# An interval constraint-based trading strategy with social sentiment for the stock market

Mingchen Li[1,2], Kun Yang[1,2,3], Wencan Lin[1,2], Yunjie Wei[1,3]* and Shouyang Wang[1,3]

*Correspondence:
weiyunjie@amss.ac.cn

[1] Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Zhongguancun East Road, # 55, Haidian District, Beijing 100190, People's Republic of China
[2] School of Mathematical Sciences, University of Chinese Academy of Sciences, Beijing 100190, People's Republic of China
[3] Center for Forecasting Science, Chinese Academy of Sciences, Beijing 100190, People's Republic of China

## Abstract

Developing effective strategies to earn excess returns in the stock market is a cutting-edge topic in the field of economics. At the same time, stock price forecasting that supports trading strategies is considered one of the most challenging tasks. Therefore, this study analyzes and extracts news media data, expert comments, social opinion data, and pandemic text data using natural language processing, and then combines the data with a deep learning model to forecast future stock price patterns based on historical stock prices. An interval constraint-based trading strategy is constructed. Using data from several typical stocks in the Chinese stock market during the COVID-19 period, the empirical studies and trading simulations show, first, that the sentiment composite index and the deep learning model can improve the accuracy of stock price forecasting. Second, the interval constraint-based trading strategy based on the proposed approach can effectively enhance returns and thus, can assist investors in decision-making.

**Keywords:** Stock price forecasting, Deep learning, Sentiment analysis, Trading strategy, COVID-19 era

## Introduction

The stock market is a crucial conduit for businesses to raise funds from investors and is a financial ecosystem connecting corporations and investors (Zhong and Enke 2019). The enormous trading volume and profitability of the stock market continues to attract investors and traders keen to employ this system to maximize their profits (Adam et al. 2016). However, the stock market has the characteristic of extreme volatility and non-stationarity, making it prone to numerous complicated shocks and games. Consequently, there are obstacles to devising solid trading methods and making profitable investment decisions (Gu and Peng 2019). Since the turn of the 20th century, a constant stream of financial institutions and researchers have been developing stock price forecasting models. With the expansion of computer technology, an increasing number of superior models, including deep learning, seek to decrease stochasticity and identify consistent trends by collecting and evaluating historical data and useful technical indicators (Salisu and Vo 2020; Liu et al. 2020).

At the beginning of the century, Hinton et al. introduced the notion of deep learning, thereby resolving the enduring impasse surrounding the arduous training of deep neural

networks (Hinton et al. 2006). Meanwhile, Bengio et al. established a robust framework for the application of deep learning in addressing language modeling challenges (Bengio et al. 2000; Khurana et al. 2023). Subsequently, deep learning and natural language processing (NLP) have gained extensive utilization across diverse domains (Lecun et al. 2015). Within the realm of finance, these techniques have been employed to facilitate financial forecasting, credit default prediction, mortgage risk estimation, and risk–return management, among other applications (Zhong and Enke 2019; Alonso Robisco and Carbó Martínez 2022; Calomiris and Mamaysky 2019; Xing et al. 2018). These technological advances along with the growth of social media have collectively driven the widespread use of unstructured data (especially news data and social media data), which has improved the predictive power of models.

Specifically, the factors that make this phenomenon noteworthy are as follows. First, the efficient market hypothesis assumes that market investors are rational and seek the greatest possible profits (Fama 1970). However, as the Dutch tulip bubble and the American Internet bubble showed, investors are not always rational (Audrino et al. 2020). According to previous studies, investor sentiment and stock returns are mutually limiting and influential. In other words, the price of stocks in the market is not only defined by the intrinsic worth of companies, but is also heavily impacted by the investing subject; that is, psychological considerations and investor behavior have a significant impact on the price decisions and movements of stocks (Bustos and Pomares-Quimbaya 2020; Liang et al. 2020). Second, social platforms or news websites, as types of digital economy presentation, are increasingly crucial avenues for consumers or investors to exchange perspectives, feelings, and knowledge. Compared to conventional data sources, these platforms' data offer the benefits of a broad user base, high socializing, high engagement, and rapid reaction times (Audrino et al. 2020; Hong et al. 2017). The efficient use of this information and its integration into research on the stock market is a highly rewarding and challenging task.

In the COVID-19 era, it is worth considering what texts should be employed for stock price analysis. In conjunction with previous studies, the following categories of data are the focus of this study: 1. news, which is a significant way for the general public to get official or more formal information through the media (Narayan 2019); 2. comments, especially from investors and practitioners, which are a synthesis of sentiment from relative professionals; 3. social media data, which reflect the collective wisdom of the general people (Teti et al. 2019); and 4. pandemic data. COVID-19 has triggered significant stock market volatility, as is common knowledge. To prevent the spread of the disease, governments across the globe implemented stricter policy controls, including limitations on labor mobility and quarantines (Salisu and Vo 2020; He et al. 2020; Li et al. 2022). This resulted in a number of challenges in global supply chains, including reduced supply and decreased demand, which discouraged investment and reduced business and consumer confidence. In this scenario, global stock markets suffer setbacks (OECD 2020). COVID-19 has had a significant impact on equity markets, and thus, pandemic-related data should also be considered. No previous study on stock price volatility has considered all of these factors.

Therefore, this study proposes a novel forecasting and interval-constraint trading approach based on deep learning and sentiment analysis for forecasting and simulating

stock price fluctuations in the COVID-19 period in conjunction with big data. First, text data from four different perspectives are collected: news media, expert comments, social opinion, and the pandemic. Second, the relevant texts are then analyzed with natural language processing techniques to provide sentiment indexes that represent a combination of official, popular, and social contexts. Third, with the support of deep learning models, we employ historical data, search engine data, and sentiment indexes to forecast stock prices, including point value forecasts and interval forecasts. Fourth, we employ a number of assessment criteria and statistical tests to test the forecasting capacity of the model. Lastly, traditional trading strategies are based on forecasts for long or short positions, which carry a lot of risk, especially when located in periods of high volatility. This study proposes combining the interval estimation algorithm to add an insurance policy to a trading strategy that can support significant improvement in trading returns under high volatility conditions.

The innovation of this study lies in three aspects. First, the data are considered comprehensively. The data form includes time-series data and textual data; the data sources include news media, stock market experts, and the general public; and the data connotation includes stock characteristics, practitioner sentiment, and pandemic information. Second, this study makes innovative use of models. We employ temporal convolutional network (TCN), an effective deep learning model, combined with interval estimation algorithms to generate a reliable forecasting framework (containing point forecasts and interval forecasts). Third, this is the first study to construct a trading strategy based on interval constraints. Interval restrictions are added to the general point forecast-based trading strategy to avoid irrational investments in high volatility periods or to generate huge returns.

The rest of the paper is structured as follows. The literature review and discussion related to this study is presented in "Literature review". The proposed methodology and related models are presented in "Methodology". "Empirical analysis" shows the experimental procedure, including the data collection, data processing, forecasting results, and trading simulations. "Conclusion" presents the conclusions. Finally, the discussion and prospects for future research are provided in "Discussion and prospects".

## Literature review

This section consists of two parts: the first introduces multiple types of models applied to stock price forecasting, including statistical models and artificial intelligence models; the second introduces the application of social media data in previous stock price forecasting studies, including search indexes and text data.

### Forecasting models for stock price

An intriguing topic in finance and forecasting research has been how to make more accurate stock price forecasts (Kumbure et al. 2022). Numerous models have been developed to describe the volatility and trend of various stock prices on different exchange platforms. These models may be categorized into two groups: traditional statistics and artificial intelligence (mainly machine learning algorithms) (Deng et al. 2022; Liu et al. 2021).

Traditional statistics models include the vector autoregressive model, autoregressive conditional heteroskedasticity model, and autoregressive integrated moving average (ARIMA) model, among others (Ribeiro Ramos 2003; Wang et al. 2022). Traditional statistical models, which are effective instruments for illuminating the inner workings of financial market functioning, are frequently utilized in stock market forecasting and analysis. For instance, Jiang et al. (2021) found that oil prices have a significant impact on stock returns in the short term by using a structural threshold vector autoregression model. However, when using traditional statistics models, the linearity or stationarity assumptions in statistical data should be satisfied, which is typically challenging when using high-frequency data from the stock market (Kumbure et al. 2022; Tang et al. 2022). Consequently, conventional statistical models have limitations in forecasting stock prices (Lin et al. 2021).

The abovementioned issues are not present in machine learning models based on artificial intelligence algorithms (Vuong et al. 2022). When complex structures, such as nonlinear high-frequency stock trading data, are present, machine learning algorithms can provide more accurate forecasts (Yun et al. 2022). Many instances of anticipating stock prices have been used in classical machine learning research. Gupta et al. (2019) examined the predictability of stock returns using the quantile random forests method. They used indicators of inequality based on consumption and income to forecast stock returns, providing a new approach for predicting stock returns. Sadaei et al. (2016) and Kao et al. (2013) applied fuzzy set and support vector regression (SVR) to forecast stock prices, respectively. They both achieved good forecasting results. For application of other classical machine learning methods in stock price forecasting, refer to Na and Kim (2021), Zhang and Lou (2021), and Shahi et al. (2020), among others.

The forecasting of stock price using machine learning techniques is now a topic of focus with rich research. Many researchers have proposed novel machine learning methods, which can achieve more robust and more accurate forecast results. Deng et al. (2022) combined a deep learning algorithm with multivariate empirical mode decomposition, and further built a multi-input and multi-output network framework to achieve multi-step forecasting of stock prices. The empirical results show this combination method can realize better prediction results. Although the combination of a machine learning algorithm with a decomposition integration method can raise forecasting accuracy, it also makes the computation more difficult. To resolve this problem, Guo et al. (2022) employed a system clustering method and particle swarm optimization to construct a decomposition and reconstruction model, which not only reduced the complexity of the algorithm, but also obtained more accurate forecasting results. Additionally, several studies have combined different machine learning techniques to overcome the drawbacks of a single technique. For example, Ghosh et al. (2022) mixed random forest and long short-term memory network (LSTM) to achieve more accurate stock price forecasting results than the single machine learning method.

TCN is a novel type of neural network improved from the one-dimensional convolutional neural network. TCN has been shown to outperform LSTM in numerous domains, including voice processing, machine translation, and time-series forecasting, while retaining the robust feature extraction capabilities of conventional convolutional neural networks (Zhu et al. 2020; Shomron and Weiser 2019).

In such a scenario, many benchmark models of the two kinds mentioned above are used in this study; they include seasonal autoregressive integrated moving average (SARIMA), exponential smoothing (ES), SVR, extreme learning machine (ELM), back propagation neural network (BPNN), LSTM, and TCN. TCN is the main model of interest owing to its benefits of higher parallelism, stable gradients, and minimal memory requirements.

### Social media and stock price forecasting

As Web 2.0 takes off, more and more investors are turning to the Internet to obtain and share real-time stock-related news (Sanford 2022). Owing to the rapid diffusion of influence through the Internet, experts' and other influential persons' written views on stocks may affect the decisions of others. The effects are dual (Maqsood et al. 2020; Gu and Kurov 2020). On the one hand, Internet user comments and event information can have a substantial effect on the price of a stock. On the other hand, sudden fluctuations in stock price may prompt the development and transmission of relevant information (e.g., government viewpoints), which may then impact public perceptions of prospective investment strategies (Shomron and Weiser 2019; Jin et al. 2020). Textual material (e.g., blogs, reviews, and status updates), online search queries (e.g., Google Trends), tags, and personal information are common forms of social media data. Social media data include individual views, ideas, and actions that affect stock market predictability and result in significant profits or losses (Bijl et al. 2016).

