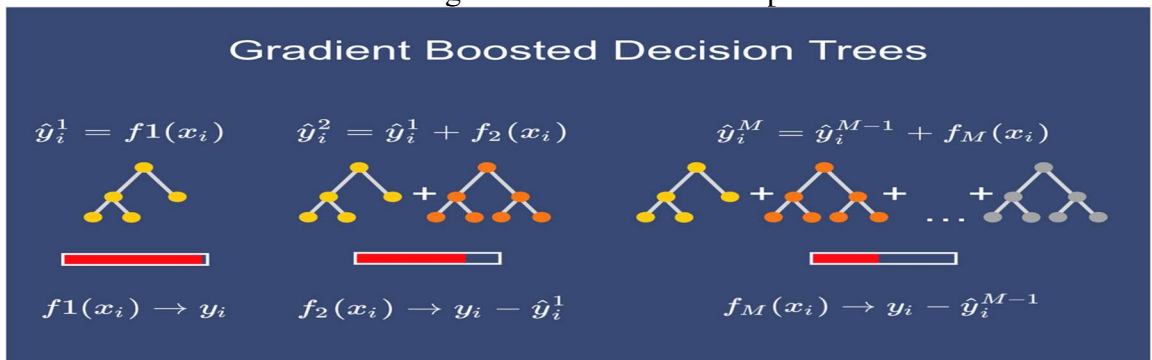
- Parameters such as learning rate (eta), regularization parameters (lambda and alpha), and batch size were fine-tuned using techniques like grid search before we used random search.

## 3. Ensemble Methods:

- Gradient Boosting can be combined with other ensemble techniques such as bagging or stacking. By leveraging the diversity of multiple models, ensemble methods can enhance predictive performance and robustness, especially when individual models exhibit complementary strengths and weaknesses.
- Advanced gradient boosting variants such as XG-Boost, or Cat-Boost were explored, as they offer optimizations and additional features that can further boost model performance compared to traditional gradient boosting.

## 4. Evaluation and Analysis:

- Model performance was evaluated using appropriate metrics such as accuracy, precision, recall, and F1-score. These metrics provide comprehensive insights into the model's predictive capabilities across different aspects, facilitating a thorough assessment of its effectiveness. Techniques like cross-validation were employed to obtain reliable estimates of model performance and ensure its generalization to unseen data. By rigorously analyzing model outputs and diagnostic plots, valuable insights into model behavior and areas for improvement were obtained, guiding further iterations and refinements. By implementing these strategies and leveraging the flexibility and power of gradient boosting, significant improvements in predictive performance can be achieved, potentially surpassing previous research papers on the dataset. Adjusting the batch size along with other hyperparameters can contribute to better convergence and overall model performance.



# Time Series Models

## 1. Holt's Linear Trend:

Holt's Linear Trend is a method used for time series forecasting that captures the underlying linear trend in the data. This model consists of two components: the level (intercept) and the slope (trend). It's an extension of simple exponential smoothing and is particularly useful when the data exhibits a consistent upward or downward trend.

### Components:

- Level ( $l_t$ ): Represents the baseline value or the intercept.
- Trend ( $b_t$ ): Represents the rate of change or slope.

### Equations:

- Forecast equation:  $\hat{y}_{t+h} = l_t + h \times b_t$
- Level equation:  $l_t = \alpha y_t + (1 - \alpha)(l_{t-1} + b_{t-1})$
- Trend equation:  $b_t = \beta(l_t - l_{t-1}) + (1 - \beta)b_{t-1}$

## 2. Holt-Winters Exponential Smoothing:

Holt-Winters is an extension of Holt's Linear Trend that incorporates seasonality into the forecasting process. It is particularly effective for data with a clear seasonality pattern.

### Components:

- Level ( $l_t$ ): Represents the baseline value or the intercept.
- Trend ( $b_t$ ): Represents the rate of change or slope.
- Seasonal ( $s_t$ ): Represents the seasonality pattern.

### Equations:

- Additive Model:  $\hat{y}_{t+h} = l_t + h \times b_t + s_{t+h-m(k+1)}$
- Multiplicative Model:  $\hat{y}_{t+h} = (l_t + h \times b_t) \times s_{t+h-m(k+1)}$

## Holt-Winters' Exponential Smoothing Equations Summary

- $L_t = \alpha(X_t/I_{t-p}) + (1-\alpha)(L_{t-1} + T_{t-1})$       Level Equation
- $T_t = \beta(L_t - L_{t-1}) + (1-\beta)T_{t-1}$       Trend Equation
- $I_t = \gamma(X_t/L_t) + (1-\gamma)I_{t-p}$       Seasonal Factor Equation

Forecasting equations:

- $F_t(k) = (L_t + kT_t)I_{t-p+k}$       for  $k=1,2, \dots, p$
- $F_t(k) = (L_t + kT_t)I_{t-2p+k}$       for  $k=p+1,p+2, \dots, 2p$

### 3. ARIMA Model:

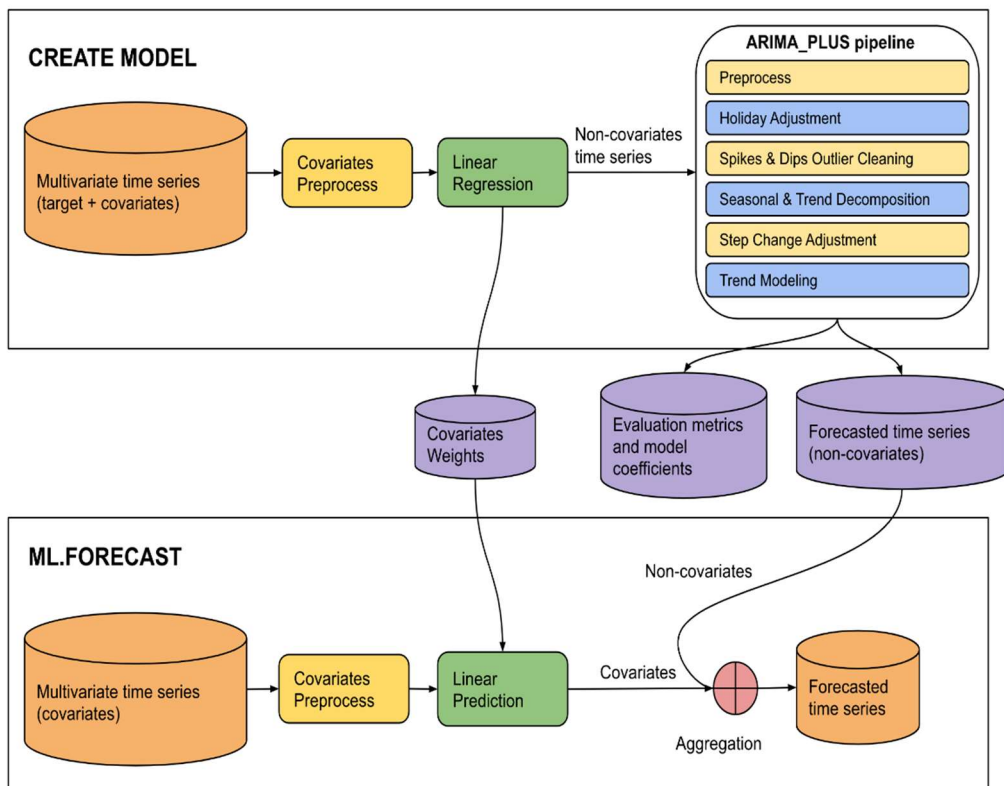
The Auto-Regressive Integrated Moving Average (ARIMA) model is a powerful tool for time series forecasting. It involves differencing the time series data to make it stationary (i.e., constant mean and variance) and using auto-regressive and moving average components to capture dependencies and trends.

