

Research on Stock Index Prediction Based on Stock Correlation Network and Deep Learning

Xueyan Li^{1,a}

¹College of Information Engineering, Nanjing University of Finance and Economics, Nanjing, China
^a350696343@qq.com

Abstract: Aiming at the problem of stock index prediction, constructing a time series stock correlation network based on the fundamentals and technology of stock index components, then using the depth map neural network to learn the hierarchical representation of stock correlation network, and obtaining the candidate prediction signal in an end-to-end way. The architecture composed of depth map neural network method and end-to-end strategy is called DIFFPOOL architecture. Taking the CSI 300 index as the research object, combining the DIFFPOOL architecture with softmax classifier, long-term and short-term memory neural network (LSTM), linear regression, and logical regression, respectively, and uses the sliding time window method to obtain the corresponding prediction accuracy of stock index. The accuracy of the combined model under the optimal parameters fluctuates in the interval [0.56, 0.62]. Ultimately, the first mock exam is based on the mean absolute error (MAE) and the root mean square error (RMSE). The regression models of DIFFPOOL and regression models are compared with LSTM, recurrent neural network (RNN), and back propagation neural network (BP). Compared with the single model, the MAE and RMSE of the combined model are smaller, 0.0061 and 0.0081, respectively. Experiments show that by aggregating the node attribute information of the stock association network hierarchically, we can dynamically capture the impact of different industry sectors on stock index price fluctuations and further improve prediction accuracy.

Keywords: stock index prediction; stock correlation network; sliding time window; graph neural network; regression prediction

1. Introduction

The stock market is a complex system with high noise, nonlinearity, non-stationarity, and chaos. The stock index time series is a comprehensive external manifestation of the inherent complexity of the stock market, which provides an essential reference for investors to formulate investment strategies. Therefore, accurate forecasting of stock indices is conducive to better monitoring and managing financial markets that are highly correlated with the stock market and provide practical guidance for investors' investment decisions.

Since Mantegna^[1] first discovered the stock market hierarchy by constructing a stock correlation network and then analyzing its correlation with the return of the S&P 500 index, more and more scholars at home and abroad have opened it. The complex network theory method was used to model and analyze the correlation between price fluctuations of stock index constituents. The construction of the stock correlation network is to treat the constituent stocks of the stock index as nodes and the price fluctuation correlation between the constituents as the edge. Considering that fully connected stock-linked networks contain a lot of noise information, to simplify the network and retain the critical information of the network, the Minimum Spanning Tree (MST)^[2] and Threshold Model (TM)^[3] are commonly used in previous studies or Planar Maximally Filtered Graph, PMFG^[4] to filter edges with weak correlation. The research in the field of stock-linked networks is mainly divided into static network research and time series network research. A stock correlation network is a category of the graph. The stock correlation network generated based on a specific time window has an implied graph label (such as the movement direction of the stock index at the next moment). The stock index prediction problem can also be formalized as a graph classification problem. Graph Classification, as an essential graph mining task, aims at classifying and predicting different types of graph data. Ying et al.^[5] proposed DIFFPOOL, a differentiable pooling module for graph classification, which can generate hierarchical representations of graphs and is integrated with various GNN models in an end-to-end manner. The empirical results show that the average accuracy of this method can be improved by 5-10% in graph classification tasks.

Although the existing graph classification research has achieved remarkable results in social network analysis, computer perspective, and compound identification, few studies have classified temporal stock correlation networks. Based on the existing graph classification algorithm, this paper applies the deep learning framework to the feature extraction of the temporal stock correlation network and the learning of candidate stock index prediction signals.

