

Joint Latent Topic Discovery and Expectation Modeling for Financial Markets

Lili Wang¹, Chenghan Huang², Chongyang Gao³, Weicheng Ma¹, and Soroush Vosoughi¹

¹ Dartmouth College, Hanover NH 03755, USA

² Millennium Management, LLC, New York NY 10022, USA

³ Northwestern University, Evanston IL 60208, USA
{lili.wang.gr,soroush}@dartmouth.edu

Abstract. In the pursuit of accurate and scalable quantitative methods for financial market analysis, the focus has shifted from individual stock models to those capturing interrelations between companies and their stocks. However, current relational stock methods are limited by their reliance on predefined stock relationships and the exclusive consideration of immediate effects. To address these limitations, we present a groundbreaking framework for financial market analysis. This approach, to our knowledge, is the first to jointly model investor expectations and automatically mine latent stock relationships. Comprehensive experiments conducted on China’s CSI 300, one of the world’s largest markets, demonstrate that our model consistently achieves an annual return exceeding 10%. This performance surpasses existing benchmarks, setting a new state-of-the-art standard in stock return prediction and multiyear trading simulations (i.e., backtesting).

Keywords: stock trend prediction · trading simulation · expectation modeling.

1 Introduction

The efficient-market hypothesis in traditional finance posits that stock prices reflect all available market information, with current prices consistently trading at their fair value [5]. Consequently, predicting future stock prices is challenging without access to new information. However, markets are often less efficient in reality [11], with stock market fluctuations driven by behavioral factors such as expectations, confidence, panic, euphoria, or herding behavior. These inefficiencies enable the use of machine learning to predict future stock movements based on historical trends.

Stock-affecting behavioral factors can be categorized into short- and long-term factors. Factors like panic, euphoria, or herding behavior are typically short-term, while subjective expectations and confidence tend to be long-term factors, only influencing stock prices imperceptibly over extended periods. These factors do not solely impact individual stocks; their effects often spread to topically related stocks, which share similarities across various explicit or latent



Fig. 1. The return of Amazon and Facebook (Meta) stocks from 2021-05-03 to 2022-07-08 with respect to their stock prices at 2021-05-03.

dimensions. Recent stock prediction works [14,21,15] utilize topic stocks to improve prediction capabilities. However, most of these methods exhibit two key limitations:

(1) Topics are typically assumed to be static and known beforehand. However, real-world topics can change and new topics may emerge. For example, during the COVID-19 pandemic, pharmaceutical companies investing in COVID vaccines (e.g., Pfizer⁴ and Moderna⁵) experienced stock price fluctuations under the new COVID topic.

(2) Only the *short-term impact* between stocks is considered, neglecting the *long-term subjective expectations*. Unlike analyst expectations, subjective expectations are based on human psychology and behavior and can be irrational. Figure 1 illustrates that Amazon and Facebook stock prices often correlate, and previous methods might reason that a significant drop in Facebook’s price would also lead to plummeting Amazon stocks. However, in the second half of 2021, Amazon’s return was lower than Facebook’s, lowering investor expectations for Amazon. Thus, when Amazon released an unremarkable financial report⁶ on February 3, 2022, its stock rose 13.5

⁴ <https://investors.pfizer.com/Investors/Stock-Info/default.aspx>

⁵ <https://investors.modernatx.com/Stock-Info/default.aspx>

⁶ https://s2.q4cdn.com/299287126/files/doc_financials/2021/q4/business_and_financial_update.pdf

In this paper, we introduce a novel attention-based framework for stock trend prediction that simultaneously discovers topical relations between stocks and models both the *short-term impact* and *long-term subjective expectations* of topically similar stocks. To the best of our knowledge, our framework is the first to:

- Model the influence of investors’ subjective expectations on stock prices.
- Automatically identify dynamic topics between stocks without making assumptions or requiring additional knowledge.

Through comprehensive experiments against 16 well-established baselines, we demonstrate that our method achieves the current state-of-the-art on the Qlib [22] quantitative investment platform.

