

Received May 29, 2021, accepted July 7, 2021, date of publication July 13, 2021, date of current version July 23, 2021.

Digital Object Identifier 10.1109/ACCESS.2021.3096825

# Stock Trend Prediction Using Candlestick Charting and Ensemble Machine Learning Techniques With a Novelty Feature Engineering Scheme

YAOHU LIN<sup>1</sup>, SHANCUN LIU<sup>1,2</sup>, HAIJUN YANG<sup>1,3</sup>, (Member, IEEE), AND HARRIS WU<sup>4</sup>

<sup>1</sup>School of Economics and Management, Beihang University, Beijing 100191, China

<sup>2</sup>Key Laboratory of Complex System Analysis, Management and Decision, Ministry of Education, Beihang University, Beijing 100191, China

<sup>3</sup>Beijing Advanced Innovation Center for Big Data and Brain Computing, Beihang University, Beijing 100191, China

<sup>4</sup>Strome College of Business, Old Dominion University, Norfolk, VA 23529, USA

Corresponding author: Haijun Yang (navy@buaa.edu.cn)

This work was supported by the National Natural Science Foundation of China under Grant 71771006 and Grant 71771008.

**ABSTRACT** Stock market forecasting is a knotty challenging task due to the highly noisy, nonparametric, complex and chaotic nature of the stock price time series. With a simple eight-trigram feature engineering scheme of the inter-day candlestick patterns, we construct a novel ensemble machine learning framework for daily stock pattern prediction, combining traditional candlestick charting with the latest artificial intelligence methods. Several machine learning techniques, including deep learning methods, are applied to stock data to predict the direction of the closing price. This framework can give a suitable machine learning prediction method for each pattern based on the trained results. The investment strategy is constructed according to the ensemble machine learning techniques. Empirical results from 2000 to 2017 of China's stock market confirm that our feature engineering has effective predictive power, with a prediction accuracy of more than 60% for some trend patterns. Various measures such as big data, feature standardization, and elimination of abnormal data can effectively solve data noise. An investment strategy based on our forecasting framework excels in both individual stock and portfolio performance theoretically. However, transaction costs have a significant impact on investment. Additional technical indicators can improve the forecast accuracy to varying degrees. Technical indicators, especially momentum indicators, can improve forecasting accuracy in most cases.

**INDEX TERMS** K-line patterns, machine learning, ensemble strategy, eight-trigram, stock forecasting.

## I. INTRODUCTION

The forecasting of the stock market is an important objective in the financial world and remains one of the most challenging problems due to the non-linear and chaotic financial nature [1], [2]. Investments in the stock market are often guided by different prediction methods which can be divided into two groups of technical analysis and fundamental analysis [3]. The fundamental analysis approach is concerned with the company which used the economic standing of the firm, employees, yearly reports, financial status, balance sheets, income reports and so on [4]. On the other hand, technical analysis, also called charting, predicts the future by studying the trends from the historical data [5]. Investors could build profitable trading strategies by using technical analysis techniques [6], [7]. Utilizing open-high-low-close

The associate editor coordinating the review of this manuscript and approving it for publication was Diego Oliva.

prices, candlestick charting can reflect not only the changing balance between supply and demand [8], but also the sentiment of the investors in the market [9].

Nti *et al.* revealed that 66% of stock market prediction documents they reviewed were based on technical analysis [3]. Traditional technical analysis is mostly limited to analyzing the candlestick charting, performing statistical analysis from past historical data to obtain the probability of prediction. For example, Caginalp and Laurent performed a statistical test including eight kinds of three-day patterns and noted that the candlestick patterns have predictive power [10]. Then Lu *et al.* examined these eight three-day patterns with three definitions of trend and four holding strategies in the DJIA component data, and found that regardless of which definition of the trend was used, eight three-day reversal patterns with a Caginalp-Laurent holding strategy were profitable [11]. Chen *et al.* gave the quantitative definitions of four pairs of two-day candlestick patterns to study their predictive

power in Chinese stock market. Results showed that these two-day candlestick patterns have different predictive capabilities [12]. Zhu *et al.* examined the effectiveness of five different candlestick reversal patterns in Chinese stock market. Statistical analysis suggested that bearish harami, and cross signals perform well in predicting head reversals for stocks of low liquidity, while bullish harami, engulfing, and piercing patterns were profitable when applied to highly liquid, small companies' stocks [13]. Lu examined the predictive power of single-day candlestick charting by using the daily data for the Taiwan stocks for the period from 4 January 1992 to 31 December 2009. Statistical results revealed that four patterns were profitable for the Taiwan stock market after transaction costs [14]. Lv *et al.* testing the predictive power of the Three Inside Up pattern and Three Inside Down pattern with the testing dataset of the K-line series data of Shanghai 180 index component stocks over the latest 10 years [15].

Andrew *et al.* provided evidence that technical analysis can be improved by using automated algorithms [16]. Recently, artificial intelligence (AI) has been applied to address the chaotic time series data [17], [18]. The intense computational use of intelligent predictive models has commonly been studied under the title of machine learning [19]. Machine learning uses historical data for parameter fitting to predict new data [20]. Many machine techniques have already been applied to forecast the stock market. For example, logistic regression (LR) and Neural Network (NNs) [21], [22], deep neural networks (DNN) [23], decision trees (DTs) [24]–[26], support vector machines (SVM) [27], k-nearest neighbors (KNN) [28], random forests (RFs) [29], [30] and long short-term memory networks (LSTMs) [31], [32] have been used to predict the stock market. Moreover, many authors try to improve the prediction ability by combining machine learning models with other methods. Ahmad *et al.* proposed a forecasting model based on chaotic mapping, firefly algorithm and support vector regression to predict stock market price. Compared with genetic algorithm-based SVR, chaotic genetic algorithm-based SVR, artificial neural networks and adaptive neuro-fuzzy inference systems, the proposed model performs best on mean squared error and mean absolute percent error [33]. Zhang *et al.* proposed a stock price predicting system by combining SVR and ensemble adaptive neuro fuzzy inference system (ENANFIS). The experimental results showed that the SVR-ENANFIS model has superior prediction performance than ENANFIS, SVR-Linear, SVR-SVR and SVR-ANN [34]. However, forecast research based on AI methods and candlestick charting is still less.

In these applications of AI methods for financial market forecasting, many studies have used technical indicators as input features. Weng *et al.* developed a financial expert system that incorporated the historical stock prices, eight kinds of technical indicators, counts and sentiment scores of published news articles, trends in Google search and Wikipedia information to predict short term stock prices [35]. Kumar *et al.* used 15 kinds of technical indicators to construct 55 input features to predict the direction of stock indices [36].

Gocken *et al.* used 44 technical indicators in their hybrid soft computing models for 1, 2, 3, 5, 7 and 10 days ahead stock price prediction [37]. Patel *et al.* selected 10 technical indicators to predict the closing price using the fusion of machine learning techniques [38]. Zhou *et al.* developed a learning architecture LR2GBDT for forecasting and trading stock indices by adding 12 technical indicators as the initial variables [39]. Bao *et al.* used 10 technical indicators in a deep learning framework where wavelet transforms (WT), stacked autoencoders (SAEs) and LSTM are combined for stock price forecasting [40]. However, the existing research is limited to using these technical indicators as input parameters and it lacks a further discussion on these technical indicators.

When comparing different machine learning effects, Weng *et al.* used four machine learning methods, including boosted regression tree (BRT), NN, SVM and RF. Results showed that BRT and RF perform the best when predicting one-day ahead stock price [35]. Krauss *et al.* deployed a statistical arbitrage strategy based on DNN, GBDT and RF to S&P 500 constituents from December 1992 to October 2015, finding that the RF outperform GBDT and DNN in their study [41]. Patel *et al.* compared four prediction models, ANN, SVM, RF and naive-Bayes. Experimental results showed that the RF outperformed the other three prediction models when the evaluation was carried out on 10 years of historical data of two stocks and two stock price indices [42]. This article attempts to construct an adaptive automatic machine learning method selection framework while the prediction effects of different machine learning methods are inconsistent in different scenarios.

The main contributions of our research are as follows: (1) This article combines traditional candlestick charting with the latest machine learning methods to enrich the research content of stock market forecasting. (2) We developed a simple eight-trigram classification following the eight-trigram scheme. We also compare different types of technical indicators in short-term stock forecasting, and further clarified the role of technical indicators in machine learning models. (3) Then an adaptive machine learning method selection framework with a novelty feature engineering scheme was proposed. Appropriate forecasting machine learning methods can be automatically selected in different candlestick charting patterns. (4) We have constructed an investment strategy based on our prediction framework. Empirical results show that this paper's strategy makes good economic returns on both individual stocks and portfolios.

The remainder of this paper is organized as follows. Section 2 outlines the design of an ensemble prediction framework using machine learning techniques. Section 3 presents the empirical results on stock data and robustness checks. Section 4 concludes the paper.

## II. METHODOLOGY

This paper proposes an adaptive prediction framework using an ensemble of machine learning models to predict the direction of the closing price, which is shown in FIGURE 1.

