# Two-stage optimization algorithm of random forest and its effectiveness in stock index prediction

**Jingdan Luo[1] and Yang Shen[1,*]**
[1] Guilin Institute of Information Technology, Guilin, Guangxi, 541000, China
Corresponding authors: (e-mail: arthur_ys@163.com).

**Abstract** A stock price index is an indicator reflecting the overall trend of the stock market, calculated through weighted averaging based on a selected sample of stocks. For investors, observing the fluctuations in stock price indices can help assess market sentiment and risk, predict future market trends, and formulate more informed investment strategies. This study introduces the particle swarm optimization (PSO) algorithm to optimize random forests, constructing a PSO-RF prediction model. Simulation experiments indicate that when the number of particles reaches 32, the model's evaluation metrics achieve optimal performance. When applying the PSO-RF algorithm to the selection of decision trees in ensemble forests, the quality of sub-forests of different sizes was evaluated using different diversity (or similarity) metrics. The PSO-RF algorithm achieved optimization effects for the random forest algorithm across all selected sub-forest sizes. Data from the CSI 500 Index from 2023 to 2027 was selected as the sample set. After validation and analysis in different experiments, the optimized random forest model demonstrated high prediction accuracy, strong stability, and good predictive performance on the CSI 500 Index across different time periods.

**Index Terms** Random Forest, Particle Swarm Optimization, Stock Index Prediction

## I. Introduction

In modern society, the economy is at the center of life, and for the economy, finance is its lifeblood[1]. Since China's reform and opening-up, government agencies and enterprises have gradually introduced various types of securities products such as stocks and bonds, which has driven China to develop and improve a systematic securities trading platform[2]-[4]. Among these, the stock price index is an indicator reflecting the overall trend of the stock market, calculated through weighted averaging based on selected stock samples[5]. For investors, observing fluctuations in stock price indices can help assess market sentiment and risk, predict future market trends, and formulate more informed investment strategies [6]-[8]. Additionally, stock price indices serve as important reference tools for financial institutions and government departments in regulating and formulating policies, thereby maintaining market stability and healthy development [9]-[11]. Stock market volatility is influenced by various complex factors, including investor sentiment, major public events, policy changes, and international economic trends [12], [13]. This results in stock price index fluctuations being more complex than those of individual stocks, exhibiting characteristics such as high noise, non-stationarity, and non-linearity [14], [15]. Therefore, to enhance the accuracy and stability of investor decision-making, it is necessary to use relevant quantitative research methods to establish mathematical models based on historical data to predict future trends [16].

For nonlinear, non-stationary, and rapidly updating financial market data, machine learning algorithms can quickly uncover more potential information in the market compared to traditional statistical analysis methods. For support vector regression models, literature [17] introduces a feature-weighted support vector machine theory for data classification, which calculates weights for different features and combines them with the K-nearest neighbor algorithm to construct a feature-weighted prediction model, thereby enabling accurate predictions of stock market indices. Literature [18] combines an investor sentiment analysis model with machine learning methods based on support vector machines. The sentiment indicators constructed can accurately reflect investors' psychological expectations in the stock market, providing reliable support for stock price predictions. Literature [19] utilizes web search technology to optimize the machine learning training dataset based on support vector regression and applies the fine-tuned model to stock prediction using time series data. This approach clearly reflects the parameter performance and associated risks of stock market fluctuations over time, thereby enhancing the accuracy of stock index predictions.

For long short-term memory networks, Literature [20] deploys long short-term memory (LSTM) networks in financial markets to perform time series predictions for stock indices, finding that compared to memory-less

classification methods, LSTM networks demonstrate higher prediction accuracy and lower exposure to risk sources. Literature [21] proposes a stock market prediction model that integrates LSTM networks and genetic algorithms, and predicts stock price indices based on the temporal characteristics of existing financial data, demonstrating high accuracy. Literature [22] indicates that the unique storage unit structure of LSTM can effectively avoid long-term dependency issues. Therefore, by combining the attention mechanism, a stock price prediction model based on financial time series data is established to provide investors with accurate stock recommendation guidance.

For extreme gradient boosting models, Reference [23] designed a stock market prediction model based on the extreme gradient boosting (XGBoost) algorithm, which improves the model's predictive effectiveness for stock price volatility trends by reformulating the prediction problem as a classification problem. Literature [24] constructs an interpretable XGBoost binary classification prediction model based on the sentiment characteristics of institutional, individual, and foreign investors, treating sentiment as an important factor in index trend prediction, thereby providing investors with valuable references. Literature [25] establishes a multi-factor stock selection model based on the XGBoost algorithm, identifying and analyzing the feature space of multi-indicator classifiers to obtain accurate predictions in complex nonlinear stock markets.

As a highly flexible ensemble algorithm that has emerged in recent years, the random forest algorithm has advantages such as being less prone to overfitting, strong noise resistance, no need for feature selection, the ability to balance errors and handle high-dimensional data, no need for data standardization, and fast training speed. Applying it to stock index prediction can provide a more precise price prediction foundation for quantitative timing strategy development compared to the aforementioned models.

This paper employs the particle swarm optimization algorithm to optimize the random forest, establishing a random forest prediction model based on the particle swarm optimization algorithm. The number of particles is determined through simulation experiments. Subsequently, the quality of sub-forests of different scales selected by the PSO-RF algorithm and their prediction results are evaluated. Data from the CSI 500 Index from 2021 to 2027 was selected as experimental data. Technical indicators were determined, and ADF and DW tests were used to examine their stationarity. Series that failed the stationarity test were subjected to differencing. Finally, the model was compared with numerous other prediction models, and its performance was analyzed from multiple perspectives.

## II. Random forest algorithm optimized based on swarm intelligence

### II. A. Random Forest Algorithm Principle

#### II. A. 1) Decision Tree

A decision tree is the most basic single classifier model and serves as the fundamental learning unit in a random forest. A decision tree is a multi-way tree composed of root nodes, intermediate nodes, and leaf nodes, with each node corresponding to an attribute of the target being studied. The training process involves progressively decomposing all samples from the root node to the leaf nodes. As each sample passes through an intermediate node, it essentially undergoes attribute classification, which is a unique process. Starting from the root node and passing through intermediate nodes, a sample can only land on a single leaf node, and the number of attributes is reduced by one. Ultimately, when all samples in an intermediate node belong to the same class, the node no longer performs further classification and automatically becomes a leaf node. By converting the path between the root node and leaf nodes in a decision tree into logical judgment rules, a judgment system composed of various judgment logics can be constructed, thereby achieving the functionality of classifying or predicting data samples.

#### II. A. 2) Bagging Algorithm

The Bagging (Bootstrap Aggregating) algorithm was proposed by Breiman. Its core principle involves using the Bootstrap sampling method to perform random sampling with replacement on the dataset. Each sample drawn yields a training dataset, and this process is repeated to generate multiple distinct training datasets. Each sample is assigned to a weak classifier for training, and the classification results are then aggregated. If the prediction is based on a classification algorithm, the weak classifier will select the classification with the most votes as the final classification [26]; if the prediction is based on a regression algorithm, the weak classifiers will calculate the arithmetic mean of their results to obtain the final model output.