2 Related Work

The stock price prediction and stock selection problems can be easily formed as a time series forecasting problem. Therefore, traditional and deep-learning-based machine learning (ML) methods, especially those for sequence learning, have been directly applied to these tasks and are widely used by investment institutions. Specifically, Qlib [22], a popular quantitative investment platform, benchmarks models based on the following ML methods: multi-layer perceptron (MLP); TabNet [1]; TCN [2]; gradient boosting models: CatBoost [12], LightGBM [8]; Recurrent Neural Network (RNN) based models: long short-term memory (LSTM) [6], gated recurrent unit (GRU) [3], DA-RNN [13], AdaRNN [4]; and attention-based models: Transformer [18], and Localformer [7]. To model the co-movement and relations among stocks, some research, such as MAN-SF [15] and STHAN-SR [14]), also adopted graph neural network methods like GCN [10] and GATs [19] to mine the correlation between different stocks.

More recent models include those specifically designed for stock trading. DoubleEnsemble [23] is an ensemble model which utilizes learning-trajectory-based sample reweighting and shuffling-based feature selection for stock prediction. ADD [16] attempt to extract clean information from noisy data to improve prediction performances. Specifically, they proposed a method for separating the inferential features from the noisy raw data to a certain degree using disentanglement, dynamic self-distillation, and data augmentation. Xu et al. assume that inter-dependencies may exist among different stocks at different time series and propose a method called IGMTF [20] to mine these relations. In their other work, they propose HIST [21], a three-step framework to mine the concept-oriented shared information and individual features among stocks.

We use most of the above-mentioned methods as baselines in our experiments.

3 Framework

3.1 Problem Definition

We formulate the stock trend prediction problem as a regression problem. Let $stock_1, stock_2, \dots, stock_n$ denote n different stocks. For each stock $stock_j$ on date i , the closing price is $price_j^i$. Given the historical information before date i , our task is to predict the one-day return $r_j^i = \frac{price_j^i - price_j^{i-1}}{price_j^{i-1}}$ for each stock j on date i . In the rest of this paper, we use r^i to denote $(r_1^i, r_2^i, \dots, r_n^i)$.

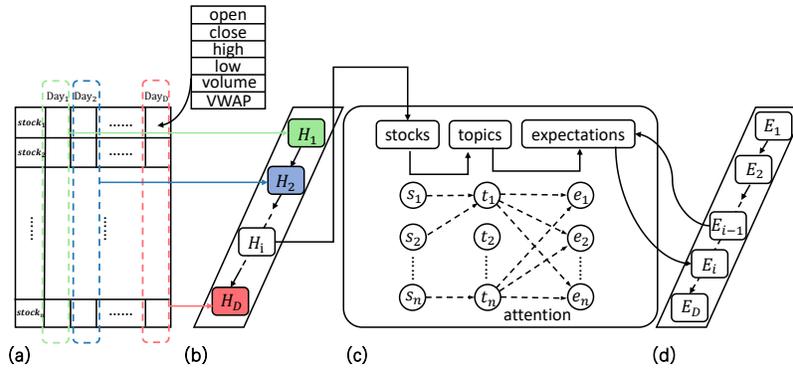


Fig. 2. Our model’s framework consists of: (a) Extracting Alpha360 features from raw data: For each stock on a given day, we combine the opening price, closing price, highest price, lowest price, trading volume, and volume-weighted average price (VWAP) into a 6-D feature vector. We then concatenate this vector with similar 6-D vectors from the preceding 59 days to form a 360-D feature vector. (b) The LSTM module processes the extracted Alpha360 features to learn temporal representations. (c) The left half of this section represents the topic module, which uses stock embeddings as input to extract latent topics. The right half illustrates the expectations module, which takes the E_{i-1} output from the expectation LSTM in part (d) as an initial embedding, employs attention with topics to update it to \hat{E}_{i-1} , and feeds it back into the expectation LSTM as input for day i . (d) The second LSTM module models the evolution of each stock’s expectation.