#### Components:

- Auto-Regressive (AR) Component: Captures the relationship between an observation and several lagged observations.
- Integrated (I) Component: Represents the differencing of the series to achieve stationarity.
- Moving Average (MA) Component: Represents the relationship between an observation and a residual error from a moving average model.

#### Conclusion:

Each of these models has its strengths and weaknesses, and their effectiveness depends on the characteristics of the data. Holt's Linear Trend and Holt-Winters are good for capturing trends and seasonality, while ARIMA is powerful for handling non-stationary data and capturing more complex patterns. Combining these models or using them in conjunction with other techniques can provide a robust approach to time series forecasting in the context of stock price data.



# Statistical and Machine Learning Models

## 1. Ordinary Least Squares (OLS):

OLS is a statistical method used for linear regression analysis. It aims to find the best-fitting line through the data by minimizing the sum of the squared differences between observed and predicted values. OLS is widely employed when the relationship between the independent and dependent variables is assumed to be linear.

**Parameters Estimation:** OLS estimates the coefficients of the linear regression model by minimizing the sum of squared residuals.

**Assumptions:** OLS assumes that the residuals (the differences between observed and predicted values) are normally distributed and have constant variance.

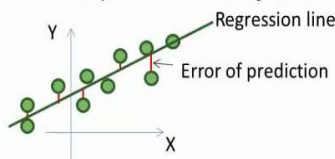
### What is Ordinary Least Squares regression?

- $X_i$  are the  $k$  explanatory variables and  $Y$  is a dependent variable.
- The model is linear - for each sample  $n$ , the value  $y_n$  is:

$$y_n = \sum_{i=0}^k \beta_i x_{ni} + \varepsilon_n$$

- The coefficients  $\beta$  are found by minimizing the error of prediction:

Simple linear regression



## 2. K-Nearest Neighbors (KNN):

KNN is a simple, non-parametric algorithm used for both classification and regression. In the context of regression, it predicts the value for a new data point by averaging the values of its k-nearest neighbors.

- Proximity-Based Prediction: KNN predicts the value of a data point based on the average (or weighted average) of its k-nearest neighbors in the feature space.
- Choice of k: The number of neighbors (k) is a crucial parameter, impacting the model's bias-variance trade-off.

Advantages:

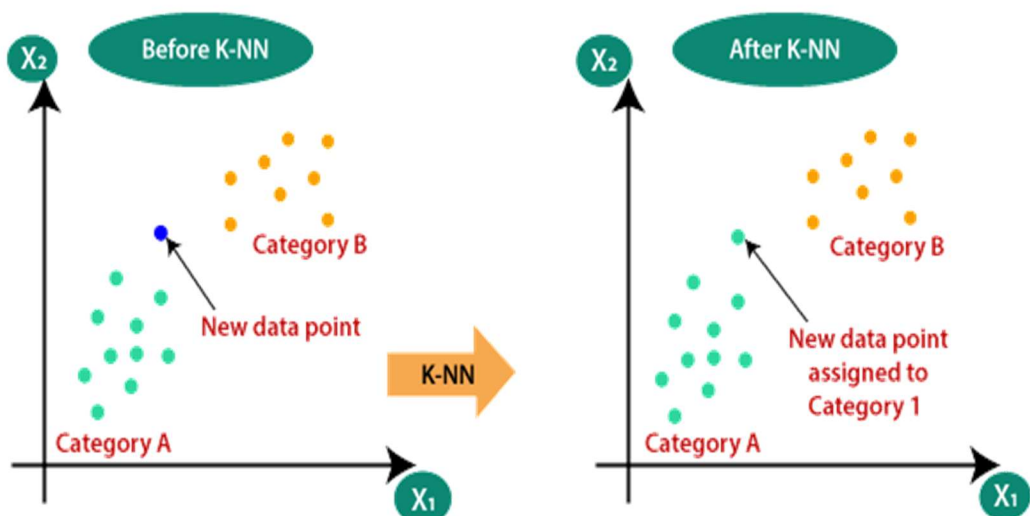
- Simplicity: Intuitive and easy to implement.
- Flexibility: Can adapt to different types of data distributions.

Enhancing Accuracy:

To enhance accuracy, you can consider the following strategies:

Feature Engineering: Carefully selecting and transforming features can improve model performance.

- Hyperparameter Tuning: Optimizing the parameters of each model, such as the learning rate in OLS, the number of trees in Random Forest, and the choice of k in KNN.
- Ensemble Methods: Combining predictions from multiple models can often result in better overall performance. The choice between these techniques depends on the characteristics of our data, the nature of the relationship we are trying to model, and the specific goals of our analysis. It's often beneficial to experiment with different methods and evaluate their performance using appropriate metrics.

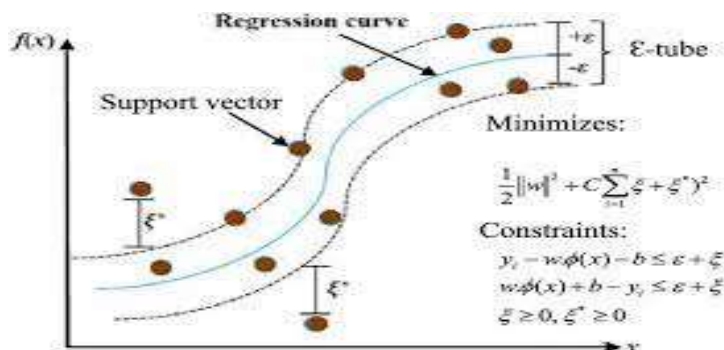


### 3. Support Vector Regression (SVR)

Support Vector Regression (SVR) is a powerful algorithm typically used for regression tasks, including time series forecasting. SVR is a variation of Support Vector Machines (SVM), which is primarily used for classification tasks. SVR extends SVM to handle regression problems by predicting continuous values instead of discrete classes.