#### II. A. 3) Random Feature Space Algorithm

In the process of constructing a random forest, there is another core algorithm proposed by Ho, known as the random feature space algorithm, also commonly referred to as the random subspace algorithm. Its most distinctive feature is the ability to randomly select feature subsets, thereby training multiple independent weak classifiers. During testing, the final classification or prediction result for the test sample is determined using a voting or arithmetic

average method. This ensemble learning strategy can improve the model's accuracy and generalization ability, reduce the risk of overfitting, and is applicable to a variety of classification and regression problems.

Similar to the Bagging algorithm, the Random Feature Space algorithm also uses bootstrap sampling to construct different weak classifiers. However, the former involves sampling training samples with replacement from the sample space, while the latter involves sampling training samples without replacement from the feature space. This is the key difference between the two.

### II. A. 4)    Random Forest Algorithm

The random forest algorithm [27] is improved on the basis of the Bagging algorithm, which combines the extraction of the training samples with return in the Bagging algorithm and the non-return extraction of attributes in the random feature space algorithm, absorbs the advantages of these two algorithms, and can search more comprehensively. This ensures that the training features in each decision tree are diverse, and can also accommodate various types and sizes of samples, and the random forest algorithm has better performance than the single Bagging algorithm and the random feature space algorithm.

(1) Random forest model

Assume that a random forest is composed of multiple decision trees $\{h(x,\theta_t), t=1,2,3,\cdots,T\}$ , where $X$ is the sample, $\theta_t$ is the random variable, and $T$ is the number of decision trees in the random forest. Then, the prediction result output by the random forest model is:

$$Y_{RF}(x) = \frac{1}{T}\sum_{i=1}^{t}\{h(x,\theta_t)\} \tag{1}$$

Among these, $Y_{RF}(x)$ represents the prediction result of the random forest algorithm for the test sample $x$ , and $h(x,\theta_t)$ represents the output based on $x$ and $\theta$ .

The random forest model has strong generalization and robustness, can handle high-dimensional data and a large number of features, and can still provide accurate prediction results for outliers in the data. It can also capture nonlinear relationships between input and output variables and provide information on the importance of each feature in the prediction model. For the nonlinear and abrupt characteristics of highway travel time, the random forest model has good applicability and stability.

(2) Random Forest Model Parameters

During the construction of the Random Forest model, the individual learners are not mutually dependent, and the number of parameters that need to be adjusted in the model is relatively fewer compared to other models, thereby accelerating the model training process. The main parameters involved include: the number of decision trees, the depth of the trees (also known as the number of layers), and the number of randomly sampled feature values.

(3) Random Forest Algorithm Process

The Random Forest algorithm process is as follows:

Step 1: Set the random forest parameters: pre-trained sample set $X$ , number of decision trees $T$ , number of feature attributes in the random feature space $m$ , and pruning threshold $\varepsilon$ .

Step 2: Sample the samples with replacement, randomly generate $T$ training sets, and select $X$ test samples from each training set.

Step 3: Select attributes for the node by randomly sampling $m$ attributes from all attribute values in the sample set as the attribute values for the node, then select the optimal attribute for splitting based on the CART algorithm's evaluation criteria.

Step 4: Train a decision tree using the training samples generated in Steps 2 and 3.

Step 5: Prune using the pruning threshold. If the number of samples in a node is less than the pruning threshold $\varepsilon$ , treat that node as a leaf node.

Step 6: If the number of decision trees obtained is less than $T$ , return to Step 2 to continue sampling and extraction. Otherwise, form a random forest using all $T$ decision trees and use the model to calculate the prediction results or final category.

### II. B. Random Forest Algorithm Based on Particle Swarm Optimization
### II. B. 1)    Particle Swarm Optimization Algorithm

The particle swarm algorithm simulates the foraging behavior of species in nature. This algorithm finds the global optimal solution to a problem through collaboration within the group. The algorithm is based on the social sharing of information, and its two core elements are the speed and position of the particles. Speed represents the direction and distance of the particle's next move, while position represents a solution to the problem being sought by the algorithm [28]. The algorithm's important parameters are as follows:

Assuming an $M$-dimensional search space, $N$ particles represent different solutions. Then:

(1) The current position of each particle $i$ is represented as (2).

$$X_{im} = (x_{i1}, x_{i2}, \ldots, x_{iM})$$ (2)

(2) The velocity of the particle determines its exploration direction and stride length; expressed as (3)

$$V_{im} = (v_{i1}, v_{i2}, \ldots, v_{iM})$$ (3)

(3) The individual optimal solution found by the particle; expressed as (4).

$$P_{im,pbest} = (p_{i1}, p_{i2}, \ldots, p_{iM})$$ (4)

(4) The optimal solution discovered jointly by all particles in the group; expressed as (5).

$$P_{m,gbest} = (p_{1,gbset}, p_{2,gbest}, \ldots, p_{M,gbest})$$ (5)

(5) The fitness value corresponding to the optimal solution associated with each particle $i$, i.e., the individual's historical optimal fitness value $f_p$ —individual historical optimal fitness value.

(6) The fitness value of the population's optimal solution, referred to as the population's historical optimal fitness value. $f_g$ —population historical optimal fitness value.

The steps of the particle swarm optimization algorithm are as follows:

Step 1: Initialization. First, determine the population size and maximum number of iterations. Next, initialize the position and velocity of all particles, and set the maximum and minimum values for these parameters.

Step 2: Calculate the initial fitness. Assign an initial fitness value to each particle as its optimal solution, and calculate the fitness function values for all particles.

Step 3: Dynamic update. Adjust the inertia weight of the particles and dynamically update their velocities and positions based on this weight.

Step 4: Fitness evaluation and update. Recalculate the fitness function values for each particle and update the individual's historical optimal position and the global optimal position of the entire population based on these values.

Step 5: Check termination conditions. Determine whether the algorithm meets the termination conditions (e.g., reaching the maximum number of iterations or the fitness function value stabilizing within a specified range). If not met, return to Step 3; otherwise, proceed to the final step.

Step 6: Output optimal parameters. When the algorithm meets the termination conditions, output the optimal solution, i.e., the optimal parameter combination for the random forest model.

### II. B. 2)    PSO-RF prediction model construction

Based on the above content, the construction process of the particle swarm optimization random forest model (PSO-RF) is as follows:

Step 1: Input the original data, preprocess the data, and divide the dataset into training and test sets.

Step 2: Select the input feature values for the random forest model.

Step 3: Set the initial velocity and coordinates of the particles, use the particle swarm algorithm to iteratively optimize the process, and adjust the movement velocity and position of each particle in the particle swarm based on this.