3.2 Overview of the Framework

The architecture of our model is shown in Figure 2. The model consists of three jointly optimized modules: temporal stock representation (which aims to extract temporal stock features), topic module (aims to discover the dynamic topics based on the extracted features), and expectation module (aims to model the subjective expectations for each stock). Below we describe each module in detail.

3.3 Temporal Stock Representation

The first step of our learning framework is to extract the Alpha360 features [22] from the raw data. The Alpha360 is a 360-D feature vector that is widely used in the quantitative investment domain. As shown in Figure 2 (a), for each stock on each day, we combine the opening price, closing price, highest price, lowest price, trading volume and volume-weighted average price (VWAP) as a 6-D feature vector and concatenate it with similar 6-D vectors from the past 59 days to get a 360-D feature vector.

To extract the temporal representation of stocks, we adopt an LSTM layer shown in Figure 2 (b)). Our framework is trained recursively by date: for each trading day i , the input is the Alpha360 features H_i of $stock_1, stock_2, \dots, stock_n$ for that day and the output of the LSTM layer is S_i , which is comprised of $s_1^i, s_2^i, \dots, s_n^i$, denoting the embeddings of each stock.

3.4 Topic Module

As mentioned before, the relations among stocks may evolve overtime, so our framework needs to be able to capture the evolution of topics and discover new topics each day. Figure 2 (c) shows the topic and expectation modules of our framework.

First, for each day i , we initialize the n topic embeddings $T_i = (t_1^i, t_2^i, \dots, t_n^i)$ using the n stock embeddings $S_i = (s_1^i, s_2^i, \dots, s_n^i)$. Then, we compute the Tanimoto coefficient (\mathcal{T}) [17] between all pairs of $t_{j_1}^i$ (topic j_1 in day i) and $s_{j_2}^i$ (stock j_2 in day i), for $\forall j_1, j_2 \in [1, n]$ with the following equation:

$$\mathcal{T}(t_{j_1}^i, s_{j_2}^i) = \frac{t_{j_1}^i s_{j_2}^i}{\|t_{j_1}^i\|^2 + \|s_{j_2}^i\|^2 - t_{j_1}^i s_{j_2}^i} \quad (1)$$

We define a function $\phi^i(s_{j_2}^i)$ that for each stock embedding, $s_{j_2}^i$, returns the most similar topic index j_1 , except for its own topic (i.e., $j_1 \neq j_2$) in date i , based on the Tanimoto coefficient:

$$\phi^i(s_{j_2}^i) = \arg \max_{j_1} \left(\mathcal{T}(t_{j_1}^i, s_{j_2}^i), j_1 \neq j_2 \right) \quad (2)$$

In the example shown in Figure 2 (c) (with the dashed lines) $\phi^i(s_1^i) = 1$, $\phi^i(s_2^i) = 1$, $\phi^i(s_n^i) = n$.

We further construct a set $valid^i$ that contains ‘‘valid’’ topics for each day i , i.e., those that are the most related to at least one stock:

$$valid^i = \left\{ x | \exists j, x = \phi^i(s_j^i) \right\} \quad (3)$$

This set denotes the topics we discovered for each day. Only if a topic $t_{j_1}^i$ is the most similar topic to at least one stock, it will be include in this set, other topics (e.g., t_2^i in Figure 2 (c)) will be excluded from the following calculations.

To update each topic embedding $t_{j_1}^i$ ($j_1 \in valid^i$), we train the fully connected layer with weight matrix W_t , bias matrix b_t and activation function \tanh to aggregate the stock embeddings using the Tanimoto coefficient:

$$t_{j_1}^i = \tanh \left(W_t \left(\sum_{\phi^i(s_{j_2}^i)=j_1} \mathcal{T}(t_{j_1}^i, s_{j_2}^i) s_{j_2}^i \right) + b_t \right) \quad (4)$$

3.5 Expectation Module

The expectations of investors change over time and our framework needs to take that into consideration. As shown in Figure 2 (d), we adopt an LSTM to model the evolving expectations of each stock. Each E_i consists of n expectation embeddings $e_1^i, e_2^i, \dots, e_n^i$, we assume that at the first timestamp, the investor's expectations are all decided by the stocks themselves, so the initial embedding $E_1 = (e_1^1, e_2^1, \dots, e_n^1)$ are initialized as the n stock embeddings $S_1 = (s_1^1, s_2^1, \dots, s_n^1)$.