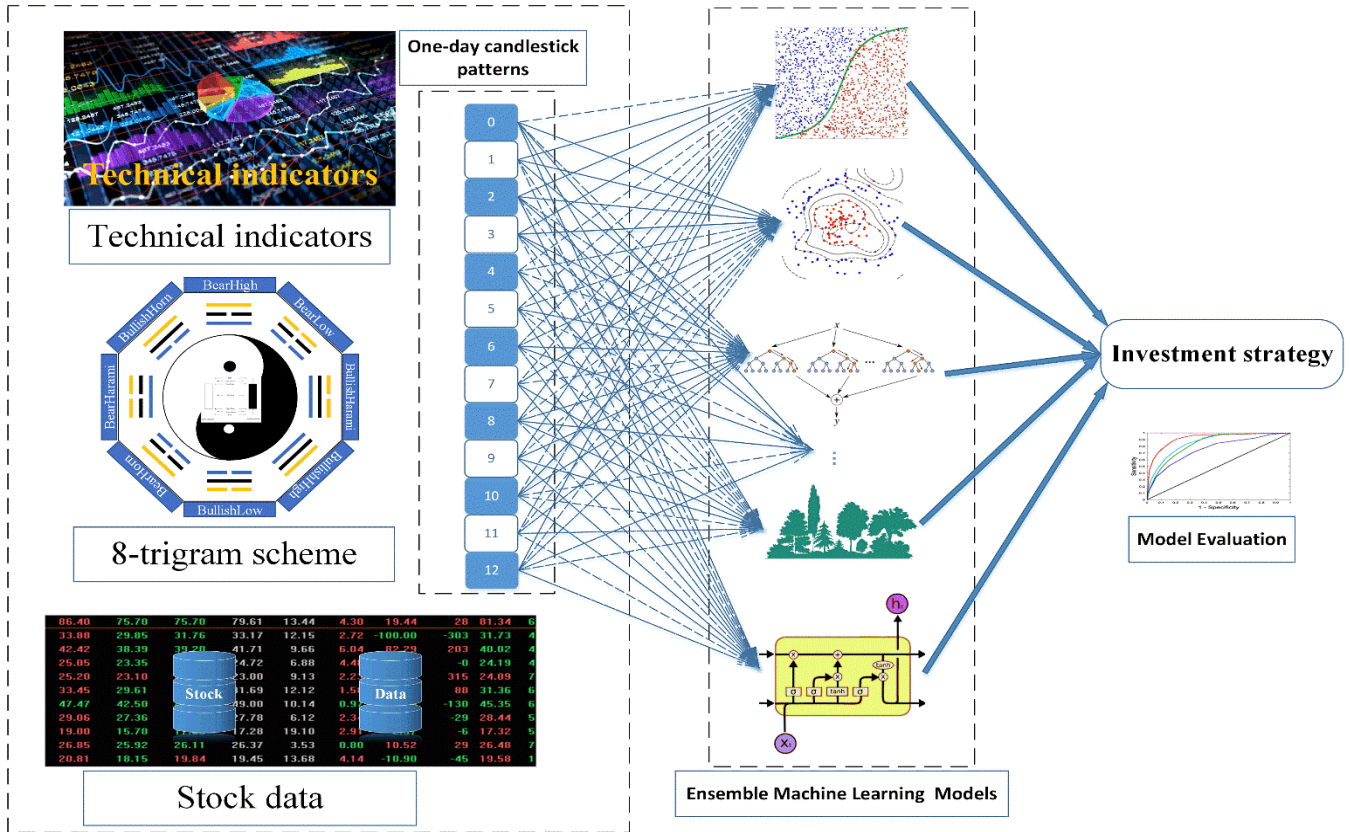


FIGURE 1. Overview of prediction framework.

First, 13 forms of one-day patterns are constructed and classified from 3,455 stocks used in this article, and then the corresponding technical indicators and eight-trigram information are calculated. Then, all feature data are passed as input to the ensemble machine learning model, which tests the prediction accuracy of each pattern. For each pattern, the machine learning method with the highest prediction will be recorded. Finally, the adaptive recommendation schedule gives corresponding stock prediction actions based on the evaluated results.

The main evaluation algorithm is shown as FIGURE 2.

### A. FEATURE ENGINEERING USING AN EIGHT-TRIGRAM SCHEME

Key to our prediction framework is a simple eight-trigram scheme to represent inter-day stock price movements based on two-day candlestick patterns. The candlestick chart, also called K-line, is drawn by close, open, high, and low, where the part between the close and open is called real body. If the asset closed higher than it opened, the body is filled with red color in Chinese stock market while represented with white or green color in European and American stock markets. If the close price is lower than open price, the body is painted with green color in Chinese stock market while filled with black or red color in European and American stock markets. And then, the candlestick pattern classification is built, which is

### Algorithm: Model Evaluation

**Input:** Patterns data which includes feature engineering data and different indicators

**Output:** *BestModel*, *F1 score*

```

0 Evaluation (features):
1   foreach p in patterns: Generate p_data of p;
2     LogisticRegression (p_data);
3     GridSearchCV of KNN (p_data);
4     GridSearchCV of SVM (p_data);
5     GridSearchCV of RF (p_data);
6     GridSearchCV of GBDT (p_data);
7     LSTM (p_data);
8     BestModel = MAXF1 (LR, KNN, SVM, RF,
9     GBDT, LSTM);
    Save the best performance model BestModel,
    and F1 score for pattern p
Output: List of best performance model, F1 score
    for each pattern
    
```

FIGURE 2. Main evaluation algorithm.

shown in FIGURE 3. We divide the candlestick patterns into 13 classes according to the candlestick basic elements: the opening, high, low, and closing prices.

Relative to yesterday's closing price, the opening position of the day reflects the accumulation of sentiments during

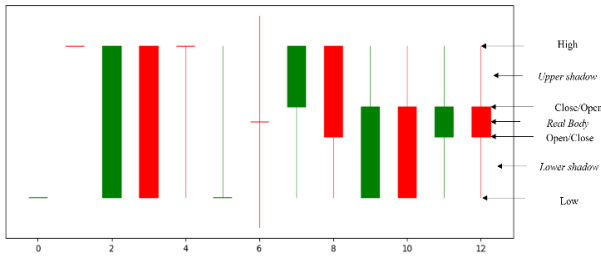


FIGURE 3. Candlestick patterns.

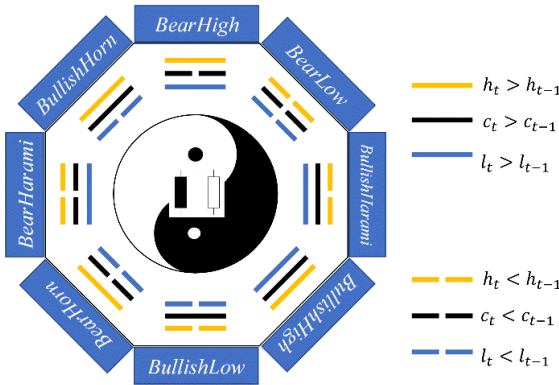


FIGURE 4. Eight-trigrams of inter-day price movement patterns.

non-trading period. We divide the inter-day price movement into eight categories based on the relative position of today’s price range and yesterday’s K-line patterns. The simple schematic diagram of eight-trigram feature engineering scheme is shown in FIGURE 4. The symbols used in eight-trigrams and the detailed expressions can be found in Appendix I.

On the other hand, trading volume is an important parameter which is not getting enough attention in academia. It is also a variable completely independent of price. The N-day moving average of volume on day  $t$  is defined by

$$MA_{n,v}(t) = \frac{1}{n} \sum_{i=0}^{n-1} V_{t-i} \quad (1)$$

The volume rate of change on day  $t$  is defined by

$$ROC(n) = \frac{V_t}{MA_{n,v}(t)} \quad (2)$$

Besides, the technical indicators could potentially have an impact on the stock price prediction [43], [44]. Four groups of technical indicators, including 21 indicators, were introduced in our research. As shown in Table 1 Detailed expressions for all indicators can be found in Appendix II.

**B. GENERATION OF TRAINING AND TESTING SETS**

We begin to extract the features after the data preprocessing process according to the feature engineering. The feature of single k-line patterns classification is generated on the basis price of Opening, High, Low and Closing price. A total of

TABLE 1. Technical indicators.

Group of indicators	Technical indicators
Overlap indicators	Moving Average (MA), Exponential Moving Average (EMA), Double Exponential Moving Average (DEMA), Kaufman's Adaptive Moving Average (KAMA), Simple Moving Average (SMA), Parabolic SAR (SAR)
Momentum indicators	Average Directional Movement Index (ADX), Price Oscillator - Absolute (APO), Balance of Power (BOP), Commodity Channel Index (CCI), Moving Average Convergence/Divergence (MACD), Money Flow Index (MFI), Momentum (MOM), Relative Strength Index (RSI)
Volume indicators	Chaikin A/D Line (AD), Chaikin Oscillator (ADOSC), On Balance Volume (OBV)
Volatility indicators	True Range (TRANGE), Average True Range (ATR), Normalized Average True Range (NATR)

13 K-line patterns was obtained after considering all the circumstances. Secondly, eight inter-day price move indicators including *BullishHorn*, *BearHorn*, *BullishHigh*, *BearHigh*, *BullishLow*, *BearLow*, *BullishHarami* and *BearHarami* are extracted. And then, the rate of volume feature which has not received enough attention is taken into consideration. Instead of simply placing the volume value directly into the forecasting model, we use the variable of the rate of change as a predictive feature. Finally, according to the indicator formula, the values of 21 indicators as prediction parameters are calculated. The research in this paper focuses on short-term forecasting, we choose 5 days or 10 days as the parameters of indicators. Parameters of 5 days for *MA*, *DEMA*, *KAMA* and 10 days for *EMA*, *SMA*, *ADX*, *APO*, *CCI*, *DX*, *MFI*, *RSI*, *ATR*, *NATR* are employed.

After all the prediction features are ready, we begin to prepare the training sets and testing sets. We divide the entire data set into two parts, 80% of which are training sets and 20% are testing sets. The corresponding result is the next day’s direction of stock price.

**C. PREDICTION MODELS**

The inference engine is introduced in this phase. Six machine learning models including Logistic Regression (LR), Support Vector Machine (SVM), k-NearestNeighbor (KNN), Random Forest (RF), Gradient Boosting Decision Tree (GBDT) and Long Short-term Memory (LSTM) are used to predict the direction of the closing price. The parameters used in these prediction models are shown as Table 2.