Step 4: Set the initial parameters of the random forest, construct decision trees based on these parameters, and generate the random forest. Evaluate the model performance using the test data and calculate its error rate.

Step 5: If the model's error rate does not meet the conditions, return to Step 3 and continue optimizing the random forest model parameters using the PSO algorithm. Once satisfactory model parameters are found or other termination conditions are met, the optimization process ends, and the PSO-RF prediction model is established, enabling stock index prediction.

### II. B. 3)    Determination of PSO-RF model parameters

In the application of particle swarm optimization algorithms, the number of particles is a critical parameter that must be determined initially. Increasing the number of particles can expand the search space and enhance the likelihood of finding the global optimal solution; however, this also means that the algorithm's runtime will increase. Generally, for conventional problems, the number of particles is recommended to be set between 20 and 40. For more complex problems, the number of particles may need to be adjusted to the range of 100 to 200. After extensive experimentation, this paper preliminarily decided to conduct simulation experiments of the particle swarm algorithm

at the values of 20, 25, 30, 35, and 40 to identify the optimal range of particle numbers. The experimental results are shown in Tables 1 and 2.

As shown in Tables 1-2, when the number of particles is set within the range of 20 to 40 and increased in increments of 5, it is observed that the model's prediction performance evaluation metrics gradually decrease as the number of particles increases. Specifically, when the number of particles is 30 and 35, the evaluation metric values are relatively close, but when the number of particles increases to 40, the evaluation metric values actually increase. This suggests that the optimal number of particles may be between 30 and 35. When the particle count reaches 32, the model's evaluation metrics achieve their optimal values, specifically 0.026, 0.085, and 0.094. Therefore, after comprehensive analysis, it is determined that the most suitable particle count for this model is 32.

Table 1: Prediction of the number of different particles (span 5)

| Number of particles | MSE | MAE | RMSE |
|---|---|---|---|
| 20 | 0.091 | 0.185 | 0.192 |
| 25 | 0.078 | 0.159 | 0.189 |
| 30 | 0.065 | 0.133 | 0.135 |
| 35 | 0.049 | 0.098 | 0.116 |
| 40 | 0.062 | 0.117 | 0.153 |

Table 2: Prediction of the number of different particles (span 1)

| Number of particles | MSE | MAE | RMSE |
|---|---|---|---|
| 30 | 0.073 | 0.132 | 0.135 |
| 31 | 0.057 | 0.115 | 0.129 |
| 32 | 0.026 | 0.085 | 0.094 |
| 33 | 0.042 | 0.093 | 0.105 |
| 34 | 0.045 | 0.096 | 0.115 |
| 35 | 0.049 | 0.098 | 0.116 |

## II. C.Experiments and Analysis of Results

### II. C. 1)  Experimental Preparation

The dataset used in the experiment is a binary classification dataset, all of which are sourced from the UCI database, with detailed descriptions as shown in Table 3. Preprocessing was performed on the dataset in the experiment: missing data were removed from the dataset; to avoid the influence of different scales between features on the experimental results, the data were standardized using the maximum-minimum method.

The random forest algorithm used in the experiment employs CART decision trees as base classifiers, with 200 decision trees and the number of candidate features at each node set to $mtry = \sqrt{m}$ is the number of features in the dataset). PSO-RF algorithm parameter settings: ant colony size $m$ = 10, pheromone importance coefficient $\alpha$ = 1, heuristic factor importance coefficient $\beta$ = 3, pheromone volatility coefficient $\rho$ = 0.2, total pheromone release quantity $Q$ = 20, maximum iteration count $t\_\max$ = 15.

Table 3: Ci data set used in the experiment

| Data set name | Sample size | Characteristic number |
|---|---|---|
| Breastcancer | 289 | 11 |
| Heart | 272 | 15 |
| Ionosphere | 355 | 36 |
| Sonar | 207 | 62 |
| ILPD | 574 | 13 |
| Australian | 694 | 16 |

The experimental code was implemented using Matlab 2014a. The experimental environment was a Win7 64-bit operating system with 4GB of memory and an Intel® Core™ i7-3520M CPU@2.90GHz.

## II. C. 2)　Verifying the quality of subforests

Evaluate the quality of the subforests selected by the PSO-RF algorithm. The similarity (or diversity) between decision trees in the ensemble forest is measured using the Kappa metric. The size M of the subforests ranges from [40, 80], increasing by increments of 10 decision trees. Experimental evaluation metrics: prediction accuracy (ACC); average similarity of the forest $\bar{\rho}$; and forest size EA. Experimental method: 10-fold cross-validation was applied to each dataset, and the average value of the sub-forest evaluation metrics was taken as the final result. The experimental results are shown in Tables 4 and 5. Symbol explanations in the tables: Original random forest $RF^{All}$; Sub-forest $RF^A$ integrated from high-performance decision trees selected from $RF^{All}$; Sub-forest $RF_M^B$ of size M selected from $RF^A$.

Table 4 shows the average size of the sub-forest $RF^A$ integrated from the decision trees with high predictive performance selected from $RF^{All}$ by the PSO-RF algorithm. It is known that the size of $RF^{All}$ in the experiment is 200. It can be seen that the size of $RF^A$ after selection is approximately half of the original size. At this point, the strength of the decision trees integrated in $RF^A$ is generally higher than that in $RF^{All}$. Table 5 shows the quality of subforests of different sizes selected by the PSO RF algorithm. It can be seen that, across the six datasets, the average similarity of $RF^A$ has increased compared to $RF^{All}$; however, the prediction accuracy of $RF^A$ has only increased on the Sonar and Ionosphere datasets, while it has decreased on the remaining datasets. Therefore, if only the strength of the decision trees in the ensemble forest is increased without considering the diversity of the trees in the ensemble forest, it may not only fail to improve the performance of the forest but may even reduce it.

Observe the quality of the sub-forest $RF_M^B$ selected by the PSO-RF algorithm in Table 5. It can be seen that the prediction accuracy of $RF_M^B$ is higher than that of $RF^{All}$ and $RF^A$ on all six datasets, and the average similarity of $RF_M^B$ is lower than that of $RF^A$ on average. For different scales of $RF_M^B$, it can be seen that except for the Heart and Ionosphere datasets, the prediction accuracy of the remaining datasets reaches its highest value when the forest scale is M=40. Among these, the Sonar dataset shows the greatest improvement, with the prediction accuracy of $RF_{40}^B$ increasing by 11.4% compared to $RF^{All}$ and by 9.07% compared to $RF^A$ The smallest improvement was observed in the Ionosphere dataset, where the prediction accuracy reached its peak at a forest size of M=70, improving by 2.45% compared to $RF^{All}$ and by 1.79% compared to $RF^A$.

Therefore, we conclude that the sub-forest $RF_M^B$ selected by the PSO-RF algorithm contains fewer trees, not only achieving higher prediction accuracy but also having a lower average similarity than $RF^A$, thus demonstrating that the PSO-RF algorithm achieves optimization of the random forest algorithm.