The expectation for one stock can also be affected by the performance of other stocks under related topics. So for day i , we take the output $E_{i-1} = (e_1^{i-1}, e_2^{i-1}, \dots, e_n^{i-1})$ of the LSTM and adopt an attention mechanism to learn the importance of each topic j_1 to the expectations:

$$\alpha(t_{j_1}^i, e_{j_2}^{i-1}) = \frac{\exp(\mathcal{T}(t_{j_1}^i, e_{j_2}^{i-1}))}{\sum_{j \in \text{valid}^i} \exp(\mathcal{T}(t_j^i, e_{j_2}^{i-1}))} \quad (5)$$

$$\hat{e}_{j_2}^{i-1} = \tanh \left(W_e^1 e_{j_2}^{i-1} + W_e^2 \left(\sum_{j_1 \in \text{valid}^i} \alpha(t_{j_1}^i, e_{j_2}^{i-1}) t_{j_1}^i \right) + b_e \right) \quad (6)$$

where $\alpha(t_{j_1}^i, e_{j_2}^{i-1})$ measures the importance of topic j_1 to the expectation of stock j_2 , and the updated $\hat{E}_{i-1} = (\hat{e}_1^{i-1}, \hat{e}_2^{i-1}, \dots, \hat{e}_n^{i-1})$ then feed back to the LSTM (d) as the input of day i .

3.6 Loss Function

The objective of our model is to predict the one-day return r of each stock. The objective relies on three components: r_{stock} , r_{topic} , and $r_{\text{expectation}}$.

The r_{stock} and $r_{\text{expectation}}$ are learnt from the temporal stock embeddings and the expectation embeddings, respectively :

$$r_{\text{stock}}^i = \tanh(W_{\text{stock}} S_i + b_{\text{stock}}) \quad (7)$$

$$r_{\text{expectation}}^i = \tanh(W_{\text{expectation}} E_i + b_{\text{expectation}}) \quad (8)$$

To learn r_{topic} , we first learn the importance of each topic to the stocks using a similar attention mechanism as the expectation module:

$$\beta(t_{j_1}^i, s_{j_2}^i) = \frac{\exp(\mathcal{T}(t_{j_1}^i, s_{j_2}^i))}{\sum_{j \in \text{valid}^i} \exp(\mathcal{T}(t_j^i, s_{j_2}^i))} \quad (9)$$

$$o_{j_2}^i = \tanh \left(W_s \left(\sum_{j_1 \in \text{valid}^i} \beta(t_{j_1}^i, s_{j_2}^i) t_{j_1}^i \right) + b_s \right) \quad (10)$$

where $\beta(t_{j_1}^i, s_{j_2}^i)$ measures the importance of topic j_1 to the expectation of stock j_2 on day i . Note that different from the expectation module which includes the term $W_e^1 e_{j_2}^{i-1}$, here $o_{j_2}^i$ measures the impact of all the topics on the stock j_2 on day i , without considering s_i . This is because s_i is already included in r_{stock} . We use O^i to denote $(o_1^i, o_2^i, \dots, o_n^i)$; r_{topic} is learnt as:

$$r_{topic}^i = \tanh(W_{topic} O^i + b_{topic}) \quad (11)$$

The predicted return \hat{r} is learnt by combining these three components:

$$\hat{r}^i = \tanh(W_{\hat{r}} r_{stock}^i + W_{\hat{r}} r_{topic}^i + W_{\hat{r}} r_{expectation}^i + b_{\hat{r}}) \quad (12)$$

The loss function of our model is defined as the mean squared error between \hat{r} and r :

$$\mathcal{L} = \frac{\sum_{i \in [1, D]} (r^i - \hat{r}^i)^\top (r^i - \hat{r}^i)}{D \cdot n} \quad (13)$$

where D corresponds to the number of trading days. Algorithm 1 shows the pseudocode of our method.