**1) LOGISTIC REGRESSION (LR)**

Logistic Regression is the most basic machine learning algorithm. The Logistic Regression model returns an equation that determines the relationship between the independent variables and the dependent variable. First, the model calculates linear functions and then converts the result into a probability. Finally, it converts the probability into a label.

TABLE 2. Parameters used in prediction models.

MLs	Parameters
LR	Regularized= L2, solver_parameter=warn, C=1.0, iteration=100, criteria=0.0001
SVM	C={1e3,5e3,1e4,5e4,1e5}, gamma={0.0001,0.0005,0.001,0.005,0.01,0.1}, optimizer= GridSearchCV cv=10
KNN	n_neighbors = range(1,10), weights = ['uniform','distance'], algorithm=['auto','ball_tree','kd_tree','brute'], leaf_size=range(1,2), optimizer= GridSearchCV cv=10
RF	n_estimators=range(10,100,5), criterion=[gini, entropy], min_samples_leaf= [2, 4, 6,50], max_depth=range(1,10), optimizer= GridSearchCV cv=10
GBDT	n_estimators=range(10,100), max_features=range(0.6,0.9), max_depth=range(1,10), optimizer= GridSearchCV cv=10
LSTM	unit=64, dropout=0.2, activation=sigmoid, loss=binary_crossentropy, optimizer=adam, epochs=20

In the empirical stage, we use L2 as a regularized parameter, specifying warn as the solver parameter which determines our optimization method for the logistic regression loss function. In terms of the termination parameters of the algorithm, the maximum number of iterations which is taken for the solvers to converge to 100 and tolerance for stopping criteria parameter is set to 0.0001.

### 2) SUPPORT VECTOR MACHINE (SVM)

In machine learning, support vector machines (SVMs) are supervised learning models with associated learning algorithms that analyze data for classification and regression analysis. The SVM training algorithm builds a model when given a set of training examples and assigns the new example to one category or another. The SVM model is to map the example as a point map in space, so that the examples of the separate categories are divided by a clear gap that is as wide as possible. Lee gave a detailed formula for the SVM's two-class problem in his 2009 article [27].

In addition to performing linear classification, SVM can effectively perform nonlinear classification using so-called kernel techniques, implicitly mapping its inputs to high-dimensional feature spaces. There are some studies using the SVM to predict the financial data [45], [46].

In the empirical stage, we get the best performance from a grid search algorithm which sets different penalty coefficients, coefficients of kernel function and degrees. Different parameter combinations produce different clustering effects.

### 3) K-NEARESTNEIGHBOR (KNN)

K nearest neighbors (KNN) is another machine learning algorithm. The k-NN algorithm looks for 'k' nearest records within the training dataset and uses the majority of the classes of the identified neighbors for classifying. Subha (2012) used k-NN to classify the stock index movement [28]. In the empirical stage, we get the best performance from a grid search algorithm which sets different neighbors, leaves and weights. Different parameter combinations produce different clustering effects.

### 4) RANDOM FOREST (RF)

Random forests are a combination of tree predictors. Each tree depends on the values of a random vector sampled independently with the same distribution for all trees in the forest. Random forests (RFs) are a nonparametric and nonlinear classification and regression algorithm [47]. Random forests not only use a subset of the training set, but also selects only a subset of the feature set when the tree is established in the decision tree. Booth (2014) used RFs to construct a n automated trading mechanism [29].

In the empirical stage, different parameter combinations may produce different classification effects. We get the best performance from a grid search algorithm which sets different leaves, depth and estimators.

### 5) GRADIENT BOOSTING DECISION TREE (GBDT)

Gradient Boosting Decision Tree (GBDT) is a popular machine learning algorithm. Friedman (2002) gave a detailed expression of the gradient descent in his research [48]. The core idea of GBDT is that the subsequent model of the sequence no longer directly predicts the predicted value of the data set, but predicts the difference between the predicted value and the true value of the previous model.

In the empirical stage, different parameter combinations may produce different classification effects. We get the best performance from a grid search algorithm which sets different features, depth and estimators.

### 6) LONG SHORT-TERM MEMORY (LSTM)

Long-short term memory is one of the recurrent neural network (RNNs) architecture [49]. Hochreiter and Schmidhuber proposed a solution by using memory cells [50] which consists of three components, including input gate, output gate and forget gate. The gates control the interactions between neighboring memory cells and the memory cell itself. The input gate controls the input state while the output gate controls the output state which is the input of other memory cells. The forget gate can choose to remember or forget its previous state.

The LSTM network used Keras, an open-source neural-network library written in Python. First, we construct a three-layer LSTM network, including two main processing layers and one output layer (dense layer). In the prediction model, we apply dropout regularization within the recurrent layer. Hereby, a fraction of the input units are randomly dropped at each update during training time to reduce the risk of overfitting and it gets better generalization. And the early stopping mechanism is also used to reduce the risk of overfitting.

In the empirical stage, the hidden neurons are set to 64 and dropout to 0.2 in the first layer and set the hidden neurons to 64 in the second layer. At the output layer, we use the sigmoid activation function to generate the classification results. This configuration yields 17,920 parameters for the first layer,

33,024 parameters for the second layer and 65 parameters for the output layer.

**D. MODEL EVALUATION**

The inference engine is introduced in this phase. We use six machine learning models to forecast the stock direction of up and down. To evaluate the performance of the prediction models, two commonly used evaluation criteria are used in this study: *Accuracy*, and *F1 score*. *Accuracy* is used to evaluate the overall classification ability of the model. The formula is as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{3}$$

*TP (True Positives)* representing model prediction is true, and the real sample is also true; *TN (True Negatives)* representing model prediction is false, and the real sample is also false; *FP (False Positives)* representing model prediction is true while the real sample is false; *FN (False Negatives)* representing model prediction is false while the real sample is true.

*Precision* is used to estimate the accuracy of positive samples in the prediction data and *recall* is used to evaluate the coverage of positive samples in the prediction data of the model. The formula is shown as (4) and (5).

$$Precision = \frac{TP}{TP + FP} \tag{4}$$

$$Recall = \frac{TP}{TP + FN} \tag{5}$$

*F1* is an indicator used in statistics to measure the accuracy of a binary model, which also considers the accuracy and recalls of the classification model. The formula is as follows:

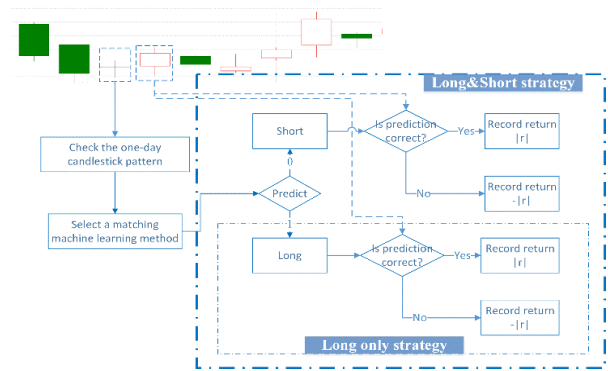
$$F1 = \frac{2 * Precision * Recall}{Precision + Recall} \tag{6}$$

And then, the evaluation progress tries to evaluate the results of the machine learning prediction model. In this study, our evaluation model not only considers the *Accuracy* index, but also considers the *Precision* and *Recall* indicators.

**E. INVESTMENT STRATEGY**

The investment strategy is constructed based on the above evaluation model. This paper considers two situations, including only long and long-short, and builds the corresponding investment strategy. This article assumes that we will invest at the closing price at time *t* and will be clear at the closing price at time *t + 1*. The flowchart of the investment strategy is shown as FIGURE 5. If the current stock trading market mechanism allows shorting, that means you can go long and short at the same time. The specific construction steps of the investment strategy are as follows:

First, the specific K-line pattern of the current stock is checked at time *t*. Next, a matching machine learning method is selected to predict the rise or fall of *t + 1* based on the above evaluation model. If the predicted result is consistent with the real result, the *t + 1* profit would be recorded. If the prediction



**FIGURE 5. Flowchart of the investment strategy.**

is wrong, the negative profit of *t + 1* would be recorded as loss value. Finally, the above steps will be repeated to calculate *t + 1*, *t + 2*, etc.

If the current stock trading market mechanism does not allow shorting, that is to say, you can only go long. The specific construction steps of the investment strategy are as follows:

First, the specific K-line mode of the current stock is checked at time *t*. Then, a matching machine learning method is selected to predict the rise or fall of *t + 1* based on the above evaluation model. We only operate the investment when the forecasting result is up. If the predicted result is consistent with the real result, we record the *t + 1* profit. Otherwise, we record the negative profit of *t + 1* as loss value. Finally, the above steps will be repeated to calculate *t + 1*, *t + 2*, etc.

**III. EMPIRICAL RESULTS**

**A. DATA AND TRAINING ENVIRONMENT**

As the world’s largest developing country and the world’s second-largest economy, China’s influence on the world is growing. The financial market of China has attracted the attention of domestic and foreign scholars and investors [12], [13], [39]. This paper selects the data of the Chinese stock market as the experimental data. The daily data of the China Stock Market from the 18-year period of 2000 to 2017 is used in this study. All the 3,455 stocks data is collected from CCER, a local data provider of China. First of all, we remove the daily data for a given stock if the trading volume is zero, which is a sign of stopped trading such as due to company reorganization. The distribution of 13 patterns in the data set is shown in Table 3. The distribution of these patterns in historical data is stable. Then, we generate feature information, which contains the date, intra-day pattern, inter-day pattern, 21 other indicator values and the next day’s closing price, for each stock on each day *t*. In order to ensure effectiveness, three rounds of training were carried out. We randomly choose 5,000 rows of daily stock data for each of 13 intra-day patterns from the database, which yields 65,000 rows of data in each round. In order to ensure the balance of classification during training, for each intra-day pattern, we choose half of the training data with

TABLE 3. Distribution of 13 patterns.