Textual data, particularly news, is a superior source of hidden information to quantitative data, because the former enables the forecasting of financial patterns with supporting evidence (Liang et al. 2020; Gu and Peng 2019). For instance, a news story about a corporation containing the terms "resignation" or "risk of default" leads the investor to anticipate a decrease in the stock price. In addition, stock market trends may be influenced by news pertaining to a variety of unforeseen events, such as terrorism, war, civil unrest, economic and political shocks, and natural disasters (Nassirtoussi et al. 2015). Similarly, Chen et al. (2014) demonstrated how information from user-generated research papers on SeekingAlpha may be utilized to anticipate earnings and stock returns. However, it is difficult for retail investors to comprehend the context of research papers completely.

Numerous social and economic effects may be forecast by Web search queries, which has attracted great interest. For example, Bijl et al. (2016) investigated the prospect of forecasting stock returns using Google Trends data and discovered that high Google search volumes are associated with negative returns. Kim et al. (2019) demonstrated that an increase in Google searches is predictive of a rise in the volatility and trading volume of the top companies listed on the Oslo Stock Exchange. Considering the national conditions of China, the Baidu index is an invaluable source for monitoring and forecasting Chinese socioeconomic activities.

According to behavioral finance theory, social network information may impact people's financial decisions to some level. To help investors understand the connection between social networks and stock prices, Liu et al. (2021) constructed daily social networks utilizing information obtained from EastyMoney, the largest social media site in China, about individuals and the stocks they followed. The empirical data indicate that

Li *et al. Financial Innovation*      (2024) 10:56

Page 6 of 31

the social network variable can greatly improve forecasting accuracy. Zhang et al. (2018) investigated characteristics relating to collective mood and perception of stock relatedness based on messages from Xueqiu, a well-known Chinese social network similar to Twitter that caters to investors, and uses nonlinear models to anticipate stock price changes. However, both EastyMoney and Xueqiu are focused sites that cater to niche audiences, ignoring the hotspots and opinions from public media.

Table 1 briefly discusses previous literature, highlighting the limitations of previous studies. First, sentiment analysis of textual data has rarely been performed jointly across multiple platforms, perspectives, and participants. Second, TCN, although applied to stock forecasting, does not combine interval estimation with point forecasting to consider the situation. Third, the double insurance trading strategy of joint point and interval forecasting has not been studied. Therefore, this study proposes a comprehensive and integrated forecasting framework and trading strategy to fill the research gaps.

## Methodology

### Temporal convolutional network

For the objective of time-series modeling, a novel algorithm that can be used to solve massively parallel computation problems in recurrent neural networks is the TCN, whose effectiveness has been verified (Bai et al. 2018; Zhu et al. 2020). Figure 1 illustrates the basic framework of the TCN model being applied to this study. The TCN model is based on a $1-D$ fully convolutional network (FCN), where each hidden layer shares the same length as the input layer. The advantages of TCN can be seen in three aspects: causal convolution, dilated convolution, and residual block.

*Causal convolutions.* First, since the convolution is causal, only the historical and present-day inputs, and not the inputs in the future, are connected to the current output. Second, the TCN only convolves the inputs at the present time and previous time since the architecture may accept a time series of any length and map it to output data of the same length.

*Dilated convolution.* The dilated convolutions can enable an exponentially large receptive field to ensure the applicability of the causal convolution on sequence tasks (Bai et al. 2018). For a $1-D$ sequence input $X \in R^n$ and a filter $f$, the dilated convolution operation on $s$th element in the sequence $X$ is defined as

$$F(X_s) = \left( X *_d f \right) = \sum_{i=0}^{k-1} f(i)X_{s-d \cdot i} \tag{1}$$

where $d$ is the dilation factor, $k$ is the filter size, and the subscript $X_{s-d \cdot i}$ denotes the direction of the past. Therefore, dilation is the same as adding a fixed step after every two consecutive filter taps. A dilated convolution becomes a regular convolution at $d = 1$. A ConvNet's receptive field is effectively expanded by larger dilation because it allows the top-level output to reflect a wider range of inputs. Figure 1a provides an example.

*Residual block.* It has been demonstrated that a residual learning framework makes network training easier and that residual blocks are useful for deep networks (Wu et al. 2021). The residual block for our baseline TCN is shown in Fig. 1b. The rectified linear unit (ReLU) is used to account for the two layers of dilated causal convolution and

Li *et al. Financial Innovation*      (2024) 10:56

Page 7 of 31

**Table 1** A brief list of selected studies

| References | Forecasted goals | Data frequency | Forecasting Method | Evaluation Indicators |
|---|---|---|---|---|
| Gupta et al. (2019) | UK stock-return | Quarterly | Quantile random forests | *p*-value |
| Kim et al. (2019) | Norway Oslo Børs OBX index | Weekly | Panel data regressions | R2 |
| Gu and Peng (2019) | Shanghai Composite Index, CITIC Securities | Weekly | TVPDM | CR |
| Sadaei et al. (2016) | TAIEX, NASDAQ, DJI, S &P 500 | Daily | Hybrid model | RMSE, MAE |
| Kao et al. (2013) | SSEC, BB, DJ, N225 | Daily | Wavelet-MAR-SVR | RMSE, MAE, MAPE, RMSPE |
| Guo et al. (2022) | SH60000 stock | Daily | EEMD-Cluster SVR-PSO-LSTM | MAE, MSE, RMSE |
| Ghosh et al. (2022) | S &P 500 | Daily | Random forests and LSTM | Sharpe ratio, standard deviation, VaR, average return |
| Na and Kim (2021) | Korean stock market data | Daily | ANN | Annualized average returns Sharpe ratio, information ratio |
| Deng et al. (2023) | Capital flow data | Daily | XGBoost, SHAP approach | F1-score, accuracy |
| Maqbool et al. (2023) | Stock data of Reliance, Tata Motors, Tata Steel and HDFC | Daily | MLP-Regressor | MAPE,F1-score, accuracy |
| Zhong and Enke (2019) | S &P 500 | Daily | ANN | MSE |
| Liang et al. (2020) | Shanghai Stock Exchange Composite Index | Daily | HAR | R2 |
| Yun et al. (2022) | Exxon Mobil stock | Daily | XGBoost | RMSE,MSE,MAE,R2 |
| Zhang et al. (2018) | A-share market stocks | Daily | SVM,MLP | ACC, AUC |
| Liu et al. (2021) | SSE 50 constituent stocks | Daily | LSTM | RMSE,MAPE |
| Maqsood et al. (2020) | Stock exchangedata for 15 companies | Daily | LR,SVR,DL | RMSE,MAE |
| This study | Five A-share stocks, SSEC | Daily | TCN/Gi-MLP | MAE, RMSE, MAPE |

| References | Considering economic variables | Considering news sentiment | Considering practitioner sentiment | Considering the context of the event (e.g., public health event, terrorist attack, etc.) |
|---|---|---|---|---|
| Gupta et al. (2019) | ✓ | × | × | × |
| Kim et al. (2019) | ✓ | ✓ | × | × |
| Gu and Peng (2019) | × | × | × | × |
| Sadaei et al. (2016) | × | × | × | × |
| Kao et al. (2013) | ✓ | × | × | × |
| Guo et al. (2022) | × | × | × | × |
| Ghosh et al. (2022) | × | × | × | × |
| Na and Kim (2021) | ✓ | × | ✓ | × |
| Deng et al. (2023) | ✓ | × | ✓ | × |
| Maqbool et al. (2023) | ✓ | ✓ | × | × |

Li *et al. Financial Innovation*      (2024) 10:56

Page 8 of 31

**Table 1** (continued)

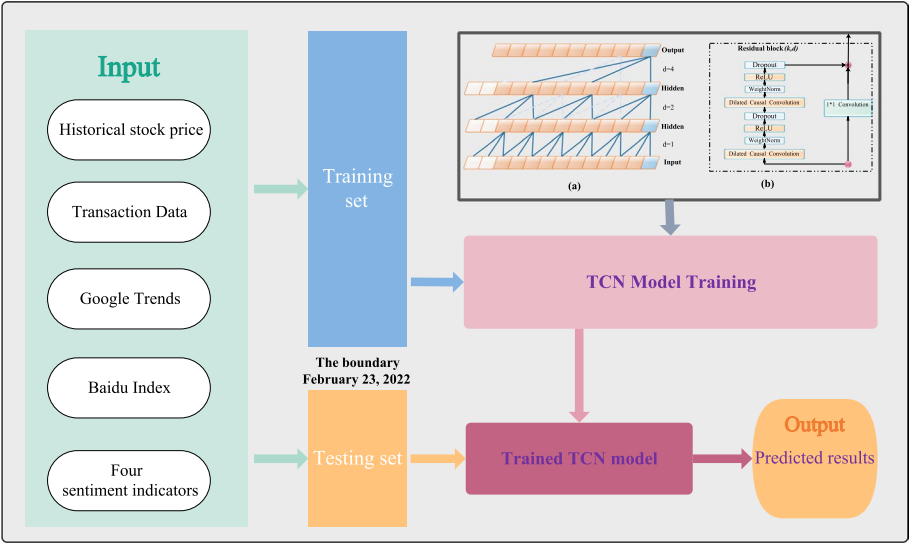| References | Considering economic variables | Considering news sentiment | Considering practitioner sentiment | Considering the context of the event (e.g., public health event, terrorist attack, etc.) |
|---|---|---|---|---|
| Zhong and Enke (2019) | ✓ | ✗ | ✗ | ✗ |
| Liang et al. (2020) | ✓ | ✓ | ✗ | ✓ |
| Yun et al. (2022) | ✓ | ✗ | ✗ | ✗ |
| Zhang et al. (2018) | ✓ | ✗ | ✓ | ✗ |
| Liu et al. (2021) | ✓ | ✗ | ✓ | ✗ |
| Maqsood et al. (2020) | ✗ | ✓ | ✗ | ✓ |
| This study | ✓ | ✓ | ✓ | ✓ |



**Fig. 1** The framework of the TCN model

non-linearity in the TCN inside a residual block. Then, weight normalization to the convolutional filters can be applied. Additionally, after each dilated convolution, a spatial dropout for regularization was implemented: during each training step, an entire channel is blank out (Poernomo and Kang 2018).

**Gi-MLP**

Multilayer perceptrons (MLP) is a widely used neural network algorithm, which can be used to find the internal relationship of high-frequency point-valued time series (PTS). Roque et al. (2007) and Sun et al. (2018) extended the MLP used for modeling PTS to traditional interval-valued time series. Furthermore, Han et al. (2012) defined the generalized random interval allows the addition of intervals until the data collection is complete, which makes interval operation easier. The generalized random interval multilayer

perceptron method (Gi-MLP) algorithm based on generalized random interval includes $N$ inputs, $M$ outputs, and $K$ hidden layers. Each hidden layer has $p_j(j = 1, \ldots, P)$ hidden nodes. For simplicity, we will introduce this algorithm with one hidden layer, namely $K = 1$.