Algorithm 1 Training pseudo-code

Input: $H = \{H_1, H_2, \dots, H_{|D|}\}$: the Alpha360 features for each trading day

Parameters:

Θ : the initialized model parameters, $epochs$: the number of training epochs, η : learning rate

Output: The predicted return \hat{r}

```

1: for  $epoch \leftarrow \{1, \dots, epochs\}$  do
2:   for  $i \leftarrow \{1, \dots, D\}$  do
3:      $S_i \leftarrow LSTM_b(H_i)$ 
4:     if  $t == 1$  then
5:        $T_i \leftarrow S_i$ 
6:        $E_i \leftarrow S_i$ 
7:     end if
8:      $\mathcal{T} \leftarrow$  Calculate Tanimoto coefficient (Eq. 1)
9:      $valid^i \leftarrow$  Calculate the valid topic set according to  $\mathcal{T}$  (Eq. 3)
10:     $T_i \leftarrow$  Aggregate information from  $S_i$  according to  $\mathcal{T}$  (Eq. 4)
11:     $\alpha \leftarrow$  Calculate the attention weight (Eq. 5)
12:     $\hat{E}_i \leftarrow$  Aggregate information from  $T_i$  according to  $\alpha$  (Eq. 6)
13:     $E_{i+1} \leftarrow LSTM_d(\hat{E}_i)$ 
14:   end for
15:   Compute the stochastic gradients of  $\Theta$  (Eq.13)
16:   Update model parameters  $\Theta$  according to learning rate  $\eta$  and gradients.
17: end for
18: return the predicted return  $\hat{r}$ 

```

3.7 Model Training

Our model is optimized by minimizing the global loss \mathcal{L} . This was done using the Adam optimizer [9]. The hyper-parameters are set as follows: the embedding size is set to 128, the learning rate is set to 0.001, the training epoch is set to 300, the dropout rate is set to 0.1. All experiments are run on a Lambda Deep Learning 2-GPU Workstation (RTX 2080) with 24GB of memory, and the random seed is set to 0 at the beginning of each experiment.

4 Experiments

4.1 Datasets

We run comprehensive evaluations of our framework on the China’s CSI 300 financial markets, from 2008 to 2022. We use the data from 01/01/2008 to 12/31/2014 as the training set, the data from 01/01/2015 to 12/31/2016 as the validation set for hyper-parameter fine-tuning, and the data from 01/01/2017 to 07/10/2022 as the test set.

4.2 Baselines

We compare our framework with a comprehensive list of 16 well-known methods which are widely used in the financial sector. These methods span six different categories and are:

- **Classic Models** - MLP, TCN [2], GATs [19]
- **Tabular Learning** - TabNet
- **Gradient Boosting Models** - CatBoost [12], LightGBM [8]
- **RNN-Based Methods**- LSTM [6], GRU [3], DA-RNN [13], AdaRNN [4]
- **Attention-Based Methods**- Transformer [18], Localformer [7]
- **Financial Prediction Methods**- DoubleEnsemble [23], ADD [16], HIST [21], IGMTF [20]

Note that although our method can mine the latent topics among stocks, the tasks in our experiments only assume access to price and volume features (opening price, closing price, highest price, lowest price, VWAP). Several recently proposed methods require additional information such as company relations [14] or social media text [15], thus these methods cannot be included as baselines.

4.3 Results

We use **stock trend prediction** and **trading simulation** for our experiments.

