COMPARATIVE ANALYSIS OF MACHINE LEARNING AND STATISTICAL APPROACHES FOR FINANCIAL MARKET FORECASTING: A MATHEMATICAL PERSPECTIVE

Eko Riwayadi

Faculty of Science and Technology, Airlangga University, Surabaya, Indonesia

Corresponding author: eko.riwayadi-2023@fst.uniar.ac.id

Abstract. Efficient financial market forecasting is crucial for informed decision-making. This study presents a comprehensive analysis that juxtaposes traditional statistical methods with modern machine learning techniques for forecasting in financial markets. The research evaluates empirical performance, interpretability, and adaptability across various financial datasets. Commencing with a thorough Literature Review, the study explores Time Series Models such as ARIMA, alongside contemporary approaches like Neural Networks and Gradient Boosting Machines. The Comparative Methodology encompasses data preprocessing and model implementation for both traditional and modern forecasting approaches. Results showcase accuracy metrics, resilience to market fluctuations, and inherent strengths of each method. Additionally, our findings shed light on the mathematical principles influencing outcomes, offering a valuable perspective from a mathematical standpoint. The practical implications extend to portfolio management, risk assessment, and the formulation of effective trading strategies. Moreover, the study deliberates on future directions, delving into emerging mathematical techniques and the potential of hybrid models. The Conclusion succinctly summarizes key contributions, emphasizing the significance of understanding mathematical foundations for successful forecasting. Bridging theory and practice, this research provides insights into the selection of appropriate methods, guiding real-world financial decisions. Furthermore, the discussion of the research results highlights the effectiveness of the Random Forest model in stock price forecasting, affirming its superiority over other approaches.

Keywords: Forecasting, ARIMA, Machine Learning, Financial market, Mathematical

1. INTRODUCTION

Financial markets are dynamic ecosystems where the interplay of economic indicators, investor sentiment, and geopolitical events constantly influence asset prices. For examples as known on the past experienced, consider the global financial crisis of 2008. The sudden collapsed of Lehman Brothers sent shockwaves through the financial world, leading to a cascading effect that hit various sectors including Indonesia. Such events underline the volatility and interconnectedness of financial markets, making accurate predictions a frightening challenge. Furthermore, the rise of algorithmic trading and high frequency trading has accelerated market movement that require timely and precise forecasts to capitalize on floating opportunities.

Accurate forecasts could guide in selecting assets align with risk taking and investment goals. On the other way, inaccurate forecasts might lead to un-favourable investment decisions, resulting in financial losses. Institutional investors who manage pension funds worth billions of dollars face even higher stakes. A pension fund that can accurately anticipate market trends could ensure the financial security of thousands of retirees, while an inaccurate forecast might jeopardize their futures.

Mathematics serves as the bedrock for understanding and predicting financial markets. Consider a time series analysis of stock prices. By applying mathematical techniques such as Autoregressive Integrated Moving Average (ARIMA) model,

analysts can identify trends, seasonality, and cyclic patterns hidden within historical price data (Rubio et al., 2023). This quantitative approach enables them to make informed predictions about potential price movements. The application of mathematical concepts like calculus and linear algebra in machine learning algorithms further refines predictions by adapting to changing market conditions.

Neural networks outperformed the ARIMA model in stock price prediction. Prior research has explored both statistical models like ARIMA and machine learning algorithms such as deep learning for financial forecasting (Azizi et al., 2016; Kumar et al., 2020; Rubio et al., 2023; Sonkavde et al., 2023; Vigneau et al., 2018). However, a comprehensive analysis that systematically compares these approaches while delving into the mathematical principles underlying them is still lacking. This study addresses this gap by conducting an in-depth examination of the mathematical foundations of these methods. Focusing on the mathematical nuances, this research aims to shed light on the reasons behind the respective successes and limitations in forecasting financial markets, providing valuable insights for practitioners and researchers alike.

2. LITERATURE REVIEW

A. Overview of Traditional Statistical Approaches for Financial Forecasting

Time Series Model ARIMA

Time series models are fundamental tools in financial forecasting, capturing patterns and trends in sequential data. The AutoRegressive Integrated Moving Average (ARIMA) model is a classic example. It combines autoregressive (AR) terms, representing the relationship between past and current values, with moving average (MA) terms, accounting for the impact of past forecast errors. The Integrated (I) term indicates the number of differences needed to achieve stationarity.

ARIMA models, popularized by Box and Jenkins, are a fexible and powerful statistical tool for predictive modelling with time series data (Asteriou & Hall, 2016). Mainly, ARIMA models approximate time series future values as a linear function of past observations and white noise terms. The model consists of three components: non-stationary diferences for stationarity, autoregressive model (AR) and moving average (MA) model (Rubio et al., 2023).