2. The Construction of Stock Correlation Network

In the stock correlation network, nodes represent constituent stocks, and the linkage between nodes represents the correlation between stocks. The stock correlation network constructed in this paper is an undirected non-entitlement network. The basic construction of the network consists of the following four steps: (1) Calculate the linkage set of constituent stocks according to the closing price sequence of the constituent stocks in the specified stock index; (2) Construct the correlation coefficient matrix between stocks; (3) The correlation coefficient matrix is used to calculate the distance matrix in the corresponding period. (4) based on the distance matrix, the minimum spanning tree (MST) or plane maximum filter graph (PMFG) algorithm is used to construct the final stock correlation network. The details are as follows:

In the time range $[t-l+1, t]$, select the closing price i sequence $C_i^{t,l} = \{c_i^{t-l+1}, \dots, c_i^t\}$ of any component stock in a specific stock index, where l represents the length of the sliding time window, as shown in Figure 1. The closing price is compared to the height of a person. It can be seen from the figure that node (m, c_i^m) is blocked by node (k, c_i^k) from seeing node (n, c_i^n) ; then there is no line between node (m, c_i^m) and node (n, c_i^n) ; otherwise, if node k is not blocked from seeing, There's a line between them.

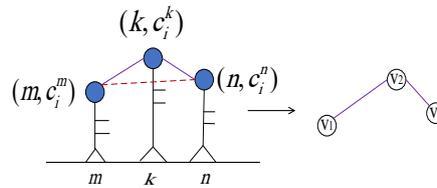


Figure 1: The visibility graph model (solid line represents connection and dotted line represents non connection).

The edge set of this component i in the time range $[t-l+1, t]$ is expressed as:

$$E_i^{t,l} = \{e_{mn}^{t,l} \mid m < k < n, c_i^n > c_i^k - (c_i^m - c_i^k) \frac{n-k}{n-m}\} \quad (1)$$

On this basis, in a specific time window $[t-l+1, t]$, the stock volatility correlation between any two component stocks i and j can be measured by the Jaccard distance of their side sets, namely:

$$d_{i,j}^{t,l} = 1 - Jaccard(E_i^{t,l}, E_j^{t,l}) = 1 - \frac{|E_i^{t,l} \cap E_j^{t,l}|}{|E_i^{t,l} \cup E_j^{t,l}|} \quad (2)$$

Because the network corresponding to the Jaccard distance matrix is a fully connected graph, the minimum spanning tree or planar pole filtering graph algorithm should be used to retain valuable edges in the network while filtering some redundant edges to form the final stock correlation network.

Table 1: Definition of constituent indicators.

| | |
|---|--|
| $ROC_i^{t,l} = (c_i^t - c_i^{t-l+1}) / c_i^{t-l+1}$ | $EMA_i^{t,l} = 2(c_i^t - EMA_i^{t-1,l}) / (l+1) + EMA_i^{t-1,l}$ |
| $RSI_i^{t,l} = 100 - 100 / (1 + RS_i^{t,l})$ | $RSI_i^{t,l} = 100 - \frac{100}{1 + RS_i^{t,l}}$ |

Finally, this paper selects four commonly used fundamental indicators (as shown in Table 1) to describe the running state of constituent stocks in time window $[t-l+1, t]$, and then constructs the

attribute graph $(G^{t,l}) = \{A^{t,l}, F^{t,l}\}$ of the stock association network in this window, where $A^{t,l} \in (0,1)^{n \times n}$ represents the adjacency matrix of the stock association network, and $F^{t,l} \in \mathbb{R}^{n \times w}$ represents the attribute matrix of the stock association network. Each row represents the fundamentals of component stock i at time t . Meanwhile, the stock correlation network $(G^{t,l})$ has an implied graph label ψ^t information, that is, the movement direction of the stock index at time $t+1$ relative to time t . Let Y_I^t represent the index price of stock index I at t , if $Y_I^t < Y_I^{t+1}$, then $\psi^t = 1$; On the contrary, $\psi^t = 0$.