Model Name	IC	ICIR	Rank IC	Rank ICIR
Transformer	0.0143±0.0024 *	0.0910±0.0180 *	0.0317±0.0024 *	0.2192±0.0190 *
TabNet	0.0286±0.0000 *	0.1975±0.0000 *	0.0367±0.0000 *	0.2798±0.0000 *
MLP	0.0267±0.0017 *	0.1845±0.0154 *	0.0362±0.0018 *	0.2681±0.0157 *
Localformer	0.0358±0.0036 *	0.2633±0.0334 *	0.0477±0.0019 *	0.3643±0.0218 *
CatBoost	0.0326±0.0000 *	0.2328±0.0000 *	0.0394±0.0000 *	0.2998±0.0000 *
DoubleEnsemble	0.0362±0.0005 *	0.2725±0.0036 *	0.0444±0.0004 *	0.3450±0.0038 *
LightGBM	0.0347±0.0000 *	0.2648±0.0000 *	0.0443±0.0000 *	0.3520±0.0000 *
TCN	0.0384±0.0015 *	0.2834±0.0164 *	0.0455±0.0012 *	0.3546±0.0077 *
ALSTM	0.0413±0.0034 *	0.3166±0.0329 *	0.0504±0.0032 *	0.3974±0.0280 *
LSTM	0.0402±0.0030 *	0.3194±0.0271 *	0.0496±0.0027 *	0.4040±0.0212 *
ADD	0.0370±0.0025 *	0.2669±0.0254 *	0.0511±0.0018 *	0.3756±0.0235 *
GRU	0.0417±0.0029 *	0.3284±0.0367 *	0.0510±0.0014 *	0.4137±0.0224 *
AdaRNN	0.0380±0.0117 *	0.2999±0.1022 *	0.0472±0.0095 *	0.3744±0.0974 *
GATs	0.0430±0.0010 *	0.3221±0.0096 *	0.0543±0.0012 *	0.4217±0.0099 *
IGMTF	0.0419±0.0004 *	0.3152±0.0055 *	0.0538±0.0014 *	0.4213±0.0171 *
HIST	0.0437±0.0012 *	0.2952±0.0108 *	0.0581±0.0013 *	0.3912±0.0096 *
Our Method	0.0489±0.0026	0.3593±0.0143	0.0605±0.0023	0.4514±0.0225

Table 1. The results of stock trend prediction on the CSI300 market from 01/01/2017 to 07/10/2022. All the results are averaged after 10 runs, and the standard deviations are shown. * corresponds to statistically significant differences between a baseline and our method ($p < 0.05$ using t-test).

Stock Trend Prediction This task aims to evaluate the ability of models to predict the future stock price trend. For each trading day i , we calculate the 1-day return \hat{r}^i of each stock based on its historical information before date i . For the results, we report the averaged information coefficient (IC), ranked information coefficient (Rank IC), information ratio of IC (ICIR), and information ratio of Rank IC (Rank ICIR). IC^i is the daily IC that measures the Pearson correlation between the predicted ratio \hat{r}^i and the ground-truth ratio r^i :

$$IC^i = \frac{(\hat{r}^i - \text{mean}(\hat{r}^i))^\top (r^i - \text{mean}(r^i))}{n \cdot \text{std}(\hat{r}^i) \cdot \text{std}(r^i)} \quad (14)$$

The IC is calculated for the average of each trading day:

$$IC = \frac{\sum_{i \in [1, D]} IC^i}{D} \quad (15)$$

The ICIR is used to show the stability of IC, which is calculated by dividing IC by its standard deviation:

$$ICIR = \frac{IC}{\text{std}(IC)} \quad (16)$$

For the calculation of Rank IC^i , we first use $R^i = \text{rank}(r^i)$, and $\hat{R}^i = \text{rank}(\hat{r}^i)$ to denote the ranks of the ground-truth and the predicted ratios, respectively:

$$\text{Rank } IC^i = \frac{(\hat{R}^i - \text{mean}(\hat{R}^i))^\top (R^i - \text{mean}(R^i))}{n \cdot \text{std}(\hat{R}^i) \cdot \text{std}(R^i)} \quad (17)$$