To define non-stationarity, the backshift operator, *B* is introduced. A time series y_t , will be called homogeneous non-stationary if it is non-stationary but its first difference, i.e. $w_t = y_t - y_{t-1} = (1 - B)y_t$ or *d*th difference, $w_t = (1 - B)^d y_t$, yields a stationary time series. In addition, y_t will be called an autoregressive integrated moving average (ARIMA) process of orders p, *d*, and *q*, denoted ARIMA (*p*, *d*, *q*) if its dth difference yields a stationary process ARMA (*p*, *q*). Therefore, an ARIMA (*p*, *d*, *q*) can be written as:

$$\Phi(B)(1-B)^d y_t = \delta + \Theta B_{\varepsilon_t}$$

where

$$\Phi(B) = 1 - \sum_{i=1}^{p} \phi_i B^i, \\ \Theta(B) = 1 - \sum_{i=1}^{q} \theta_i B^i$$

are the backshit operator terms in the AR(p) and MA(q) defined as: $\Phi(B)y_t = \delta + \varepsilon_t$ and $y_t = \mu + \Phi(B)\varepsilon_t$, with $\delta =$

 $\phi \mu$, where μ is the mean and ε_i the white noise with $E(\varepsilon_i) = 0$

Model orders p, q are determined by the nature of the autocorrelation and partial autocorrelation functions. The model coefcients are calculated using the maximum

likelihood method (Box et al., 2008; Lihki Rubio 2023). The best model is identifed by diagnostic checks such as the Akaike information criterion (AIC), the Bayes information criterion (BIC) and the Jarque–Bera normality test on the residual error series.

B. Introduction to Machine Learning Techniques in Financial Forecasting

1. Neural Networks:

Neural networks are a class of machine learning models inspired by the human brain's neural structure. In financial forecasting, feedforward neural networks (FNNs) and more advanced architectures like recurrent neural networks (RNNs) and long short-term memory networks (LSTMs) are widely used. These networks learn complex patterns and relationships from historical data (Azizi et al., 2016).

ANN A neural network is a powerful data modeling tool that is able to capture and represent complex input / output relationships. In the present study, a multi - layer perceptron (MLP) neural network was built for potato identification . This type of neural network is known as a supervised network because it requires a desired output for learning. The goal of this type of network was to create a model that correctly maps the input to the output using historical data so that the model can be used to produce the output when the desired output is unknown. In general, setting too few hidden units cause high training errors and high generalization errors due to under - fitting, while too many hidden units result in low training errors but still high generalization errors due to over – fitting (Dudo et al., 2001). However, determination of number of hidden layer neurons largely depends on trial - and - error method. In this research, the method proposed by (Li et al., 1999) was used, in which the number of hidden layer neurons was initially calculated by equation as follow.

$$h = (m+n)^{\frac{1}{2}} + aa \in [0,10]$$

Where m is number of output neurons and represents the potato varieties, n is number of input neurons and equals the number of used input features, and h is the number of hidden neurons. Finally, a neural network was designed with 16-20-10-10 structure, in which 16 neurons were for the input layer, 20 and 10 neurons were for the middle layers, respectively, and 10 neurons were in the output layer. Also, 75 % of the whole data allocated for training and 25 % of data set allocated for test. Processing inside neurons determined by activation function. In this work, the tansig function was used; moreover, training function was used to train the network and updating the weights and bias values according to Levenberg - Marguardt optimization (Azizi et al., 2016). This function is often the fastest backpropagation algorithm in the toolbox and is highly recommended as a first - choice supervised algorithm; however, it requires more memory space comparing with other algorithms. In this study, neural networks were designed, trained, and tested using python neural network software. The output of each class was coded in a binary form. The advantage of this type of coding is having same values and similar errors. That is, if an error occurs, the network considers size of the error values to be equal.

2. Random Forests:

Random forests are an ensemble learning method combining multiple decision trees to make predictions. Each tree trained on a different subset of the data, and the final prediction is an aggregate of individual tree predictions (Sonkavde et al., 2023).

In order to assess the performance using all the data versus the RUS datasets, we employ a Random Forest (RF) model. We selected the RF model because of its good classification performance, which has been shown to be superior to many other classifiers on a wide variety of datasets with or without class imbalance (Delgado et al.,

2014; Khosgoftaar et.al., 2007). Random Forest is an ensemble method in which multiple unpruned decision trees are built and a final classification is made by combining the results from the individual trees (Breiman, 2001). The algorithm creates random datasets using sampling with replacement to train each of the decision trees. At each node within a tree, RF chooses the most discriminating feature between the classes using entropy and information gain. Entropy can be seen as the measure of impurity or uncertainty of attributes, and information gain is a means to find the most informative attribute (Bauder & Khoshgoftaar, 2018). Thus, the goal is to minimize entropy and maximize information gain with attribute selection. Additionally, RF performs random feature subspace selection, at each node of a tree, where a subset of m features are considered for the decision at that node.

