

# Long Time Spans Stock Trend Prediction based on LSTM

Xiaozhao Tu , Lin Zhang

School of Shanghai Maritime University, Shanghai 201306, China

## Abstract

**The inherent volatility of stock markets around the world makes the task of prediction challenging. Therefore, forecasting and diffusion modeling underling many difficulties in the field of stock trend forecasting. To minimize the error rate of the prediction is to minimize the risk of the investment. In the current work, this problem was used as a directional prediction exercise with predicted direction as profit and loss. A model based on the comparison of classifying stock closing price up and down N days earlier was realized. LSTM(Long Short-Term Memory), an algorithm suitable for processing and predicting data with relatively long intervals and delays in time series. Based on LSTM, a new model was developed, and according to experiment results, this model was improved in each data set by comparing with other existing models. Compared with other traditional models, the innovation point of this model lies in simple extraction of characteristics and attention to a up and down prediction relative to the past in the future. Due to high prediction accuracy, this model is of great practical value.**

## Keywords

**long short-term memory, stock market predict, Long time spans, classification.**

## 1. Introduction

For a long time, it was believed that changes in the prices of stocks is not forecastable. Among them are the random swimming theory proposed by Malkiel and Fama in 1970 [1] and the effective market hypothesis proposed by Jensen in 1978 [2]. In addition, Carpenter et al[3] believes that the stock market is affected by many complicated factors, such as economic and business status, national policies, and investor's personal investment sentiment. Therefore, the high degree of uncertainty in the stock market makes the trend forecast of stock prices very difficult. However, how to estimate this uncertain risk has always been the research direction of the industry. Therefore, forecasting stock prices is a very challenging job, and reducing the risk of forecasting has a very substantial value for achieving high return returns.

With the globalization of the economy and the development of some information and technology tools, many people want to earn excess returns in the stock market relative to other stable investment methods. As we all know, the stock market price series is inherently dynamic, non-parametric, chaotic and noisy and inherent risks of investment. The stock market price movement is considered to be a stochastic volatility process, and this volatility is more pronounced in the short term. An in-depth understanding of recent share price movements would help to minimise the risk of such volatility. Investors are more willing to buy stocks that are expected to increase in the future, and sell stocks that are expected to fall in the future. Therefore, improving the accuracy of the stock market price trend forecast can maximize capital gains and minimize losses.

The forecasting methods of the stock market can be divided into two categories, one is fundamental analysis and the other is technical analysis. The fundamental approach is to analyze the factors that affect the intrinsic value of the company. While the technical method is to obtain the future trend of stock prices by analyzing the price of stock history. Machine learning and neural networks have been used to predict price changes in the stock market for

a long time. However, because many financial market factors are directly or indirectly intertwined, the trend of traditional machine learning to predict the trend of stock prices is a bit weak. In recent years, with the advent of some auxiliary predictive financial time series tools, neural networks have performed very well in predicting stock market trends [4,5,6,7,8]. Stock price changes are non-linear and unstable. The prediction of stock price is not only related to the information of current time window but also influenced by the information of earlier time. Unlike traditional machine learning methods, circular neural networks establish connections between hidden layers, which preserves the ability to remember previous data [9]. Recurrent neural network processes pre-correlated data through its unique memory function, which is very suitable for predicting time series [10,11]. Long short-term memory networks (lstm), as an improvement of recurrent neural networks, are often used in other fields such as natural language processing, time series prediction, etc [12,13]. Selectively filter information through the "gate" structure of lstm, so that this model can extract more useful information from the training historical data.

This paper introduces the use of lstm to predict the rise and fall of a company's stock price relative to the past, so as to obtain the accurate measurement of price changes. In this paper, we used stock trend forecasting as a classification problem. The class label of each sample is up or down by its relative to the previous N-day trading window. In our analysis, we conducted experiments on the comparison time windows of  $N = 3, 5, 10, 15, 30$  and 60 days respectively. Our goal is to design an intelligent model that uses deep learning techniques to learn from stock market data and predict the direction of stock price changes at the close of the daily stock market. In the study, we compared some predictions of accuracy and loss to some classical models, such as Random Forest [14], Support Vector Machine (svm), and Gradient Boosting Decision Tree (gbdt), which proved that we have certain advantages in prediction accuracy and loss. Therefore, when it comes to the ability to predict the closing trend of stocks, this paper is very helpful in the formulation of investment strategies for individuals or companies in the stock market.

## 2. Related Work

### 2.1. Stock Market Forecast

Stock market forecasts are known to be a very challenging job, as stock markets are full of instability and high levels of uncertainty. But some empirical studies have shown that stock market forecasts are possible [15,16,17]. Previous research has often used statistical or machine learning techniques to predict future price trends in the stock market. The traditional stock market forecasting technology is based on statistical methods to establish models through linear processes. For example, differential integrated moving average autoregressive (arima) model, the autoregressive conditional heteroscedasticity (arch) model, and generalized autoregressive conditional heteroskedasticity (garch) model are widely used in the prediction of financial time series data [18,19,20,21]. But this purely statistically based predictive model does not work well; and they have their own limitations because they require more historical data to satisfy statistical assumptions, such as data satisfying the normal distribution assumption.