3. Hierarchical Representation and Graph Classification Algorithm of Stock Correlation Network

The existing GNN model only transmits information through the edges of the graph and cannot infer and aggregate information in a hierarchical manner, and embeds all nodes together for global pooling, ignoring any hierarchical structure that may exist in the graph. Ying et al.^[5] proposed DIFFPOOL, a micropoolable operation module whose core idea is to delaminate aggregate graph nodes through a differential module, to construct a deep and multilayer GNN^[6] model. DIFFPOOL maps nodes to collections of clusters based on learned embeddedness and learns a differentiable soft distribution at each level of GNN assignment. By hierarchically stacking multiple GNN layers, deep GNNS can be established for graph sorting tasks. The overall architecture of DIFFPOOL is shown in Figure 2. The message transmission mode of Graph Convolution Network (GCN)^[7] is used in the single-layer GNN module to learn useful Graph classification representation. Given a stock correlation network $(G^{t,l}, \psi^t)$, the initial node embedding matrix is $H^{(0)} = F^{t,l}$, and the message propagation function of GCN can be expressed as:

$$H^{(k)} = ReLU(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} H^{(k-1)} W^{(k-1)} H^{(k-1)}) \quad (3)$$

where \tilde{A} is the adjacency matrix of the stock correlation network with self-loop, \tilde{D} is the degree matrix, and $W^{(k)} \in \mathbb{R}^{w \times w}$ is the trainable weight matrix. The single-layer GCN module can generate the final output node embedding matrix $M = H^{(k)}$ by iterating k times.

Because the DIFFPOOL architecture does not require knowledge of the message passing function of the single-layer GNN module, the outgoing direct extraction of the single-layer GCN module through k iterations can be denoted as $M = GCN(A^{t,l}, F^{t,l})$. DIFFPOOL's goal is to define a general, end-to-end, differentiable strategy. Enables multiple GNN models to be stacked hierarchically. Given the output $M = GCN(A^{t,l}, F^{t,l})$ of a GCN module and the adjacency matrix $A^{t,l} \in (0,1)^{m \times m}$ of the stock correlation network, DIFFPOOL tries to find a way to get a new containing $M' \in \mathbb{R}^{m \times l}$. The coarse graph can be used as the input to the next layer of GCN. Repeat d times to get a series of coarse graphs. DIFFPOOL aims to learn a pooling strategy, which can be generalized to graphs with different nodes and edges and adapted to different graph structures during inference.

The membership matrix of the class cluster learned at the d layer is defined as $S^{(d)} \in \mathbb{R}^{n_{d+1} \times n_d}$ ($n_d > n_{d+1}$), n_d represents the number of nodes at the d layer, n_{d+1} represents the number of nodes at the $d+1$ layer, and $S^{(d)}$ can be understood as the probability from each node at the d layer to each node (class cluster) at the $d+1$ layer. The figure coarsening method is based on $(A^{(d+1)}, V^{(d+1)}) = DIFFPOOL(A^{(d)}, M^{(d)})$, where $A^{(d+1)}$ represents the connection strength between the class clusters at layer $d+1$. The following formula is used to update:

$$A^{(d+1)} = S^{(d)T} A^{(d)} S^{(d)} \in \mathbb{R}^{n_{d+1} \times n_{d+1}} \quad (4)$$

$V^{(d+1)}$ represents the embedded matrix of the class clusters at layer $d+1$, and its update strategy is as follows:

$$\mathbf{V}^{(d+1)} = \mathbf{S}^{(d)\top} \mathbf{M}^{(d)} \in \mathbb{R}^{n_{d+1} \times l} \quad (5)$$

DIFFPOOL adopts two different GNN modules to learn the class cluster membership matrix of each layer as $\mathbf{S}^{(d)}$ and node embedding matrix $\mathbf{M}^{(d)}$:

$$\mathbf{M}^{(d)} = GNN_{d,embed}(\mathbf{A}^{(d)}, \mathbf{V}^{(d)}) \quad (6)$$

$$\mathbf{S}^{(d)} = \text{soft max}\left(GNN_{d,pool}\left(\mathbf{A}^{(d)}, \mathbf{V}^{(d)}\right)\right) \quad (7)$$