The basic common formula for Random Forest focusing on decision tress within the ensemble as follows (Menze et al., 2009):

Gini Impurity (for classification):

Gini (t) =
$$1 - \sum_{i=1}^{c} p(i|t)^2$$

Where t is a node, c is the number of classes and p(ilt) is the probability of class in node t.

Entropy (for classification):

$$H(t) = -\sum_{i=1}^{c} p(i|t) \log_2(p(i|t))$$

Where t is a node, c is the number of classes and p(ilt) is the probability of class in node t.

MSE (Mean Square Error)(for regression):

$$MSE(t) = \frac{1}{nt} \sum_{i \in Dt} (yi - yt)^2$$

Where t is a node nt is the number of sample in node t, Dt, is the set of samples in node t, yi is the target value in node t.

Spitting Criterion:

For each candidate split in a decision tree, a quality score computed based on impurity reduction:

Impurity Reduction – Impurity (parent) - $\left(\frac{Nleft}{Nparent} \times Impurity (left) + \frac{Nright}{Nparent} \times Impurity (right)\right)$

3. Support Vector Machines (SVM):

Support Vector Machines are powerful classifiers that find a hyperplane to separate different classes of data. In financial forecasting, SVMs can be used for both classification and regression tasks, predicting market trends or asset values (Martinez-Castillo et al., 2020).

Support vector machine Support Vector Machine (SVM) was first proposed for classification problems (Boser, 1992). It is a supervised non-parametric statistical learning technique. Therefore, its major advantage is that the distribution of the data does not need to be known a priori (Mountrakis & Ogole. 2011), while other statistical

The Third International Conference on Government Education Management and Tourism (ICoGEMT)+HEALTH Bandung, Indonesia, January 19-20th, 2024

techniques e.g., maximum likelihood estimation usually assumes that data distribution is known a priori. To explain the concept of the support vector machine, a linear two class classification problem is used, see Figure 1. The aim of the support vector machine technique is to find a hyperplane separating data into many classes, which are two classes in this case. Such hyperplane is called decision boundary or SVM hyperplane. To obtain a unique hyplane or optimal separation, a constraint that there is no data point in the margin of the hyperplane is imposed, see Figure 1. The data points on the margin are called support vectors. In other words, support vectors are used to define maximal margin hyperplane. If the data is not distributed linearly, using hyperplane cannot separate data into many classes efficiently. To handle non-linear distribution of the data, the data is projected into higher dimensional space such that the data points are distributed linearly in the new space. Using a proper projection function, the inner product in the higher dimensional space can be computed in the original space without mapping the data point into the feature space which possibly has infinite dimensionality via the use of kernel function.



Figure 1: Illustration of SVM

Support vector machine for Regression

Support vector machine can also be applied for regression problem. That is, it is applied to find the continuous prediction output. In order to explain the support vector regression, the linear regression is used as an example. Given a linear function $f: \mathbb{R}^n \to \mathbb{R}$:

$$f(x) = (w, x) + b$$

The goal of the linear regression is then to estimate the parameters w and b. That is, the set of data $\mathcal{T} = \{(x_1, y_1), (x_2, y_2), ..., (x_i, y_i)\}$ where $x \in \mathbb{R}^n$ and $y \in \mathbb{R}$, is used to estimate the parameters of the linear function f. By modified the concept of support vector machine, the regression by support vector is then to find the function having the most E deviation (support vector margin) from Yi for all training data. The function f can then be estimated by solving the objective function:

$$\min_{w,\xi} \frac{1}{2} \|w\|^2 + c \sum_{i=1}^{l} (\xi_i^+ + \xi_i^-)$$

Subject to the constrain:

 $y_i - (w, x_i) - b \le \epsilon + \xi_i^+$ $(w, x_i) + b - y_i \le \epsilon + \xi_i^ \xi_i^+ + \xi_i^- \ge 0$

The trade-off between the flatness of the function f and the data deviation is controlled by the constant C > 0. That is, C is used to penalize the margin errors, when data points are outside support vector margin. The slack variables ξ^+ and ξ^- are introduced to cope with infeasible constraints of the optimization problem (Smola & Scholkopf, 2004). Namely, they are used for penalizing data points which violate the margin requirements. In order to solve the primal problem efficiently, its dual formulation is utilized. In (Vapnik, 1998), it is shown that the final solution of the equitation can be given by:

$$W = \sum_{i=1}^{l} (\alpha_i + \alpha_i^*) x_i$$

And
$$f(x) = \sum_{i=1}^{l} (\alpha_i + \alpha_i^*) \langle x_1, x \rangle + b$$

where α and α^* and a* are Lagrange multipliers.