Table 4: The size of the selected Subforest PSO-RF of the experimental data set

| Data set | Breastcancer | Heart | Ionosphere | Sonar | ILPD | Australian |
|---|---|---|---|---|---|---|
| EA | 105.5 | 106.6 | 107.4 | 100.2 | 103.0 | 109.3 |

Table 5: The quality of the subforests of different sizes

| | Breastcancer | | Heart | | Ionosphere | |
|---|---|---|---|---|---|---|
| | ACC | $\bar{\rho}$ | ACC | $\bar{\rho}$ | ACC | $\bar{\rho}$ |
| $RF^{All}$ | 0.8110 | 0.3379 | 0.8544 | 0.4695 | 0.9537 | 0.6951 |
| $RF^A$ | 0.7964 | 0.3951 | 0.8507 | 0.5504 | 0.9594 | 0.7726 |
| $RF_{40}^B$ | 0.8398 | 0.3771 | 0.8989 | 0.5435 | 0.9708 | 0.7622 |
| $RF_{50}^B$ | 0.8325 | 0.3912 | 0.9026 | 0.5382 | 0.9737 | 0.7632 |
| $RF_{60}^B$ | 0.8253 | 0.3833 | 0.8878 | 0.5463 | 0.9737 | 0.7666 |
| $RF_{70}^B$ | 0.8217 | 0.3826 | 0.8915 | 0.5440 | 0.9765 | 0.7683 |
| $RF_{80}^B$ | 0.8181 | 0.3902 | 0.8878 | 0.5451 | 0.9708 | 0.7681 |
| | Sonar | | ILPD | | Australian | |
| $RF^{All}$ | 0.8295 | 0.3032 | 0.7092 | 0.2101 | 0.8797 | 0.6046 |
| $RF^A$ | 0.8491 | 0.3847 | 0.7041 | 0.2297 | 0.8783 | 0.6724 |

| | | | | | | |
|---|---|---|---|---|---|---|
| $RF_{40}^{B}$ | 0.9355 | 0.3801 | 0.7663 | 0.2263 | 0.9087 | 0.6603 |
| $RF_{50}^{B}$ | 0.9210 | 0.3742 | 0.7628 | 0.2268 | 0.9044 | 0.6642 |
| $RF_{60}^{B}$ | 0.9160 | 0.3768 | 0.7507 | 0.2274 | 0.9015 | 0.6671 |
| $RF_{70}^{B}$ | 0.9114 | 0.3797 | 0.7508 | 0.2271 | 0.8986 | 0.6683 |
| $RF_{80}^{B}$ | 0.8969 | 0.3834 | 0.7438 | 0.2284 | 0.8913 | 0.6677 |

## II. C. 3) Verification of results obtained using different measurement methods

This chapter introduces three different measurement methods: Kappa, Q-statistic, and dissimilarity measure. In the PSO-RF algorithm, these three methods are used to measure the similarity (or diversity) between decision trees in the ensemble forest. The prediction accuracy (ACC) is used as the evaluation metric for the experiment, with the remaining settings consistent with those in the previous section. Figure 1 shows the average prediction accuracy of the selected sub-forests $RF_M^B$ using the PSO-RF algorithm with the three different measurement methods. In the figure, the horizontal axis of the six bar charts starts at the prediction accuracy of $RF^{All}$ corresponding to the dataset or a higher value. It is evident that, across the six datasets, the PSO-RF algorithm effectively measures the similarity (or diversity) between decision trees in the ensemble forest using the three methods, achieving optimization of the random forest algorithm.
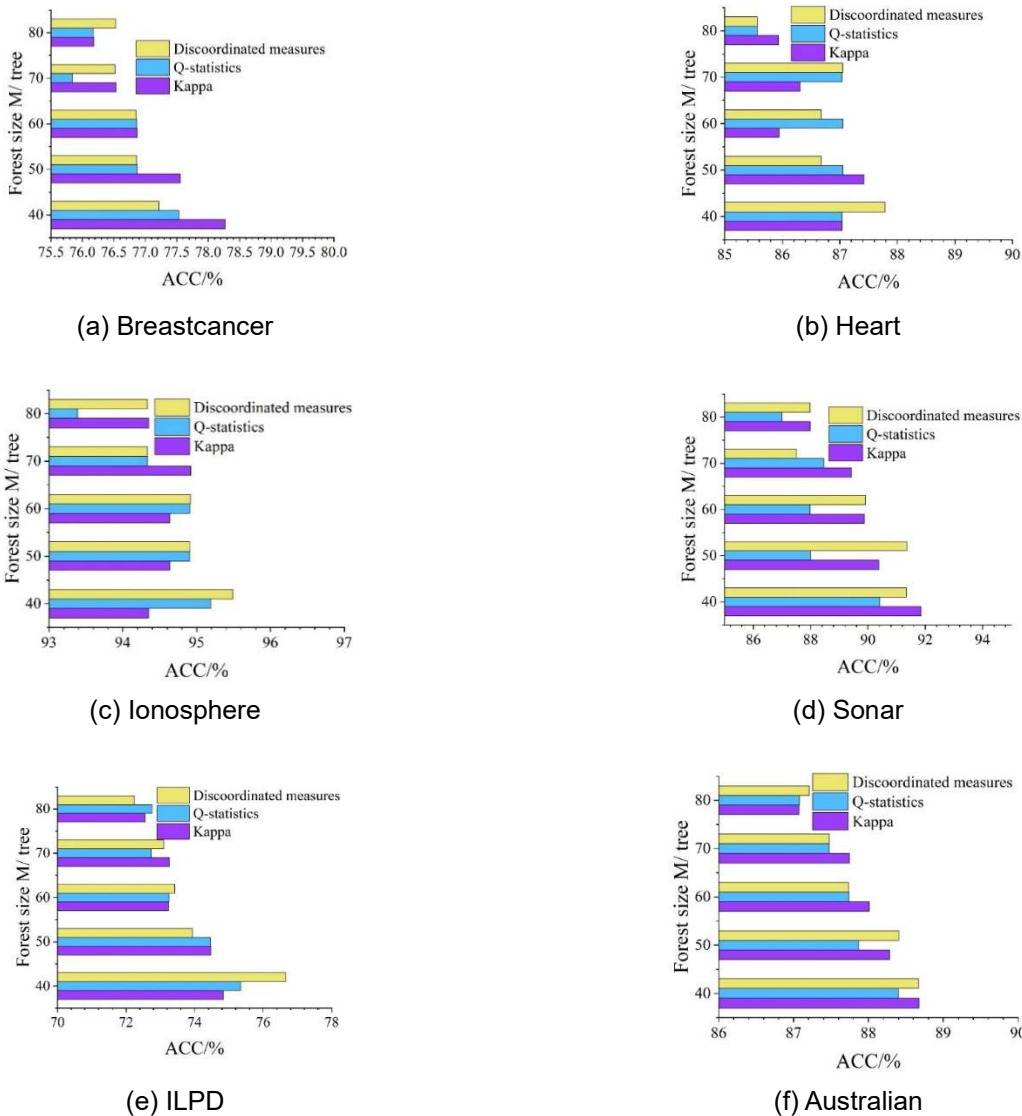


(a) Breastcancer

(b) Heart

(c) Ionosphere

(d) Sonar

(e) ILPD

(f) Australian

Figure 1: Mean prediction accuracy

## III. Stock index prediction based on improved random forests

### III. A. Data Selection and Explanation

The experiment selected commonly used technical indicators in the stock market as the basis for analysis and used an improved random forest model to train the CSI 500 index dataset.