Where $GNN_{d,pool}$ represents pooled GNN at the d layer, whose output dimension corresponds to the maximum number of predefined class clusters at the d layer, at the penultimate layer of the DIFFPOOL framework, $\mathbf{S}^{(D-1)}$ is set as an all-1 vector; all nodes are divided into the same class cluster at the last layer d . Thus, DIFFPOOL generates an embedded vector $\mathbf{X}^{(l)} \in \mathbb{R}^l$ that corresponds to the original stock association network. The whole system can be trained end to end using stochastic gradient descent. Due to training $GNN_{d,pool}$, in essence, belongs to a nonconvex optimization problem, drawing lessons from DIFFPOOL thought, the use of the auxiliary link to predict target training pooling Geri weis-corbley, namely in the first d layer to minimize the loss function is as follows:

$$L_L P^{(d)} = \|\mathbf{A}^{(d)}, \mathbf{S}^{(d)} \mathbf{S}^{(d)\top}\|_F \quad (8)$$

Where, $\|\cdot\|_F$ stands for Frobenius norm. Considering that the node's class cluster membership should be close to the one-hot vector, the relationship between each class cluster or subgraph is clearly defined. Therefore, at each layer, the following loss function is introduced to standardize the entropy of the membership matrix of class clusters:

$$L_E = \frac{1}{n} \sum_{i=1}^n H\left(\mathbf{S}_i^{(d)}\right) \quad (9)$$

Where H is the entropy function and $\mathbf{S}_i^{(d)}$ is the i line of $\mathbf{S}^{(d)}$.

Compared with the traditional stock index prediction method, which extracted the prediction factors directly from the original stock index time series, this study extracted the nonlinear characteristics of the stock correlation network hierarchically based on the micropoolable framework, effectively aggregated the technical indicators of different class clusters, and captured the interaction of different class clusters based on the hierarchical clustering results of stock correlation network. Then automatically learn the prediction variables suitable for the short-term trend prediction of the stock index.

4. Experimental results and analysis

4.1. Experimental Data

The CSI 300 Index is a composite index jointly launched by the Shanghai and Shenzhen Stock Exchanges on April 8, 2005, to reflect the overall trend of China's A-share market. Because its industry market coverage is relatively high, and the main component weight is relatively dispersed, it can effectively prevent the possible stock index operation behavior in the financial market. Therefore, choosing the CSI 300 index as the empirical research object of the stock index time series is of typical significance.

The component stocks of the CSI 300 Index are adjusted every six months. To ensure the steadiness of the 300 component stocks and the reliability of the experimental results, this study selects the component stocks of the CSI 300 Index from December 14, 2020, to June 14, 2021, as the research object. When constructing the stock correlation network, the length of the sliding window is set to 15. The Jaccard distance of the corresponding edge set measures the stock price fluctuation correlation between any two component stocks. Finally, the distance matrix is trimmed by the plane maximum filter graph algorithm (PMFG) to generate the stock correlation network. Figure 2 shows the stock correlation network generated on January 25, 2021. The whole network contains several clusters. The connections among stocks in each cluster are relatively dense, while the connections among each cluster are relatively

sparse.

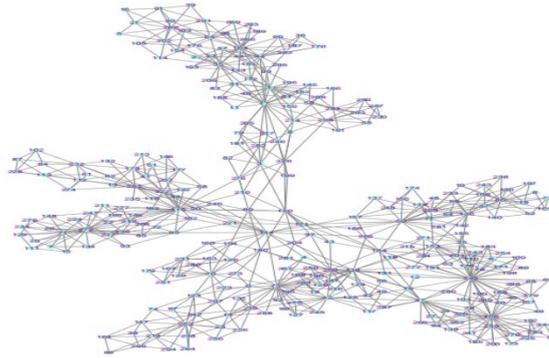


Figure 2: Shanghai Shenzhen 300 Index stock correlation Network.

4.2. Analysis of Stock Correlation Network

In this paper, nodes with degree greater than 20 are selected to analyze the feasibility of constructing the stock correlation network method, as shown in Table 2.