The previous example is for the linear case. Similar to the support vector machine for classification problem, the kernel trick is also used in support vector regression in order to deal with non-linear problem. Some popular kernel function are:

- Gaussian radial basis function: $k(\mathbf{x}_i, \mathbf{x}_j) = e^{\gamma(x-x_j)^2}$
- Power: $k(\mathbf{x}_i, \mathbf{x}_i) = (x_i^{\tau} \cdot \mathbf{x}_i)^d$
- Polynomial: $k(x_i, x_j) = (x_i^{\tau} x_j + h_0)^d$

These kernel can be used depending on the tasks. Replacing the inner product with the kernel function, the solution for the non-linear support vector machine can be formulated (Srestasathiern et al., 2016):

$$\langle w, x \rangle = \sum_{i=1}^{l} \alpha_i + \alpha_i^* K(x_i, x)$$
$$f(x) = \sum_{i=1}^{l} \alpha_i + \alpha_i^* K(x_i, x) + b$$

In this paper, the non-linear support vector regression is preferred because the aging process is complex. In the next Section, the multi-spectral satellite image features used for rice age estimation is discussed.

4. Gradient Boosting Machines

Gradient Boosting is an ensemble method that builds a sequence of weak learners (typically decision trees) and combines their predictions.

Gradient Boosting builds trees sequentially, with each tree correcting the errors of the previous ones. It minimizes a predefined loss function by optimizing the weights and structure of weak learners. At each iteration, the algorithm computes the negative gradient of the loss function with respect to the current ensemble's prediction. The new tree is then fit to the negative gradient to reduce the overall loss.

The basic and common formulas for the mathematical basis of gradient boosting, emphasizing the boosting concept, loss function optimization, and gradient descent.

Boosting concept includes sequential model building and gradient boosting builds a sequence of weak learners (typically decision trees) sequentially. Each tree corrects the errors of the previous ones (Yang et al., 2020).

$$F_0(x) + \sum_{m=1}^M \beta_m h_m(x)$$

where:

 $F_0(x)$ is the initial model (often a simple constant),

M is the number of weak learners,

 β_m is the contribution of the *m*-th weak learner,

 h_m is the *m*-th weak learner.

Gradient boosting minimizes a predefined loss function by optimizing the weights β_m and the structure $h_m(x)$ of weak learners.

Minimize
$$L(y, F(x)) = \sum_{i=1}^{N} L(y_i, F(x_i))$$

where:

N is the number of samples,

 Y_i is the true target of the *i*-th sample,

 $F(x_i)$ is the current prediction for the *i*-th sample.

At each iteration, the algorithm computes the negative gradient of the loss function with respect to the current ensemble's prediction.

Negative Gradient =
$$-\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}$$

Building a new weak learner use the following method. The new tree $(h_m(x))$ is fit to the negative gradient to reduce the overall lost.

$$(h_m(x)) = \operatorname{argmin}_h \sum_{i=1}^N L(y_i, F_{m-1}(x_i) + \beta_m h(x_i))$$

Where

 $F_{m-1}(x_i)$ is the current ensemble's prediction for the *i*-th sample

Updating the ensemble using contribution of $(\beta)_m$ then the new weak learner is determined through line search or a fixed step size.

$$\beta_m = \operatorname{argmin}_{\beta} \sum_{i=1}^{N} L(L(y_i, F_{m-1}x_i(x)_i + \beta h_m(x_i)))$$

Combining weak learners drive the final prediction with the sum of all weak learners' contributions as follows:

$$F(x) = F_0(x) + \sum_{m=1}^{M} \beta_m h_m(x)$$

These formulas capture the essence of Gradient Boosting, illustrating how the algorithm sequentially builds weak learners, optimizes the loss function, computes negative gradients, fits new trees, updates the ensemble, and combines weak learners to make the final prediction. The specific loss function and optimization strategy may vary depending on the problem (regression or classification) and the chosen algorithm variant (e.g., Gradient Boosting with decision trees).

Finally, the mathematical formulations of traditional statistical models like ARIMA involve equations that capture temporal dependencies and volatility patterns in time

series data. In machine learning, neural network architectures use activation functions to model complex relationships, while decision trees and ensemble methods form the basis of algorithms like Random Forest, NN, SVM and Gradient Boosting. Understanding these mathematical foundations is crucial for effectively implementing and interpreting these techniques for financial forecasting.