### III. A. 1) Technical Specifications

Technical indicators are calculated based on stock index time series data and can be used to measure the profitability, solvency, and operational capabilities of the constituent stocks included in the CSI 500 Index. The technical indicators selected below are relatively widely used by investors and investment institutions at present, covering fundamental analysis indicators and dynamic indicators, and can assess the trend of stock index data changes from multiple perspectives. The specific definitions of technical indicators are shown in Table 6.

Table 6: Technical index information list

| Serial number | Technical index | For short |
|---|---|---|
| 1 | Market value | CMV |
| 2 | Return on equity | ROE |
| 3 | Net profit | ROA |
| 4 | Mobility ratio | CR |
| 5 | P/e | PE |
| 6 | Market net | PB |
| 7 | Earnings per share | EPS |
| 8 | Dividend yield | DY |
| 9 | Relative strength and strength index | RSI |
| 10 | The exponential smoothness is the average | MACD |
| 11 | Momentum line index | MOM |

### III. A. 2) Experimental Data

Select data from the CSI 500 Index to form the training sample set and test sample set. The top 500 stocks by total market capitalization constitute the sample stocks of the CSI 500 Index, reflecting the comprehensive changes in stock value of small and medium-sized companies. The CSI 500 Index differs from others in its constituent stock composition, as its sample stocks primarily consist of small and medium-cap stocks with strong growth potential and significant development prospects.

The experimental data spans from January 4, 2023, to December 31, 2027. Technical indicator data from the past three years is used to form the experimental sample space, with each trading day's data serving as a single sample for calculation. The data from January 4, 2023, to December 31, 2027, is used as the training sample set for model training and parameter adjustment, while the data from January 4, 2027, to December 31, 2027, is used as the test data set to evaluate the model's predictive performance.

### III. A. 3) Evaluation Indicators

(1) Mean Absolute Error (MAD)

The Mean Absolute Error (MAD) is the average of the absolute values of the differences between the sample label values and the model output values.

$$MAD = \frac{\sum_{i=1}^{n} |Y_i - F_i|}{n} \tag{6}$$

(2) Mean square error

Mean square error (MSE) is calculated by taking the average of the squared differences between the sample label values and the model output values.

$$MSE = \frac{\sum_{i=1}^{n} (Y_i - F_i)^2}{n} \tag{7}$$

(3) Mean Absolute Percentage Error (MAPE)

MAPE is calculated by first summing the absolute values of the differences between the sample label values and the model output values, divided by the model output values, and then taking the average. MAPE reflects the relative magnitude of prediction errors and eliminates the influence of measurement units.

$$MAPE = \frac{\sum_{i=1}^{n}\left(\frac{|Y_i - F_i|}{Y_i} \times 100\right)}{n} \tag{8}$$

### III. A. 4)　Modeling Process

Using an improved random forest model to train a financial time series dataset, the main process can be divided into two parts: preprocessing of financial time series data and construction of a regression random forest model. In the preprocessing of financial time series data, the sequence data must first undergo a stationarity test. If the test is passed, a pure randomness test (i.e., white noise test) can be conducted. If the test is not passed, differential processing must be performed, followed by another test to determine whether the sequence is a stationary non-white noise sequence. Sequences that pass the test can then be input into the improved random forest model for training and prediction. The flowchart for stock index prediction modeling is shown in Figure 2.
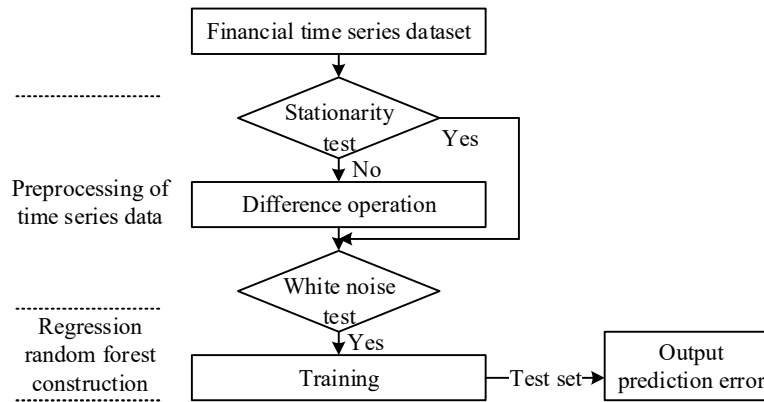


Figure 2: Index prediction modeling flow chart

### III. B.　Preprocessing of stock index data

### III. B. 1)　Stability test

(1) Unit root test

　　The ADF test [29] is an extension of the DF test, which is more suitable for small sample tests and can avoid the impact of high-order lag correlation in time series. Therefore, the ADF test is selected for the following tests.

　　The DW test is used to test the first-order correlation of time series residual values. Assuming that the residuals are $e_t$, the correlation equation for each residual is expressed as $e_i = \rho e_{i-1} + v_i$.

　　The results of the unit root test for the CSI 500 are shown in Table 7. Among the original indicators, only the ADF test p-value for the CSI 500's DY indicator is less than the significance level of 0.05, allowing us to reject the null hypothesis and conclude that the CSI 500's DY indicator series is stationary. The p-values of the ADF tests for the remaining indicators are all much greater than 0.05, meaning that the series did not pass the unit root test and the null hypothesis that the time series contains a unit root cannot be rejected. Therefore, the original series is not stationary and requires further differencing. We performed first-order differencing on the indicators and then conducted ADF tests again. As shown in the table, the differenced indicator data all passed the ADF test, meaning that the differenced indicator series is stationary.