Table 2: Stocks with point degree greater than 20 and their industries.

| id | Node degree | stock | industry |
|-----|-------------|-------------------------|------------------------------|
| 1 | 21 | Ping An Bank | finance |
| 17 | 24 | Tianmao Group | Chemical raw material |
| 47 | 21 | Yunda Stock | Electronic information |
| 75 | 26 | Rongsheng Petrochemical | fossilization |
| 106 | 30 | Xinwei communication | communication |
| 115 | 27 | Kangtai biology | Pharmaceutical manufacturing |
| 138 | 23 | Oriental wealth | securities |
| 249 | 35 | China Life Insurance | securities |
| 272 | 24 | Bank of China | finance |

Table 2 shows that the maximum node degree in the stock correlation network is 35, and the listed stock nodes are all in the center of the stock correlation network. These nine stocks belong to the core industries of the Chinese economy. Finance, securities, and pharmaceutical manufacturing are the lifeblood of the Chinese economy.

4.3. Experimental Setup and Analysis

First, the stock index prediction problem is regarded as a supervised graph classification problem to determine the optimal parameters of the model. Due to the small number of nodes in the stock correlation network, the number of iterations of GCN and pooling layers in the DIFFPOOL model is set to 2, according to the suggestion in the literature. In the experiment, SHSZ stock index samples from December 14, 2020, to June 14, 2021 were selected to create a hyper reference network: test set size:[5,6,7,8,9,10,11,12], the training set size:[20,25,30,35,40,45,50,55], hidden layer dimension: [20,40,60,80,100,120,140,160], node coarsening rate:[0.1,0.15,0.2,0.25,0.3,0.35,0.4,0.45], number of iterations :[5,10,15,20,25,30,35]. It is not difficult to find that in the dynamic training stock correlation network hierarchical representation model, the model's investment win rate increases first and then decreases with the increase of the test set and training set size. When hidden layer dimensions continue to increase, the prediction model of investment conferences shows a relatively stable upward trend, and when the parameters of over 140, investment promotion is more gentle. The influence of node coarsening rate on the model prediction performance is further analyzed. This parameter directly affects the resolution of the hierarchical clustering cluster structure of the stock correlation network: a low coarsening rate will produce a large number of small clusters. Otherwise, a small number of large class clusters will be produced. Based on the parameter analysis results in Figure 3 (d), the coarsening node rate is set as 30% in this paper. Finally, by observing Figure 4 (e), it can be found that the increase in model iteration times cannot significantly improve the model's prediction accuracy, and a higher number of iterations may cause an overfitting phenomenon.

To sum up, in building the stock correlation network, the length of the sliding window is set to 15,

the share price volatility correlation between any two components with the corresponding visual network edge set Jaccard similarity measure and finally uses, the plane great filter (PMFG) methods to cut the similarity matrix, generate stock correlation network. When training the stock correlation network hierarchical representation model, the GNN model for DIFFPOOL is based on the graph convolutional neural network GCN architecture, and its iteration number is set to 2. One DIFFPOOL layer is applied after every two GCN layers. A total of two DIFFPOOL layers were used on each stock index dataset. At the same time, two different sparse GraphSAGE models^[8] were used to calculate the cluster membership matrix of each layer as S^d and node embedding matrix Z^d . Outside the sample TEST process, the dynamic training mode to the DIFFPOOL model to generate candidate index prediction signal: test set size set to 8, the training set size train set to 30, hidden layer dimensions set to 140, node coarsening rate set to 30%, the number of iterations is set to 5.