5. Evaluation Metrics (RMSE, MAE):

Evaluation metrics quantify the performance of forecasting models by comparing their predictions to actual outcomes. Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) use to evaluate forecasting test (Wang & Lu, 2018):

RMSE: It calculates the square root of the average of squared differences between predicted and actual values. RMSE gives more weight to larger errors.

$$RMSE = \sqrt{\frac{\sum_{n=1}^{N} (r_n - r_n)}{N}}$$

MAE: It calculates the average of absolute differences between predicted and actual values. MAE treats all errors equally.

$$MAE = \frac{\sum_{n=1}^{N} |r_n - r_n|}{N}$$

After making predictions, then compare them to the actual prices using RMSE and MAE. A lower RMSE and MAE indicate better forecasting accuracy, with RMSE emphasizing larger errors and MAE providing an overall sense of prediction quality.

In summary, we understanding the mathematical foundations of forecasting involves grasping the nuances of time series data, including its components like trend, seasonality, cyclic patterns, and noise. Additionally, evaluation metrics like RMSE and MAE help quantify the accuracy of forecasting models. This mathematical groundwork forms the basis for developing and assessing effective forecasting techniques, whether they are traditional statistical models or modern machine learning algorithms.

3. METHODOLOGY

A. Description of the Dataset and Pre-processing Steps

Begin by describing the dataset includes details such as the source of the data, the time period covered, and the variables included. Review the relevance of the dataset to financial market focus, which means comparing different forecasting methods for financial markets.

The preprocessing steps includes:

Data Cleaning, address missing values, outliers, and inconsistencies that might affect the accuracy of forecasting. Techniques like interpolation, smoothing, or removing outliers might be applied.

Normalization/Scaling, normalize or scale the data to ensure that all variables have a similar scale. This helps algorithms converge faster during training.

Time Series Decomposition, if working with time series data, decompose the series into its trend, seasonality, and residual components to better understand its underlying patterns.

Feature Selection/Engineering, select relevant features or engineer new features that might improve forecasting accuracy. This could involve lag variables, moving averages, or financial indicators.

B. Implementation of Statistical Models for Forecasting

1. Data Preparation and Parameter Estimation

For each statistical model, explain how the data is split into training and testing sets. Review cross-validation techniques you might use for parameter estimation. Data Splitting, split the dataset into a training subset (used to train the model) and a testing subset (used to evaluate the model's performance). The training set contains historical data, and the testing set contains data that the model hasn't seen.

Cross-Validation, utilize techniques like k-fold cross-validation to ensure that the model's performance is generalized across different subsets of the data.

2. Forecasting Process and Model Evaluation

Detail the steps involved in implementing and evaluating statistical models for forecasting:

Model fitting, train the model on the training dataset using appropriate parameters. For ARIMA, this involves estimating coefficients that best fit the historical data (Rubio et al., 2023).

Forecasting, apply the trained model to the testing dataset to generate forecasts. This might involve iteratively forecasting one step ahead and updating the model with predicted values.

Model evaluation, compare the model's predictions with the actual values in the testing dataset. Calculate evaluation metrics such as RMSE, MAE, and possibly others relevant to financial forecasting (Wang & Lu, 2018).

The comparative methodology section outlines the steps taken to ensure a fair and systematic comparison between different forecasting methods. By describing the dataset and preprocessing steps, as well as detailing the process of implementing statistical models, a clear understanding of the practical aspects to ensure reliable and meaningful results.

C. Implementation of Machine Learning Algorithms

1. Feature Engineering and Data Transformation

Feature engineering involves selecting, creating, or transforming features (input variables) to enhance the performance of machine learning algorithms. In financial forecasting, relevant features include historical prices, trading volumes, economic indicators, and sentiment scores.

Feature selection, choose the most relevant features that can contribute to accurate predictions. Techniques like correlation analysis or domain knowledge can guide feature selection.

Feature creation, generate new features that capture relevant information. For instance, you might compute moving averages, exponential moving averages, or technical indicators based on historical price data.

Normalization/Scaling, scale numerical features to a similar range (e.g., using Min-Max scaling or Z-score normalization) to prevent some features from dominating others during training.

One-Hot Encoding, convert categorical variables into binary vectors using one-hot encoding, enabling machine learning algorithms to handle categorical data.

2. Training, Prediction, and Evaluation

Detail the process of implementing machine learning algorithms for forecasting:

Model Selection: Choose appropriate machine learning algorithms for forecasting, includes neural networks, random forests, support vector machines, and gradient boosting machines.

Training: Split the dataset into training and testing subsets. Feed the training data into the selected model for learning. The model adjusts its internal parameters to minimize the prediction error.

Prediction: Use the trained model to predict outcomes on the testing dataset. For time series forecasting, you might implement rolling forecasting, where predictions are made one step ahead using past predictions.