Because the stock market is considered a nonlinear, non-parametric dynamic system, then a more flexible method of learning complex dimensions is essential. Machine learning and deep learning have great advantages in this respect, because these methods can extract the nonlinear relationship between data without knowing the prior knowledge of the input data. The following literature aims to focus on the use of machine learning and neural network prediction methods in various stock data sets to demonstrate the value of our views from the perspective of companies or investors. The use of predictive algorithms to determine future

trends in stock market prices (Widom, 1995; Hellstrom & Holmstrom, 1998) [24, 25] is a way to improve predictive power and facilitates the reassessment of efficient market hypotheses and diffusion models (Saha, Routh, & Goswami, 2014)[26]. Because computer algorithms can build more complex dynamic economic systems, this also exacerbates the controversy over whether stock prices are fully predictable. In addition, since individual behavior and reactions play an important role in determining stock turnover and price volatility, some studies have begun to involve lexical analysis of news articles (Kim, Jeong, & Ghani, 2014) [29].

In our research, we used very common daily stock data. These data sets include the daily opening price, closing price, highest price, lowest price, daily trading volume of individual stocks. Features are derived using historical inventory data and information from external variables mentioned earlier in this section. In the article by Dai and Zhang (2013) [29], the data used for analysis is the closing price of 3M's stock. The data set contains a total of 1471 data points for stock day trading data from September 1, 2008 to August 11, 2013. A variety of predictive models were used in their experiments, such as logistic regression models, secondary analysis discriminant models, and SVM models. These models predict stocks' movements over the next day based on a given data sample, and they also forecast trends for several consecutive days. The accuracy of the continuous day prediction model is 44.52% ~ 58.2%. In their article, they came to the conclusion that there is reason to believe that the US stock market is semi-strong, that is, neither fundamental analysis nor technical analysis can be used to achieve better returns. However, the model performed quite well in the long-term forecast. When the predicted trading window reached 44 days, the SVM model achieved the best accuracy of 79.3%. In 1998, Sadd compared three neural network models, delay, recursive and probabilistic neural networks, using a conjugate gradient training method and using multi-stream extended Kalman filter delay neural network and recurrent neural network to train stock prediction model. RNN shows better predictive performance than other models. Chen et al. (2005) used neural networks, TS fuzzy systems and hierarchical model systems to verify the effectiveness of hybrid models, and used particle swarm search [31] algorithm to optimize the parameters of each model. Recently, with the excellent performance of deep learning on various classification problems, everyone tried to use the deep learning technology to predict the stock market. Deep learning techniques have significant success in many predictive tasks because they can automatically extract useful features during the learning process [32, 33]. Chong et al. (2017) predict the future stock market trends by analyzing the effects of three unsupervised feature extraction methods [34]. Sezer et al. (2017) proposed a stock trading system based on deep neural networks [35]. In other studies, artificial neural networks (ANN) have very good prediction accuracy even in the case of complex variable relationships. Boonpeng and Jeatrakul (2016) [36] proposed one-to-one and one-to-one signal prediction models for predictive trading or positions, and compared their models with traditional neural network models. They tested on the dataset of the Thai stock market and found that one-to-many models performed better than one-to-one models, with an average prediction accuracy of 72.5%. In Qiu and Song (2016) [36], an optimized neural network model using genetic algorithms was used to predict the trend of the stock market. The predictions of the two different feature sets on a data set (Tokyo Stock Exchange) are significantly different, with one prediction accuracy of 61.87% and the other up to 81.27%.

## 2.2. RNN for Time Series Prediction

Most ANNs, including multi-layer perceptrons, can only learn spatial patterns from independent time inputs and outputs. RNN has certain advantages over other traditional ANNs, because the memory mechanism of RNN can lead to the time pattern in the data. Due to this

characteristic, RNN is often used for time series analysis. In recent years, with the rapid development of deep learning technology, RNN has been widely used in various fields, such as natural language processing, speech recognition, computer vision, and serialized data [37, 38]. In addition, some studies have used RNN and LSTM networks to achieve satisfactory results in financial time series prediction. Lin et al. (2009) used RNN to predict the closing price of the next trading day. They used the Hurst index to select the initial transient and selected the sub-sequence with the strongest predictive ability during training [39]. Wei and Cheng (2012) propose a method for detecting key technical indicators of stock market forecasting using comprehensive feature selection [40]. They use stepwise regression and decision trees to reduce the dimensions of financial data. Fischer and Krauss (2018) [41] used the LSTM network to predict the trend of S&P 500 constituents from 1992 to 2015. They compared test results on random forests, deep neural networks, and logistic regression models. The LSTM model has obvious advantages over other models. The experiment found that the LSTM model network is suitable for the financial sector.

Therefore, the focus of this paper is to implement LSTM and discuss its advantages over traditional machine learning on stock data sets. These models are trained in 3, 5, 10, 15, 30, 60 days of different contrast time windows, and the predictions are very different. Most of our work is concentrated on the comparison time window of 10 ~ 60 days. Since most previous studies have tended to use classifiers for unsmoothed time series data, these models cannot learn from data sets when it comes to predictions of long-running windows. In our proposed method, we first use exponential smoothing to preprocess the data, then calculate the increase and decrease of the price, and then the classification algorithm is used to predict the rise and fall of the stock price.

The entire article is organized as follows: The third part introduces the methods we use, the processing of data, and the extraction of features. The fourth part introduces our experimental results and simple model comparison analysis. The fifth part introduces the summary of the experiment. Appendix A further expands the results. Appendix B gives our data sources.

### 3. Methods and Analysis

In our experiment, after the exponential smoothing of the data of each dimension of the stock, the feature extraction is performed. The extracted feature indicators provide information on the direction of the stock market in the future. These extracted indicators will be used to train as feature of the classifier. In this part we mainly describe data processing, feature extraction, classification algorithms, and some parameters of the model.