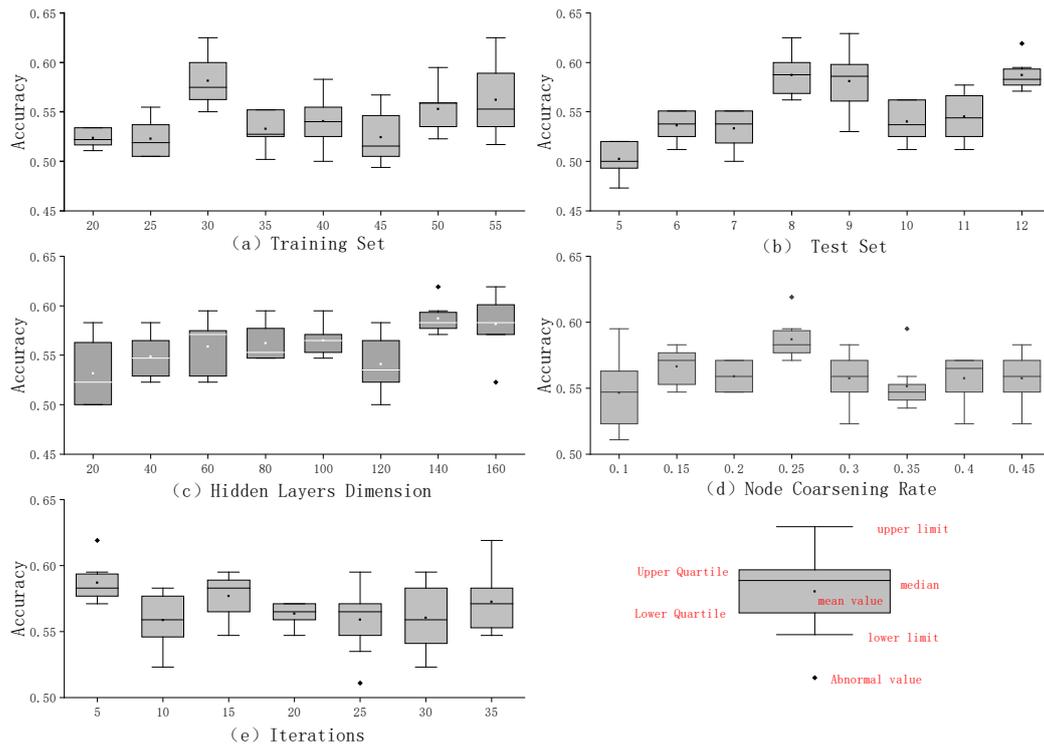


Figure 3: Parameter analysis of DIFFPOOL architecture and softmax classifier combination.

The above experiments combine DIFFPOOL architecture and SOFTMAX classifier for prediction research. In this paper, linear regression, logistic regression, LSTM long and short-term neural networks^[9], and DIFFPOOL architecture are combined, respectively, to reflect the feasibility of combining stock correlation network and deep learning. Specific combinations are as follows:

- 1) DP-LNR: Stock correlation network regression model based on DIFFPOOL and linear regression.
- 2) DP-LGR: Stock correlation network regression model based on DIFFPOOL and logistic regression.
- 3) DP-SFTM: Stock association, network classification model, based on DIFFPOOL and SOFTMAX classifier.
- 4) DP-LSTM: Stock association network classification model based on DIFFPOOL and LSTM Long and Short Term memory network.

The prediction accuracy of the DIFFPOOL architecture and regression model is based on the same sign of the actual stock index return rate and predicted stock index return rate. The same sign represents the consistency of the actual stock index return rate and predicted stock index return rate, while the different sign represents their inconsistency. The prediction accuracy of the DIFFPOOL architecture and classification model is based on the classification prediction of stock association network tags. The experiment conducted prediction analysis based on the above parameter Settings. All the same, parameters were evaluated eight times, and the maximum accuracy was selected as the value under the parameter setting. The results are shown in Figure 4.