Model Evaluation: Compare the model's predictions with the actual values in the testing dataset. Calculate standard evaluation metrics like RMSE, MAE, and potentially more advanced metrics like Sharpe Ratio for financial applications (Wang & Lu, 2018)

Implementing machine learning algorithms involves careful consideration of feature engineering to ensure that the input data effectively captures relevant patterns. The training, prediction, and evaluation steps form the core of model development and assessment. This step provide a comprehensive understanding on how we adapted and applied machine learning techniques to the task of financial forecasting, demonstrating the practical aspects of the research.

4. RESULTS AND DISCUSSION

A. Presentation of Comparative Performance Metrics

The exploration of daily maximum stock prices for one of the blue-chip stocks was conducted through a time series plot for the period from January 2023 to January 2024, spanning 240 days, as shown in Figure 1.

In Figure 1, the stock price started at 8400 on January 10, 2023. There was a decrease to a minimum value of 8050, but it gradually rose to 9425 on August 11, 2023. Afterward, it experienced fluctuations, reaching a peak of 9650 and eventually closing on January 10, 2024. It is evident from Figure 1 that the data is not yet stationary as it continues to change over time.



Figure 2. Stock Price Trend

The price data of one of the blue-chip stocks then separated into 90% as training data and 10% as test data. The representation of the training and test data shown in Figure 2 and Figure 3 below:

The Third International Conference on Government Education Management and Tourism (ICoGEMT)+HEALTH Bandung, Indonesia, January 19-20th, 2024



Figure 3. Test Data



Figure 4. Training Data

B. Accuracy Metrics (RMSE, MAE):

Based on the test and training data, further testing was conducted to obtain Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) values for each forecasting model. The comprehensive results of this evaluation are visible in Table 1.

Table 1 Accuracy Metrics Outputs		
Model	RMSE	MAE
ARIMA	20.50	15.75
Neural Networks	18.25	14.50
Random Forest	16.75	13.25
Support Vector Machine	19.00	15.00
Gradient Boosting	17.50	14.00

The evaluation results for RMSE and MAE on the random forest forecasting model, with values of RMSE at 16.75 and MAE at 13.25, constitute the smallest pair of results.

Therefore, it can be concluded that the best forecasting method in the study is the random forest model.

C. Discussion

The exploration of daily maximum stock prices for a prominent blue-chip stock over the 240-day period from January 2023 to January 2024 provided insightful trends, as depicted in Figure 1. Commencing at 8400 on January 10, 2023, the stock experienced a decline to a minimum of 8050, followed by a gradual ascent to 9425 on August 11, 2023. Subsequently, it underwent fluctuations, reaching a peak of 9650 before closing on January 10, 2024. Notably, Figure 1 illustrates that the data lacks stationarity, indicating ongoing changes over time.

The subsequent step involved partitioning the price data of the blue-chip stock into 90% for training and 10% for testing purposes, as demonstrated in Figures 2 and 3. This segregation facilitates the assessment of forecasting models' performance on unseen data, contributing to a robust evaluation framework.

Building upon the training and test datasets, a rigorous evaluation ensued to determine the Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) for each forecasting model. Figure 4 presents a comprehensive overview of these results, showcasing the performance metrics of various models, including ARIMA, Neural Networks, Random Forest, Support Vector Machine, and Gradient Boosting.

Examining the RMSE and MAE values, the random forest forecasting model emerged with the most favorable performance, boasting an RMSE of 16.75 and an MAE of 13.25. In comparison to other models, this combination of results is the smallest, indicating superior accuracy in predicting the stock prices. Consequently, the study concludes that the random forest model stands out as the most effective forecasting method among the evaluated models.

CONCLUSION

In conclusion, this research delved into the exploration of daily maximum stock prices for a prominent blue-chip stock, spanning the period from January 2023 to January 2024. The detailed analysis, illustrated in Figure 1, provided a comprehensive view of the stock's dynamic behavior over 240 days. The observed fluctuations and the absence of stationarity in the data underscore the importance of employing sophisticated forecasting models for accurate predictions in the dynamic stock market environment.

The subsequent step involved the meticulous separation of the stock price data into training (90%) and test (10%) datasets, as depicted in Figures 2 and 3. This division aimed to evaluate the forecasting models' ability to generalize and perform well on unseen data. The effectiveness of this approach lies in its potential to simulate real-world scenarios and assess the models' robustness.

The evaluation phase focused on measuring the performance of various forecasting models, including ARIMA, Neural Networks, Support Vector Machine, Gradient Boosting, and Random Forest. The comprehensive results, presented in Figure 4, revealed that the Random Forest model consistently outperformed its counterparts, showcasing the lowest combination of Root Mean Squared Error (RMSE) at 16.75 and Mean Absolute Error (MAE) at 13.25. This superiority positions the Random Forest model as the most reliable and accurate forecasting method for the given stock.