#### 3.1. Data Processing

Exponential smoothing is used to give greater weight to recent observations, thereby reducing the weight of past observations. The exponential smoothing statistic for time series  $Y$  can be recursively calculated as:

$$S_0 = Y_0$$

$$\text{for } t > 0, S_t = \beta * Y_t + (1 - \beta) * S_{t-1}$$

$\beta$  is the smoothing factor and  $0 < \beta < 1$ , the larger the value of  $\beta$ , the smaller the degree of smoothing. The smoothed data is equal to the original data when  $\beta = 1$ .  $S_t$  is the smoothed data. By smoothing, the noise of random variables can be removed, making it easier for the model to accurately identify the stock price trend. After calculating the smoothed data, the feature index is calculated based on the data, and a feature matrix is formed. Our goal is to predict the rise or fall of the  $i$ -day in the future relative to the previous  $d$ -day ( $i$  takes 1 in our experiment):

$$\text{target}_i = \text{sign}(\text{close}_i - \text{close}_{i-d})$$

$close_{i-d}$  represents the closing price of  $d$  days prior to the predicted  $i$ -day;  $close_i$  represents the predicted closing price of the  $i$ th day. When  $target_i = +1$ , it indicates that the predicted date has a positive change with respect to the price before  $d$  days, and when  $target_i = 0$ , it indicates that the predicted date has a negative change with respect to the price before  $d$  days.  $target_i$  is the target value of the  $i$ -th row feature matrix.

### 3.2. Feature Extraction

In our solution, we collected a lot of data, but only consider the closing price of the stock. So our data format is (data, closing price). In addition to the original indicators of the data (daily closing price, daily opening price, daily high price, daily low price, daily trading volume), this paper also calculated the following indicators from these data:

Relative strength index (RSI):RSI is based on the balance of supply and demand in the stock market. By comparing the rise and fall of the price of a single stock or the rise and fall of the entire market index over a period of time, the strength of the trading force in the market is analyzed. To judge the trend of the future market.

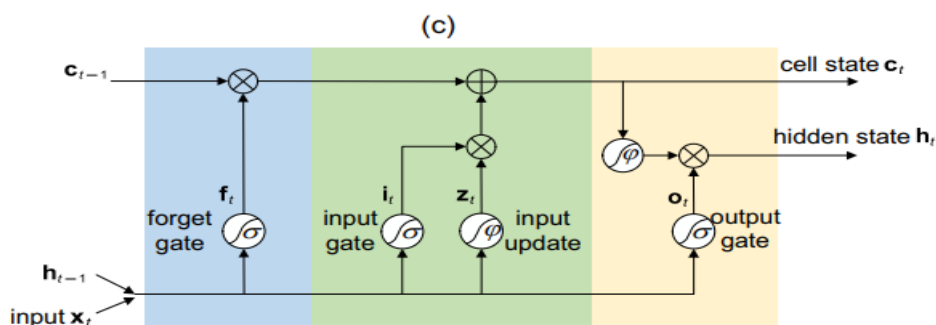
Williams percentage range (W%R): The William indicator mainly judges the overbought and oversold phenomenon of the stock market by analyzing the relationship between the highest price, the lowest price and the closing price of the stock price over a period of time, and predicts the short-term trend of the stock price. It mainly uses the oscillation point to reflect the overbought and oversold behavior of the market, analyzes the comparison of the strength of both sides, and proposes effective signals to judge the trend of short-term behavior in the market.

Moving average convergence divergence (MACD):The MACD is characterized by the separation and aggregation of short- and medium-term fast and slow moving averages, plus double smoothing to determine the timing and signals of buying and selling.

Price rate of change (PROC):A technical indicator that reflects the percentage change between the current price and the trading window price.

On balance volume (OBV):This technical indicator is to quantify the volume of transactions, in line with the stock price trend, from the price changes and the increase and decrease of the volume, speculate the market atmosphere. The theoretical basis of OBV is that the change of market price must have the volume of cooperation, the price rises and the volume does not rise and fall accordingly, then the original trend of market price is difficult to continue.

### 3.3. Machine Learning and Deep Learning Algorithms



**Figure 1:** An illustration of LSTM memory block

Long Short-Term Memory networks, usually just called “LSTMs”, are a special RNNs that are suitable for learning long-term dependencies . The key part that enhances LSTMs’capability to model long-term dependencies is a component called memory block. As illustrated in Fig.1, the memory block is a recurrently connected subnet that contains functional modules called

the memory cell and gates. The memory cell is in charge of remembering the temporal state of the neural network and the gates formed by multiplicative units are responsible for controlling the pattern of information flow. According to the corresponding practical functionalities, these gates are classified as input gates, output gates and forget gates. Input gates control how much new information flows into the memory cell, while forget gates control how much information of the memory cell still remains in the current memory cell through recurrent connection, and output gates control how much information is used to compute the output activation of the memory block and further flows into the rest of the neural network. Before going through the details of LSTM, some simple yet useful activation functions need to be reviewed. The sigmoid function  $\sigma(x) = 1/(1 + e^{-x})$  and the tanh function  $\psi(x) = 2\sigma(2x) - 1$  are commonly used as the activation function in ANNs. The domain of both functions is the real number field, but the return value for the sigmoid function ranges from 0 to 1, while the tanh function ranges from -1 to 1. Fig.1 explains in detail how LSTM works. The first step is to decide what kind of information will be removed from the memory cell state, which is implemented by a sigmoid layer (i.e., the forget gate). The next step is to decide what new information will be stored in the memory cell state. This operation can be divided into two steps. First, a sigmoid layer (i.e., the input gate) determines what will be updated, and a tanh layer creates a vector of new candidate values  $z_t$  that can be added to the memory cell state, where the subscript  $t$  denotes the current moment. Next, these two parts are combined to trigger an update to the memory cell state. To update the old memory cell state  $c_{t-1}$  into the new memory cell state  $c_t$ , we can first multiply the corresponding elements of  $c_{t-1}$  and the output of forget gate layer (i.e.  $f_t$ ), which is just like the oblivion mechanism in the human brain, and then add  $i_t * z_t$ , where  $i_t$  denotes the output of input gate and  $*$  denotes element-wise multiplication. The last step is to decide what to output, which is realized by element-wise multiplication between the value obtained from a tanh function of  $c_t$  and the output of a sigmoid layer (i.e., the output gate),  $o_t$ . Through the cooperation between the memory cell and the gates, LSTM is endowed with a powerful ability to predict time series with long-term dependences.