Table 7: Comparison of unit root test of CSI 500

|  | Perform an ADF test for p values | DW statistics | The p value of the ADF test after difference is one | Differenced DW statistics |
|---|---|---|---|---|
| CMV | 0.8709 | 1.9672 | 0.0005 | 2.0031 |
| ROE | 0.5259 | 2.0599 | 0.0065 | 2.051 |
| ROA | 0.7845 | 1.9398 | 0.0100 | 2.0001 |
| CR | 0.2484 | 1.9882 | 0.0138 | 2.0018 |
| PE | 0.6763 | 1.8702 | 0.0084 | 2.0036 |
| PB | 0.3558 | 1.9998 | 0.0003 | 2.0019 |
| EPS | 0.3603 | 1.9932 | 0.0022 | 2.0014 |
| DY | 0.0512 | 2.0351 | - | - |

| RSI | 0.777 | 1.9818 | 0.0006 | 2.0023 |
| MACD | 0.7792 | 1.9681 | 0.0004 | 2.0024 |
| MOM | 0.7731 | 1.9808 | 0.0075 | 2.0024 |

(2) Auxiliary verification

For the technical indicators CMV, ROE, and ROA in the CSI 500, the autocorrelation coefficient and partial autocorrelation coefficient with a maximum lag period of 3 periods are calculated, and the autocorrelation plot and partial autocorrelation plot are shown in Figure 3. The autocorrelation coefficient plots in the figure are all first-order truncated, and the partial autocorrelation plots of the differenced technical indicator sequences all exhibit a tailing characteristic. Therefore, it can be verified that the time series data of each technical indicator after differencing is stationary and can be further analyzed.
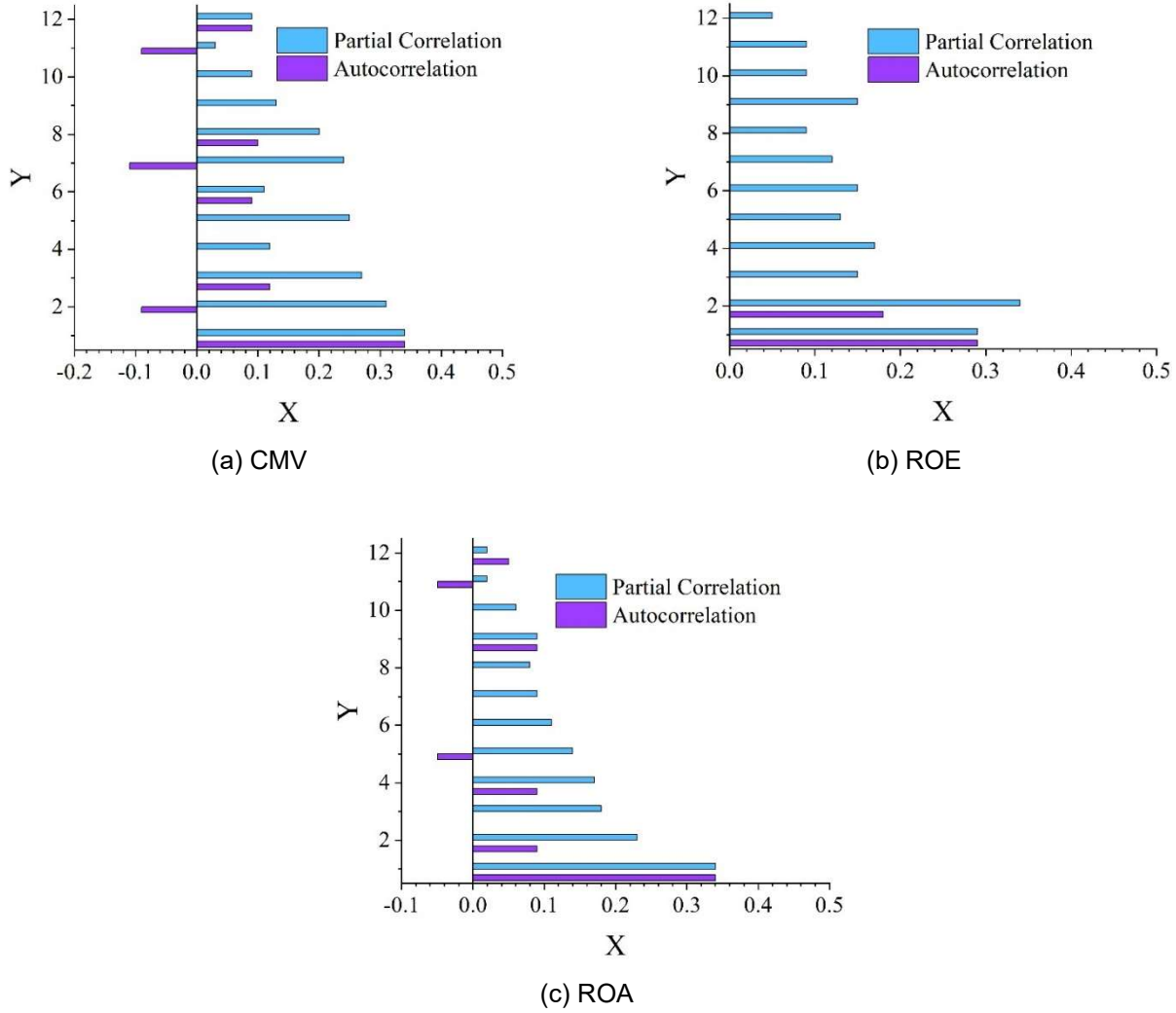


(a) CMV

(b) ROE

(c) ROA

Figure 3: Correlation and partial correlation diagram of CSI 500 index

### III. B. 2) Pure randomness test

Pure randomness tests generally select the $Q_{BP}$ statistic and the $Q_{LB}$ statistic. The $Q_{LB}$ statistic is a modified version of the $Q_{BP}$ statistic, and $Q_{LB}$ is more suitable for small sample tests. Therefore, based on the characteristics of the dataset, the following experiment selects the $Q_{LB}$ statistic as the test statistic. The formula for calculating the $Q_{LB}$ statistic is as follows:

$$Q_{LB} = n(n+2)\sum_{k=1}^{h}\frac{\hat{\rho}_k^2}{n-k} \tag{9}$$

where $n$ is the sample size, $h$ is the number of sequences, $\hat{\rho}_k^2$ is the correlation coefficient of the $k$ th lag of the sample, and the statistic follows a chi-square distribution with $h$ degrees of freedom. At a given significance level $\alpha$, the rejection region is $Q_{LB} > \chi^2_{1-\alpha,h}$, i.e., the sequence is a white noise sequence at this point.

The maximum lag order of 1–12 for the $Q_{LB}$ statistic after calculating the difference for each indicator in the CSI 500 is shown in Table 8, along with the mean and range of the statistic. As shown in the table, the p-values for the randomness test are all less than 0.05 when retaining four decimal places, so they are not listed in the table. Since all p-values are less than 0.05, the null hypothesis can be rejected, meaning the sequence is not white noise. The data of each index sequence is not purely random, and further model training can be conducted to extract effective information and relevant trends from the sequence.