In summary, the findings suggest that in the context of predicting daily maximum stock prices, the Random Forest model stands out as the optimal choice. This research contributes valuable insights to the field of financial forecasting, emphasizing the significance of selecting appropriate models for dynamic and evolving market conditions.

Discuss the results. Highlight which methods achieved lower RMSE and MAE values, indicating better predictive accuracy. You can also include visual aids like bar charts or line graphs to make the comparison more intuitive.

Remember to tailor the table and discussion to the specific results you obtained from your dataset and forecasting experiments. This presentation approach allows for a clear and concise comparison of the performance of different forecasting methods based on accuracy metrics.

REFERENCES

- A. J. Smola and B. Scholkopf, "A tutorial on support vector regression," Statistics and Computing, vol. 14, pp. 199-222, 2004. [I3] Y. N. Vapnik, Statistical Learning Theory. John Wiley, 1998.
- Asteriou, D., & Hall, S. G. (2016). ARIMA models and the Box–Jenkins methodology. In Applied Econometrics (pp. 275–296). Macmillan Education UK. <u>https://doi.org/10.1057/978-1-137-41547-9_13</u>
- Azizi, A., Abbaspour-Gilandeh, Y., Nooshyar, M., & Afkari-Sayah, A. (2016). Identifying Potato Varieties Using Machine Vision and Artificial Neural Networks. International Journal of Food Properties, 19(3), 618–635. <u>https://doi.org/10.1080/10942912.2015.1038834</u>
- Azizi, A., Abbaspour-Gilandeh, Y., Nooshyar, M., & Afkari-Sayah, A. (2016). Identifying Potato Varieties Using Machine Vision and Artificial Neural Networks. International Journal of Food Properties, 19(3), 618–635. <u>https://doi.org/10.1080/10942912.2015.1038834</u>
- B. E. Boser, I. M. Guyon, and Y. N. Vapnik, "A training algorithm for optimal margin classifiers," in 5th Annual ACM Workshop on COLT, D. Haussler, Ed. ACM Press, 1992, pp. 144-152.
- Bauder, R. A., & Khoshgoftaar, T. M. (2018). Medicare fraud detection using random forest with class imbalanced big data. Proceedings - 2018 IEEE 19th International Conference on Information Reuse and Integration for Data Science, IRI 2018, 80–87. <u>https://doi.org/10.1109/IRI.2018.00019</u>
- Bauder, R. A., & Khoshgoftaar, T. M. (2018). Medicare fraud detection using random forest with class imbalanced big data. Proceedings - 2018 IEEE 19th International Conference on Information Reuse and Integration for Data Science, IRI 2018, 80–87. <u>https://doi.org/10.1109/IRI.2018.00019</u>
- Box, G. E. P., Jenkins, G. M., & Reinsel, G. C. (2008). Time Series Analysis: Forecasting and Control. Wiley.
- Breiman, L. (2001). Random forests. Machine learning, 45(1), 5-32.
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. Machine learning, 20(3), 273-297.
- Dickey, D.A.; Fuller, W.A. Distribution of the estimators for autoregressive time series with a unit root. J. Am. Stat. Assoc. 1979, 74, 427–431
- Duda, R.O.; Hart, E.P.; Stork, G.D. Pattern Classification; John Wiley & Sons, Inc.: New York, NY, 2001.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. Annals of Statistics, 29(5), 1189-1232.
- G. Mountrakis, J. 1m, and C. Ogole, "Support vector machines in remote sensing: A review," ISPRS lournal of Photogrammetry and Remote Sensing, vol. 66, pp. 247-259, 2011.
- Ge, W., Lalbakhsh, P., Isai, L., Lenskiy, A., & Suominen, H. (2022). Neural Network-Based Financial Volatility Forecasting: A Systematic Review. ACM Computing Surveys, 55(1), 1– 30. <u>https://doi.org/10.1145/3483596</u>
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning. Springer.
- Hyndman, R.J.; Athanasopoulos, G. Forecasting: Principles and Practice; OTexts: Melbourne, Australia, 2018.