The input is $N$ generalized random interval data, namely $x_i = [x_i^L, x_i^R], i = 1, \ldots, N$, where $x_i^L$ and $x_i^R$ represent the left and right endpoints of the interval respectively. The value of the $j$th hidden node is obtained from the linear combination of these $N$ generalized random intervals and two trend interval terms. The specific form is:

$$L_j = \beta_j^l[1, 1] + \beta_j^r[-\frac{1}{2}, \frac{1}{2}] + \sum_{i=1}^{N} \alpha_{ji} x_i, \quad j = 1, \ldots, P, \tag{2}$$

where $\beta_j^l, \beta_j^r, \alpha_{ji}$ are constant parameters. The $\beta_j^l$ is the trend item representing the overall level of the interval, while the coefficient $\beta_j^r$ represents the trend term of interval radius fluctuation. $F$ is an active function, which is usually chosen as a hyperbolic tangent function and sigmoid function. We can obtain the outputs of the hidden layer after using the active function:

$$H_j = [H_j^L, H_j^R] = F(L_j) = [F(L_j^L), F(L_j^R)], \quad j = 1, \ldots, P. \tag{3}$$

After obtaining the interval value of the hidden layer node, the $i$th interval value of the output layer is constructed as the linear combination of the node value of the hidden layer and two types of trend terms. It has the following form:

$$\hat{Y}_i = \zeta_i^l[1, 1] + \zeta_i^r[-\frac{1}{2}, \frac{1}{2}] + \sum_{j=1}^{P} \theta_{ij} H_j, \quad i = 1, \ldots, M, \tag{4}$$

where the meaning of $\zeta_i^l, \zeta_i^r, \theta_{ij}$ is the same as (2). Figure 2 illustrates the basic framework of the Gi-MLP applied to this study.
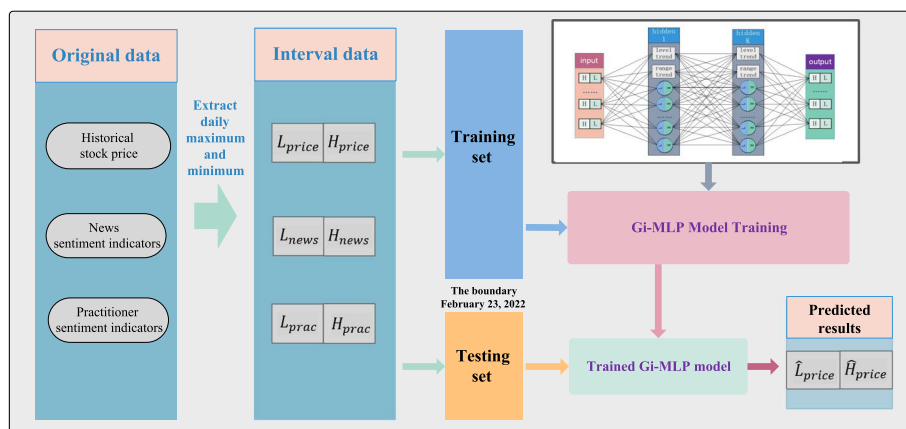


**Fig. 2** The structure diagram of Gi-MLP

**Table 2** The parameters of the forecasting models

| Model | Parameters | Determination approach | Values |
|---|---|---|---|
| SARIMA | Seasonal period | Preset | 5 |
| | AR | Partial autocorrelation function | [0,5] |
| | I | Augmented dickeye-fuller test | 0 or 1 |
| | Ma | Autocorrelation function | [0,5] |
| SVR | Regularization coefficient c | Grid search | [1,300] |
| | Kernel parameter g | Grid search | $[2^{-5}, 2^5]$ |
| ELM& BPNN | Input dimension | Preset | – |
| | Number of hidden layer nodes | Trial and error approach | 24 |
| | Output dimension | Preset | 1 |
| | Maximum of epochs | Preset | 100 |
| | Number of parameters (BPNN) | – | 361 |
| LSTM | Input dimension | Preset | – |
| | Number of hidden layer nodes | Trial and error approach | 24 |
| | Output dimension | Preset | 1 |
| | Maximum of epochs | Preset | 100 |
| | Number of parameters | – | 3673 |
| TCN | Input dimension | Preset | – |
| | Nb_filters | Trial and error approach | 32 |
| | Kernel_size | Trial and error approach | 2 |
| | Output dimension | Preset | 1 |
| | Maximum of epochs | Preset | 100 |
| | Number of parameters | – | 24225 |

**Table 3** Running time of each model

| | SARIMA | SVR | ELM | LSTM | TCN | The proposed approach |
|---|---|---|---|---|---|---|
| Running time (s) | 0.041 | 0.061 | 0.036 | 4.372 | 9.697 | 9.723 |
| Running time (s) (Including sentiment) | 0.041 | 0.062 | 0.036 | 4.789 | 9.931 | 10.024 |

### Benchmark models and parameter setting

In this study, we used seven benchmark models, namely, SARIMA, ES, SVR, ELM, BPNN, LSTM and TCN, for stock price forecasting, and the parameter settings of these models are listed in Table 2, including the determination methods. Numpy, Pandas, Tensorflow, and Keras packages are used in python3 for model training and testing. SVR is implemented using the libsvm toolbox in MATLAB 2018b.

Furthermore, it is imperative for researchers to consider the operational efficiency of neural network models. Hence, the number of parameters of the neural network is presented in Table 2. The running time of each neural network is presented in Table 3. The recorded time values are obtained by averaging the results from 10 separate runs. It is observed that the neural networks, including the proposed model, exhibit longer running times compared to conventional machine learning models and time-series models. However, the overall time required by the neural networks remains within acceptable limits when compared to the forecasting of daily data. The central processing unit (CPU)

employed in this research is the Intel(R) Core(TM) i9-10900K CPU operating at a frequency of 3.70 GHz. It is accompanied by a random access memory (RAM) capacity of 32.0 GB. Additionally, the graphics processing unit (GPU) utilized is the NVIDIA GeForce RTX 3090.

### The framework

The methodology combining natural language processing and deep learning forecasting is constructed with four parts: data collection, data processing, empirical forecasting, and trading simulations, as shown in Fig. 3.

*Stage 1*, data collection. We collect stock prices as the forecast target. We also collect transaction-related feature data, search engine data, news media reports, expert comments, public opinion, and pandemic-related text data.
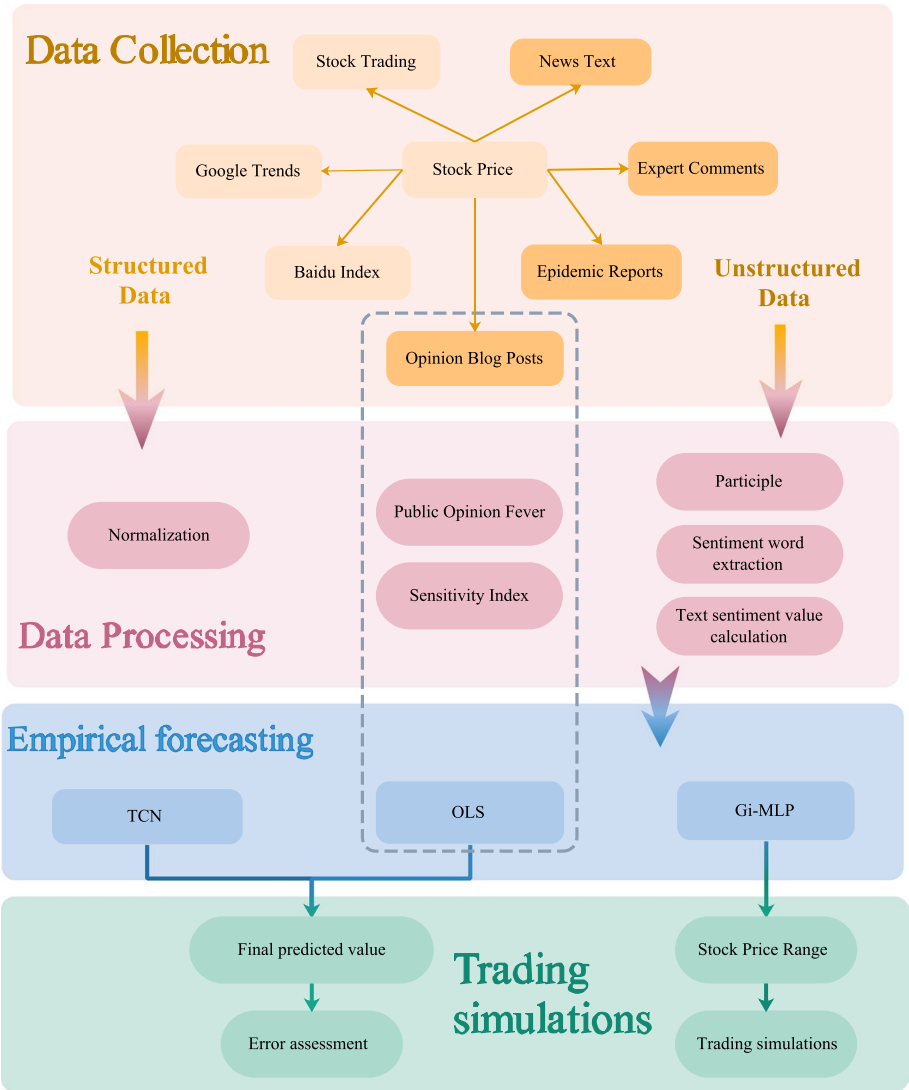


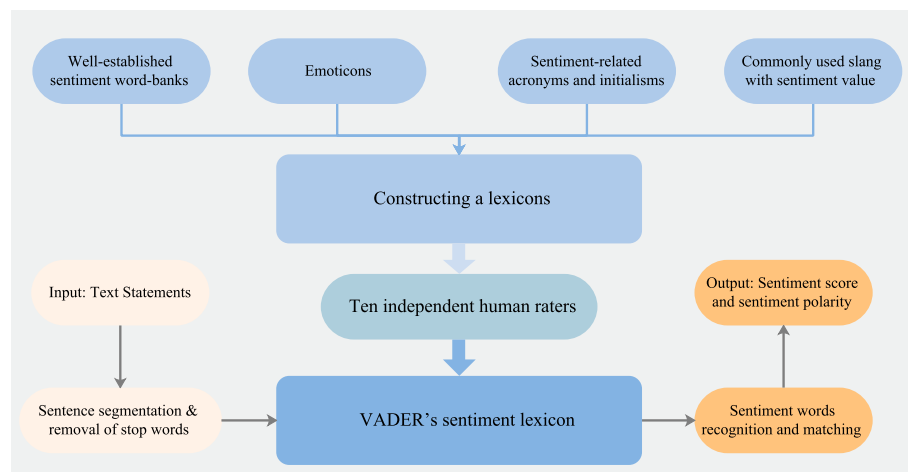**Fig. 3** The framework of the research

**Fig. 4** The framework of natural language processing

*Stage 2*, data processing. We standardize the structured data, such as time series, to eliminate the dimension (unit); meanwhile, we extract the sentiment from the sentiment-rich text information and extract the public opinion fever and sensitivity index from the text comments containing public opinion. The framework of natural language processing used in this study is shown in Fig. 4. Some detailed descriptions are placed "Data processing" in section.

*Stage 3*, empirical forecasting. We use three effective models to forecast the corresponding series. Finally, the results of point forecasting and interval forecasting can be obtained. We will also use error evaluation metrics and statistical tests to evaluate the results.

*Stage 4*, trading simulations. We will combine the results of point forecasts and interval forecasts to innovatively propose trading strategies that will guide investors' decisions with insurance recommendations for great returns.

## Empirical analysis

### Data collection

In the COVID-19 era, scholars from a variety of fields have concentrated on emerging markets since they have become the essential contributor to global economic growth. Numerous researchers have examined the characteristics and performance of developing economies, such as innovation creation, international trade, etc (Gu and Peng 2019; Kang 2018). As the world's most largest emerging economy, China has led the world in economic growth since the turn of the 21st century, and its stock market performance is notable (Shen et al. 2017). Therefore, the subject of this study is the performance of stock prices in the Chinese stock market.

This research focuses on the five stocks with the biggest market capitalization in their respective sectors among the most popular sectors of the Chinese stock market. They are Kweichou Maotai (600519) in the wine category, Hengrui Pharmaceutical (600276) in the pharmaceutical category, Zhongxing Telecom Equipment (000063) in the technology category, Shanghai Airport (600009) in the logistics category, and Industrial and Commercial Bank of China (601398) in the banking category. Stock trading data are

obtained from The Wind Database, including daily high and low stock prices, opening and closing prices, trading volume, and turnover, etc. The data are captured from January 1, 2020 to August 29, 2022.