Patterns	No.	Ratio	Ratio1	Ratio2
0	6,946	0.10%	0.10%	0.06%
1	15,268	0.21%	0.23%	0.13%
2	56,477	0.77%	0.79%	0.70%
3	57,910	0.79%	0.83%	0.63%
4	12,549	0.17%	0.19%	0.10%
5	7,501	0.10%	0.12%	0.05%
6	170,234	2.33%	2.34%	2.31%
7	630,482	8.63%	8.65%	8.56%
8	247,028	3.38%	3.53%	2.72%
9	186,261	2.55%	2.59%	2.39%
10	563,830	7.72%	7.60%	8.27%
11	2,442,444	33.45%	33.15%	34.77%
12	2,905,003	39.78%	39.89%	39.31%

The statistical data cycle is from Jan 1, 2000 to Dec 31, 2017. The No. column refers to the total number of occurrences of the specified pattern. The Ratio column indicates the proportion of the pattern in the total data set. Ratio1 column indicates the proportion of the pattern in the data set from

rising prices (closing price lower than next day's) and half of the training data with falling prices. Finally, the average is obtained based on three rounds of results.

**B. MODEL COMPARISON AND EVALUATION**

To clarify the role of technical indicators in machine learning models, the 21 technical indicators are divided into four groups, including Overlap indicators, Momentum indicators, Volume indicators and Volatility indicators. And then each

pattern containing feature engineering information is put into the ensemble machine learning prediction framework that contains six machine-learning models for training. 80% of the data as training data and the remaining 20% as test data were set to verify the predictive validity of the model.

First of all, we train the machine learning models without indicators. 13 patterns and 6 machine learning models resulting in a total of 78 predictions shows the effectiveness of feature engineering, which is shown in the Fig. a. of FIGURE 6. The abscissa represents the 13 kinds of one-day candlestick patterns and the ordinate represents the *F1 score*. 49 of the 78 prediction models exceeded the random walk probability. Pattern 0 and pattern 1, which indicate that the maximum price limit has been reached, show strong long-short signals. In this case, pattern 4 and pattern 5 show significant prediction effects, and the *F1 score* value can reach about 0.57. In all prediction models, SVM, RF and GBDT showed good predictions for all patterns, while LR, KNN and LSTM only performed well in individual patterns.

Fig. c. of FIGURE 6 shows the prediction results of the introduction of the Overlap indicators which contain 6 indicators of MA, DEMA, EMA, KAMA, SMA and SAR. 60 of the 78 predictions exceeded the random walk probability. And after the introduction of this indicator group, the predicted maximum *F1 score* of each pattern are all improved slightly. Among them, pattern 1, which indicates a board daily limit, shows a strong long signal, the *F1 score* is more than 0.80. This Overlap indicator group has a significant effect on the

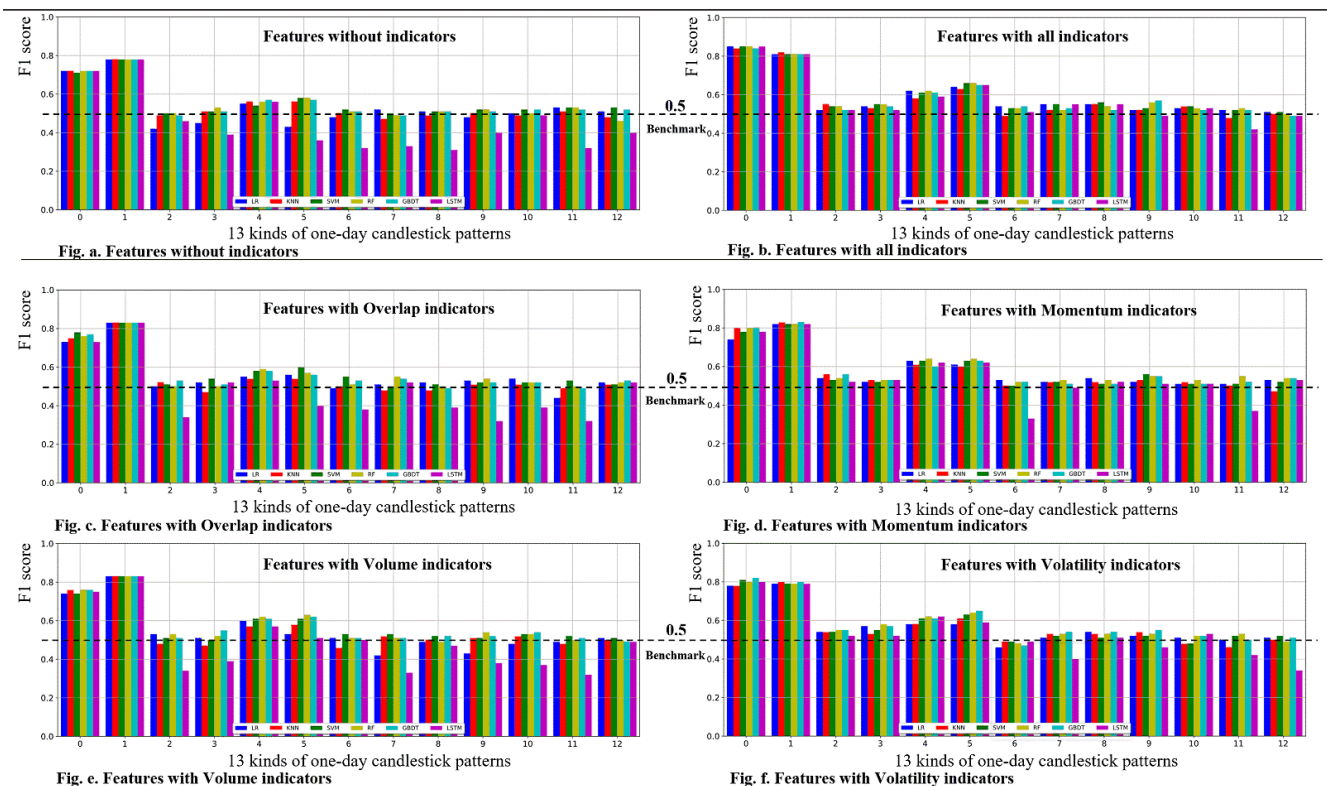


FIGURE 6. Forecast performance of different combinations of indicator groups.

TABLE 4. Best performance parameters.

Parameters	One-day candlestick patterns												
	0	1	2	3	4	5	6	7	8	9	10	11	12
ML Models	RF	KNN	GBDT	RF	RF	SVM	SVM	GBDT	SVM	GBDT	SVM	RF	RF
Indicators	All	Mom	Mom	Volatility	Mom	All	MA	All	All	All	All	Mom	Mom

improvement of the predictive ability of the SVM model, and the predictive ability of the LR model has also been improved to a certain extent.

Fig. d. of FIGURE 6 shows the forecast results of the introduction of the Momentum indicators which contain 9 indicators of ADX, APO, BOP, CCI, DX, MACD, MFI, MOM and RSI. 73 of the 78 predictions exceeded the probability of random walk. And after the introduction of this indicator, the predicted maximum *F1 score* of each pattern are all improved. Among them, the *F1 score* of pattern 0 and pattern 1 exceeds 0.80. The prediction effect of most prediction patterns has been improved after the introduction of these indicators, reflecting the obvious momentum characteristics in short-term prediction. The prediction ability of RF has been significantly improved, and the prediction effects of pattern 4 and pattern 5 have been significantly improved which reach about 0.62.

Fig. e. of FIGURE 6 shows the predicted results of the introduction of the Volume indicators which contain 3 indicators of AD, ADOSC and OBV. 57 of the 78 predictions exceeded the random walk probability. And after the introduction of this indicator, the predicted maximum *F1 score* of each pattern are slightly improved. Among them, only pattern 4 and pattern 5 show good prediction effects, with a *F1 score* about 0.62. RF and GBDT perform well in all prediction modes. From this we can find that the characteristic project has already included the information of short-term volume change, and the introduction of more volume characteristics cannot significantly improve the forecast level.

Fig. f. of FIGURE 6 shows the results of the introduction of the Volatility indicators which contain 3 indicators of ATR, NATR and TRANGE. 63 of the 78 predictions also exceeded the random walk probability. And after the introduction of the indicator, the highest prediction *F1 score* of each pattern are all improved, and different prediction models have a clear distinction between different pattern prediction effects. After the introduction of this indicator, all the patterns have been improved slightly. The overall prediction results reflect the obvious volatility characteristics in the short term.

In order to test whether more technical indicators can improve the forecast level, we introduce all the technical indicators into our forecasting framework. Fig. b. of FIGURE 6 shows the results of the introduction of all indicators which including 21 indicators. 69 of the 78 predictions exceeded the random walk probability. And after the introduction all indicators, the predicted maximum *F1 score* of each mode are improved. Among them, pattern 0 and pattern 1, that is, a board daily limit shows a strong signal, the *F1 score* is more than 0.80. Pattern 4 and pattern 5 still show good prediction

results, *F1 score* is more than 0.62. These indicators have a significant effect on the prediction ability of the LR model, which means that more parameters can improve the linear fitting effect and verify that the stock market has complex nonlinear characteristics.