Table 8: LB statistical range

| Index | Certificate 500 |
|---|---|
| CMV | 198.232±5.83 |
| ROE | 117.162±1.724 |
| ROA | 165.441±2.901 |
| CR | 199.159±9.874 |
| PE | 171.215±2.572 |
| PB | 179.643±2.875 |
| EPS | 195.738±3.214 |
| DY | 169.223±0.146 |
| RSI | 199.674±5.173 |
| MACD | 192.291±5.011 |
| MOM | 199.183±5.343 |

### III. C. Empirical Analysis

The loss of the PSO-RF stock index prediction model during the training process is shown in Figure 4, and the results of the model's actual values and predicted values on the test set are shown in Figure 5. As can be seen from Figure 4, during the first few rounds of training, the RMSE and loss showed an upward trend, but during the last few rounds of training, they began to show a downward trend and finally slowly leveled off. In Figure 5, the purple line represents the actual closing price data, and the blue line represents the predicted values from the PSO-RF model. It can be observed that in some samples, the predicted values are slightly higher than the actual values, in some samples they are slightly lower, and in some samples they are closely aligned with the actual values. Overall, the fitting performance is superior to that of the RF model, indicating that the stock price index prediction model proposed in this paper, which utilizes particle swarm optimization to enhance the random forest algorithm, demonstrates excellent fitting performance.



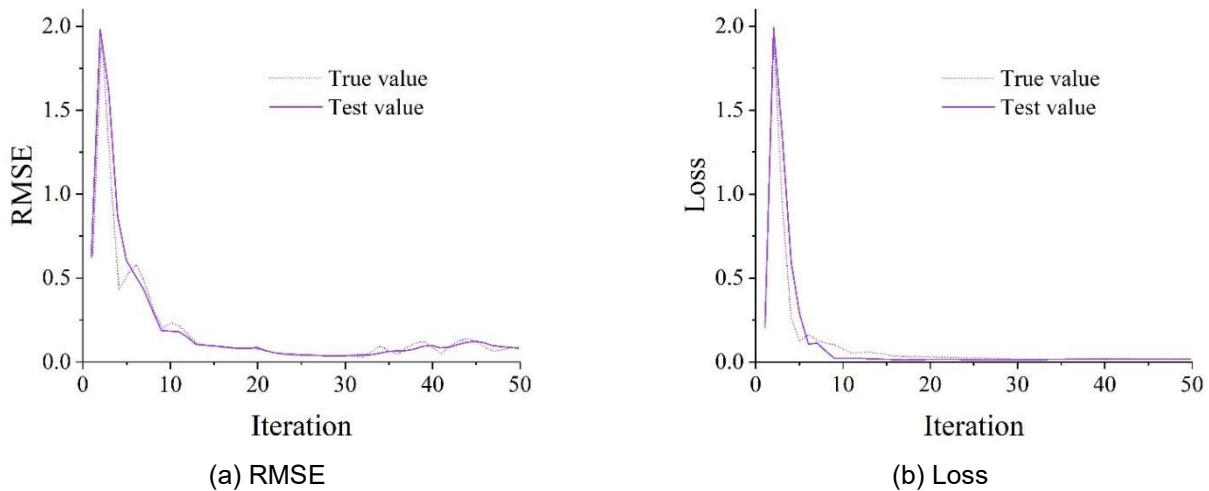(a) RMSE　　　　　　　　　　　　(b) Loss

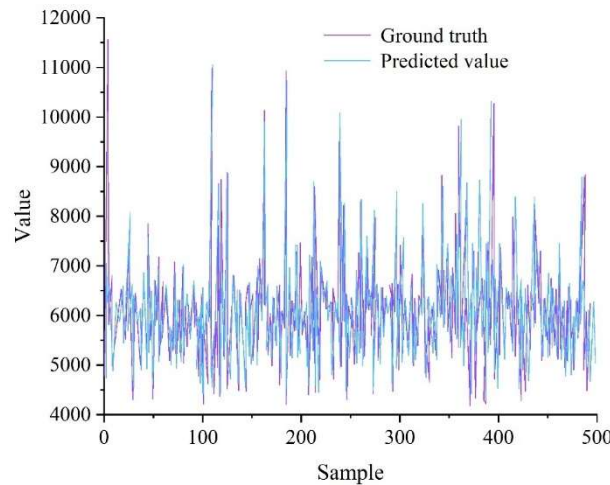Figure 4: Loss diagram during training

Figure 5: PSO-RF model prediction result figure

The performance of the PSO-RF stock index prediction model was validated on datasets from different time periods of the CSI 500 Stock Price Index, as shown in Table 9. During the period from January 4, 2023, to December 31, 2024, the PSO-RF prediction model achieved a $R^2$ value of the PSO-RF prediction model's evaluation metrics was 98.94%, slightly higher than the $R^2$ value over the four-year period from January 4, 2021, to December 21, 2027. However, its MAE and RMSE values were higher than those of the latter. Overall, the six-year data prediction results were closer to the actual values. Over the four-year period from January 4, 2021, to December 31, 2024, the PSO-RF prediction model's evaluation metric $R^2$ value of the PSO-RF prediction model's evaluation metrics was 99.47%, which is 0.51 percentage points higher than the $R^2$ value for the four-year period. Both the MAE and RMSE values are significantly lower than those of the latter, indicating that the PSO-RF prediction model performs optimally in the four-year period. However, overall, it demonstrates good performance on the CSI 500 Index dataset.

Table 9: PSO-RF prediction effect in different time periods

| Time | $R^2$ | RMSE | MAE |
|------|-------|------|-----|
| Two years | 98.94% | 178.42 | 151.07 |
| Four years | 99.47% | 67.39 | 55.22 |
| Six years | 98.86% | 125.41 | 91.43 |

Table 10 presents the evaluation metrics results for several common regression-based stock index prediction models. Among the single prediction models—LR, SVR, GRU, LSTM, and TCG—the GRU model has the highest $R^2$ value, the lowest RMSE and MAE values, and the best prediction performance; the SVR model, on the other hand, has the worst prediction performance; while the LR, LSTM, and TCG prediction models fall between the two, with the time-convolutional network (TCG) having a higher accuracy rate. The GRU network model is an effective variant of the LSTM network, with fewer parameters while retaining the long-short-term memory functionality of the LSTM network. This demonstrates that the LSTM network and its optimized models perform well in financial time series prediction. In contrast, the LR and SVR models have lower prediction accuracy compared to the LSTM and GRU models.

For the CNN-LSTM and WOA-LSTM models, the accuracy of the proposed prediction model is the highest, while the CNN-LSTM model performs the worst, with an $R^2$ value of 93.66%, an RMSE value of 239.22, and an MAE value of 179.22. The GRU model has an $R^2$ value of 98.12%, an RMSE value of 141.24, and an MAE value of 90.72, which are respectively 0.81 percentage points lower, 16.8 higher, and 0.46 lower than the PSO-RF stock index prediction model. Overall, the PSO-RF model demonstrates superior prediction performance and achieves the best fit with stock index data among the aforementioned single and combined models.