### 3.4. Data Sets and Parameter Settings

The research data for this experiment comes from the total 1763 trading days of the S&P500 from January 4, 2010 to December 30, 2016. Each sample contains the lowest, highest, opening, closing and trading volume of the trading day. The entire data set is divided into an 80% training set and a 20% test set, with 10% of the training set used for cross-validation. We randomly selected 8 stocks to present our experimental results. The stocks selected do not strictly consider their background or their economic impact on society. These companies come from different industries, such as software industry (AAPL), power industry (ETR), clothing industry (NKE), etc. It is to test the stability of our algorithm model.

In all of our experiments, we used the following as settings for the preprocessing and classification:

1. Smoothing rate  $\beta=0.05$
2. The comparison window is also called span was varied 3, 5, 10, 15, 30, and 60 days.
3. The number of hidden layer neurons in LSTM is 256 in two layers. The activation function of LSTM layer is dropout function, 50% random unit deletion probability, and the number of all connected layer neurons is 32, using l1 weight loss constraint, output layer For the two-category task, optimize with the adam optimizer. The parameters of the LSTM model are batch\_size=64, epochs=40, the input dimension of the data (the number of features of the data).

4.The depth of the random forest is 8, and the number of base classifiers is 100. The Gini index is used as an indicator of node division.

#### 4. Experimental Result

The stocks we choose come from a variety of industries, and the company's choices are random without any human factors. To test the robustness of our model, we considered several indicators for the two-division problem: accuracy, precision, recall (also known as sensitivity), F1 score (harmonic average of accuracy and recall), and the area under the curve (AUC) of the ROC curve and Brier score (mean squared error of loss).

In order to demonstrate the results of the experiment, we chose AAPL and AMZN as representative samples, and the following presentations are all results on the test set.

**Table 1:** Experimental results of LSTM

stock	span	accuracy	precision	recall	f1_score	brier_score	auc
AAPL	3	0.72	0.72	0.72	0.72	0.28	0.72
	5	0.76	0.76	0.76	0.76	0.24	0.76
	10	0.87	0.87	0.87	0.87	0.13	0.87
	15	0.91	0.91	0.91	0.91	0.09	0.91
	30	0.91	0.91	0.91	0.91	0.09	0.91
	60	0.91	0.91	0.92	0.91	0.09	0.92
AMZN	3	0.75	0.74	0.74	0.74	0.25	0.74
	5	0.81	0.8	0.8	0.8	0.19	0.8
	10	0.87	0.86	0.87	0.87	0.13	0.87
	15	0.88	0.87	0.86	0.87	0.12	0.86
	30	0.94	0.93	0.93	0.93	0.06	0.93
	60	0.96	0.94	0.94	0.94	0.04	0.94

**Table 2:** Experimental results of the Random forest

stock	span	accuracy	precision	recall	f1_score	brier_score	auc
AAPL	3	0.55	0.55	0.55	0.55	0.45	0.55
	5	0.65	0.65	0.65	0.65	0.35	0.65
	10	0.75	0.75	0.75	0.75	0.25	0.75
	15	0.81	0.82	0.81	0.81	0.19	0.81
	30	0.77	0.77	0.77	0.77	0.23	0.77
	60	0.78	0.78	0.78	0.78	0.22	0.78
AMZN	3	0.52	0.52	0.52	0.52	0.48	0.52
	5	0.64	0.62	0.59	0.59	0.36	0.59
	10	0.75	0.77	0.71	0.72	0.25	0.71
	15	0.83	0.82	0.81	0.81	0.17	0.81
	30	0.89	0.88	0.84	0.86	0.11	0.84
	60	0.94	0.93	0.92	0.92	0.06	0.92

### 4.1. LSTM Experimental Results

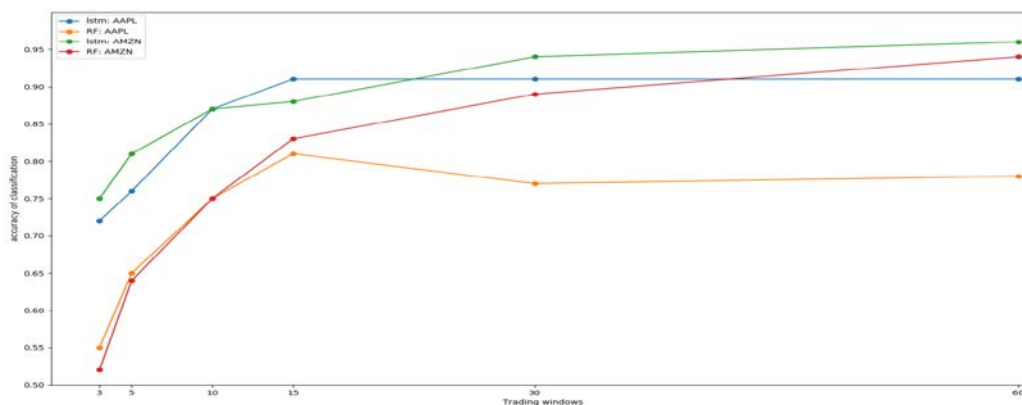
The sample classification results of the AAPL and AMZN stock data sets are shown in Table 1. Span represents the contrast window. In general, as the comparison window increases, the accuracy of the prediction and the f1\_score are also increasing.