From the above results, we can see that RF and GBDT have good predictive ability in most cases in short-term prediction. KNN only showed relatively good predictive power in the first six patterns. Although SVM will take too long time to process big data, the prediction level is still significant in some cases. The advantage of the deep learning model LSTM in this scenario is not fully reflected. In addition, we can see that the increase in the number of parameters can increase the level of linear prediction, which also reflects the complexity and diversity of financial markets from another perspective. In most cases, an increase in the number of indicators can increase the level of prediction. However, in the short-term forecast, in the prediction of some models, using a smaller number of momentum indicators and volatility indicators can achieve satisfactory results. The best performance results is shown in Table 4.

C. ROBUSTNESS CHECKS

In order to test the validity of the prediction framework introduced in this paper, we selected the daily data of Shanghai and Shenzhen 300 Index constituent stocks of the China Stock Market from the 18-year period of 2000 to 2017 for testing. At the same time, to ensure the validity of the data, we exclude stocks that are no longer constituents in the last five years. The resulting sample consists of 168 stocks, 553,028 rows of data during that period. We use the best performance parameters of each pattern to verify the predictive validity of the model.

The validity of the prediction results is shown in FIGURE 7. The upper part is the *F1 score* and the lower part is the accuracy rate. From the results, we can see that all the pattern prediction accuracy is greater than 52%, and the *F1 score* is basically more than 50% after using our prediction framework. Among them, pattern 0, pattern 1, pattern 3, pattern 4, pattern 5, and pattern 8 have obvious prediction effects.

D. EMPIRICAL RESULT OF INVESTMENT STRATEGY

Based on the prediction framework of this paper, we construct a trading strategy as follows:

1. Identify the k-line pattern at day t.
2. Use the prediction framework to predict rise or fall of stock closing price at day t + 1.



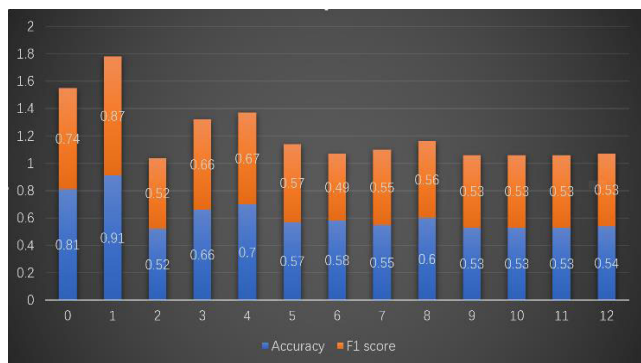


FIGURE 7. Robustness test.

3. If the prediction is correct, as validated by the empirical data, record a profit in the amount of price change (assuming short transactions are allowed and therefore profit can be made even when stock falls).

4. If the prediction is incorrect, record a loss in the amount of price change.

We construct a long-only strategy as a limited version to above, to record a profit only when the stock rises and the prediction correctly predicts the rise. These strategies can be executed at an X%-confidence, that is, be executed only when the confidence of prediction exceeds X%.

FIGURE 8 shows the forecast results of random stock ‘000001.SZ’. From the figure, we can see that before the 2008 financial crisis, a trading strategy based on the forecasting framework performs better than holding the stock itself. After the financial crisis, the prediction-based strategy has a smaller retraction and will soon outperform the market. However, the prediction-based strategy showed greater volatility during the 2015 stock market crash. The predicted maximum drawdown is 71.4%, which is less than 77.5% of the original stock. And the predicted Sharpe Ratio is 0.31, which is bigger than 0.25 of the original stock. The predicted Sortino Ratio is 0.0348, which is bigger than 0.033 of the original stock. Explain that there are fewer risks and greater benefits based on our forecasting framework. However, after considering the transaction costs, the profits are significantly reduced. FIGURE 8 shows the profitability under different transaction costs of 0.1%, 0.2%, and 0.3%.

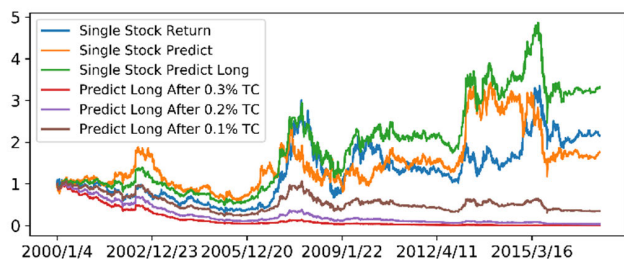


FIGURE 8. Return on single stock using prediction strategies.

Then, we build an equal-weighted portfolio based on the constituents of CSI 300. FIGURE 9 shows the historical

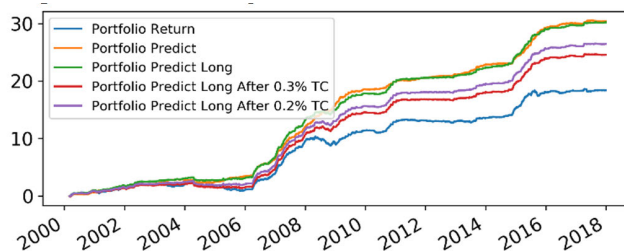


FIGURE 9. Return on CSI 300 using prediction strategies with 55% confidence threshold.

revenue and forecast for the portfolio. From the figure, we can see that the predicted results are significantly better than the portfolio itself. The predicted maximum drawdown is 43.7%, which is less than 63.5% of the original stock. And the predicted Sharpe Ratio is 0.62, which is bigger than 0.35 of the original portfolio, the predicted Sortino Ratio is 0.808, which is bigger than 0.368 of the original portfolio supporting that there are fewer risks and greater benefits based on our forecasting framework. After considering the transaction costs, although the return has declined, it is still possible to obtain excess return, which is higher than the investment portfolio itself. The figure shows the return on the portfolio considering 0.2% and 0.3% transaction cost.

FIGURE 10 shows portfolio return based on the constituents of Shanghai Composite Index. The index component covers 1,380 stocks, and we selected the 2016-2017 data outside the sample for testing to verify the effectiveness of our forecasting framework.

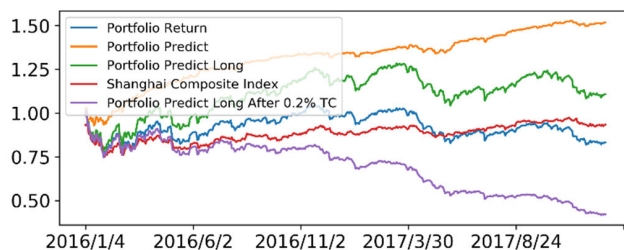


FIGURE 10. Return performance of portfolio from Shanghai composite index.

We can see that the return rate based on our research method is better than the market performance from the figures. The effect will be more prominent if we can short-sell stocks. However, after considering the transaction costs, the profit is significantly reduced. Table 5 shows the finance performance.

TABLE 5. Finance performance.

	Portfolio	Predict1	Predict2	Predict 3	Index
Max drawdown	26.1%	9.7%	24.1%	59.04%	21.0%
Sharpe Ratio	-0.36	1.65	0.24	-1.76	-0.08
Sortino Ratio	-0.023	0.174	0.029	-0.139	-0.009

Predict 1 means predict the portfolio including long and short. Predict 2 means the portfolio only go long. Predict 3 means the portfolio only go long after the 0.2% transaction cost. The index means Shanghai Composite Index.

IV. CONCLUSION

This paper develops an ensemble machine learning prediction model that automatically selects appropriate prediction methods for each daily k-line pattern. The empirical results show that the forecasting framework of this paper has predictive power, and the investment strategy based on the forecasting model can generate superior returns.

This study makes contributions into four aspects. Firstly, this article combines traditional candlestick charting with the latest artificial intelligence methods to enrich the forecasting research of the stock market. By studying the prediction effects of all 13 one-day candlestick patterns under different machine learning methods, we combine traditional technical analysis methods with AI technology. We also concluded that certain candlestick patterns, for example, pattern 4 and pattern 5, have apparent predictive effects in the stock market.

Secondly, in feature engineering, an eight-trigram classification of two-day k-line patterns is developed, in addition to the 13 patterns of daily k-line patterns and volume change features. The simple eight-trigram classification follows the eight-trigram scheme, or Bagua, a key concept in Taoism cosmology. The eight-trigram classification provides a simple set of features based on opening, closing, high and low prices of two consecutive trading days. To improve the forecast level, we introduced four sets of technical indicators: overlap, momentum, volume, and volatility as auxiliary feature variables. We find that the momentum indicators are significantly better via empirical testing than other indicators in short-term forecasting. Additional technical indicators can improve forecasting in most cases. However, in the short-term forecast, in predicting some patterns, a smaller number of momentum or volatility indicators can achieve satisfactory results.

Thirdly, we introduce a framework for assembling multiple machine prediction models to select the optimal prediction method for different feature modes. The ensemble model includes six commonly-used effective prediction models (RF, GBDT, LR, KNN, SVM, LSTM) and optimizes the parameters of each model. In the empirical study, we find that RF and GBDT have a good predictive ability for short-term prediction in most cases. The prediction level of LR needs to be improved by adding features. KNN and SVM only fit in some patterns. The advantage of the deep learning model LSTM in this scenario is not fully reflected.

Finally, based on the prediction results of this paper, we have constructed an investment strategy. The empirical results show that this paper’s strategy makes good economic returns on both individual stocks and portfolios theoretically. This also shows that through big data, multiple rounds of training, feature standardization, etc., the prediction of results is effective. The predicted maximum drawdown, Sharpe Ratio, and Sortino Ratio of this investment strategy are better than buying and holding the original stock. However, the transaction costs have a significant impact on actual transactions. In actual investment, other factors need to be considered to obtain excess returns.

The forecasting framework of this paper has predictive power, although it is difficult to profit from certain patterns due to the stop-trading rules of the Chinese market. One of the machine learning methods, the SVM method, is not suitable for predicting large-scale stock data. We intend to incorporate more suitable machine learning methods for prediction, such as reinforcement learning methods, into the ensemble model. Furthermore, we plan to utilize additional predictive factors, such as major news events and market sentiment, to improve forecasting results in the future.