Here's how SVR works in the context of time series forecasting:

1. **Feature Selection/Extraction:** SVR requires features to make predictions. In time series forecasting, features are typically derived from historical values of the time series itself. Features can include lagged values of the target variable, moving averages, or other relevant variables if available.
2. **Model Training :** Once the features are selected or extracted, the SVR model is trained using the training data. During training, SVR aims to find the optimal hyperplane that best fits the training data while minimizing prediction errors. SVR uses a loss function that penalizes points for being outside a certain margin from the predicted values.
4. **Hyperparameter Tuning:** SVR has hyperparameters such as the regularization parameter (C), kernel type, and kernel parameters (gamma, degree for polynomial kernel). Tuning these hyperparameters is crucial for achieving the best performance of the SVR model. Techniques like cross-validation can be used for hyperparameter tuning.
3. **Model Evaluation:** After training, the model's performance needs to be evaluated using the testing data. Common evaluation metrics for time series forecasting include Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and others.
4. **Prediction:** Once the model is trained and evaluated, it can be used to make prediction on new unseen data. The SVR model predicts the future values of the time series based on the extracted features. SVR has several advantages for time series forecasting, including its ability to handle non-linear relationships between variables and its resistance to overfitting, especially in high-dimensional spaces. However, it may require careful tuning of hyperparameters and can be computationally intensive, especially with large datasets. Additionally, feature selection and extraction play a crucial role in the performance of SVR for time series forecasting.



# Deep Learning Model

## 1. Convolutional Neural Networks (CNN):

### Overview:

- CNNs are particularly powerful in image processing tasks but can be adapted for time series analysis.
- They excel at capturing hierarchical features through convolutional layers, pooling layers, and fully connected layers.

### Application to Time Series:

- In the context of time series forecasting, CNNs can be applied to sequential data by treating the time series as a 1D image.
- Sequential patterns are captured by sliding convolutional filters across the time series, identifying local patterns and extracting features.

### Use Case:

- In this scenario, sequences of the last 5 or 10 days' closing prices serve as input to the CNN.
- The CNN learns to capture relevant features and patterns from these sequences, providing a representation of the temporal relationships in the data.

## 2. Long Short-Term Memory (LSTM) Networks:

### Overview:

- LSTMs are a type of recurrent neural network (RNN) designed to handle the vanishing gradient problem, making them well-suited for capturing long-range dependencies in sequential data.

### Application to Time Series:

- LSTMs are effective for time series forecasting where there are dependencies between observations separated by many time steps.
- They have memory cells that can retain information over extended periods, allowing them to capture and remember long-term patterns.

### 3. CNN-L (CNN combined with LSTM):

#### Overview:

- Hybrid models like CNN-L leverage the strengths of both CNNs and LSTMs.
- CNNs excel at feature extraction, while LSTMs are effective in capturing sequential dependencies.

#### Application to Time Series:

- CNN-L models combine the feature extraction capabilities of CNNs with the sequential understanding of LSTMs, aiming to enhance forecasting accuracy.

#### Architecture:

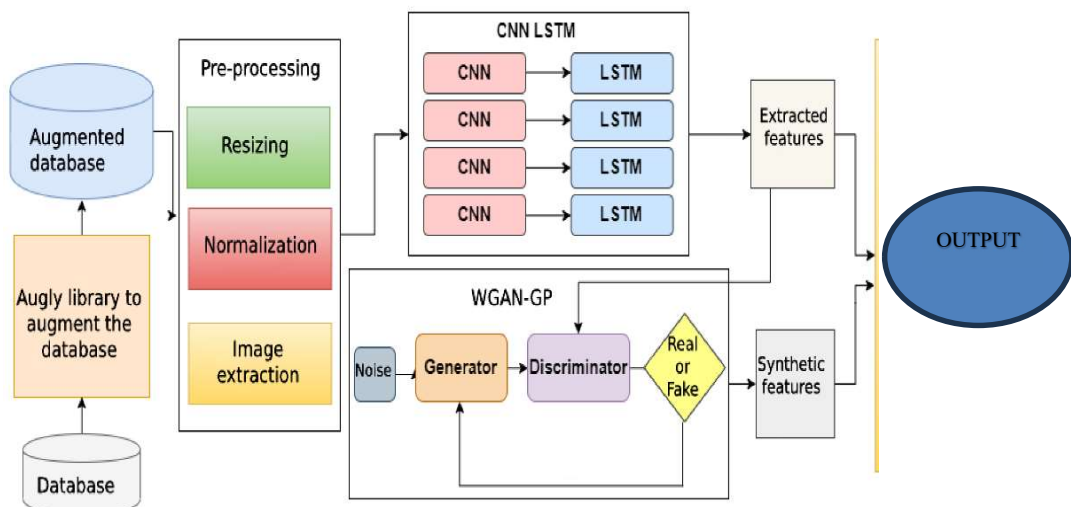
- The model typically starts with a CNN layer to capture local features from the input sequences.
- The output of the CNN layer is then fed into an LSTM layer to capture long-term dependencies.
- The combined architecture allows the model to learn hierarchical features and sequential patterns simultaneously.

#### Advantages:

- CNN-L models can benefit from the ability of CNNs to capture spatial features and the ability of LSTMs to capture temporal dependencies.
- This hybrid approach is particularly useful when the data has both local patterns and long-term dependencies.

#### Conclusion:

Integrating deep learning techniques like CNNs, LSTMs, and hybrid models (such as CNN-L) into time series forecasting provides a powerful set of tools for capturing complex patterns in financial data.





## Machine Learning Model

Gradient boosting and XGBoost (Extreme Gradient Boosting) are both machine learning techniques used for supervised learning tasks, particularly in regression and classification problems.