### 4.2. Random Forests Experimental Results

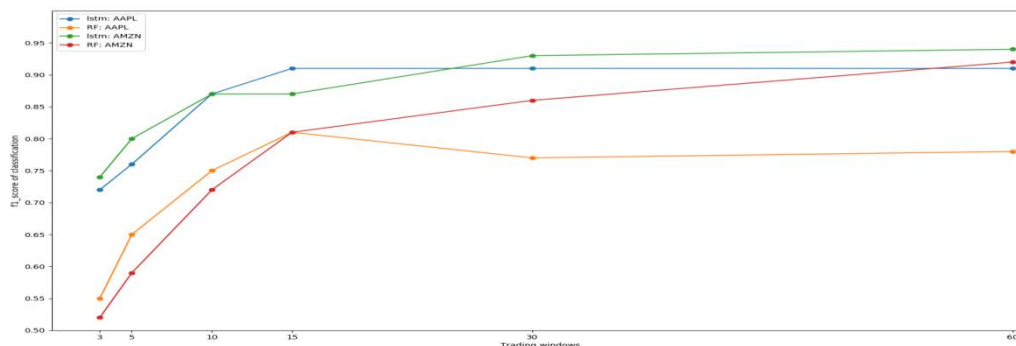
The sample classification results of the AAPL and AMZN stock data sets are shown in Table 2. In general, the accuracy of the incremental prediction of the comparison window and the f1\_score are also increasing. But significantly worse than the results of the various evaluation dimensions of the LSTM prediction. Especially for the prediction evaluation indicators of low contrast windows, LSTM has a better prediction effect.

### 4.3. Comparison of Prediction Results

In this section, we compare the predictive performance of LSTM and Random Forests from multiple dimensions. Finally, we implemented some classic classifiers and compared their average accuracy on different data sets.

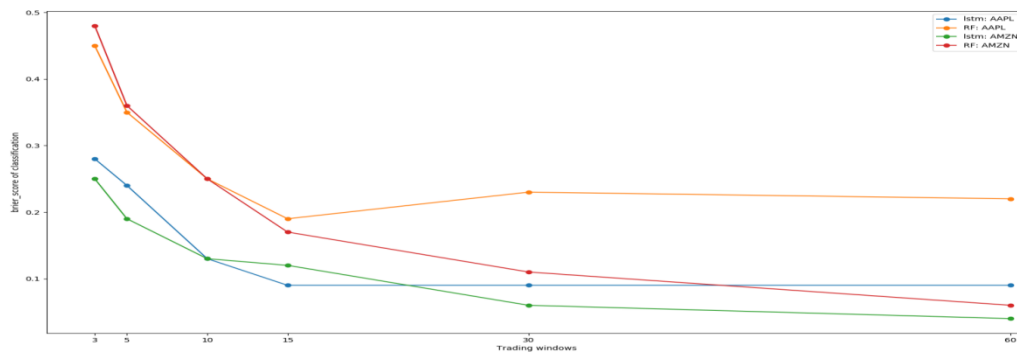


**Figure 2:** Prediction accuracy of random forest and LSTM in different contrast windows. The performance of LSTM in the high contrast window tends to be stable and the prediction accuracy is quite high, but the accuracy of random forest prediction in high contrast windows is not very stable, even a slight decrease in the AAPL data set. And LSTM is significantly better than the random forest prediction accuracy in the low contrast window.



**Figure 3:** Random forest and LSTM predictive f1\_score performance in different contrast windows, using two data sets(AAPL,AMZN). F1\_score increases in the comparison window range of 3-15 with the increase of the contrast window, and tends to be stable in the high contrast window.

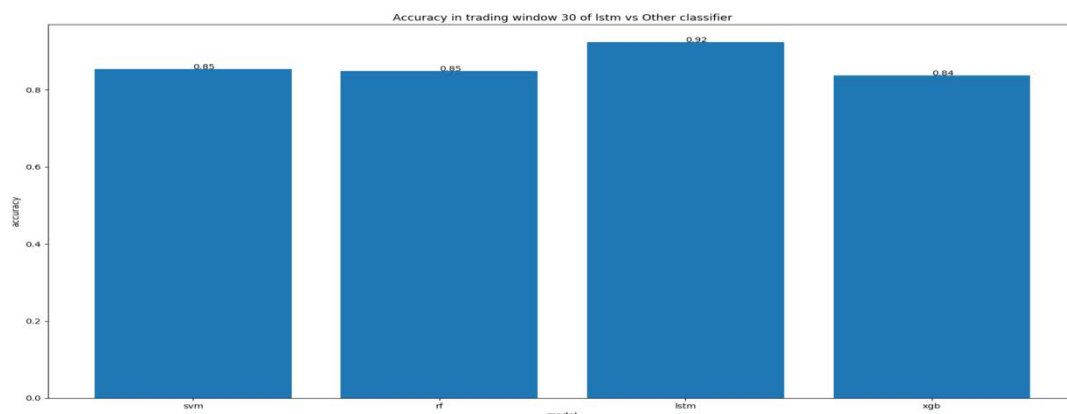




**Figure 4:** Random forest and LSTM predictive brier\_score (loss mean squared error) performance in different contrast windows, using two data sets. The Brier scores of the random forest and LSTM models decrease as the prediction contrast window increases.