Model Name	Annualized Return	Max Drawdown	Information Ratio
Transformer	0.0069±0.0181 *	-0.2131±0.0868 *	0.0753±0.2138 *
TabNet	0.0719±0.0000 *	-0.1139±0.0000	0.8155±0.0000 *
MLP	0.0441±0.0153 *	-0.1512±0.0375 *	0.5163±0.1882 *
Localformer	0.0498±0.0228 *	-0.1268±0.0235	0.6194±0.2843 *
CatBoost	0.0585±0.0013 *	-0.1364±0.0051	0.7270±0.0162 *
DoubleEnsemble	0.0642±0.0112 *	-0.0900±0.0103 *	0.8234±0.1398 *
LightGBM	0.0707±0.0000 *	-0.0835±0.0000 *	0.9487±0.0000 *
TCN	0.0781±0.0203 *	-0.0849±0.0151 *	1.0205±0.2350 *
ALSTM	0.0777±0.0220 *	-0.1031±0.0204	1.0226±0.2859 *
LSTM	0.0826±0.0242 *	-0.0908±0.0132 *	1.0706±0.2771 *
ADD	0.0759±0.0178 *	-0.0939±0.0237	0.9471±0.2101 *
GRU	0.0815±0.0258 *	-0.0917±0.0270 *	1.0826±0.3671
AdaRNN	0.0619±0.0589 *	-0.1392±0.1622	0.8439±0.7172
GATs	0.0886±0.0115 *	-0.1022±0.0184	1.1524±0.1469 *
IGMTF	0.0903±0.0095 *	-0.0986±0.0174	1.1825±0.1035
HIST	0.0854±0.0119 *	-0.0919±0.0152 *	1.0879±0.1504 *
Our Method	0.1063±0.0187	-0.1191±0.0301	1.3315±0.2169

Table 2. The results of trading simulation on the CSI300 market from 01/01/2017 to 07/10/2022. All the results are averaged after 10 runs, and the standard deviations are shown. * corresponds to statistically significant differences between a baseline and our method ($p < 0.05$ using t-test).

The Rank IC and Rank ICIR are calculated similarly as before:

$$\text{Rank IC} = \frac{\sum_{i \in [1, D]} \text{Rank IC}^i}{D} \quad (18)$$

$$\text{Rank ICIR} = \frac{\text{Rank IC}}{\text{std}(\text{Rank IC})} \quad (19)$$

The results of the stock trend prediction task on the test set of the China CSI300 market (01/01/2017 to 07/10/2022) are shown in Table 1. Our method significantly outperforms all the 16 baselines across all four metrics (IC, ICIR, Rank IC, and Rank ICIR) with around 10% enhancement over the second-place model for each metric. These results indicate the importance of modeling expectations and dynamic topics in financial market analysis. It is also interesting to note that the traditional RNN-based methods (such as GRU and LSTM) achieve similar or even better results compared to the models specifically designed for financial analysis (such as ADD, IGMTF, and DoubleEnsemble). This may be attributed to the low signal-to-noise ratio in the financial market since the simpler models may be more robust to noise. These observations further demonstrate the hardness of this task.

Trading Simulation In quantitative investment, "backtesting" refers to applying a trading strategy to historical data, simulating trading, and measuring the return of the strategy. For this task, we employ the top- k dropout strategy

for each method, reporting the **annualized return**⁷ (the geometric average of money earned by an investment strategy each year over a given time period), **max drawdown**⁸ (maximum observed loss from a peak to a trough), and the **information ratio**⁹ (ratio of returns above the returns of the CSI300 benchmark). The top- k dropout strategy is a straightforward quantitative investment approach: for each trading day, we hold k stocks, sell d stocks with the worst predicted 1-day return, and buy d unheld stocks with the best-predicted 1-day return. In our experiments, k is set to 50, and d is set to 5. The trading simulation task results on the test set of the China CSI300 market are displayed in Table 2. Our method surpasses all 16 baselines in annualized return and information ratio. To improve the stability of profitability, future research could explore modifications designed to reduce the max drawdown of our approach.

5 Conclusion

In this paper, we introduce a novel framework for stock trend prediction, suitable for quantitative analysis of financial markets and stock selection. To the best of our knowledge, our method is the first to consider (1) investors’ subjective expectations, and (2) automatically mined dynamic topics that do not require additional knowledge. Through experiments on 16 baselines using the CSI 300 market, we demonstrate that our model achieves a stable annual return above 10%, outperforming all existing baselines and attaining the current state-of-the-art results for stock trend prediction and trading simulation tasks.