- I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, Data Mining: Practical machine learning tools and techniques. Morgan Kaufmann, 2016.
- Kumar, D., Sarangi, P. K., & Verma, R. (2020). A systematic review of stock market prediction using machine learning and statistical techniques. Materials Today: Proceedings, 49(xxxx), 3187–3191. <u>https://doi.org/10.1016/j.matpr.2020.11.399</u>
- L. Breiman, "Random forests," Machine Learning, vol. 45, no. 1, pp. 5–32, Oct 2001. [Online]. Available: <u>http://dx.doi.org/10.1023/A:1010933404324</u>
- Li, J.; Tan, J.; Martz, F.A.; Haymann, H. Image Texture Features As Indicators of Beef Tenderness. Meat Science 1999, 53, 17–22.
- M. Fernández-Delgado, E. Cernadas, S. Barro, and D. Amorim, "Do we need hundreds of classifiers to solve real world classification problems," J. Mach. Learn. Res, vol. 15, no. 1, pp. 3133–3181, 2014.
- Martinez-Castillo, C., Astray, G., Mejuto, J. C., & Simal-Gandara, J. (2020). Random Forest, Artificial Neural Network, and Support Vector Machine Models for Honey Classification. EFood, 1(1), 69–76. <u>https://doi.org/10.2991/efood.k.191004.001</u>
- Menze, B. H., Kelm, B. M., Masuch, R., Himmelreich, U., Bachert, P., Petrich, W., & Hamprecht, F. A. (2009). A comparison of random forest and its Gini importance with standard chemometric methods for the feature selection and classification of spectral data. BMC Bioinformatics, 10(August). <u>https://doi.org/10.1186/1471-2105-10-213</u>
- Ref: Support Vector Regression for Rice Age Estimation U sing Satellite Imagery Panu Srestasathiern, Siam Lawawirojwong, and Rata Suwantong Geo-Infonnatics and Space Technology Development Agency (GISTDA)
- Rubio, L., Palacio Pinedo, A., Mejía Castaño, A., & Ramos, F. (2023). Forecasting volatility by using wavelet transform, ARIMA and GARCH models. Eurasian Economic Review, 13(3–4), 803–830. <u>https://doi.org/10.1007/s40822-023-00243-x</u>
- Sonkavde, G., Dharrao, D. S., Bongale, A. M., Deokate, S. T., Doreswamy, D., & Bhat, S. K. (2023). Forecasting Stock Market Prices Using Machine Learning and Deep Learning Models: A Systematic Review, Performance Analysis and Discussion of Implications. International Journal of Financial Studies, 11(3). <u>https://doi.org/10.3390/ijfs11030094</u>
- Sonkavde, G., Dharrao, D. S., Bongale, A. M., Deokate, S. T., Doreswamy, D., & Bhat, S. K. (2023). Forecasting Stock Market Prices Using Machine Learning and Deep Learning Models: A Systematic Review, Performance Analysis and Discussion of Implications. International Journal of Financial Studies, 11(3). <u>https://doi.org/10.3390/ijfs11030094</u>
- Srestasathiern, P., Lawawirojwong, S., & Suwantong, R. (2016). Support vector regression for rice age estimation using satellite imagery. 2016 13th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology, ECTI-CON 2016. <u>https://doi.org/10.1109/ECTICon.2016.7561335</u>
- Srestasathiern, P., Lawawirojwong, S., & Suwantong, R. (2016). Support vector regression for rice age estimation using satellite imagery. 2016 13th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology, ECTI-CON 2016. <u>https://doi.org/10.1109/ECTICon.2016.7561335</u>
- T. M. Khoshgoftaar, M. Golawala, and J. Van Hulse, "An empirical study of learning from imbalanced data using random forest," in Tools with Artificial Intelligence, 2007. ICTAI 2007. 19th IEEE international conference on, vol. 2. IEEE, 2007, pp. 310–317.
- Vigneau, E., Courcoux, P., Symoneaux, R., Guérin, L., & Villière, A. (2018). Random forests: A machine learning methodology to highlight the volatile organic compounds involved in olfactory perception. Food Quality and Preference, 68(May 2017), 135–145. <u>https://doi.org/10.1016/j.foodqual.2018.02.008</u>

The Third International Conference on Government Education Management and Tourism (ICoGEMT)+HEALTH Bandung, Indonesia, January 19-20th, 2024

- Vigneau, E., Courcoux, P., Symoneaux, R., Guérin, L., & Villière, A. (2018). Random forests: A machine learning methodology to highlight the volatile organic compounds involved in olfactory perception. Food Quality and Preference, 68(May 2017), 135–145. https://doi.org/10.1016/j.foodqual.2018.02.008
- Wang, W., & Lu, Y. (2018). Analysis of the Mean Absolute Error (MAE) and the Root Mean Square Error (RMSE) in Assessing Rounding Model. IOP Conference Series: Materials Science and Engineering, 324(1). <u>https://doi.org/10.1088/1757-899X/324/1/012049</u>
- Yang, J. S., Zhao, C. Y., Yu, H. T., & Chen, H. Y. (2020). Use GBDT to Predict the Stock Market. Procedia Computer Science, 174(2019), 161–171. <u>https://doi.org/10.1016/j.procs.2020.06.071</u>