Internet data, such as search engine data and text data, are frequently employed in research. In this study, two search indexes—Baidu Index and Google Trends—are employed. Meanwhile, four types of textual data are collected for this research: news data (from Tencent News), practitioner commentary data (from the stock bar forum of Oriental Fortune), public opinion data (from SinaPublic Opinion Communication), and pandemic data (from the China National Health Commission). News data, commentary data and pandemic data are collected from January 1, 2020 to August 29, 2022; opinion text data are collected from March 1, 2021 to August 29, 2022. Public opinion data provide the most extensive information; however, because of their accessibility, they cannot be maintained as long as other data. As a result, we further incorporate these data into the model after ordinary least squares manipulation. To present a fuller picture of

**Table 4** The description of the data

|  | Data | Description | Time range | Data source |
|---|---|---|---|---|
| Dependent variable | Stock price history data | Closing prices of the five stocks (¥) | January 1, 2020 to August 29, 2022 | https://www.wind.com.cn/ |
| Characteristics associated with the stock price | Opening price | Opening price of the five stocks (¥) | January 1, 2020 to August 29, 2022 | https://www.wind.com.cn/ |
|  | Highest Price | Daily highest price of the five stocks (¥) | January 1, 2020 to August 29, 2022 | https://www.wind.com.cn/ |
|  | Lowest Price | Daily lowest price of the five stocks (¥) | January 1, 2020 to August 29, 2022 | https://www.wind.com.cn/ |
|  | Turnover | Turnover of the five stocks (¥) | January 1, 2020 to August 29, 2022 | https://www.wind.com.cn/ |
|  | Transaction volume | Transaction volume of the five stocks (¥) | January 1, 2020 to August 29, 2022 | https://www.wind.com.cn/ |
| Search Engine Data | Baidu Index | The most popular search engine in China. Selected keywords: The names and ticker symbols of the five stocks | January 1, 2020 to August 29, 2022 | https://index.baidu.com/ |
|  | Google Trends | The most popular search engine worldwide. Selected keywords: The names and ticker symbols of the five stocks | January 1, 2020 to August 29, 2022 | https://trends.google.com/ |
| Text Data | News data | Stock-related news texts | January 1, 2020 to August 29, 2022 | https://news.qq.com/ |
|  | Practitioner commentary data | Text of comments from practitioners and investors | January 1, 2020 to August 29, 2022 | http://guba.eastmoney.com/ |
|  | Pandemic data | Official news reports on the outbreak | January 1, 2020 to August 29, 2022 | http://www.nhc.gov.cn/ |
|  | Public opinion data | Content of discussion among the public | March 1, 2021 to August 29, 2022 | https://yqt.midu.com/ |

the overall trend of public opinion in society, we also collect data from Baidu index and Google Trends to use as predictive features.

Table 4 provides a description of the data. The experimental process is discussed in the following subsections after the data collection is complete. We use the stock price of Hengrui Pharmaceutical as an example, because the development of pharmaceuticals was a core issue during the pandemic, and the fluctuation of the company's stock price reflected this context. February 23, 2022 was adopted as the boundary to divide the data into the training set and a testing set; that is, the training set used data from January 1, 2020 to February 22, 2022 for model training and the testing set used data from February 23, 2022 to August 29, 2022 for empirical testing. The experimental results of other groups of data are shown in the Appendix (i.e., Table 14 and Table 15).

### Data processing

The structured data in this research, including stock price history data, stock trading data, and search engine data, are normalized to eliminate the dimension (unit) and accelerate gradient descent, as indicated in Eq. (5).

$$X_{\mathrm{minmax}} = \frac{X - X_{\mathrm{min}}}{X_{\mathrm{max}} - X_{\mathrm{min}}} \tag{5}$$

where $X_{\mathrm{max}}$ represents the maximum value in the sequence and $X_{\mathrm{min}}$ represents the minimum value in the sequence. $X$ represents the data before normalization, and $X_{\mathrm{minmax}}$ represents the data after normalization.

For text data, this study uses natural language processing techniques for analysis. Specifically, Valence Aware Dictionary and sEntiment Reasoner (VADER) is used to determine the sentiment score of each day's text. It is worth noting that VADER is a lexicon and rule-based sentiment analysis tool specifically attuned to sentiments expressed in social media (Audrino et al. 2020; Shahi et al. 2020). As depicted in Fig. 4, the building of the lexicon includes well-established sentiment word-banks (e.g., Linguistic Inquiry and Word Count, LIWC; and the General Inquirer, GI), emoticons, sentiment-related acronyms and initialisms, commonly used slang with sentiment value, and the use of independent reviewers (Hutto and Gilbert 2014). If one day has many texts, we average them, and the sentiment score of the day is the average of the sentiment ratings of all texts. Additionally, to prevent the influence of extreme values, the daily sentiment interval is determined by the 25% to 75% percentile of the sentiment score of all comments each day. The results of the sentiment scores of these texts are shown in Fig. 5. Meanwhile, we provide word clouds for both types of text data (showing the high frequency words), as shown in Fig. 6, where the language is Chinese, because our context is the Chinese stock market. As shown in Table 5, we mark the terms that occur most frequently.

Figure 5 demonstrates that practitioner comment data has more sentiment variation and that the public influences individual sentiments. News media, however, have more severe sentiments (positive emotions tend to be 1 and negative emotions tend to be 0) and are predominantly positive, indicating that news is always presented with a bias. In addition, Table 6 shows that the news focuses on broad issues such as scientific research, experiments, and research duration. Practitioners' worries tend to center on personal interest views, such as stock price fluctuations and direct stock appraisal. As a result,
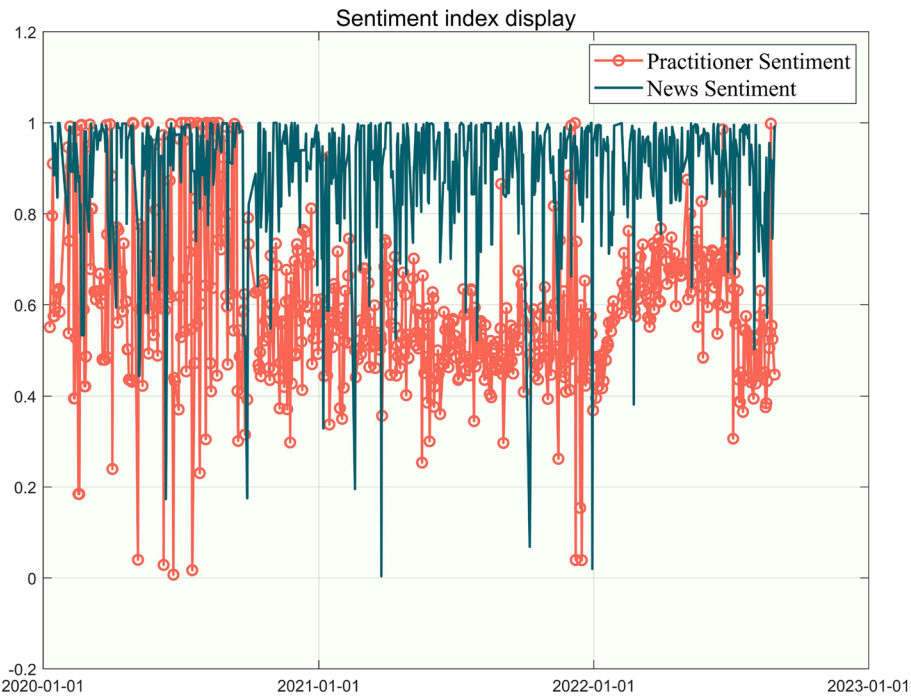
**Fig. 5** The sentiment

while performing a sentiment sequence research, there is variation in the sentiments of different groups, which supports the use of text material from various viewpoints in this study.

**Evaluation criteria**

This study employs two error assessment criteria based on data values, mean absolute error (MAE) and root mean square error (RMSE), along with one error evaluation criterion based on error percentages, mean absolute percentage error (MAPE), to properly assess forecasting performance (Li et al. 2023; Zhao et al. 2022). They can be computed numerically as follows:

$$MAE = \frac{1}{N} \sum_{t=1}^{N} (\hat{y}_t - y_t)^2 \tag{6}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (\hat{y}_t - y_t)^2} \tag{7}$$

$$MAPE = \frac{1}{N} \sum_{t=1}^{N} \left| \frac{y_t - \hat{y}_t}{y_t} \right| \tag{8}$$

where $\hat{y}_t$ and $y_t$ are the forecast stock price and the actual stock price at time $t$, respectively, and $N$ refers to the number of samples in the test set.
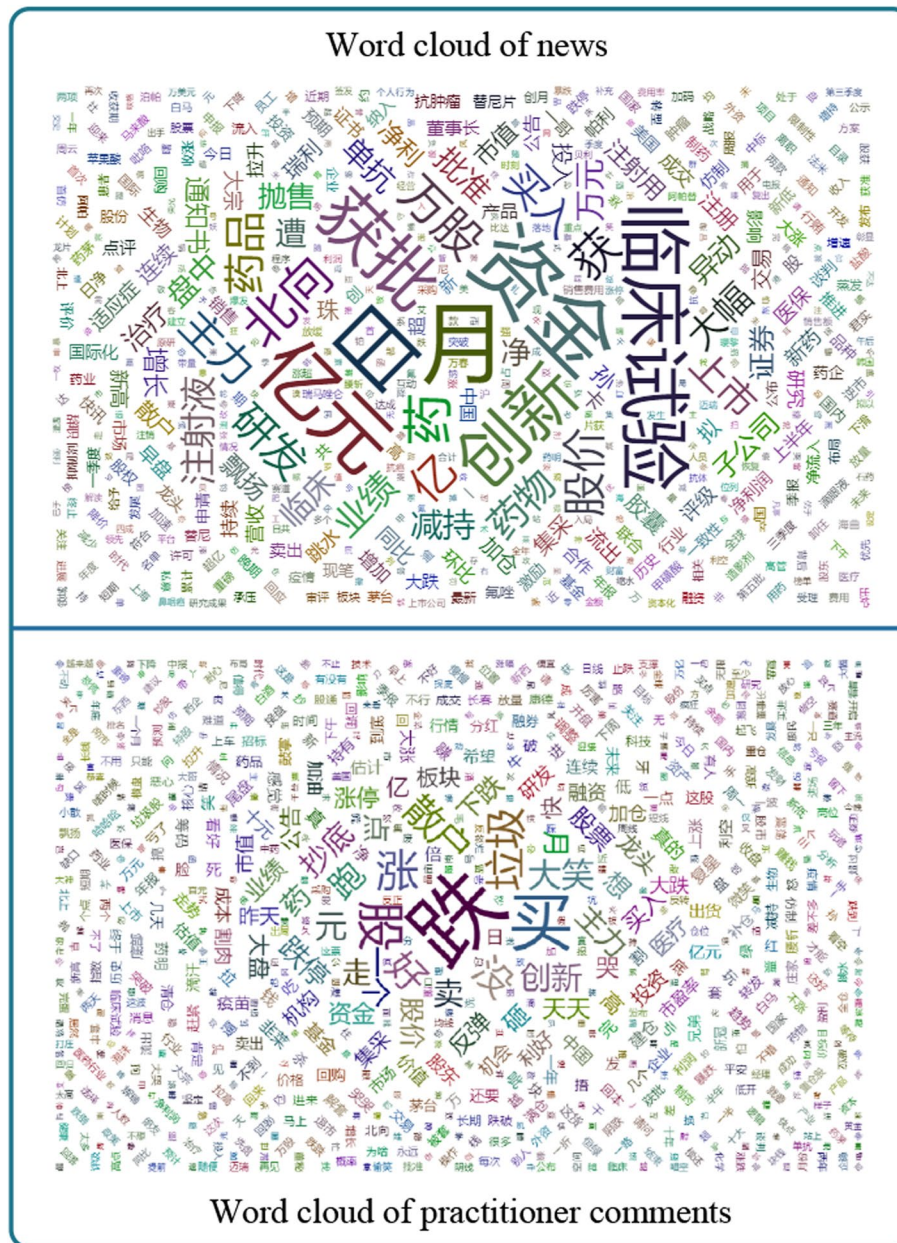
**Fig. 6** The word clouds

Then, this study employs the Wilcoxon signed-rank (WSR) test, the superior predictive ability (SPA) test, and the Kolmogorov-Smirnov predictive accuracy (KSPA) test to assess the model's predictive power from a statistical perspective in addition to the error evaluation index (Zhao and Itti 2018; Hansen 2005). The null hypothesis for the WSR test is the following loss differential series:

$$l(t) = f(e_A(t)) - f(e_B(t)) \qquad (9)$$

where $f(*)$ represents the loss function (MSE in this study), $e_A(t)$ and $e_B(t)$ represent the forecasting error series of models A and B, respectively.