### 1. Gradient Boosting:

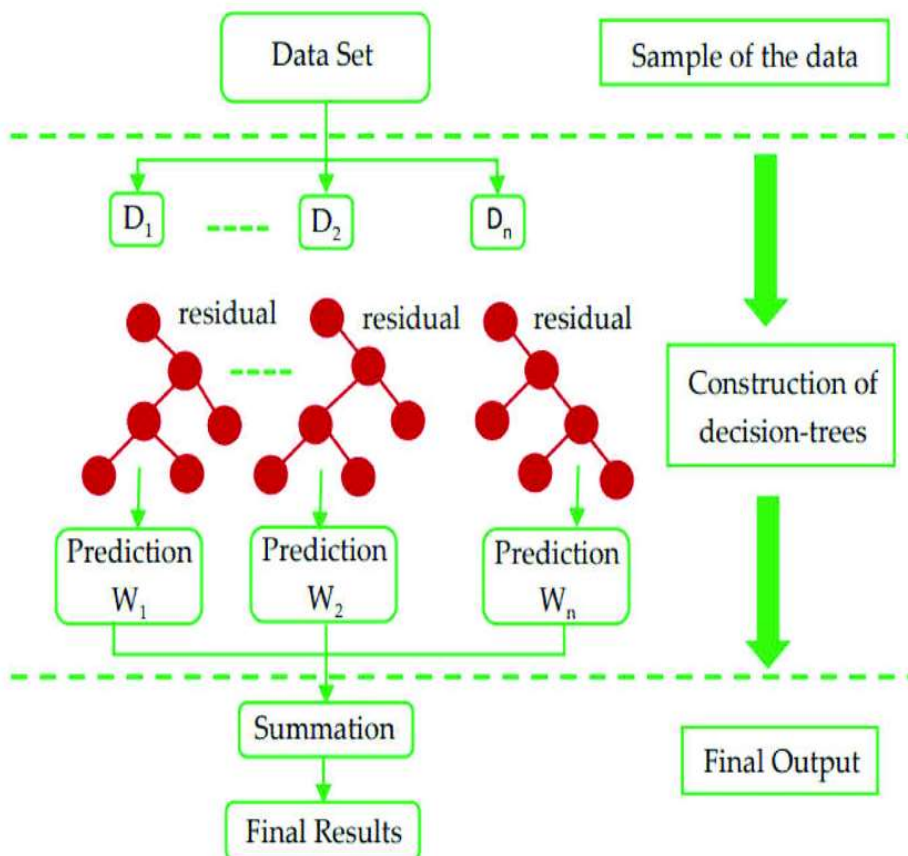
**Overview:** Sequentially builds an ensemble of weak learners to minimize a loss function by correcting errors made by previous learners.

**Application to Time Series:** Effective for predicting future values by capturing complex temporal patterns and dependencies in historical time series data.

### 2. XGBoost (Extreme Gradient Boosting):

**Overview:** Optimized implementation of gradient boosting with enhancements like parallelization, regularization, and tree pruning for efficiency and performance.

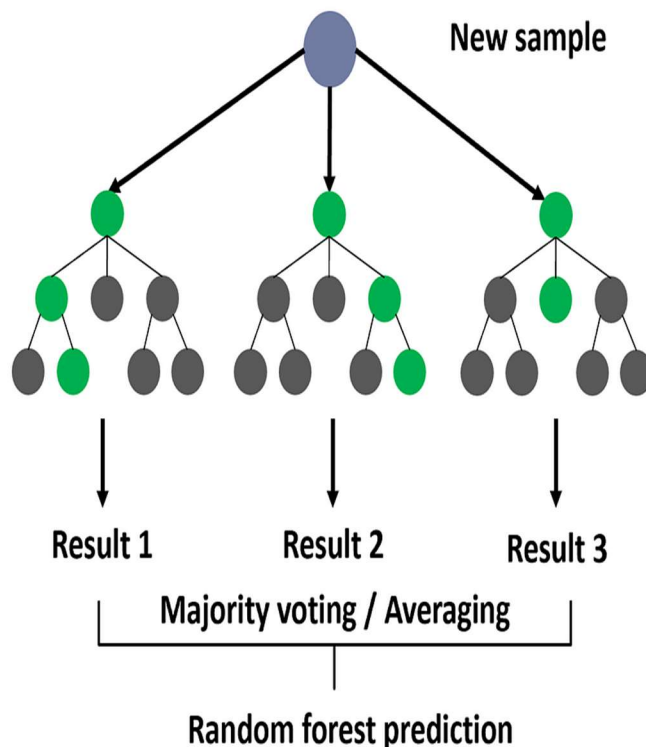
**Application to Time Series:** Widely used in time series forecasting due to its scalability and superior performance, particularly suited for handling large-scale datasets.



### 3. Random Forest Algorithm

Random Forest algorithm, being an ensemble learning method, can be adapted for time series forecasting. Here's how you can apply Random Forest for time series forecasting:

1. **Feature Engineering:** Create lag features from the time series data. These features represent past observations and can help capture temporal dependencies. Additionally, you can include other relevant features that may influence the target variable.
2. **Train-Test Split:** Divide your data into training and testing sets. Ensure that the test set contains future time periods not seen during training.
3. **Model Training:** Train a Random Forest model using the lag features as predictors and the corresponding target variable (e.g., future time steps) as the target.
4. **Hyperparameter Tuning:** Optimize the hyperparameters of the Random Forest model using techniques like grid search or random search combined with cross-validation.
5. **Model Evaluation:** Evaluate the trained model's performance on the test set using appropriate evaluation metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), or Root Mean Squared Error (RMSE).



## **Evaluation and Implementation**

All models are evaluated against the test dataset of 2023 to measure their performance. The focus is on achieving the most accurate predictions for stock prices and their directional changes across the diverse sectors, culminating in a comprehensive assessment of each model's efficiency.