Table 10: Different prediction model regression effect

| Prediction model | $R^2$ | RMSE | MAE |
|---|---|---|---|
| LR | 93.74% | 237.48 | 91.58 |
| SVR | 91.25% | 311.36 | 253.34 |
| GRU | 98.12% | 141.24 | 90.72 |
| LSTM | 96.47% | 189.75 | 129.54 |
| TCG | 97.69% | 162.11 | 103.63 |
| CNN-LSTM | 93.66% | 239.22 | 179.22 |
| WOA-LSTM | 98.52% | 160.36 | 131.46 |
| PSO-RF | 98.93% | 124.45 | 91.18 |

## IV. Conclusion

The study proposes a random forest model based on the particle swarm optimization (PSO) algorithm, adjusting key parameters within the random forest model and setting the number of particles to 32. Experimental validation demonstrates that the PSO-RF algorithm reduces the forest size while enhancing prediction accuracy, achieving the desired optimization effects for the random forest algorithm. Using the optimized random forest model to conduct empirical analysis on the CSI 500 index data, the ADF test and DW statistic method were employed to examine the stationarity of the selected indicators. The results showed that most indicators were greater than 0.05, failing to reject the null hypothesis, necessitating differential processing. Those that passed the test were used as feature inputs. The empirical results indicate that the index prediction model proposed in this paper achieves excellent results in processing the CSI 500 time series data.

## References

[1] Yang, R., Yu, L., Zhao, Y., Yu, H., Xu, G., Wu, Y., & Liu, Z. (2020). Big data analytics for financial Market volatility forecast based on support vector machine. International Journal of Information Management, 50, 452-462.

[2] Nti, I. K., Adekoya, A. F., & Weyori, B. A. (2020). A comprehensive evaluation of ensemble learning for stock-market prediction. Journal of Big Data, 7(1), 20.

[3] Yu, Y., & Kim, Y. J. (2019). Two-dimensional attention-based LSTM model for stock index prediction. Journal of Information Processing Systems, 15(5), 1231-1242.

[4] Shah, D., Isah, H., & Zulkernine, F. (2019). Stock market analysis: A review and taxonomy of prediction techniques. International Journal of Financial Studies, 7(2), 26.

[5] Singh, R., & Srivastava, S. (2017). Stock prediction using deep learning. Multimedia Tools and Applications, 76, 18569-18584.

[6] Chiang, W. C., Enke, D., Wu, T., & Wang, R. (2016). An adaptive stock index trading decision support system. Expert Systems with Applications, 59, 195-207.

[7] Baek, Y., & Kim, H. Y. (2018). ModAugNet: A new forecasting framework for stock market index value with an overfitting prevention LSTM module and a prediction LSTM module. Expert Systems with Applications, 113, 457-480.

[8] Qiu, M., & Song, Y. (2016). Predicting the direction of stock market index movement using an optimized artificial neural network model. PloS one, 11(5), e0155133.

[9] Kyoung-Sook, M. O. O. N., & Hongjoong, K. I. M. (2019). Performance of deep learning in prediction of stock market volatility. Economic Computation & Economic Cybernetics Studies & Research, 53(2).

[10] Gandhmal, D. P., & Kumar, K. (2019). Systematic analysis and review of stock market prediction techniques. Computer Science Review, 34, 100190.

[11] Patel, J., Shah, S., Thakkar, P., & Kotecha, K. (2015). Predicting stock and stock price index movement using trend deterministic data preparation and machine learning techniques. Expert systems with applications, 42(1), 259-268.

[12] Liu, X., Guo, J., Wang, H., & Zhang, F. (2022). Prediction of stock market index based on ISSA-BP neural network. Expert Systems with Applications, 204, 117604.

[13] Jiang, M., Liu, J., Zhang, L., & Liu, C. (2020). An improved Stacking framework for stock index prediction by leveraging tree-based ensemble models and deep learning algorithms. Physica A: Statistical Mechanics and its Applications, 541, 122272.

[14] Sheta, A. F., Ahmed, S. E. M., & Faris, H. (2015). A comparison between regression, artificial neural networks and support vector machines for predicting stock market index. Soft Computing, 7(8), 2.

[15] Wang, C., Chen, Y., Zhang, S., & Zhang, Q. (2022). Stock market index prediction using deep Transformer model. Expert Systems with Applications, 208, 118128.

[16] Gao, P., Zhang, R., & Yang, X. (2020). The application of stock index price prediction with neural network. Mathematical and Computational Applications, 25(3), 53.

[17] Chen, Y., & Hao, Y. (2017). A feature weighted support vector machine and K-nearest neighbor algorithm for stock market indices prediction. Expert Systems with Applications, 80, 340-355.

[18] Ren, R., Wu, D. D., & Liu, T. (2018). Forecasting stock market movement direction using sentiment analysis and support vector machine. IEEE Systems Journal, 13(1), 760-770.

[19] Dash, R. K., Nguyen, T. N., Cengiz, K., & Sharma, A. (2023). Fine-tuned support vector regression model for stock predictions. Neural Computing and Applications, 35(32), 23295-23309.

[20] Fischer, T., & Krauss, C. (2018). Deep learning with long short-term memory networks for financial market predictions. European journal of operational research, 270(2), 654-669.

[21] Chung, H., & Shin, K. S. (2018). Genetic algorithm-optimized long short-term memory network for stock market prediction. Sustainability, 10(10), 3765.

[22] Qiu, J., Wang, B., & Zhou, C. (2020). Forecasting stock prices with long-short term memory neural network based on attention mechanism. PloS one, 15(1), e0227222.

[23] Dey, S., Kumar, Y., Saha, S., & Basak, S. (2016). Forecasting to Classification: Predicting the direction of stock market price using Xtreme Gradient Boosting. PESIT South Campus, 1, 1-10.

[24] Deng, S., Huang, X., Zhu, Y., Su, Z., Fu, Z., & Shimada, T. (2023). Stock index direction forecasting using an explainable eXtreme Gradient Boosting and investor sentiments. The North American Journal of Economics and Finance, 64, 101848.

[25] Zhang, X., & Chen, W. (2019, July). Stock selection based on extreme gradient boosting. In 2019 Chinese Control Conference (CCC) (pp. 8926-8931). IEEE.

[26] Jiaru Wang & Zhenhan Chen. (2024). Gaussian Process Regression Model Based on Multi-Kernel Learning and Bagging Algorithm. Academic Journal of Mathematical Sciences,5(3),

[27] Jason Gasper & Jennifer Cahalan. (2025). Utilizing the random forest algorithm and interpretable machine learning to inform post-stratification of commercial fisheries data. Fisheries Research,281,107253-107253.

[28] Chander Diwaker,Vijay Hasanpuri,Yonis Gulzar & Bhanu Sharma .(2025). Optimizing MapReduce efficiency and reducing complexity with enhanced particle Swarm Optimization (MR-MPSO). Swarm and Evolutionary Computation,95,101917-101917.

[29] Somak Maitra & Dimitris N. Politis. (2024). Prepivoted Augmented Dickey-Fuller Test with Bootstrap-Assisted Lag Length Selection. Stats,7(4),1226-1243.