APPENDIX I  
DEFINITION OF EIGHT-TRIGRAM

Following table shows the symbols used in eight trigrams:

TABLE 6. Symbols in eight trigrams.

Symbols	Description
$h$	$h_t$ represents the highest price at time $t$ , $h_{t-1}$ represents the highest price at time $t-1$
$c$	$c_t$ represents the closing price at time $t$ , $c_{t-1}$ represents the closing price at time $t-1$
$l$	$l_t$ represents the lowest price at time $t$ , $l_{t-1}$ represents the lowest price at time $t-1$

*BullishHorn* reflects that the oscillations of the day exceeded the previous cycle and reached a new high and a new low, reflecting the strong characteristics of an oscillation. The *BullishHorn* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t > h_{t-1}$
- $l_t < l_{t-1}$
- $c_t > c_{t-1}$

where  $h_t, l_t, c_t$  represent highest price, lowest price and closing price at time  $t$ .

*BearHorn* reflects that the oscillations of the day exceeded the previous cycle and reached a new high and a new low, but the closing price was lower than the previous cycle, reflecting the weak characteristics of an oscillation. The *BearHorn* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t > h_{t-1}$
- $l_t < l_{t-1}$
- $c_t < c_{t-1}$

*BullishHigh* shows that the candlestick creates a new high price and the lowest price is higher than the previous period, reflecting a rising strong feature. The *BullishHigh* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t > h_{t-1}$
- $l_t > l_{t-1}$
- $c_t > c_{t-1}$

*BearHigh* reflects a weak rising feature. The candlestick creates a new high price while the closing price is lower than the previous period. *BearHigh* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t > h_{t-1}$
- $l_t > l_{t-1}$
- $c_t < c_{t-1}$

*BullishLow* contains the information that the stock is getting stronger. The stock hit a new low, but the closing price exceeds the previous cycle. The *BullishLow* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t < h_{t-1}$
- $l_t < l_{t-1}$
- $c_t > c_{t-1}$

*BearLow* reflects the weak characteristics. The stock hit a new low and the closing price is lower than the previous periods. The *BearLow* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t < h_{t-1}$
- $l_t < l_{t-1}$
- $c_t < c_{t-1}$

*BullishHarami* shows the process of energy accumulation. The stock's amplitude is within the range of the previous period and the closing price is higher than the previous period. The *BullishHarami* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t < h_{t-1}$
- $l_t > l_{t-1}$
- $c_t > c_{t-1}$

*BearHarami* reflects another process of energy accumulation. The stock's amplitude is within the range of the previous period while the closing price is lower than the previous period. *BearHarami* candlestick at time  $t$  should fulfill the following three conditions:

- $h_t < h_{t-1}$
- $l_t > l_{t-1}$
- $c_t < c_{t-1}$

## APPENDIX II DEFINITION OF TECHNICAL INDICATORS

TABLE 7. Classification of technical indicators.

Overlap indicators	MA, DEMA, EMA, KAMA, SMA, SAR
Momentum indicators	ADX, APO, BOP, CCI, DX, MACD, MFI, MOM, RSI
Volume indicators	AD, ADOSC, OBV
Volatility indicators	ATR, NATR, TRANGE

### A. OVERLAP INDICATORS

#### 1) MOVING AVERAGE (MA)

$$MA(t) = \frac{1}{n} \sum_{i=0}^{n-1} C_{t-i}$$

where  $n$  refers to the time interval, and  $C$  is the close price.

#### 2) EXPONENTIAL MOVING AVERAGE (EMA)

The Exponential Moving Average is a staple of technical analysis and is used in countless technical indicators.

$$EMA(t) = \frac{2}{n+1} C_t + \frac{n-1}{n+1} EMA(t-1)$$

#### 3) DOUBLE EXPONENTIAL MOVING AVERAGE (DEMA)

The DEMA is a smoothing indicator with less lag than a straight exponential moving average.

$$DEMA(t) = 2 * EMA(C_t) - EMA(EMA(C_t))$$

#### 4) KAUFMAN'S ADAPTATIVE MOVING AVERAGE (KAMA)

The KAMA automatically increases EMA's smoothing during weak trends and during ranging trends.

$$ER(t) = \frac{Abs(C_t - C_{t-n})}{\sum_{i=t-n}^t Abs(C_i - C_{i-1})}$$

$$sc(t) = (ER(t) * \left( \frac{2}{n1+1} - \frac{2}{n2+1} \right) + \frac{2}{n2+1})^2$$

$$KAMA(t) = sc(t) * (C_t - KAMA(t-1)) + KAMA(t-1)$$

where  $n1$  refers to the fast period,  $n2$  refers to the slow period and  $Abs$  indicates absolute value.

#### 5) SIMPLE MOVING AVERAGE (SMA)

Moving Averages are used to smooth the data in an array to help eliminate noise and identify trends.

$$SMA(t) = \frac{1}{n} (m * MA(t) + (n-m) * SMA(t-1))$$

where  $m$  represents the weight.

#### 6) PARABOLIC SAR (SAR)

The Parabolic SAR calculates a trailing stop.

$$SAR(t) = SAR(t-1) + af_t * (xp_{t-1} - SAR(t-1))$$

where  $af$  is acceleration factor and  $xp$  is the extreme point.

### B. VOLUME INDICATORS

Directional Movement Index (+DI and -DI) The +DI is the percentage of the true range that is up. The -DI is the percentage of the true range that is down.

$$\Delta H = H_{t-1} - H_t$$

$$\Delta L = L_t - L_{t-1}$$

where  $H$  refers to the high price and  $L$  refers to the low price. The calculation logic of  $DI$  is as follows:

If  $(\Delta H < 0 \text{ and } \Delta L < 0)$  or  $\Delta H = \Delta L$  then

$$plusDM = 0$$

$$minusDM = 0$$

If  $\Delta H > \Delta L$  then

$$plusDM = \Delta H$$

$$minusDM = 0$$

If  $\Delta L > \Delta H$  then

$$plusDM = 0$$

$$minusDM = \Delta L$$

Then

$$\begin{aligned}
 & plusDMsum(t) \\
 &= plusDMsum(t-1) - \frac{plusDMsum(t-1)}{n} \\
 & \quad + plusDM \\
 minusDMsum(t) \\
 &= minusDMsum(t-1) - \frac{minusDMsum(t-1)}{n} \\
 & \quad + minusDM \\
 TR(t) &= H_t - L_t \\
 TRsum(t) \\
 &= TRsum(t-1) - \frac{TRsum(t-1)}{n} + TR(t) \\
 + DI(t) &= 100 * \frac{plusDMsum(t)}{TRsum(t)} \\
 - DI(t) &= 100 * \frac{minusDMsum(t)}{TRsum(t)}
 \end{aligned}$$

Directional Movement Index (DX) The DX is usually smoothed with a moving average.

$$DX(t) = \frac{(+DI(t)) - (-DI(t))}{(+DI(t)) + (-DI(t))}$$

### 1) AVERAGE DIRECTIONAL MOVEMENT INDEX (ADX)

The ADX is a Welles Wilder style moving average of the Directional Movement Index (DX).

$$ADX(t) = \frac{ADX(t-1) * (n-1) + DX(t)}{n}$$

### 2) PRICE OSCILLATOR-ABSOLUTE (APO)

The Price Oscillator shows the difference between two moving averages.

$$APO(t) = MA(t1) - MA(t2)$$

where  $t1$  is the slow-moving average and  $t2$  is the fast-moving average.

### 3) BALANCE OF POWER (BOP)

BOP attempts to measure the strength of buyers vs. sellers by assessing the ability of each to push price to an extreme level. BOP calculates raw values for each bar as:

$$BOP(t) = \frac{C_t - O_t}{H_t - L_t}$$

where  $O$  refers to the open price.

### 4) COMMODITY CHANNEL INDEX (CCI)

The CCI is designed to detect beginning and ending market trends.

$$CCI(t) = \frac{\frac{H_t + L_t + C_t}{3} - MA(n)}{0.015 * \frac{1}{n} \sum_{i=t-n}^n MA(i) - C_t}$$

### 5) MOVING AVERAGE CONVERGENCE/DIVERGENCE (MACD)

The Moving Average Convergence Divergence (MACD) is the difference between two Exponential Moving Averages.

$$\begin{aligned}
 shortMA(t) &= 0.15 * C_t + 0.85 * shortMA(t-1) \\
 longMA(t) &= 0.075 * C_t + 0.925 * longMA(t-1) \\
 MACD(t) &= 0.15 * C_t + 0.85 * shortMA(t-1)
 \end{aligned}$$

### 6) MONEY FLOW INDEX (MFI)

The Money Flow Index calculates the ratio of money flowing into and out of a security. T

$$\begin{aligned}
 typicalPrice(t) &= \frac{H_t + L_t + C_t}{3} \\
 moneyFlow(t) &= typicalPrice(t) * V_t
 \end{aligned}$$

The calculation logic of MFI is as follows:

If  $typicalPrice(t) > typicalPrice(t-1)$

$$\begin{aligned}
 positiveMoneyFlow(t) &= positiveMoneyFlow(t-1) \\
 & \quad + moneyFlow(t)
 \end{aligned}$$

else

$$\begin{aligned}
 negativeMoneyFlow(t) &= negativeMoneyFlow(t-1) \\
 & \quad + moneyFlow(t)
 \end{aligned}$$

then

$$\begin{aligned}
 moneyRatio(t) &= \frac{\sum_{t-n}^t positiveMoneyFlow(t)}{\sum_{t-n}^t negativeMoneyFlow(t)} \\
 MFI(t) &= 100 - \frac{100}{1 + moneyRatio(t)}
 \end{aligned}$$