**Figure 5:** Compare the prediction accuracy of LSTM and several popular models in the industry with a comparison window of 15. We forecast all data sets using different models and calculate the average accuracy of each model under the same data set. It can be seen from the figure that the average accuracy of the LSTM model is as high as 0.89, which is much higher than the accuracy of the second-ranked random forest.



**Figure 6:** Comparing the prediction accuracy of LSTM and several popular models in the industry with a comparison window of 30. As the contrast window increases, the prediction accuracy of each model increases to a certain extent, and the accuracy of the SVM model increases up to 7.59%. The accuracy of the LSTM model is still far superior to other models.



**Figure 7:** According to our model (the data is not evenly distributed in the comparison window, the paper uses span=15 as the comparison window). The prediction results on the test set are used to buy and sell AMZN stocks. The indication (red means sell, green means buy). The x-axis represents the number of days on the test set.

## 5. Conclusion

In this paper, we use the LSTM algorithm to construct our prediction model and get satisfactory results. The model is robust in predicting stock movements. The robustness of the model was evaluated by calculating parameters such as accuracy, precision, recall rate and F1-score. Through comparison and analysis, the validity of the model is proved, and its performance is better than the model discussed in the literature survey. In addition, a novelty of current work is the choice of technical indicators and their application as features.

Our models can be used to design new trading strategies or execute stock portfolio management to change stock trading based on trend forecasts. Compared with the existing forecasting methods, the model can predict stock volatility more accurately, thus reducing the investment risk of the stock market to a lower level. In the future, we can optimize this model to fit the stock trend forecast for high-frequency data trading. We also recommend using some swarm intelligence algorithms to assist our model to find the optimal parameters to achieve better predictions.

## 6. Appendix A. Results

In this appendix, we elaborate the experimental results achieved by performing the experiments (Tables 3,4)

**Table 3:** LSTM Results: results of long short term memory network implemented on stocks of AAPL, AMZN, FB, ETR, and MSI,NKE.

stock	span	accuracy	precision	recall	f1_score	brier_score	auc
AAPL	3	0.72	0.72	0.72	0.72	0.28	0.72
	5	0.76	0.76	0.76	0.76	0.24	0.76
	10	0.87	0.87	0.87	0.87	0.13	0.87
	15	0.91	0.91	0.91	0.91	0.09	0.91
	30	0.91	0.91	0.91	0.91	0.09	0.91
	60	0.91	0.91	0.92	0.91	0.09	0.92
AMZN	3	0.75	0.74	0.74	0.74	0.25	0.74
	5	0.81	0.8	0.8	0.8	0.19	0.8
	10	0.87	0.86	0.87	0.87	0.13	0.87
	15	0.88	0.87	0.86	0.87	0.12	0.86
	30	0.94	0.93	0.93	0.93	0.06	0.93
	60	0.96	0.94	0.94	0.94	0.04	0.94
FB	3	0.7	0.7	0.7	0.7	0.3	0.7
	5	0.76	0.75	0.76	0.75	0.25	0.76
	10	0.83	0.83	0.83	0.83	0.17	0.83
	15	0.87	0.88	0.86	0.87	0.13	0.86
	30	0.94	0.94	0.92	0.93	0.06	0.92
	60	0.95	0.93	0.92	0.92	0.05	0.92
ETR	3	0.74	0.73	0.73	0.73	0.26	0.73
	5	0.77	0.79	0.77	0.77	0.23	0.77
	10	0.84	0.84	0.84	0.84	0.16	0.84
	15	0.87	0.87	0.87	0.87	0.13	0.87
	30	0.95	0.95	0.95	0.95	0.05	0.95
	60	0.93	0.93	0.93	0.93	0.07	0.93
MSI	3	0.73	0.73	0.73	0.73	0.27	0.73
	5	0.82	0.81	0.82	0.82	0.18	0.82
	10	0.87	0.87	0.87	0.87	0.13	0.87
	15	0.89	0.89	0.89	0.89	0.11	0.89
	30	0.96	0.95	0.95	0.95	0.04	0.95
	60	0.96	0.96	0.96	0.96	0.04	0.96
NKE	3	0.72	0.73	0.72	0.72	0.28	0.72
	5	0.79	0.79	0.79	0.79	0.21	0.79
	10	0.84	0.84	0.85	0.84	0.16	0.85
	15	0.9	0.9	0.91	0.9	0.1	0.91
	30	0.89	0.89	0.89	0.89	0.11	0.89
	60	0.96	0.95	0.96	0.96	0.04	0.96