**Table 5** High frequency words list

| News sentiment | | | Practitioner sentiment | | |
|---|---|---|---|---|---|
| Chinese words | English translation | Words frequency | Chinese words | English translation | Words frequency |
| 月 | Month | 419 | 跌 | Fall | 1669 |
| 日 | Day | 392 | 买 | Buy | 1186 |
| 资金 | Funding | 391 | 股 | Stock | 896 |
| 亿元 | 0.1 billion | 351 | 涨 | Up | 780 |
| 临床试验 | Clinical trials | 342 | 垃圾 | Garbage | 740 |
| 创新 | Innovation | 321 | 好 | Excellent | 654 |
| 获批 | Receive approval | 313 | 散户 | Individual stock investors | 646 |
| 药 | Pharmaceuticals | 223 | 没 | No more | 620 |

SPA is a statistical test method for evaluating point forecasting's greater predictive power. Specifically, it can determine whether the target model's accuracy performance is better than that of other benchmark models. The conclusion that the benchmark is the best is this null hypothesis:

$$H_0 : \max_i E[L_i] \geq E[L_{bm}] \tag{10}$$

where $L_i$ represents MSE of the $i$th model, and $L_{bm}$ is the MSE of the benchmark model. The KSPA is a statistical test that aids in determining the predictive accuracy of the two models. Its advantage is that it determines not only the predictive distribution of the two models, but also whether the models have minimal random error. The test is not affected by any autocorrelation in the forecasting errors (Hassani and Silva 2015; Fan et al. 2022a, b).

**Forecasting results**

In the first place, the performance of seven single forecasting models (i.e., SARIMA, ES, SVR, ELM, BPNN, LSTM, and TCN) is comprehensively examined in this study to illustrate the predictive performance of different models and the justification for using TCN in this study. Specifically, the performance of single forecasting models, including SARIMA, ES, SVR, ELM, BPNN, LSTM and TCN, is compared using the Hengrui Pharmaceutical stock price dataset, with MAE, RMSE, and MAPE used as error assessment criteria. Table 6 details the forecasting accuracy and ranking of various forecasting methodologies.

When evaluating the performance of seven single forecasting models, the most notable observations that emerged from the data comparison are as follows. First, the model suggested in this study (i.e., TCN) outperforms other models in forecasting the stock price of Hengrui Pharmaceutical. Second, based on the overall forecasting findings, the performance of LSTM, a neural network model, is comparable to that of TCN; SVR, ELM and BPNN perform somewhat worse. Third, the SARIMA model and ES differs in performance from the five methods discussed above, yielding somewhat unsatisfactory forecasting results using the research dataset analyzed in this study. This coincides with the findings of some previous studies on the accuracy of machine learning and statistical methods.

**Table 6** Forecasting error and ranking of single models

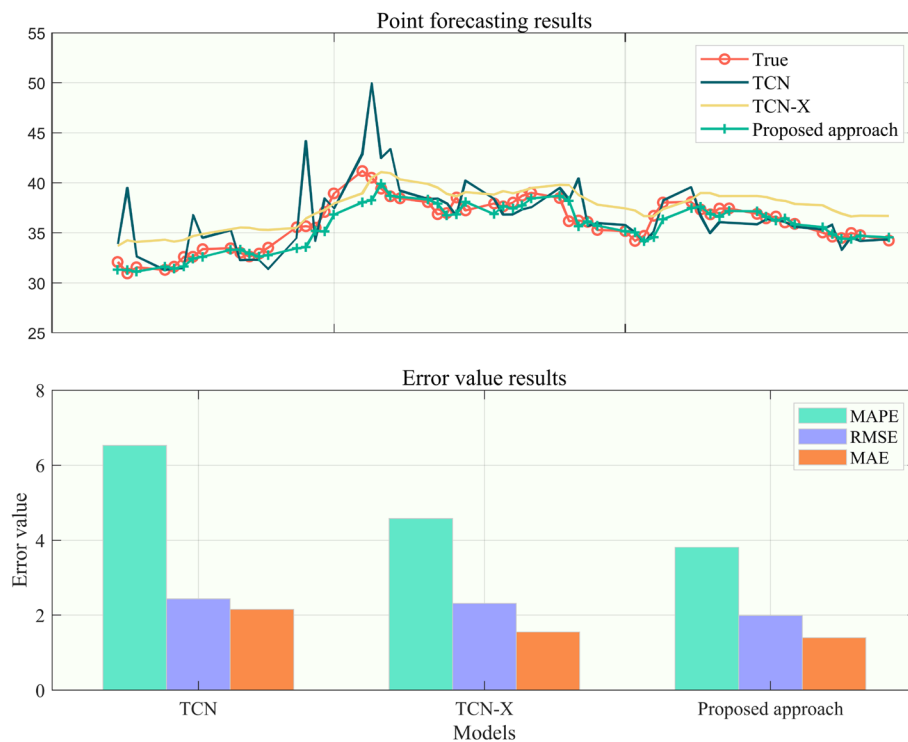| Indicators | SARIMA | ES | SVR | ELM | BPNN | LSTM | TCN |
|---|---|---|---|---|---|---|---|
| MAE | 4.17 | 4.09 | 3.44 | 3.13 | 3.07 | 2.43 | **2.15** |
| RMSE | 4.34 | 4.26 | 3.63 | 3.35 | 3.24 | 2.71 | **2.44** |
| MAPE (%) | 12.3 | 12.07 | 10.15 | 9.24 | 9.08 | 7.34 | **6.53** |
| Rank | 7 | 6 | 5 | 4 | 3 | 2 | **1** |

Bold represents the smallest error or best performance

**Table 7** Forecasting error and ranking of models with sentiment index

| Indicators | SARIMA-X | ES-X | SVR-X | ELM-X | BPNN-X | LSTM-X | TCN-X | Proposed approach |
|---|---|---|---|---|---|---|---|---|
| MAE | 2.2 | 2.22 | 2.15 | 2.37 | 2.01 | 1.8 | 1.55 | **1.4** |
| RMSE | 2.6 | 2.46 | 2.39 | 2.97 | 2.27 | 2.11 | 2.31 | **2** |
| MAPE (%) | 6.56 | 6.54 | 6.35 | 6.97 | 5.92 | 5.42 | 4.58 | **3.82** |
| Rank | 7 | 6 | 5 | 8 | 4 | 3 | 2 | **1** |

Bold represents the smallest error or best performance

Based on the single-model forecasting, we match each of the seven models with the extracted sentiment indexes for further assessment and give them the corresponding names SARIMA-X, ES-X, SVR-X, ELM-X, BPNN-X, LSTM-X, and TCN-X. These models incorporate stock related characteristics, search indexes, news sentiment, and practitioner sentiment to compare and assess the proposed approach's forecast accuracy and robustness (with the inclusion of short-length social opinion information). Table 7



**Fig. 7** Point forecasting results

displays the assessment findings for each measure, whereas Fig. 7 displays the anticipated fitted curves and comparison graphs.

Regarding these models' forecasting errors (see Table 7), this table is informative in several ways. First, the models that used the sentiment index all had lower error values compared to the single model, and MAE values for the models decreased by 47.24%, 45.81%, 37.50%, 24.28%, 34.80%, 35.00%, and 27.91%, respectively. Second, all machine learning models except ELM-X beat SARIMA-X, indicating that SARIMA is applicable and is not inferior to artificial intelligence approaches in certain forecasting scenarios. Third, the proposed approach has the lowest error value and ranks first within the experimental control group, showing that it has the highest predictive accuracy. These outcomes demonstrate the market's sensitivity to sentiment. The incorporation of sentiment from multiple groups has the ability to properly predict stock values, and a more inclusive group sentiment produces greater outcomes. Furthermore, it is noteworthy to highlight that neural network models exhibit a higher number of parameters, as evidenced by the data presented in Tables 2 and 3. Consequently, this leads to longer runtimes, necessitating a heightened focus on optimizing the computational efficiency of the model, particularly when dealing with larger datasets.

The Wilcoxon signed-rank test and the SPA test are employed in this study to statistically assess the differences in the models' forecasting abilities and to confirm the validity of the aforementioned results. The results with *p*-values are given in Table 8. In each table, the *p*-value in panel A denotes the significant difference between the two comparison models, and the *p*-value in panel B denotes the existence of a significant superiority link between the two comparison models. For example, in Panel A, the *p*-value in row 2, column 3 is 0.000, meaning the test rejects at a 99% confidence level the null hypothesis (i.e., there is a significant difference between SARIMA-X and SVR-X models).

**Table 8** Results of the Wilcoxon signed rank test and the SPA test

| | SARIMA-X | ES-X | SVR-X | ELM-X | BPNN-X | LSTM-X | TCN-X | Proposed approach |
|---|---|---|---|---|---|---|---|---|
| *Panel A: The Wilcoxon signed rank test* | | | | | | | | |
| SARIMA-X | | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| ES-X | | | 0.0003 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| SVR-X | | | | 0.0004 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| ELM-X | | | | | 0.0000 | 0.0015 | 0.0000 | 0.0000 |
| BPNN-X | | | | | | 0.0000 | 0.0000 | 0.0000 |
| LSTM-X | | | | | | | 0.0007 | 0.0000 |
| TCN-X | | | | | | | | 0.0000 |
| *Panel B: The SPA test* | | | | | | | | |
| SARIMA-X | | 0.0000 | 0.0000 | 1.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| ES-X | 1.0000 | | 0.0000 | 1.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| SVR-X | 1.0000 | 1.0000 | | 1.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| ELM-X | 0.0000 | 0.0000 | 0.0000 | | 0.0000 | 0.0040 | 0.0000 | 0.0000 |
| BPNN-X | 1.0000 | 1.0000 | 1.0000 | 1.0000 | | 0.0000 | 0.0000 | 0.0000 |
| LSTM-X | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | | 0.0000 | 0.0000 |
| TCN-X | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | | 0.0000 |
| Proposed approach | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | |

Additionally, in Panel B, the *p-* value in the second column of row 3 is 1.000, indicating the test rejects the null hypothesis at a 99% confidence level (i.e., in forecasting performance, SVR-X beats the SARIMA-X model).

Table 8 statistically compares each models' forecasting performance for stock price data. Some striking observations are as follows. First, when the proposed approach, is considered the benchmark model for Panels A and B, the *p*-values for both tests are 0.000 and 1.000 in all cases, respectively, indicating that the proposed approach significantly outperforms all other comparative models at the 99.9% confidence level. Second, the machine learning models generally outperform SARIMA-X in addition to ELM-X, with TCN performing the best, which validates the correctness of our choice of this model. Third, the ELM-X does not perform as well as SARIMA-X, indicating that SARIMA can still achieve satisfactory results in some scenarios. Fourth, among machine learning models, LSTM with multilayer neural network structure outperforms classical SVR and ELM, also indicating its capability in time-series analysis.

The results of the KSPA test also support the conclusion that the proposed approach has the best predictive performance (all corresponding p-values are less than 0.01). Figure 9 shows the KSPA error distribution and the empirical cumulative distribution function (c.d.f.) in this study. The proposed model better describes the random deviation; in other words, it has better prediction performance and higher prediction accuracy.