### 7) MOMENTUM (MOM)

The Momentum is a measurement of the acceleration and deceleration of prices.

$$MOM(t) = C_t - C_{t-n}$$

### 8) RELATIVE STRENGTH INDEX (RSI)

The Relative Strength Index (RSI) calculates a ratio of the recent upward price movements to the absolute price movement. The RSI is generated as follows:

If  $C_t > C_{t-1}$

$$\begin{aligned}
 up(t) &= C_t - C_{t-1} \\
 dn(t) &= 0
 \end{aligned}$$

else

$$\begin{aligned}
 up(t) &= 0 \\
 dn(t) &= C_{t-1} - C_t
 \end{aligned}$$

then

$$\begin{aligned}
 upAvg(t) &= \frac{(n-1) upAvg(t-1) + up(t)}{n} \\
 dnAvg(t) &= \frac{(n-1) dnAvg(t-1) + dn(t)}{n} \\
 RSI(t) &= 100 * \frac{upAvg(t)}{upAvg(t) + dnAvg(t)}
 \end{aligned}$$

### C. VOLUME INDICATORS

#### 1) CHAIKIN A/D LINE (AD)

The AD Line is calculated as follows:

$$Clv(t) = \frac{2 * C_t - H_t - C_t}{H_t - L_t}$$

$$AD(t) = AD(t-1) + V_t * Clv(t)$$

#### 2) CHAIKIN OSCILLATOR (ADOSC)

The Chaikin Oscillator is essentially a momentum of the Accumulation/Distribution Line.

$$ADOSC(t) = EMA(AD(n1)) - EMA(AD(n2))$$

where *EMA* is the exponential moving average, *n1* represents the fast period, *n2* represents the slow period.

#### 3) ON BALANCE VOLUME (OBV)

The OBV is a cumulative total of the up and down volume. The calculation logic of *OBV* is as follows:

If  $C_t > C_{t-1}$

$$OBV(t) = OBV(t-1) + V_t$$

If  $C_t < C_{t-1}$

$$OBV(t) = OBV(t-1) - V_t$$

Else

$$OBV(t) = OBV(t-1)$$

where *V* refers to the volume.

### D. VOLATILITY INDICATORS

#### 1) TRUE RANGE (TRANGE)

The True Range is a base calculation that is used to determine the normal trading range of a stock or commodity.

$$TR(t) = \max(H_t, C_{t-1}) - \min(L_t, C_{t-1})$$

#### 2) AVERAGE TRUE RANGE (ATR)

The ATR is a moving average of the True Range.

$$ATR(t) = \frac{1}{n} \sum_{i=1}^n TR(t-i+1)$$

#### 3) NORMALIZED AVERAGE TRUE RANGE (NATR)

The *NATR* is the normalized of *ATR*.

### REFERENCES

- [1] M. Kumar and M. Thenmozhi, "Forecasting stock index returns using ARIMA-SVM, ARIMA-ANN, and ARIMA-random forest hybrid models," *Int. J. Banking, Account. Financ.*, vol. 5, no. 3, pp. 284–308, 2014, doi: [10.1504/IJBAAF.2014.064307](https://doi.org/10.1504/IJBAAF.2014.064307).
- [2] P. C. S. Bezerra and P. H. M. Albuquerque, "Volatility forecasting via SVR-GARCH with mixture of Gaussian kernels," *Comput. Manage. Sci.*, vol. 14, no. 2, pp. 179–196, Apr. 2017, doi: [10.1007/s10287-016-0267-0](https://doi.org/10.1007/s10287-016-0267-0).
- [3] I. K. Nti, A. F. Adekoya, and B. A. Weyori, "A systematic review of fundamental and technical analysis of stock market predictions," *Artif. Intell. Rev.*, vol. 53, pp. 1–51, Aug. 2019.
- [4] A. Ghaznavi, M. Aliyari, and M. R. Mohammadi, "Predicting stock price changes of tehran artmis company using radial basis function neural networks," *Int. Res. J. App. Basic Sci.*, vol. 10, no. 8, p. 972, 2016.
- [5] E. Ahmadi, M. Jasemi, L. Monplaisir, M. A. Nabavi, A. Mahmoodi, and P. A. Jam, "New efficient hybrid candlestick technical analysis model for stock market timing on the basis of the support vector machine and heuristic algorithms of imperialist competition and genetic," *Expert Syst. Appl.*, vol. 94, pp. 21–31, Mar. 2018.
- [6] C. L. Osler, "Currency orders and exchange rate dynamics: An explanation for the predictive success of technical analysis," *J. Finance*, vol. 58, no. 5, pp. 1791–1819, Oct. 2003, doi: [10.1111/1540-6261.00588](https://doi.org/10.1111/1540-6261.00588).
- [7] Y. Zhu and G. Zhou, "Technical analysis: An asset allocation perspective on the use of moving averages," *J. Financial Econ.*, vol. 92, no. 3, pp. 519–544, 2009, doi: [10.1016/j.jfineco.2008.07.002](https://doi.org/10.1016/j.jfineco.2008.07.002).
- [8] H. Bessembinder and K. Chan, "Market efficiency and the returns to technical analysis," *Financial Manage.*, vol. 27, no. 2, p. 5, 1998, doi: [10.2307/3666289](https://doi.org/10.2307/3666289).
- [9] B. R. Marshall, M. R. Young, and L. C. Rose, "Candlestick technical trading strategies: Can they create value for investors?" *J. Banking Finance*, vol. 30, no. 8, pp. 2303–2323, Aug. 2006, doi: [10.1016/j.jbankfin.2005.08.001](https://doi.org/10.1016/j.jbankfin.2005.08.001).
- [10] G. Caginalp and H. Laurent, "The predictive power of price patterns," *Appl. Math. Finance*, vol. 5, nos. 3–4, pp. 181–205, Sep. 1998, doi: [10.1080/135048698334637](https://doi.org/10.1080/135048698334637).
- [11] T.-H. Lu, Y.-C. Chen, and Y.-C. Hsu, "Trend definition or holding strategy: What determines the profitability of candlestick charting?" *J. Banking Finance*, vol. 61, pp. 172–183, Dec. 2015, doi: [10.1016/j.jbankfin.2015.09.009](https://doi.org/10.1016/j.jbankfin.2015.09.009).
- [12] S. Chen, S. Bao, and Y. Zhou, "The predictive power of Japanese candlestick charting in Chinese stock market," *Phys. A, Stat. Mech. Appl.*, vol. 457, pp. 148–165, Sep. 2016, doi: [10.1016/j.physa.2016.03.081](https://doi.org/10.1016/j.physa.2016.03.081).
- [13] M. Zhu, S. Atri, and E. Yegen, "Are candlestick trading strategies effective in certain stocks with distinct features?" *Pacific-Basin Finance J.*, vol. 37, pp. 116–127, Apr. 2016, doi: [10.1016/j.pacfin.2015.10.007](https://doi.org/10.1016/j.pacfin.2015.10.007).
- [14] T.-H. Lu, "The profitability of candlestick charting in the taiwan stock market," *Pacific-Basin Finance J.*, vol. 26, pp. 65–78, Jan. 2014, doi: [10.1016/j.pacfin.2013.10.006](https://doi.org/10.1016/j.pacfin.2013.10.006).
- [15] L. Tao, Y. Hao, H. Yijie, and S. Chunfeng, "K-line patterns' predictive power analysis using the methods of similarity match and clustering," *Math. Probl. Eng.*, vol. 2017, May 2017, Art. no. 3096917, doi: [10.1155/2017/3096917](https://doi.org/10.1155/2017/3096917).
- [16] A. W. Lo, H. Mamaysky, and J. Wang, "Foundations of technical analysis: Computational algorithms, statistical inference, and empirical implementation," *J. Finance*, vol. 55, no. 4, pp. 1705–1765, Aug. 2000, doi: [10.1111/0022-1082.00265](https://doi.org/10.1111/0022-1082.00265).
- [17] D. Yan, Q. Zhou, J. Wang, and N. Zhang, "Bayesian regularization neural network based on artificial intelligence optimisation," *Int. J. Prod. Res.*, vol. 55, no. 8, pp. 2266–2287, Apr. 2017, doi: [10.1080/00207543.2016.1237785](https://doi.org/10.1080/00207543.2016.1237785).
- [18] J.-J. Wang, J.-Z. Wang, Z.-G. Zhang, and S.-P. Guo, "Stock index forecasting based on a hybrid model," *Omega*, vol. 40, pp. 758–766, Dec. 2012, doi: [10.1016/j.omega.2011.07.008](https://doi.org/10.1016/j.omega.2011.07.008).
- [19] B. M. Henrique, V. A. Sobreiro, and H. Kimura, "Literature review: Machine learning techniques applied to financial market prediction," *Expert Syst. Appl.*, vol. 124, pp. 226–251, Jun. 2019, doi: [10.1016/j.eswa.2019.01.012](https://doi.org/10.1016/j.eswa.2019.01.012).
- [20] Y. Xiao, J. Xiao, F. Lu, and S. Wang, "Ensemble ANNs-PSO-GA approach for day-ahead stock e-exchange prices forecasting," *Int. J. Comput. Intell. Syst.*, vol. 7, no. 2, pp. 272–290, 2014, doi: [10.1080/18756891.2013.864472](https://doi.org/10.1080/18756891.2013.864472).
- [21] D. Brownstone, "Using percentage accuracy to measure neural network predictions in stock market movements," *Neurocomputing*, vol. 10, no. 3, pp. 237–250, 1996, doi: [10.1016/0925-2312\(95\)00052-6](https://doi.org/10.1016/0925-2312(95)00052-6).
- [22] F. Kamalov, "Forecasting significant stock price changes using neural networks," *Neural Comput. Appl.*, vol. 32, no. 23, pp. 17655–17667, Dec. 2020, doi: [10.1007/s00521-020-04942-3](https://doi.org/10.1007/s00521-020-04942-3).
- [23] P. Yu and X. Yan, "Stock price prediction based on deep neural networks," *Neural Comput. Appl.*, vol. 32, no. 6, pp. 1609–1628, Mar. 2020, doi: [10.1007/s00521-019-04212-x](https://doi.org/10.1007/s00521-019-04212-x).
- [24] M.-C. Wu, S.-Y. Lin, and C.-H. Lin, "An effective application of decision tree to stock trading," *Expert Syst. Appl.*, vol. 31, no. 2, pp. 270–274, Aug. 2006, doi: [10.1016/j.eswa.2005.09.026](https://doi.org/10.1016/j.eswa.2005.09.026).
- [25] J. Dopke, U. Fritsche, and C. Pierdzioch, "Predicting recessions with boosted regression trees," *Int. J. Forecasting*, vol. 33, no. 4, pp. 745–759, 2017.
- [26] S. Barak, A. Arjmand, and S. Ortobelli, "Fusion of multiple diverse predictors in stock market," *Inf. Fusion*, vol. 36, pp. 90–102, Jul. 2017.