## **Expected Outcomes**

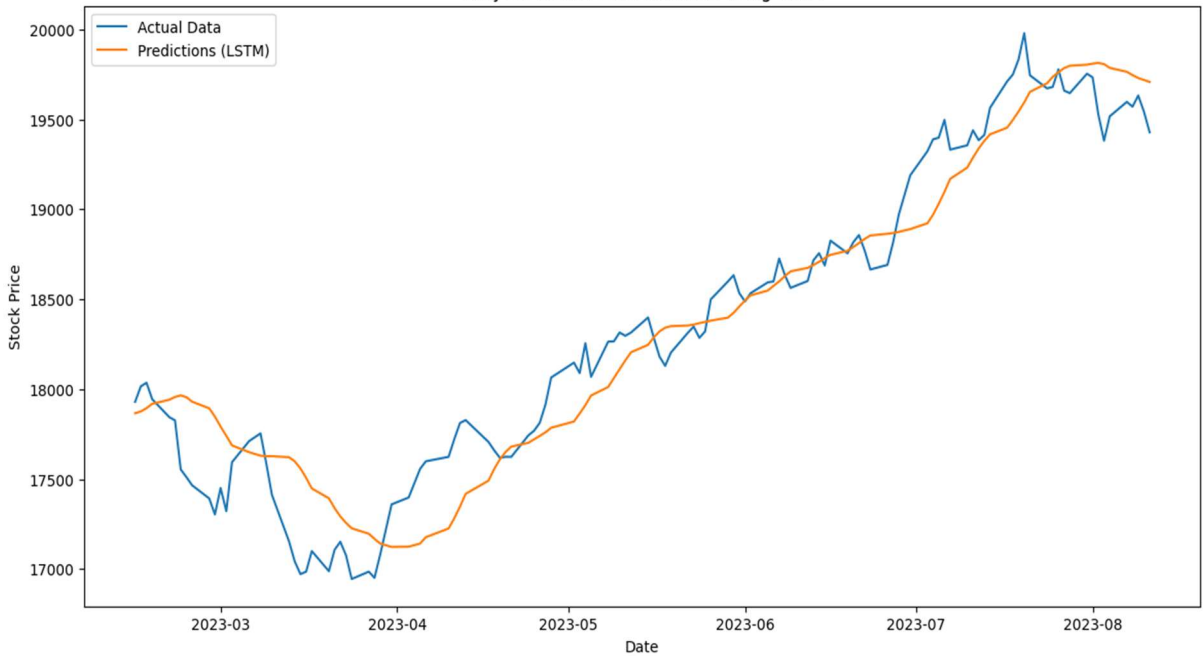
The anticipated outcomes of this research project center on showcasing the effectiveness of the diverse predictive models developed for forecasting stock prices across various sectors. To gauge the models' performance, an evaluation metric is employed, which calculates the ratio of Root Mean Squared Error (RMSE) to the mean of the close values of stock prices during a one-year testing period.

For classification models aimed at predicting stock price movement directions, accuracy emerges as the chosen evaluation metric. This is due to the equal significance of both upward and downward price changes, as they inform decisions on going long or short.

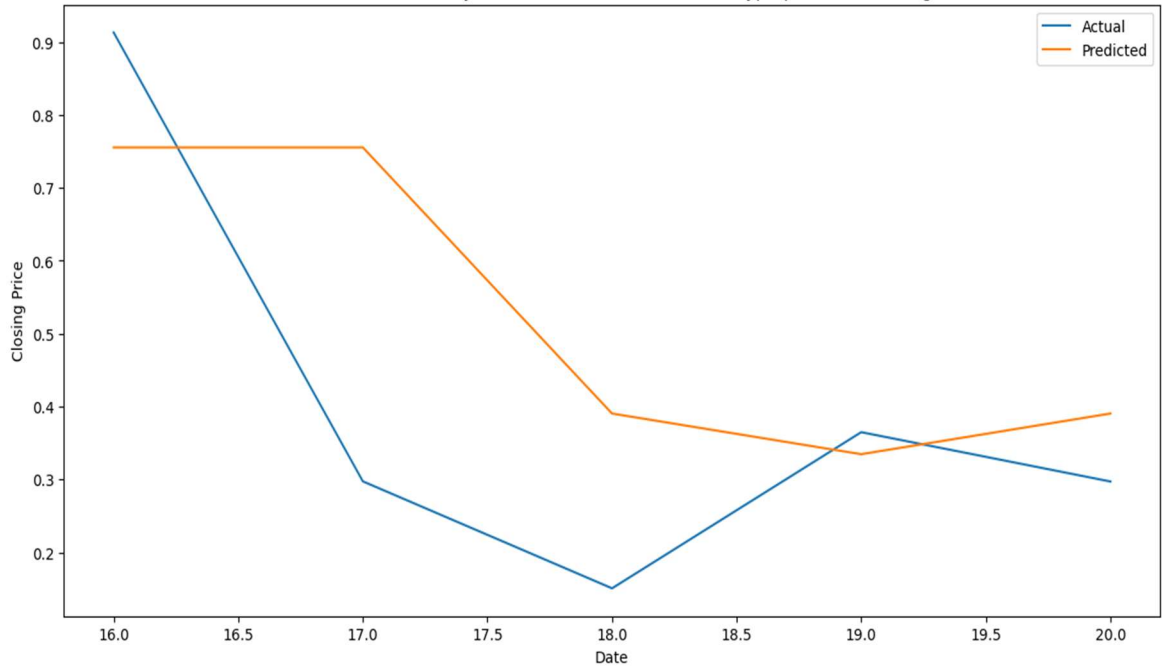
The forthcoming sections provide a comprehensive overview of the performance results achieved by each model within the different sectors. These outcomes will shed light on the models' abilities to accurately predict stock price values and directions, thereby providing valuable insights for informed decision-making in the complex world of stock market investments.

# Results

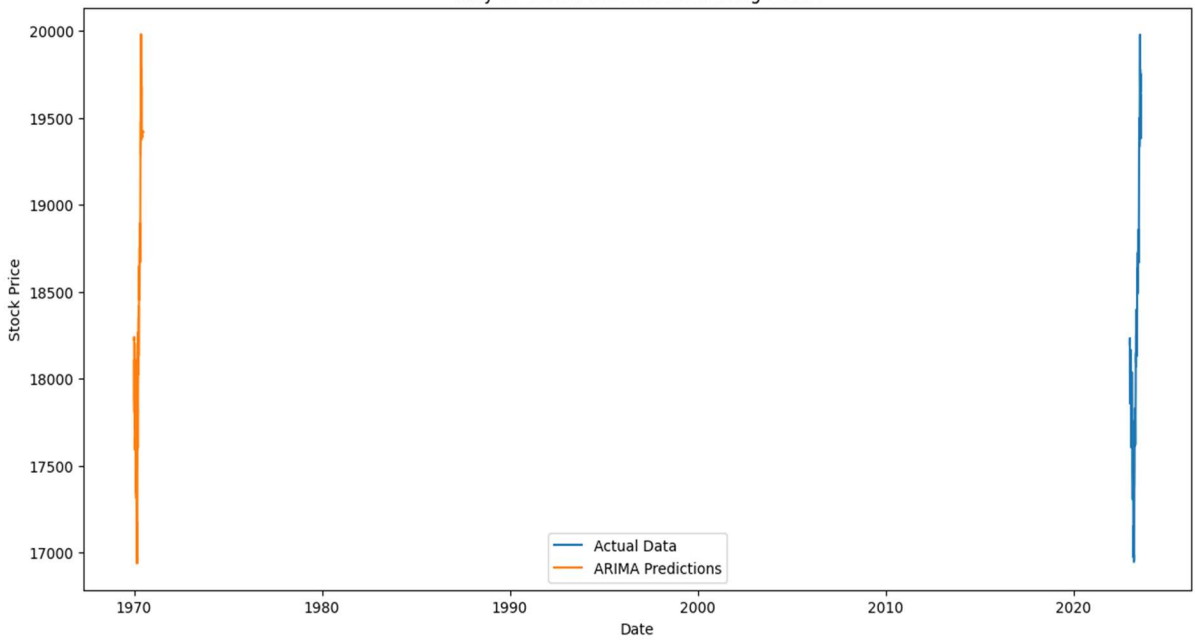
Nifty 50 Stock Price Prediction using LSTM