After collecting the closing values from the point forecasting models, we utilize the lowest daily stock price value as the left endpoint of the interval and the maximum value as the right endpoint of the interval, and combine the 25% and 75% quantile intervals of the experts sentiment and news sentiment as inputs to the Gi-MLP model. Following the interval forecasting model, the expected values of the daily high and low prices of the stocks are then produced. The Gi-MLP model's



**Fig. 8** Interval forecasting results

training and test sets are divided in a manner consistent with the point-value fore-casting model, and the rolling forecasting method is used to make predictions one step ahead of time in the forecasting window. The forecasting results are shown in Fig. 8 alongside the actual results.

The following points can be observed from Fig. 8. First, the daily high and low stock prices may be maintained in the forecast outputs of the Gi-MLP model. The anticipated lowest price might be higher than the forecasting maximum price if point-value regression models (e.g., SVR) are used to forecast the high and low stock values independently, but the Gi-MLP model overcomes this issue because it makes use of the interval's entire information. Second, a broader range of anticipated high and low prices is indicated when the actual volatility is strong, which depicts the general trend of the stock and indicates the change in volatility of the stock. Third,



a) Distribution of errors



(b) Empirical cumulative distribution functions (c.d.f.)

**Fig. 9** Distribution of errors and empirical cumulative distribution function of error (c.d.f.)

there is little deviation between the Gi-MLP model's forecast high and low stock prices and the actual predicted values, with an average MAPE of less than 6%.

### Trading simulations

In previous experiments, we examine the forecasting framework using error evaluation metrics, appropriate statistical tests, and interval estimates; nevertheless, the precision of a forecast is not synonymous with the real return. In the practice of stock trading, the objective of investors is to develop a profitable strategy. Therefore, we use a long-short trading strategy to show the profitability of the proposed approach. The practical strategy is to go "short" on the position when the forecast return is below zero, and to go "long" on the position when the forecast return is above zero. The "short" and "long" stock price position is defined as selling and buying the stock at each respective current price. When the forecast return is zero, we maintained our position.

We compute the relevant return estimates based on the projected closing prices from June 8, 2022 to August 29, 2022, with the expected return on day $t$ determined as follows:

$$\hat{R}_t = \frac{\hat{C}_t - \hat{C}_{t-1}}{\hat{C}_{t-1}}, \tag{11}$$

where $\hat{C}_t$ represents the closing price on day $t$ and uses $R_t$ to denote the true rate of return on day $t$. Based on the estimated rate of return, the following investment strategy can be developed.

Scheme 1.1: Invest based on the forecast return for the following day. If the next day's return is positive, long the position at the next day's opening, and if the next day's return is negative, short the position at the next day's opening, achieving a positive return as long as $\hat{R}_{t+1}$ is the same as $R_{t+1}$; however, if the two variables have opposite signs, the investment will be damaged.

Scheme 1.2: Transaction fee is 0.5‰ in each transaction based on Scheme 1.1.

Scheme 2.1: To mitigate the loss from a wrong investment, we use the expected high and low stock values as insurance. When the expected closing price is between the high and low prices of the stock predicted using Gi-MLP, that is, when $\hat{Y}_{t+1}^L \leq \hat{C}_{t+1} \leq \hat{Y}_{t+1}^R$ is satisfied, the transaction stated in Scheme 1.1 is executed the next day. Otherwise, the stock is believed to be more volatile the next day, and the return forecasting is erroneous, therefore no deal is executed. Figure 10 gives an illustration of such a judgment (the marked points will not be involved in the transaction).

Scheme 2.2: Transaction fee is 0.5‰ in each transaction based on Scheme 2.1.

According to the aforementioned investment scenario, we simulated the real stock trading from June 9, 2022 to August 29, 2022, a total of 58 days, and the outcomes of Hengrui Pharmaceutical's simulation are displayed in Table 9. As the table shows, the trading strategy proposed in this study reduces the number of transactions overall owing to its double insurance structure (a greater percentage of transactions that may deplete returns). Therefore, this strategy is particularly suitable for periods of high volatility, such as large stock rises and falls, which can guide investors to make fewer wrong transactions and helps investors make rational judgments, enhance risk management, and thus, obtain high returns.

**Fig. 10** Strategy judgment curve

**Table 9** Simulated transaction results

| Scheme 1.1 | | | Scheme2.1 | | | Comparison | | |
|---|---|---|---|---|---|---|---|---|
| Add transaction cost | | | | | | | | |
| #(T) | #(+) | #(−) | #(T) | #(+) | #(−) | #(T) | #(+) | #(−) |
| 58 | 34 | 24 | 39 | 22 | 17 | 19↓ | 12↓ | 7↓ |
| CR:7.22% | | | CR:26.52% | | | CR:19.30%↑ | | |
| **Scheme 1.2** | | | **Scheme2.2** | | | **Comparison** | | |
| Without transaction cost | | | | | | | | |
| #(T) | #(+) | #(−) | #(T) | #(+) | #(−) | #(T) | #(+) | #(−) |
| 58 | 34 | 24 | 39 | 22 | 17 | 19↓ | 12↓ | 7↓ |
| CR:4.32% | | | CR:24.57% | | | CR:20.25%↑ | | |

1. #(T) indicates the number of transactions among all the simulation trading times

2. #(+) and #(−) represent the number of transactions that bring positive returns and the number of transactions that bring negative returns, respectively

3. CR is the cumulative rate of return

4. ↓ and ↑ respectively represent the increase or decrease of the latter scheme compared with the former one

Meanwhile, we give the simulated trading results for different models, as shown in Table 10. The results show that, our proposed approach generates the best rate of return after considering interval restrictions. The statistical precision and profitability of our strategy can be explained from two perspectives. First, a variety of variables impact the volatility of stock prices. Our suggested strategy considers sentiments from numerous platforms and extracts information using an advanced TCN to drive price forecasting. Second, interval estimation can provide an insurance cover for trading. By not depending solely on point forecast results, the frequency of trading errors is decreased, while the confidence level of the interval is examined to prevent unwanted losses.

**Table 10** Simulated transaction results for different models

|  | SARIMA-X | ES-X | SVR-X | ELM-X | BPNN-X | LSTM-X | TCN-X | Propose approach |
|---|---|---|---|---|---|---|---|---|
| Return (%) | 6.21 | − 1.41 | 1.88 | − 2.09 | 8.58 | 2.59 | 8.08 | 24.57 |
| Mean of return (%) | 0.1041 | − 0.0253 | 0.0322 | -0.0365 | 0.1421 | 0.0442 | 0.1346 | 0.3803 |
| Standard deviation of return (%) | 2.43 | 2.42 | 2.42 | 2.41 | 2.39 | 2.40 | 2.40 | 2.13 |
| Sharpe Ratio | 0.5277 | – | 0.4418 | – | 0.3711 | 0.4613 | 0.5753 | 1.2030 |
| Maximum of drawback (%) | 11.15 | 12.65 | 15.40 | 12.63 | 9.58 | 14.73 | 10.69 | 6.27 |

## Robustness analysis

To test the robustness of the proposed method, two aspects are considered for further experiments in this subsection. The first is the prediction of the stock index. Stock index forecasting is also of great interest to financial managers and researchers as a hot topic. Second, trading simulations for longer time periods. Since this study prioritizes the medium stock price fluctuations and simulated trading in the stock market during the pandemic period, the simulation period involved is relatively short, so a set of trading simulations for a longer time period will be conducted. Specifically, we will supplement the experiments of SSE index forecasting with Hengrui Pharmaceutica trading simulation, involving dates from 2017-2022, for a total of five years. We take the first 80% as the training set and the last 20% as the test set.

The results regarding SSE predictions are shown in Table 11. Conclusions similar to those of previous experiments can be drawn: firstly, TCN has an absolute performance advantage in the single prediction model comparison. Second, machine learning models overall due to time-series models with econometric models, such as SARIMA and ES. Third, when sentiment-rich textual information is added, the overall prediction performance is improved in both cases. Fourth, the proposed model has the lowest prediction error value, that is, the best prediction performance. The forecast results and simulated trading results regarding HR stock are shown in Tables 12 and 13. It can be clearly observed that the proposed model still achieves satisfactory results under a year-long back-testing period. The performance of the corresponding benchmark model is roughly similar to previous findings.

**Table 11** Forecasting error of SSE

|  | SARIMA | ES | SVR | ELM | BPNN | LSTM | TCN |
|---|---|---|---|---|---|---|---|
| MAE | 4.4051 | 4.2063 | 3.5062 | 3.3520 | 2.4389 | 2.8545 | 2.2961 |
| RMSE | 5.7473 | 5.4348 | 4.3640 | 3.5306 | 2.7734 | 3.2071 | 2.8030 |
| MAPE (%) | 12.5363 | 11.9055 | 9.4956 | 8.2887 | 6.3553 | 6.4840 | 6.1561 |
| **SARIMA-X** | **EX-X** | **SVR-X** | **ELM-X** | **BPNN-X** | **LSTM-X** | **TCN-X** | **Proposed approach** |
| 2.9649 | 2.4770 | 1.9997 | 2.1578 | 2.0351 | 1.8782 | 1.9747 | 1.4862 |
| 3.1638 | 2.6957 | 2.3961 | 2.6034 | 2.2772 | 2.2573 | 2.2243 | 1.7914 |
| 7.2819 | 6.0349 | 5.0831 | 5.6970 | 4.9135 | 4.7919 | 4.7718 | 3.4967 |

**Table 12** Forecasting errors of hengrui pharmaceutical

|  | SARIMA | ES | SVR | ELM | BPNN | LSTM | TCN |
|---|---|---|---|---|---|---|---|
| MAE | 244.2574 | 275.7138 | 235.2861 | 136.1548 | 147.3738 | 134.6134 | 100.9706 |
| RMSE | 286.4579 | 322.6071 | 276.2383 | 160.7180 | 174.2493 | 158.7485 | 118.5251 |
| MAPE | 6.9976 | 7.8992 | 6.7404 | 3.9053 | 4.2247 | 3.8616 | 2.9048 |
| **SARIMA-X** | **EX-X** | **SVR-X** | **ELM-X** | **BPNN-X** | **LSTM-X** | **TCN-X** | **Proposed approach** |
| 71.0749 | 73.6239 | 61.0251 | 68.2109 | 61.0816 | 51.4318 | 41.8059 | 40.1873 |
| 82.6532 | 85.8654 | 71.6140 | 79.6781 | 71.6719 | 60.8306 | 50.2892 | 47.6682 |
| 2.0601 | 2.1313 | 1.7743 | 1.9780 | 1.7759 | 1.5017 | 1.2293 | 1.1824 |

**Table 13** Simulated transaction results for Robustness analysis

|  | SARIMA-X | ES-X | SVR-X | ELM-X | BPNN-X | LSTM-X | TCN-X | Propose approach |
|---|---|---|---|---|---|---|---|---|
| Return (%) | 17.78 | 21.58 | 21.47 | 25.26 | 22.24 | 18.94 | 29.49 | 37.10 |
| Mean of return (%) | 0.2831 | 0.3382 | 0.3227 | 0.3892 | 0.3471 | 0.2999 | 0.4470 | 0.5465 |
| Standard deviation of return (%) | 2.32 | 2.31 | 1.38 | 1.36 | 2.31 | 1.39 | 2.31 | 1.91 |
| Sharpe ratio | 0.8667 | 0.9992 | 1.6073 | 1.9310 | 1.0227 | 1.5095 | 1.3130 | 2.0330 |
| Maximum of drawback (%) | 23.94 | 30.24 | 25.78 | 24.90 | 23.95 | 28.07 | 24.50 | 20.13 |

When considered together, the results of the two sets of experiments in this section (considering index forecasting with extended back-testing period) support the conclusion drawn from previous experiments that the study of stock prices, and thus financial management decisions, can be effectively performed using the forecasting methodology and simulated trading strategy, combining multi-platform textual information proposed in this study.