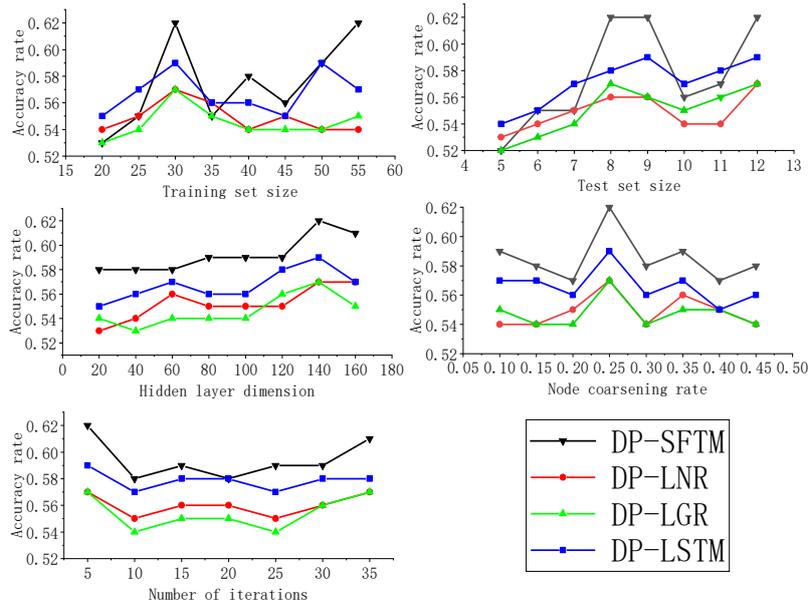


Figure 4: Combined prediction of diffpool framework and four models.

It can be seen from Figure 4 that compared with DP-LNR and DP-LGR, DP-LSTM and DP-SFTM have a better effect on stock index prediction, so the combined prediction effect of DIFFPOOL architecture and classification model is better than that of DIFFPOOL architecture and regression model.

To test the prediction effect of DIFFPOOL architecture and regression model combination, this paper studied the return rate of the CSI 300 index and selected LSTM, RNN^[10], BP^[11], SVR^[12], and other typical deep learning, neural network, machine learning and time series models to compare and analyze with DP-LNR. The parameter values of these five reference models are optimized and determined on the same training data set, then applied to the subsequent test data set. In this experiment, the training set size is 74, and the test set size is 30. The experimental results are shown in Figure 5. The date corresponding to the time in the figure is from April 28, 2021, to June 11, 2021.

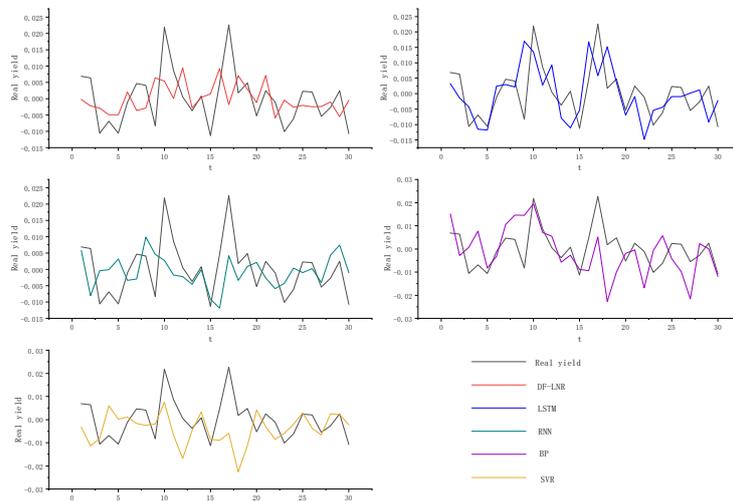


Figure 5: Regression prediction of each model.

In order to enable DP-LNR to be compared with other models, Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) are selected as unified and objective evaluation criteria in this paper.

The comparison results are shown in Table 3.