- [27] M.-C. Lee, "Using support vector machine with a hybrid feature selection method to the stock trend prediction," *Expert Syst. Appl.*, vol. 36, no. 8, pp. 10896–10904, Oct. 2009, doi: [10.1016/j.eswa.2009.02.038](https://doi.org/10.1016/j.eswa.2009.02.038).
- [28] M. V. Subha and S. T. Nambi, "Classification of stock index movement using k-nearest neighbours (k-NN) algorithm," *WSEAS Trans. Inf. Sci. Appl.*, vol. 9, no. 9, pp. 261–270, 2012.
- [29] A. Booth, E. Gerding, and F. McGroarty, "Automated trading with performance weighted random forests and seasonality," *Expert Syst. Appl.*, vol. 41, no. 8, pp. 3651–3661, Jun. 2014, doi: [10.1016/j.eswa.2013.12.009](https://doi.org/10.1016/j.eswa.2013.12.009).
- [30] C. Lohrmann and P. Luukka, "Classification of intraday S&P500 returns with a random forest," *Int. J. Forecasting*, vol. 35, no. 1, pp. 390–407, Jan. 2019.
- [31] T. Fischer and C. Krauss, "Deep learning with long short-term memory networks for financial market predictions," *Eur. J. Oper. Res.*, vol. 270, no. 2, pp. 654–669, 2018, doi: [10.1016/j.ejor.2017.11.054](https://doi.org/10.1016/j.ejor.2017.11.054).
- [32] A. Mundra, S. Mundra, V. K. Verma, and J. S. Srivastava, "A deep learning based hybrid framework for stock price prediction," *J. Intell. Fuzzy Syst.*, vol. 38, no. 5, pp. 5949–5956, May 2020.
- [33] A. Kazem, E. Sharifi, F. K. Hussain, M. Saberi, and O. K. Hussain, "Support vector regression with chaos-based firefly algorithm for stock market price forecasting," *Appl. Soft Comput.*, vol. 13, no. 2, pp. 947–958, 2013.
- [34] J. Zhang, L. Li, and W. Chen, "Predicting stock price using two-stage machine learning techniques," *Comput. Econ.*, vol. 57, pp. 1237–1261, 2021, doi: [10.1007/s10614-020-10013-5](https://doi.org/10.1007/s10614-020-10013-5).
- [35] B. Weng, L. Lu, X. Wang, F. M. Megahed, and W. Martinez, "Predicting short-term stock prices using ensemble methods and online data sources," *Expert Syst. Appl.*, vol. 112, pp. 258–273, Dec. 2018.
- [36] D. Kumar, S. S. Meghwani, and M. Thakur, "Proximal support vector machine based hybrid prediction models for trend forecasting in financial markets," *J. Comput. Sci.*, vol. 17, pp. 1–13, Nov. 2016.
- [37] M. Göçken, M. Özçalıcı, A. Boru, and A. T. Dosdoğru, "Stock price prediction using hybrid soft computing models incorporating parameter tuning and input variable selection," *Neural Comput. Appl.*, vol. 31, no. 2, pp. 577–592, Feb. 2019.
- [38] J. Patel, S. Shah, P. Thakkar, and K. Kotecha, "Predicting stock market index using fusion of machine learning techniques," *Expert Syst. Appl.*, vol. 42, no. 4, pp. 2162–2172, Mar. 2015, doi: [10.1016/j.eswa.2014.10.031](https://doi.org/10.1016/j.eswa.2014.10.031).
- [39] F. Zhou, Q. Zhang, D. Sornette, and L. Jiang, "Cascading logistic regression onto gradient boosted decision trees for forecasting and trading stock indices," *Appl. Soft Comput.*, vol. 84, Nov. 2019, Art. no. 105747, doi: [10.1016/j.asoc.2019.105747](https://doi.org/10.1016/j.asoc.2019.105747).
- [40] W. Bao, J. Yue, and Y. Rao, "A deep learning framework for financial time series using stacked autoencoders and long-short term memory," *PLoS ONE*, vol. 12, no. 7, Jul. 2017, Art. no. e0180944, doi: [10.1371/journal.pone.0180944](https://doi.org/10.1371/journal.pone.0180944).
- [41] C. Krauss, X. A. Do, and N. Huck, "Deep neural networks, gradient-boosted trees, random forests: Statistical arbitrage on the S&P 500," *Eur. J. Oper. Res.*, vol. 259, no. 2, pp. 689–702, Jun. 2017, doi: [10.1016/j.ejor.2016.10.031](https://doi.org/10.1016/j.ejor.2016.10.031).
- [42] J. Patel, S. Shah, P. Thakkar, and K. Kotecha, "Predicting stock and stock price index movement using trend deterministic data preparation and machine learning techniques," *Expert Syst. Appl.*, vol. 42, no. 1, pp. 259–268, 2015.
- [43] K.-J. Kim and I. Han, "Genetic algorithms approach to feature discretization in artificial neural networks for the prediction of stock price index," *Expert Syst. Appl.*, vol. 19, no. 2, pp. 125–132, Aug. 2000, doi: [10.1016/S0957-4174\(00\)00027-0](https://doi.org/10.1016/S0957-4174(00)00027-0).
- [44] C.-F. Tsai and Y.-C. Hsiao, "Combining multiple feature selection methods for stock prediction: Union, intersection, and multi-intersection approaches," *Decis. Support Syst.*, vol. 50, no. 1, pp. 258–269, 2010, doi: [10.1016/j.dss.2010.08.028](https://doi.org/10.1016/j.dss.2010.08.028).
- [45] H. Li, L.-Y. Hong, Y.-C. Mo, B.-Z. Zhu, and P.-C. Chang, "Restructuring performance prediction with a rebalanced and clustered support vector machine," *J. Forecasting*, vol. 37, no. 4, pp. 437–456, Jul. 2018, doi: [10.1002/for.2512](https://doi.org/10.1002/for.2512).
- [46] L. F. S. Vilela, R. C. Leme, C. A. M. Pinheiro, and O. A. S. Carpinteiro, "Forecasting financial series using clustering methods and support vector regression," *Artif. Intell. Rev.*, vol. 52, no. 2, pp. 743–773, Aug. 2019, doi: [10.1007/s10462-018-9663-x](https://doi.org/10.1007/s10462-018-9663-x).
- [47] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001, doi: [10.1023/A:1010933404324](https://doi.org/10.1023/A:1010933404324).
- [48] J. H. Friedman, "Stochastic gradient boosting," *Comput. Statist. Data Anal.*, vol. 38, no. 4, pp. 367–378, 2002, doi: [10.1016/S0167-9473\(01\)00065-2](https://doi.org/10.1016/S0167-9473(01)00065-2).
- [49] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2013, pp. 6645–6649, doi: [10.1109/ICASSP.2013.6638947](https://doi.org/10.1109/ICASSP.2013.6638947).
- [50] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.



**YAOHU LIN** was born in Fujian, China, in April 1984. He received the B.S. and M.S. degrees in computer science and engineering from Beihang University, China, where he is currently pursuing the Ph.D. degree with the School of Economics and Management. His main research interests include machine learning, deep learning, and relation classification.



**SHANCUN LIU** was born in Hebei, China, in October 1964. He received the B.S. degree in mathematics from Lanzhou University, the M.S. degree in fundamental mathematics from Wuhan University, and the Ph.D. degree from Beihang University. He is currently a Full Professor with Beihang University. He has published more than 100 articles on different journals, such as, *Applied Economics Letters*, *Economics Letters*, *Journal of Systems Science and Complexity*, and *Finance Research Letters*.



**HAIJUN YANG** (Member, IEEE) was born in Tianjin, China, in January 1970. He received the B.S. and M.S. degrees in mathematics and management science from Nankai University, China, and the Ph.D. degree from Tianjin University.

He is currently a Full Professor with Beihang University. He is also a Fulbright Scholar. He has published more than 30 articles on different journals, such as, *Journal of Evolutionary Economics*, *Entropy*, and *International Journal of Information Technology & Decision Making*.

Dr. Yang is a member of AEA.



**HARRIS WU** was born in Shenyang, China. He received the Ph.D. degree from University of Michigan.

He is currently a Full Professor with the Department of Information Technology and Decision Sciences, College of Business and Public Administration, Old Dominion University. He has published some articles on *Decision Support Systems*, *Information and Management*, and *Applied Soft Computing Journal*. His research interests

include data analytics, social media, cybersecurity, text mining, enterprise information systems, and system integration.

...