## Conclusion

Investor sentiment is closely related to stock price volatility. This study proposes a deep learning approach based on sentiment indexes, which includes a framework for both prediction and trading strategies. The method addresses the key problem of price prediction in the stock market by analyzing the sentiment expressions of different groups, guiding the interpretation of stock price fluctuations. Specifically, this paper selects news texts, practitioner comments, public opinion texts, and pandemic reports to depict the sentiment orientations of different groups from macro media and micro individuals' perspectives. This sentiment information is then integrated into the prediction model to enhance the forecasting accuracy. Using the deep learning model to predict stock prices, this study proposes an innovative trading strategy based on the prediction framework and interval estimation.

In the empirical analysis, this study conducts stock price forecasting and simulated trading experiments on five stocks, with Hengrui Pharmaceutica as a representative case. The forecasting results demonstrate that the proposed approach outperforms other benchmark models, yielding more robust stock price predictions. Furthermore, after

conducting robustness tests, the proposed prediction method continues to perform well. This finding indicates that incorporating comprehensive integration and sentiment indexes from multiple groups significantly reduces prediction errors and generates more robust prediction results. The simulated trading results based on the prediction demonstrate that the proposed trading strategy with interval constraints improves the return on investment in stocks, such as Hengrui Pharmaceutica, compared to the benchmark strategy, which relies solely on point forecasting. This strategy effectively guides investors to make prudent and accurate trades, achieving greater returns on investment at lower trading costs.

The literature about forecasting stock prices and developing trading and investment models is extensive, as mentioned in Literature review. In comparison to these studies, this paper introduces an innovative approach by utilizing interval forecasting results as constraints to reduce investment risk in point-based trading strategies. This idea has not been addressed in previous research literature. Through simulated trading, this study demonstrates that return on investment can be improved not only by enhancing forecasting accuracy but also by reducing the number of failed trades through the application of interval constraints. This interval-constrained trading strategy provides a feasible and convenient investment solution to achieve higher returns with lower transaction costs, offering new and instructive avenues for future research.

## Discussion and prospects

The proposed approach makes three significant contributions. First, it considers multiple sources of information and various data types. The research dataset includes variables related to the stock market as well as the opinions of practitioners, which are valuable for both forecasting and simulated trading. Second, the study utilizes a deep learning model called TCN for forecasting and achieves outstanding results. This highlights the effectiveness and reliability of deep learning models in forecasting stock prices. Third, the study incorporates interval forecasting and applies the results to simulated trading. This approach helps practitioners to mitigate losses substantially during market fluctuations and plays a vital role in guiding decision-making among managers and practitioners.

The practical significance of this study lies in providing an approach based on deep learning and sentiment analysis to help investors accurately predict stock price fluctuations and make effective investment decisions. By capturing investor sentiment and market expectations, this method can assist investors in making informed buying and selling decisions in the stock market, controlling investment risks, and maximizing return on investment. Furthermore, by combining the forecasting results with interval estimation, investors can better assess the balance between risk and return and formulate rational asset allocation strategies. This research is of great importance to individual investors, asset management companies, and other financial institutions, demonstrating the application potential of deep learning and sentiment analysis in investment decision-making, and providing new ideas and methods for related research and practice.

This research holds managerial and societal significance as it utilizes textual data to capture the impact of social events and investor sentiment on the stock market. By analyzing news, comments, social media, and pandemic-related data, this study provides insights into the dynamics of the stock market during significant events like the COVID-19 pandemic. This understanding contributes to maintaining market stability and supporting investment decision-making. By considering a range of integrated data sources, including sentiment analysis, this research enhances our ability to assess market sentiment and investor behavior. This can enable us to formulate effective risk management strategies, thereby contributing to the overall stability and efficiency of financial markets. Furthermore, by combining interval forecasting and trading strategies, this research provides investors with a framework to navigate market fluctuations and make informed investment decisions. Ultimately, the findings of this study have practical implications for investors, financial institutions, and policymakers pursuing sustainable and profitable investment strategies.

However, the study still has certain limitations. First, the unstructured data used in this study are all text, and in the era of big data, images and videos should also be taken into account in the future. Second, as this study has yielded these outstanding findings while focusing on emerging economic markets, the effectiveness of the proposed approach should be evaluated further in other markets. Third, the volatility of financial markets related with stocks (e.g., gold, crude oil) may also play a role in stock price volatility and these should be incorporated effectively into the features. We will investigate these promising issues soon.

## Appendix

See Tables 14, 15.

**Table 14** Performance of benchmark models with four Chinese stock data (without sentiment index)

|  |  | SARIMA | ES | SVR | ELM | BPNN | LSTM | TCN |
|---|---|---|---|---|---|---|---|---|
| Kweichou Maotai | MAE | 66.48 | 80.55 | 67.49 | 82.87 | 66.38 | 50.70 | **42.60** |
|  | RMSE | 81.70 | 97.51 | 79.90 | 96.82 | 81.61 | 62.61 | **56.80** |
|  | MAPE | 3.51 | 4.31 | 3.63 | 4.55 | 3.50 | 2.70 | **2.36** |
| Shanghai Airport | MAE | 1.99 | 1.81 | 1.81 | 1.60 | 1.38 | 1.38 | **1.30** |
|  | RMSE | 2.39 | 2.21 | 2.21 | 2.10 | 1.75 | 1.81 | **1.70** |
|  | MAPE | 3.97 | 3.60 | 3.60 | 3.12 | 2.72 | 2.67 | **2.52** |
| Commercial Bank of China | MAE | 1.19 | 1.13 | 1.14 | 0.96 | 0.95 | 0.95 | **0.97** |
|  | RMSE | 1.52 | 1.45 | 1.47 | 1.35 | 1.22 | 1.29 | **1.24** |
|  | MAPE | 3.05 | 2.87 | 2.91 | 2.37 | 2.43 | 2.35 | **2.35** |
| Zhongxing Telecom Equipment | MAE | 2.53 | 2.35 | 2.11 | 2.71 | 2.02 | 1.88 | **1.67** |
|  | RMSE | 2.61 | 2.44 | 2.20 | 2.87 | 2.13 | 2.01 | **2.09** |
|  | MAPE | 10.43 | 9.75 | 8.75 | 11.20 | 8.42 | 7.85 | **6.94** |

Bold represents the smallest error or best performance

**Table 15** Performance of the models with four Chinese stock data (with sentiment index)

|  |  | SARIMA-X | ES-X | SVR-X | ELM-X | BPNN-X | LSTM-X | TCN-X | Proposed approach |
|---|---|---|---|---|---|---|---|---|---|
| Kweichou Maotai | MAE | 49.84 | 50.77 | 53.70 | 46.72 | 46.30 | 46.34 | 40.85 | **37.68** |
|  | RMSE | 59.17 | 59.62 | 62.82 | 57.21 | 55.22 | 59.54 | 49.40 | **43.33** |
|  | MAPE | 2.73 | 2.71 | 2.96 | 2.57 | 2.48 | 2.51 | 2.23 | **1.94** |
| Shanghai Airport | MAE | 1.45 | 1.35 | 1.36 | 1.52 | 1.28 | 1.21 | 1.14 | **1.17** |
|  | RMSE | 1.83 | 1.71 | 1.72 | 1.92 | 1.64 | 1.61 | 1.52 | **1.51** |
|  | MAPE | 2.87 | 2.66 | 2.67 | 2.98 | 2.51 | 2.35 | 2.21 | **2.16** |
| Commercial Bank of China | MAE | 1.06 | 1.10 | 1.08 | 0.89 | 0.89 | 0.88 | 0.84 | **0.75** |
|  | RMSE | 1.35 | 1.43 | 1.43 | 1.24 | 1.17 | 1.22 | 1.17 | **0.95** |
|  | MAPE | 2.71 | 2.80 | 2.74 | 2.22 | 2.27 | 2.20 | 2.11 | **2.06** |
| Zhongxing Telecom Equipment | MAE | 1.90 | 1.82 | 1.45 | 2.06 | 1.41 | 1.86 | 1.37 | **0.98** |
|  | RMSE | 2.01 | 1.92 | 1.57 | 2.24 | 1.51 | 1.99 | 1.57 | **1.32** |
|  | MAPE | 7.87 | 7.56 | 6.00 | 8.53 | 5.87 | 7.76 | 5.66 | **1.97** |

Bold represents the smallest error or best performance

**Abbreviations**

| | |
|---|---|
| ARCH | Autoregressive conditional heteroskedasticity model |
| ARIMA | Autoregressive integrated moving average model |
| CR | Cumulative rate of return |
| ELM | Extreme learning machine |
| FCN | Fully-convolutional network |
| GARCH | Generalized autoregressive conditional heteroscedasticity model |
| Gi-MLP | Generalized random interval multilayer perceptron method |
| LSTM | Long short-term memory network |
| MAE | Mean absolute error |
| MAPE | Mean absolute percentage error |
| MLP | Multilayer perceptrons |
| MSE | Mean square error |
| MS-VAR | Markov-switching vector autoregressive model |
| NLP | Natural language processing |
| RMSE | Root mean square error |
| SARIMA | Seasonal autoregressive integrated moving average |
| SPA | Superior predictive ability test |
| SVR | Support vector regression |
| TCN | Temporal convolutional network |
| TVAR | Threshold vector autoregression model |
| VAR | Vector autoregressive model |
| WSR | Wilcoxon signed-rank test |

## Declarations

**Competing interests**
The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

Adam K, Marcet A, Nicolini JP, Adam K, Marcet A, Nicolini JP (2016) Stock market volatility and learning. J Finance 71(1):33–82. https://doi.org/10.1111/JOFI.2016.71.ISSUE-1

Alonso Robisco A, Carbó Martínez JM (2022) Measuring the model risk-adjusted performance of machine learning algorithms in credit default prediction. Financ Innov 8(1):1–35. https://doi.org/10.1186/S40854-022-00366-1/FIGURES/9

Audrino F, Sigrist F, Ballinari D (2020) The impact of sentiment and attention measures on stock market volatility. Int J Forecast 36(2):334–357. https://doi.org/10.1016/J.IJFORECAST.2019.05.010

Bai S, Kolter JZ, Koltun V. (2018) An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. arXiv preprint arXiv:1803.01271

Bengio Y, Ducharme R, Vincent P (2000) A neural probabilistic language model. Adv Neural Inf Process Syst 13

Bijl L, Kringhaug G, Molnár P, Sandvik E (2016) Google searches and stock returns. Int Rev Financ Anal 45:150–156

Bustos O, Pomares-Quimbaya A (2020) Stock market movement forecast: a systematic review. Expert Syst Appl. https://doi.org/10.1016/J.ESWA.2020.113464

Calomiris CW, Mamaysky H (2019) How news and its context drive risk and returns around the world. J Financ Econ 133(2):299–336. https://doi.org/10.1016/J.JFINECO.2018.11.009

Chen H, De P, Hu YJ, Hwang B-H (2014) Wisdom of crowds: the value of stock opinions transmitted through social media. Rev Financ Stud 27(5):1367–1403

Deng C, Huang Y, Hasan N, Bao Y (2022) Multi-step-ahead stock price index forecasting using long short-term memory model with multivariate empirical mode decomposition. Inf Sci 607:297–321

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. N Am J Econ Finance 64:101848. https://doi.org/10.1016/J.NAJEF.2022.101848

Fama E (1970) Efficient market hypothesis: a review of theory and empirical work. J Finance 25(2)

Fan GF, Zhang LZ, Yu M, Hong WC, Dong SQ (2022) Applications of random forest in multivariable response surface for short-term load forecasting. Int J Electr Power Energy Syst 139(January):108073. https://doi.org/10.1016/j.ijepes.2022.108073

Fan GF, Peng LL, Hong WC (2022) Short-term load forecasting based on empirical wavelet transform and random forest. Electr Eng 104(6):4433–4449. https://doi.org/10.1007/S00202-022-01628-Y/FIGURES/14

Ghosh P, Neufeld A, Sahoo JK (2022) Forecasting directional movements of stock prices for intraday trading using lstm and random forests. Financ Res Lett 46:102280. https://doi.org/10.1016/j.frl.2021.102280