Table 3: Comparison and effect analysis of different models.

| Evaluation index | DF-LNR | LSTM | RNN | BP | SVR |
|------------------|--------|------|------|------|------|
| MAE/% | 0.61 | 0.68 | 0.68 | 0.70 | 0.86 |
| RMSE/% | 0.81 | 0.85 | 0.87 | 0.87 | 1.08 |

As seen from Table 3, DF-LNR has the best performance in both mean absolute error and root mean square error compared with other models. Based on RNN neural network, the LSTM model increases the filtering of past states, and the accuracy of RMSE is slightly higher than that of RNN prediction. The accuracy of MAE and RMSE is slightly higher than BP because RNN considers the order relation in time. Compared with the neural network method, the machine learning SVR model and the traditional time series ARIMA model is insufficient in the prediction effect.

5. Conclusions

This paper takes the stocks in the Shanghai and Shenzhen 300 Index as the research object, constructs a temporal stock association network based on a viewable graph and planar maximum filtered graph, formalizes the stock index prediction problem into a graph classification problem-oriented to the stock association network, and then introduces a depth map neural network to learn the hierarchical representation of the stock association network, and generates candidate prediction signals in an end-to-end manner. The DIFFPOOL architecture is combined with the SOFTMAX classifier LSTM short-term and short-term neural memory networks, linear regression, and logical regression are combined to obtain the corresponding accuracy of stock index prediction. Finally, this paper compares the DIFFPOOL framework and regression model combination with regression models such as LSTM, RNN, and BP. The experiment shows that the combined model performs best on MAE and MRSE compared to other regression models, fully demonstrating the method's effectiveness in predicting stock indexes.

Acknowledgements

This article is supported by the Jiangsu Postgraduate Research and Innovation Program (No. KYCX21_1533).

References

- [1] Mantegna R N. Hierarchical Structure in Financial Markets [J]. *The European Physical Journal B*, 1999, 11(1): 193-197.
- [2] Miccich S, Bonanno G, Mantegna R N, et al. Degree stability of a minimum spanning tree of price return and volatility[J]. *Physica A: Statistical Mechanics and its Applications*, 2003, 324(1-2):66-73.
- [3] Namaki A, Shirazi A H, Raei R, et al. Network analysis of a financial market based on genuine correlation and threshold method[J]. *Physica A: Statistical Mechanics and its Applications*, 2011, 390 21-22: 3835-3841.
- [4] Tumminello M, Aste T, Di Matteo T, et al. A tool for filtering information in complex systems[J]. *Proceedings of the National Academy of Sciences*, 2005, 102 30: 10421-10426.
- [5] Ying R, You J, Morris C, et al. Hierarchical graph representation learning with differentiable pooling [C]. *Annual Conference on Neural Information Processing Systems NeurIPS*, 2018, 4805-4815.
- [6] Wu Z, Pan S, Chen F, et al. A Comprehensive Survey on Graph Neural Networks [J]. *IEEE Transactions on Neural Networks and Learning Systems*, 2021, 32(1): 4-24.
- [7] Zhou M, Xu W, Zhang W, et al. Leverage knowledge graph and GCN for fine-grained-level clickbait detection[J]. *World Wide Web*, 2022, 25(3):1243-1258.
- [8] Hamilton W L, Ying R, Leskovec J. Inductive representation learning on large graphs[C]. *Annual Conference on Neural Information Processing Systems NeurIPS*, 2017, 1024-1034
- [9] Li Q, Tan J, Wang J, et al. A Multimodal Event-driven LSTM Model for Stock Prediction Using Online News [J]. *IEEE Transactions on Knowledge and Data Engineering*, 2020, PP (99):1-1.
- [10] Dong-Ha, Shin, Kwang-Ho, Choi, et al. Deep Learning Model for Prediction Rate Improvement of Stock Price Using RNN and LSTM[J]. *The Journal of Korean Institute of Information Technology*, 2017, 15(10):9-16.
- [11] Zhang D, Lou S. The application research of neural network and BP algorithm in stock price pattern classification and prediction – ScienceDirect [J]. *Future Generation Computer Systems*, 2021, 115:872-879.
- [12] Ghanbari M, Goldani M. Support Vector Regression Parameters Optimization using Golden Sine Algorithm and its application in stock market [J]. 2021.