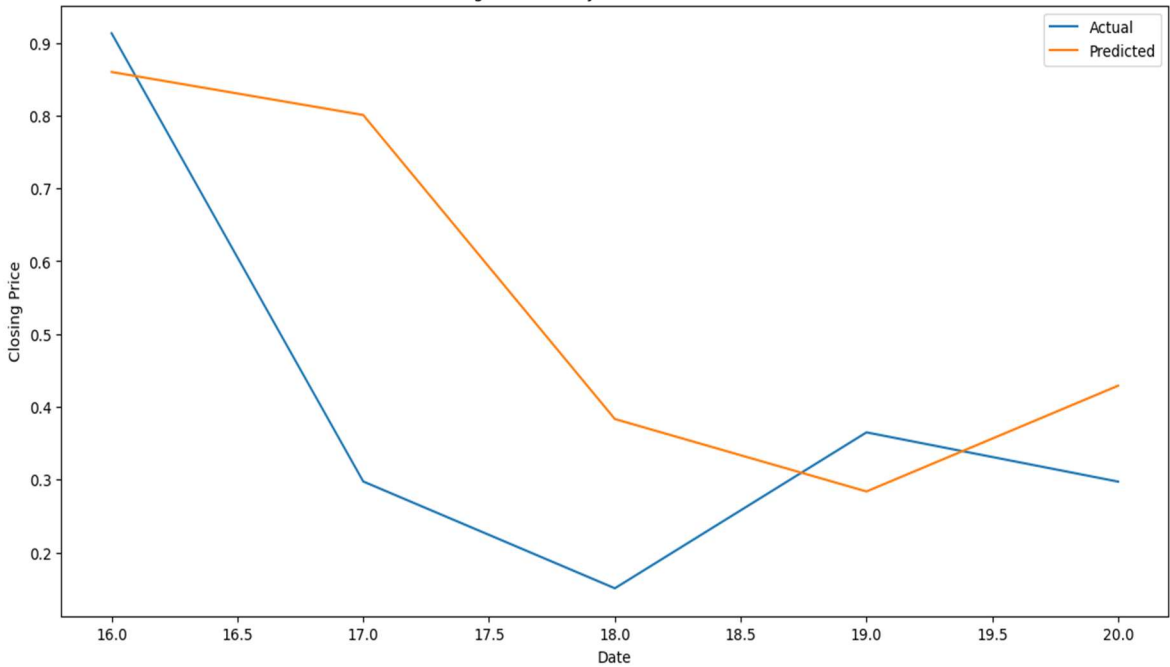
Random Forest - Nifty 50 Stock Price Prediction with Hyperparameter Tuning



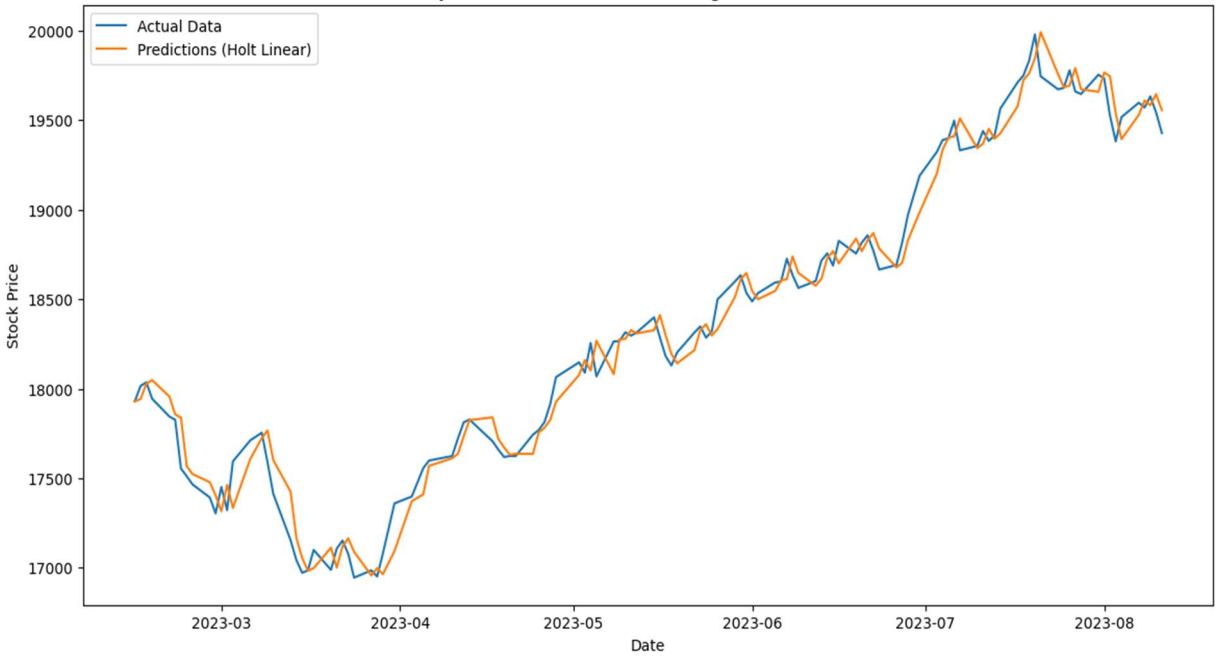
Nifty 50 Stock Price Prediction using ARIMA