Future work could explore modifications to decrease the max drawdown of our method, resulting in more stable profitability. Additionally, since expectations are influenced by external factors such as financial reports or discussions on social media, future research could investigate incorporating this information into our model.

References

1. Arik, S.Ö., Pfister, T.: Tabnet: Attentive interpretable tabular learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 6679–6687 (2021)
2. Bai, S., Kolter, J.Z., Koltun, V.: An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. arXiv preprint arXiv:1803.01271 (2018)
3. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078 (2014)
4. Du, Y., Wang, J., Feng, W., Pan, S., Qin, T., Xu, R., Wang, C.: Adarnn: Adaptive learning and forecasting of time series. In: Proceedings of the 30th ACM International Conference on Information & Knowledge Management. pp. 402–411 (2021)

⁷ <https://www.investopedia.com/terms/a/annualized-rate.asp>

⁸ <https://www.investopedia.com/terms/m/maximum-drawdown-mdd.asp>

⁹ <https://www.investopedia.com/terms/i/informationratio.asp>

5. Fama, E.F.: The behavior of stock-market prices. *The journal of Business* **38**(1), 34–105 (1965)
6. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9**(8), 1735–1780 (1997)
7. Jiang, J., Kim, J.B., Luo, Y., Zhang, K., Kim, S.: Adamct: Adaptive mixture of cnn-transformer for sequential recommendation. *arXiv preprint arXiv:2205.08776* (2022)
8. Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., Liu, T.Y.: Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems* **30** (2017)
9. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
10. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016)
11. Poterba, J.M., Summers, L.H.: Mean reversion in stock prices: Evidence and implications. *Journal of financial economics* **22**(1), 27–59 (1988)
12. Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A.V., Gulin, A.: Catboost: unbiased boosting with categorical features. *Advances in neural information processing systems* **31** (2018)
13. Qin, Y., Song, D., Chen, H., Cheng, W., Jiang, G., Cottrell, G.: A dual-stage attention-based recurrent neural network for time series prediction. *arXiv preprint arXiv:1704.02971* (2017)
14. Sawhney, R., Agarwal, S., Wadhwa, A., Derr, T., Shah, R.R.: Stock selection via spatiotemporal hypergraph attention network: A learning to rank approach. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 35, pp. 497–504 (2021)
15. Sawhney, R., Agarwal, S., Wadhwa, A., Shah, R.: Deep attentive learning for stock movement prediction from social media text and company correlations. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. pp. 8415–8426 (2020)
16. Tang, H., Wu, L., Liu, W., Bian, J.: Add: Augmented disentanglement distillation framework for improving stock trend forecasting. *arXiv preprint arXiv:2012.06289* (2020)
17. Tanimoto, T.T.: *Elementary mathematical theory of classification and prediction* (1958)
18. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
19. Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y.: Graph attention networks. *arXiv preprint arXiv:1710.10903* (2017)
20. Xu, W., Liu, W., Bian, J., Yin, J., Liu, T.Y.: Instance-wise graph-based framework for multivariate time series forecasting. *arXiv preprint arXiv:2109.06489* (2021)
21. Xu, W., Liu, W., Wang, L., Xia, Y., Bian, J., Yin, J., Liu, T.Y.: Hist: A graph-based framework for stock trend forecasting via mining concept-oriented shared information. *arXiv preprint arXiv:2110.13716* (2021)
22. Yang, X., Liu, W., Zhou, D., Bian, J., Liu, T.Y.: Qlib: An ai-oriented quantitative investment platform. *arXiv preprint arXiv:2009.11189* (2020)
23. Zhang, C., Li, Y., Chen, X., Jin, Y., Tang, P., Li, J.: Doubleensemble: A new ensemble method based on sample reweighting and feature selection for financial data analysis. In: *2020 IEEE International Conference on Data Mining (ICDM)*. pp. 781–790. IEEE (2020)