**Table 4:** RF Results: results of random forest implemented on stocks of AAPL,AMZN, FB, ETR, and MSI,NKE.

stock	span	accuracy	precision	recall	f1_score	brier_score	auc
AAPL	3	0.55	0.55	0.55	0.55	0.45	0.55
	5	0.65	0.65	0.65	0.65	0.35	0.65
	10	0.75	0.75	0.75	0.75	0.25	0.75
	15	0.81	0.82	0.81	0.81	0.19	0.81
	30	0.77	0.77	0.77	0.77	0.23	0.77
	60	0.78	0.78	0.78	0.78	0.22	0.78
AMZN	3	0.52	0.52	0.52	0.52	0.48	0.52
	5	0.64	0.62	0.59	0.59	0.36	0.59
	10	0.75	0.77	0.71	0.72	0.25	0.71
	15	0.83	0.82	0.81	0.81	0.17	0.81
	30	0.89	0.88	0.84	0.86	0.11	0.84
	60	0.94	0.93	0.92	0.92	0.06	0.92
FB	3	0.59	0.62	0.59	0.56	0.41	0.59
	5	0.72	0.72	0.71	0.71	0.28	0.71
	10	0.7	0.69	0.69	0.69	0.3	0.69
	15	0.76	0.77	0.73	0.74	0.24	0.73
	30	0.88	0.88	0.85	0.86	0.12	0.85
	60	0.94	0.96	0.88	0.91	0.06	0.88
ETR	3	0.56	0.55	0.56	0.55	0.44	0.56
	5	0.57	0.57	0.57	0.56	0.43	0.57
	10	0.73	0.73	0.72	0.72	0.27	0.72
	15	0.84	0.84	0.83	0.83	0.16	0.83
	30	0.91	0.92	0.91	0.91	0.09	0.91
	60	0.85	0.87	0.84	0.84	0.15	0.84
MSI	3	0.5	0.48	0.48	0.47	0.5	0.48
	5	0.58	0.57	0.53	0.48	0.42	0.53
	10	0.75	0.78	0.73	0.73	0.25	0.73
	15	0.81	0.81	0.8	0.81	0.19	0.8
	30	0.92	0.91	0.92	0.92	0.08	0.92
	60	0.94	0.95	0.91	0.93	0.06	0.91
NKE	3	0.56	0.59	0.58	0.55	0.44	0.58
	5	0.55	0.59	0.57	0.54	0.45	0.57
	10	0.68	0.7	0.7	0.68	0.32	0.7
	15	0.79	0.79	0.79	0.79	0.21	0.79
	30	0.79	0.78	0.78	0.78	0.21	0.78
	60	0.93	0.94	0.91	0.93	0.07	0.91

## 7. Appendix B. Supplementary Data

Supplementary data associated with this article can be found, in the online version, <https://www.kaggle.com/camnugent/sandp500/download>.

### References

- [1] Malkiel, B. G., & Fama, E. F. (1970). Efficient capital markets: A review of theory and empirical work. *The Journal of Finance*, 25(2), 383–417.
- [2] Jensen, M. C. (1978). Some anomalous evidence regarding market efficiency. *Journal of Financial Economics*. 6(2), 95–101.
- [3] G. A. Carpenter, S. Grossberg, N. Markuzon, J. H. Reynolds, and D. B. Rosen, “Artmap: a neural network architecture for incremental learning supervised learning of analog multidimensional maps,” *IEEE Transactions in Neural Networks*, vol. 3, no. 5, pp.698–713, 1992.
- [4] E. Guresen, G. Kayakutlu, T.U. Daim, Using artificial neural network models in stock market index prediction, *Expert Syst. Appl.* 38 (8) (2011) 10389–10397.
- [5] T.J. Hsieh, H.F. Hsiao, W.C. Yeh, Forecasting stock markets using wavelet transforms and recurrent neural networks: An integrated system based on artificial bee colony algorithm, *Appl. Soft Comput. J.* 11 (2) (2011) 2510–2525.
- [6] K.J. Kim, I. Han, Genetic algorithms approach to feature discretization in artificial neural networks for the prediction of stock price index, *Expert Syst. Appl.* 19 (2) (2000) 125–132.
- [7] J. Wang, J. Wang, Forecasting stock market indexes using principle component analysis and stochastic time effective neural networks, *Neurocomputing* 156 (C) (2015) 68–78.
- [8] D. Pradeepkumar, V. Ravi, Forecasting financial time series volatility using particle swarm optimization trained quantile regression neural network, *Appl. Soft Comput.* 58 (2017) 35–52.
- [9] T. Lin, B.G. Horne, C.L. Giles, How embedded memory in recurrent neural network architectures helps learning long-term temporal dependencies, *Neural Netw.* 11 (5) (1998) 861–868.
- [10] P.A. Chen, L.C. Chang, F.J. Chang, Reinforced recurrent neural networks for multi-step-ahead flood forecasts, *J. Hydrol.* 497 (2013) 71–79.
- [11] T. Guo, Z. Xu, X. Yao, H. Chen, K. Aberer, K. Funaya, Robust online time series prediction with recurrent neural networks, in: *IEEE International Conference on Data Science and Advanced Analytics*, vol. 9, IEEE, 2016, pp. 816–825.
- [12] M. Sundermeyer, R. Schlüter, H. Ney, LSTM neural networks for language modeling, in: *Interspeech*, 2012, pp. 601–608.
- [13] M. Sundermeyer, H. Ney, R. Schlüter, From feedforward to recurrent lstm neural networks for language modeling, *IEEE/ACM Trans. Audio Speech Lang. Process.* 23 (3) (2015) 517–529.
- [14] Suryoday Basak, Saibai kar, Snehanshu Saha, Predicting the direction of stock market prices using tree-based Classifiers, *North American Journal of Economics and Finance* 47 (2019) 552–567.
- [15] Tay, F.E.; Cao, L. Application of support vector machines in financial time series forecasting. *Omega* 2001, 29, 309–317.
- [16] Kim, J.H.; Shamsuddin, A.; Lim, K.P. Stock return predictability and the adaptive markets hypothesis: Evidence from century-long US data. *J. Empir. Finan.* 2011, 18, 868–879.
- [17] Kumar, D.A.; Murugan, S. Performance analysis of Indian stock market index using neural network time series model. In *Proceedings of the International Conference on Pattern Recognition, Informatics and Mobile Engineering*, Salem, India, 21–22 February 2013; pp. 72–78.
- [18] Armano, G.; Marchesi, M.; Murru, A. A hybrid genetic-neural architecture for stock indexes forecasting. *Inf. Sci.* 2005, 170, 3–33.
- [19] Rao, J.N.K.; Box, G.E.P.; Jenkins, G.M. Time Series Analysis Forecasting and Control. *Econometrica* 1972, 40, 970.