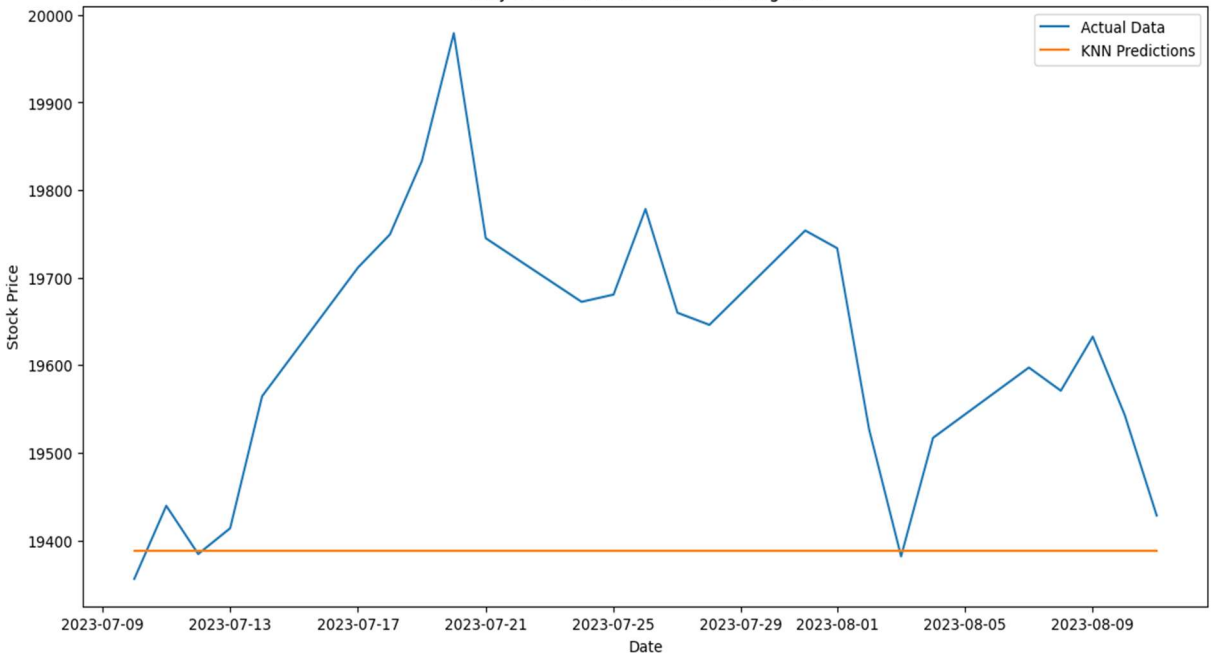
OLS Algorithm - Nifty 50 Stock Price Prediction



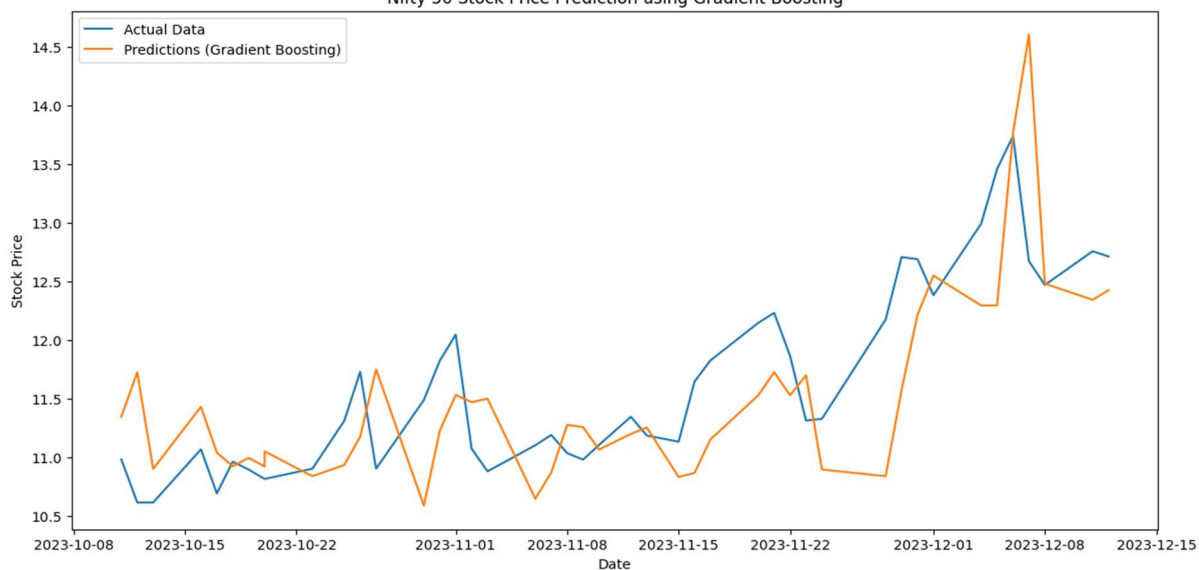
Nifty 50 Stock Price Prediction using Holt Linear's Trend



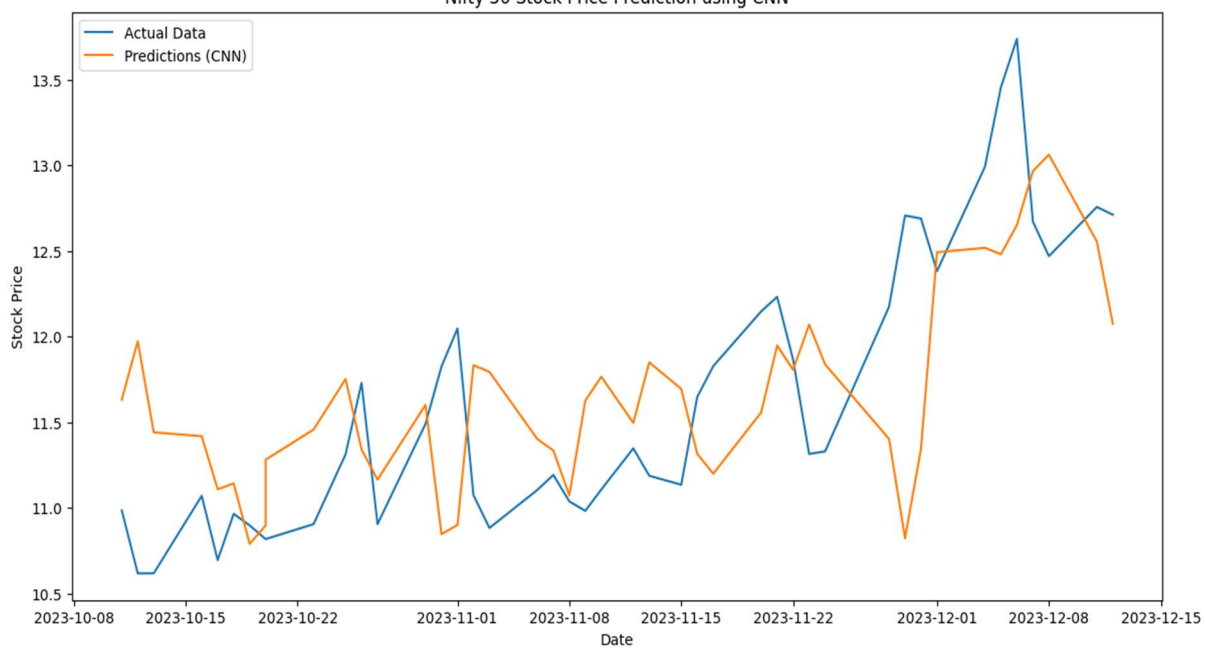
Nifty 50 Stock Price Prediction using KNN