Gu C, Kurov A (2020) Informational role of social media: evidence from Twitter sentiment. J Bank Finance. https://doi.org/10.1016/J.JBANKFIN.2020.105969

Gu W, Peng Y (2019) Forecasting the market return direction based on a time-varying probability density model. Technol Forecast Soc Chang 148(August):119726. https://doi.org/10.1016/j.techfore.2019.119726

Guo Y, Guo J, Sun B, Bai J, Chen Y (2022) A new decomposition ensemble model for stock price forecasting based on system clustering and particle swarm optimization. Appl Soft Comput. https://doi.org/10.1016/j.asoc.2022.109726

Gupta R, Pierdzioch C, Vivian AJ, Wohar ME (2019) The predictive value of inequality measures for stock returns: an analysis of long-span uk data using quantile random forests. Financ Res Lett 29:315–322. https://doi.org/10.1016/j.frl.2018.08.013

Han A, Hong Y, Wang S. (2012) Autoregressive conditional models for interval-valued time series data. In: The 3rd international conference on singular spectrum analysis and its applications, p 27

Hansen PR (2005) A test for superior predictive ability. J Bus Econ Stat 23(4):365–380. https://doi.org/10.1198/073500105000000063

Hassani H, Silva ES (2015) A Kolmogorov–Smirnov based test for comparing the predictive accuracy of two sets of forecasts. Econometrics 3(3):590–609. https://doi.org/10.3390/econometrics3030590

He P, Sun Y, Zhang Y, Li T (2020) COVID-19's impact on stock prices across different sectors-An event study based on the Chinese stock market. Emerg Mark Financ Trade 56(10):2198–2212. https://doi.org/10.1080/1540496X.2020.1785865

Hinton GE, Osindero S, Teh YW (2006) A fast learning algorithm for deep belief nets. Neural Comput 18(7):1527–1554. https://doi.org/10.1162/NECO.2006.18.7.1527

Hong H, Xu D, Wang GA, Fan W (2017) Understanding the determinants of online review helpfulness: a meta-analytic investigation. Decis Support Syst 102:1–11. https://doi.org/10.1016/J.DSS.2017.06.007

Hutto CJ, Gilbert E. (2014) VADER: A parsimonious rule-based model for sentiment analysis of social media text. In: Proceedings of the 8th international conference on weblogs and social media, ICWSM 2014, 216–225. https://doi.org/10.1609/icwsm.v8i1.14550

Jiang Y, Wang G-J, Ma C, Yang X (2021) Do credit conditions matter for the impact of oil price shocks on stock returns? Evidence from a structural threshold var model. Int Rev Econ Finance 72:1–15. https://doi.org/10.1016/j.iref.2020.10.019

Jin Z, Guo K, Sun Y, Lai L, Liao Z (2020) The industrial asymmetry of the stock price prediction with investor sentiment: based on the comparison of predictive effects with SVR. J Forecast 39(7):1166–1178. https://doi.org/10.1002/FOR.2681

Kang Y (2018) Regulatory institutions, natural resource endowment and location choice of emerging-market FDI: A dynamic panel data analysis. J Multinatl Financ Manag 45:1–14. https://doi.org/10.1016/J.MULFIN.2018.04.003

Kao L-J, Chiu C-C, Lu C-J, Chang C-H (2013) A hybrid approach by integrating wavelet-based feature extraction with mars and svr for stock index forecasting. Decis Support Syst 54(3):1228–1244. https://doi.org/10.1016/j.dss.2012.11.012

Khurana D, Koli A, Khatter K, Singh S (2023) Natural language processing: state of the art, current trends and challenges. Multimed Tools Appl 82(3):3713–3744. https://doi.org/10.1007/S11042-022-13428-4/FIGURES/3. arXiv:1708.05148

Kim N, Lučivjanská K, Molnár P, Villa R (2019) Google searches and stock market activity: evidence from Norway. Financ Res Lett 28:208–220

Kumbure MM, Lohrmann C, Luukka P, Porras J (2022) Machine learning techniques and data for stock market forecasting: a literature review. Expert Syst Appl. https://doi.org/10.1016/J.ESWA.2022.116659

Lecun Y, Bengio Y, Hinton G (2015) Deep learning. Nature 521(7553):436–444. https://doi.org/10.1038/nature14539

Li M, Zhang C, Wang S, Sun S (2022) Multi-scale analysis-driven tourism forecasting: insights from the peri-COVID-19. Curr Issues Tour. https://doi.org/10.1080/13683500.2022.2144151

Li M, Cheng Z, Lin W, Wei Y, Wang S (2023) What can be learned from the historical trend of crude oil prices? An ensemble approach for crude oil price forecasting. Energy Econ 123:106736. https://doi.org/10.1016/J.ENECO.2023.106736

Liang C, Tang L, Li Y, Wei Y (2020) Which sentiment index is more informative to forecast stock market volatility? Evidence from China. Int Rev Financ Anal. https://doi.org/10.1016/J.IRFA.2020.101552

Lin Y, Yan Y, Xu J, Liao Y, Ma F (2021) Forecasting stock index price using the ceemdan-lstm model. N Am J Econ Finance 57:101421. https://doi.org/10.1016/j.najef.2021.101421

Liu Y, Yang C, Huang K, Gui W (2020) Non-ferrous metals price forecasting based on variational mode decomposition and LSTM network. Knowl Based Syst. https://doi.org/10.1016/J.KNOSYS.2019.105006

Liu J, Zhang Z, Yan L, Wen F (2021) Forecasting the volatility of EUA futures with economic policy uncertainty using the GARCH-MIDAS model. Financ Innov. https://doi.org/10.1186/S40854-021-00292-8

Liu K, Zhou J, Dong D (2021) Improving stock price prediction using the long short-term memory model combined with online social networks. J Behav Exp Financ 30:100507

Maqbool J, Aggarwal P, Kaur R, Mittal A, Ganaie IA (2023) Stock prediction by integrating sentiment scores of financial news and MLP-regressor: a machine learning approach. Procedia Comput Sci 218:1067–1078. https://doi.org/10.1016/J.PROCS.2023.01.086

Maqsood H, Mehmood I, Maqsood M, Yasir M, Afzal S, Aadil F, Selim MM, Muhammad K (2020) A local and global event sentiment based efficient stock exchange forecasting using deep learning. Int J Inf Manage 50:432–451. https://doi.org/10.1016/J.IJINFOMGT.2019.07.011

Na H, Kim S (2021) Predicting stock prices based on informed traders' activities using deep neural networks. Econ Lett 204:109917. https://doi.org/10.1016/j.econlet.2021.109917

Narayan PK (2019) Can stale oil price news predict stock returns? Energy Econ 83:430–444. https://doi.org/10.1016/J.ENECO.2019.07.022

Nassirtoussi AK, Aghabozorgi S, Wah TY, Ngo DCL (2015) Text mining of news-headlines for forex market prediction: a multi-layer dimension reduction algorithm with semantics and sentiment. Expert Syst Appl 42(1):306–324

OECD: Coronavirus: the world economy at risk (2020). http://www.oecd.org/

Poernomo A, Kang DK (2018) Biased dropout and crossmap dropout: learning towards effective dropout regularization in convolutional neural network. Neural Netw 104:60–67. https://doi.org/10.1016/J.NEUNET.2018.03.016

Ribeiro Ramos FF (2003) Forecasts of market shares from var and bvar models: a comparison of their accuracy. Int J Forecast 19(1):95–110

Roque AMS, Maté C, Arroyo J, Sarabia Á (2007) Imlp: applying multi-layer perceptrons to interval-valued data. Neural Process Lett 25(2):157–169

Sadaei HJ, Enayatifar R, Lee MH, Mahmud M (2016) A hybrid model based on differential fuzzy logic relationships and imperialist competitive algorithm for stock market forecasting. Appl Soft Comput 40:132–149. https://doi.org/10.1016/j.asoc.2015.11.026

Salisu AA, Vo XV (2020) Predicting stock returns in the presence of COVID-19 pandemic: the role of health news. Int Rev Financ Anal. https://doi.org/10.1016/J.IRFA.2020.101546

Sanford A (2022) Does perception matter in asset pricing? Modeling volatility jumps using Twitter-based sentiment indices. J Behav Financ 23(3):262–280. https://doi.org/10.1080/15427560.2020.1866573

Shahi TB, Shrestha A, Neupane A, Guo W (2020) Stock price forecasting with deep learning: a comparative study. Mathematics. https://doi.org/10.3390/MATH8091441

Shen D, Zhang Y, Xiong X, Zhang W (2017) Baidu index and predictability of Chinese stock returns. Financ Innov. https://doi.org/10.1186/S40854-017-0053-1

Shomron G, Weiser U (2019) Spatial correlation and value prediction in convolutional neural networks. IEEE Computer Architec Lett 18(1):10–13. https://doi.org/10.1109/LCA.2018.2890236. arXiv:1807.10598

Sun S, Sun Y, Wang S, Wei Y (2018) Interval decomposition ensemble approach for crude oil price forecasting. Energy Econ 76:274–287

Tang L, Li J, Du H, Li L, Wu J, Wang S (2022) Big data in forecasting research: a literature review. Big Data Res. https://doi.org/10.1016/J.BDR.2021.100289

Teti E, Dallocchio M, Aniasi A (2019) The relationship between twitter and stock prices. Evidence from the US technology industry. Technol Forecast Soc Change 149:119747. https://doi.org/10.1016/j.techfore.2019.119747

Vuong PH, Dat TT, Mai TK, Uyen PH, Bao PT (2022) Stock-price forecasting based on XGBoost and LSTM. Comput Syst Sci Eng 40(1):237–246. https://doi.org/10.32604/CSSE.2022.017685

Wang X, Kang Y, Hyndman RJ, Li F (2022) Distributed arima models for ultra-long time series. Int J Forecast. https://doi.org/10.1016/j.ijforecast.2022.05.001

Wu Z, Pan S, Chen F, Long G, Zhang C, Yu PS (2021) A comprehensive survey on graph neural networks. IEEE Trans Neural Netw Learn Syst 32(1):4–24. https://doi.org/10.1109/TNNLS.2020.2978386. arXiv:1901.00596

Xing FZ, Cambria E, Welsch RE (2018) Natural language based financial forecasting: a survey. Artif Intell Rev 50(1):49–73. https://doi.org/10.1007/S10462-017-9588-9/FIGURES/5

Yun KK, Yoon SW, Won D (2022) Interpretable stock price forecasting model using genetic algorithm-machine learning regressions and best feature subset selection. Expert Syst Appl 213:118803. https://doi.org/10.1016/j.eswa.2022.118803

Zhang D, Lou S (2021) The application research of neural network and bp algorithm in stock price pattern classification and prediction. Futur Gener Comput Syst 115:872–879. https://doi.org/10.1016/j.future.2020.10.009

Zhang X, Shi J, Wang D, Fang B (2018) Exploiting investors social network for stock prediction in china's market. J Comput Sci 28:294–303

Zhao J, Itti L (2018) shapeDTW: shape dynamic time warping. Pattern Recognit 74:171–184. https://doi.org/10.1016/J.PATCOG.2017.09.020. arXiv:1606.01601

Zhao E, Sun S, Wang S (2022) New developments in wind energy forecasting with artificial intelligence and big data: a scientometric insight. Data Sci Manag 5(2):84–95. https://doi.org/10.1016/J.DSM.2022.05.002

Zhong X, Enke D (2019) Predicting the daily return direction of the stock market using hybrid machine learning algorithms. Financ Innov. https://doi.org/10.1186/S40854-019-0138-0

Zhong X, Enke D (2019) Predicting the daily return direction of the stock market using hybrid machine learning algorithms. Financ Innov 5(1):1–20. https://doi.org/10.1186/S40854-019-0138-0/TABLES/10

Zhu R, Liao W, Wang Y (2020) Short-term prediction for wind power based on temporal convolutional network. Energy Rep 6:424–429. https://doi.org/10.1016/j.egyr.2020.11.219

**Publisher's Note**

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.