- [20] Engle, R.F. Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation. *Econometrica* 1982, 50, 987–1007.
- [21] Karolyi, G.A. A multivariate GARCH model of international transmissions of stock returns and volatility: The case of the United States and Canada. *J. Bus. Econ. Stat.* 1995, 13, 11–25.
- [22] Franses, P.H.; Van Dijk, D. Forecasting stock market volatility using (nonlinear) GARCH models. *J. Forecast.* 1996, 15, 229–235.
- [23] Abu-Mostafa, Y.S.; Atiya, A.F. Introduction to financial forecasting. *Appl. Intell.* 1996, 6, 205–213.
- [24] Widom, J. (1995). Research problems in data warehousing. Proceedings of the fourth international conference on information and knowledge management. CIKM '95 (pp.25–30). New York, NY, USA: ACM. <http://dx.doi.org/10.1145/221270.221319>.
- [25] Hellstrom, T. & Holmstrom, K. (1998). Predictable Patterns in Stock Returns. Technical Report Series IMA-TOM-1997-09.
- [26] Saha, S., Routh, S., & Goswami, B. (2014). Modeling vanilla option prices: A simulation study by an implicit method. *Journal of Advances in Mathematics.* 6(1),834–848.
- [27] Maymin, P. (2012). Music and the market: Song and stock volatility. *The North American Journal of Economics and Finance*, 23(1), 70–85. <http://dx.doi.org/10.1016/j.najef.2011.11.004>.
- [28] Khanal, A. R., & Mishra, A. K. (2017). Stock price reactions to stock dividend announcements: A case from a sluggish economic period. *The North American Journal of Economics and Finance*, 42, 338–345. <http://dx.doi.org/10.1016/j.najef.2017.08.002>.
- [29] Dai, Y., & Zhang, Y. (2013). Machine learning in stock price trend forecasting. Stanford University [http://cs229.stanford.edu/proj2013/DaiZhang-Machine Learning In Stock Price Trend Forecasting.pdf](http://cs229.stanford.edu/proj2013/DaiZhang-Machine%20Learning%20In%20Stock%20Price%20Trend%20Forecasting.pdf).
- [30] Saad, E.W.; Prokhorov, D.V.; Wunsch, D.C. Comparative study of stock trend prediction using time delay, recurrent and probabilistic neural networks. *IEEE Trans. Neural Netw.* 1998, 9, 1456–1470.
- [31] Chen, Y.; Abraham, A.; Yang, J.; Yang, B. Hybrid methods for stock index modeling. In Proceedings of the International Conference on Fuzzy Systems and Knowledge Discovery, Changsha, China, 27–29 August 2005; Springer: Berlin/Heidelberg, Germany; pp. 1067–1070.
- [32] Guo, Y.; Liu, Y.; Oerlemans, A.; Lao, S.; Wu, S.; Lew, M.S. Deep learning for visual understanding: A review. *Neurocomputing* 2016, 187, 27–48.
- [33] Lee, J.; Jang, D.; Park, S. Deep Learning-Based Corporate Performance Prediction Model Considering Technical Capability. *Sustainability* 2017, 9, 899.
- [34] Chong, E.; Han, C.; Park, F.C. Deep learning networks for stock market analysis and prediction: Methodology, data representations, and case studies. *Expert Syst. Appl.* 2017, 83, 187–205.
- [35] Sezer, O.B.; Ozbayoglu, M.; Dogdu, E. A Deep Neural-Network Based Stock Trading System Based on Evolutionary Optimized Technical Analysis Parameters. *Procedia Comput. Sci.* 2017, 114, 473–480.
- [36] Boonpeng, S., & Jeatrakul, P. (2016). Decision support system for investing in stock market by using OAA-neural network. In: 8th International Conference on Advanced Computational Intelligence Chiang Mai, Thailand.
- [37] Qiu, M., & Song, U. (2016). Predicting the direction of stock market index movement using an optimized artificial neural network model. *PLoS One*, 11(5).
- [38] Brocki, Ł.; Marasek, K. Deep belief neural networks and bidirectional long-short term memory hybrid for speech recognition. *Arch. Acoust.* 2015, 40, 191–195.
- [39] Donahue, J.; Anne Hendricks, L.; Guadarrama, S.; Rohrbach, M.; Venugopalan, S.; Saenko, K.; Darrell, T. Long-term recurrent convolutional networks for visual recognition and description. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 7–12 June 2015; pp. 2625–2634.
- [40] Lin, X.; Yang, Z.; Song, Y. Short-term stock price prediction based on echo state networks. *Expert Syst. Appl.* 2009, 36, 7313–7317.

- [41] Wei, L.Y.; Cheng, C.H. A hybrid recurrent neural networks model based on synthesis features to forecast the Taiwan stock market. *Int. J. Innov. Comput. Inf. Control* 2012, 8, 5559–5571.
- [42] Fischer, T.; Krauss, C. Deep learning with long short-term memory networks for financial market predictions. *Eur. J. Oper. Res.* 2018, 270, 654–669.