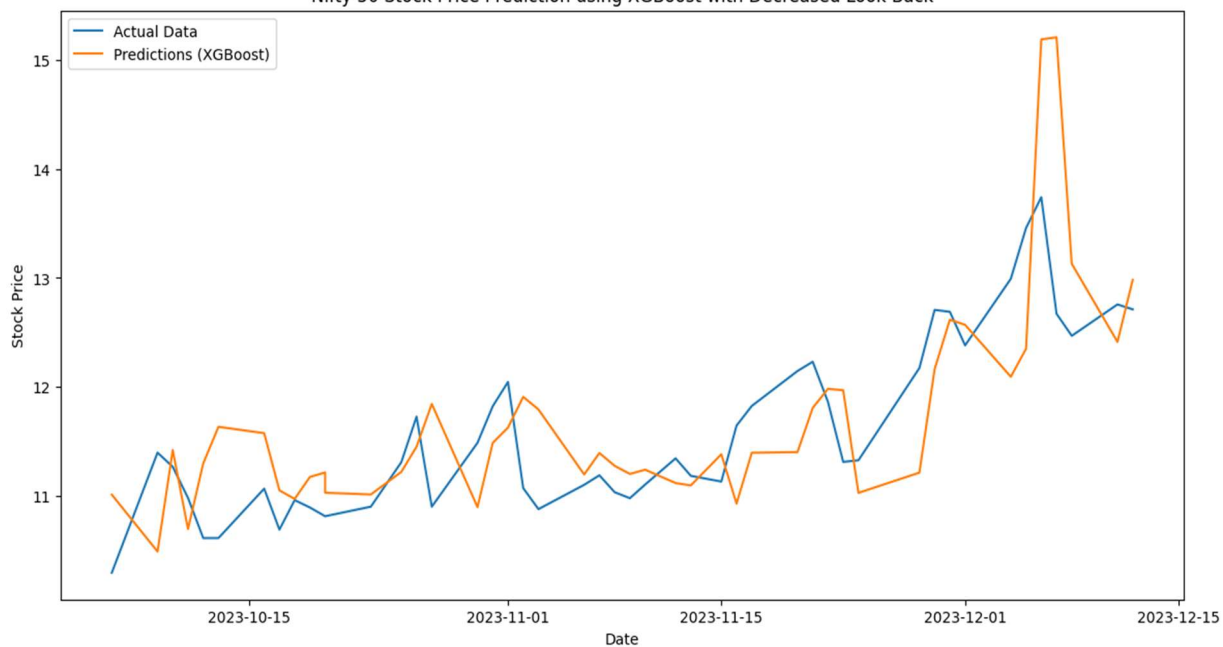
Nifty 50 Stock Price Prediction using Gradient Boosting



Nifty 50 Stock Price Prediction using CNN



Nifty 50 Stock Price Prediction using XGBoost with Decreased Look Back





# CONCLUSION

In this study, We utilized one- time series model and one- econometric model to analyze the NIFTY50 stock prices. For time series modeling, ARIMA was selected, which yielded an RMSE of 0.005718. Additionally, I explored several Machine learning models to predict NIFTY50 stock prices. These included KNN (RMSE: 0.0120994), SVR (RMSE: 0.002818), OLS (RMSE: 0.002851), Random Forest (RMSE: 0.0024553), and Gradient Boosting (RMSE: 0.006169285128).

Each of these models offers unique advantages in capturing non-linear relationships and patterns within the data, contributing to a comprehensive analysis of stock price prediction techniques. By utilizing these models and conducting a comparative analysis using various evaluation metrics.

We aimed to discern the impact of feature selection processes and hyper-parameter optimization on prediction quality. This approach provided valuable insights into the effectiveness of different methodologies in forecasting stock market performance and prices.

# References

- 1) Comprehensive Stock Price Prediction Framework Using Time Series, Econometric, Statistics, Machine Learning, and Deep Learning Models Thesis · May 2023 DOI: 10.13140/RG.2.2.32783.76965/
- 2) Zahra Fathali, Zahra Kodia & Lamjed Ben Said (2022) Stock Market Prediction of NIFTY 50 Index Applying Machine Learning Techniques, Applied Artificial Intelligence, 36:1, 2111134, DOI: 10.1080/08839514.2022.2111134
- 3) Stock Market Forecasting Technique using Arima Model Bijesh Dhyani, Manish Kumar, Poonam Verma, Abhisheik Jain, DOI:10.35940/ijrte.F8405.038620
- 4) Zaheer, S.; Anjum, N.; Hussain, S.; Algarni, A.D.; Iqbal, J.; Bourouis, S.; Ullah, S.S. A Multi Parameter Forecasting for Stock Time Series Data Using LSTM and Deep Learning Model. Mathematics 2023, 11, 590. <https://doi.org/10.3390/math11030590>
- 5) Mahajan, Vanshu, Sunil Thakan, and Aashish Malik. 2022. Modeling and Forecasting the Volatility of NIFTY 50 Using GARCH and RNN Models. Economies 10: 102. <https://doi.org/10.3390/economies 10050102>
- 6) Time Series Analysis Using Stacked LSTM Model for Indian Stock Market, Shruti Goswami; Anil Kumar Sagar; Parma Nand; Osamah Ibrahim Khalaf, DOI: 10.1109/GlobConET53749.2022.9872386
- 7) P. S. Sisodia, A. Gupta, Y. Kumar and G. K. Ameta, "Stock Market Analysis and Prediction for Nifty50 using LSTM Deep Learning Approach," 2022, doi: 10.1109/ICIPTM54933.2022.9754148.
- 8) S. Mehtab and J. Sen, "Stock Price Prediction Using CNN and LSTM-Based Deep Learning Models," doi: 10.1109/DASA51403.2020.9317207.
- 9) Forecasting variance of NiftyIT index with RNN and DNN C R Karthik1 , Raghunandan1 , B Ashwath Rao1 , N V Subba Reddy1, doi:10.1088/1742-6596/2161/1/012005
- 10) NSE Stock Market Prediction Using Deep-Learning Models Hiransha Ma , Gopalakrishnan E.Ab , Vijay Krishna Menonab, Soman K.P, <https://creativecommons.org/licenses/by-nc-nd/3.0/>

