# Railway Freight Volume Prediction Based on Spatiotemporal Graph Convolutional Neural Network

Huan Dai[1,a,*]

[1]School of Transportation and Logistics, Southwest Jiaotong University, Chengdu, 610031, China
[a]daihuan4550@163.com
*Corresponding author

**Abstract:** *Accurate railway freight volume prediction can effectively support the dynamic adjustment of freight station organization, enhance the service level and competitiveness of railway transportation. Considering the impact of railway network relationships on freight volume, this paper proposes a railway freight volume prediction method based on Spatiotemporal Graph Convolutional Neural Network (STGNN). The spatial convolution module adopts relational graph convolution to explore and integrate spatial characteristics of the railway physical network, inter-station relationships based on freight volume, and service relationships based on operational plans. In the time series module, a multi-layer Gated Recurrent Unit (GRU) is used for multi-step freight volume prediction of freight station groups. Using freight volume data along the Xiang-Yue section as the research object, the prediction results under different step sizes are compared, showing that STGNN significantly outperforms the baseline models.*

**Keywords:** *Railway freight volume prediction; Relational graph convolution; Spatiotemporal graph convolutional neural network*

## 1. Introduction

With rapid economic development and the continuous expansion of the railway network, China's railway freight volume has shown a generally upward trend year by year. Accurate railway freight volume prediction is crucial for railway transportation management, resource optimization, and efficiency improvement. Freight volume prediction not only uses historical data to forecast future trends but also accurately captures dynamic changes in freight demand in an uncertain market environment, guiding actual transportation planning and decision-making.

In recent years, with the development and application of big data technology, the collection and storage of freight volume data have become more convenient, enabling the use of machine learning and deep learning in freight volume research. Researchers such as Wu Wei[1], Xie Jianwen[2], and Wang Xifu[3] have built BP neural networks for volume prediction. Zhang Guandong proposed a prediction method based on multidimensional long short-term memory networks(LSTM), forecasting freight volumes for railways, highways, and civil aviation[4].

The models mentioned above mainly focus on mining temporal sequence features and do not consider the spatial correlations in traffic flow prediction problems. The combination of graph convolutional networks (GCN) and recurrent neural networks (RNN) has shown strong capabilities in handling and predicting spatio-temporal data, accurately capturing spatial topologies and predicting complex graph-structured data. From a broad perspective, traffic data is a type of graph-structured data, making GCNs more suitable for calculating traffic data than traditional convolutional neural networks. Liu Qidong and others proposed a time-aware Transformer model for passenger flow prediction[5]. Xu Li et al.[6] built a GC-STGCN model that combines GCN and gated recurrent neural networks for traffic flow prediction. Li[7] and other researchers introduced DCRNN to integrate spatial and temporal correlations, achieving accurate traffic flow prediction. It can be seen that spatio-temporal correlation prediction models have achieved good results in the field of traffic flow.

This paper comprehensively considers the freight conditions of railway stations and the relationships among them, introducing GCN to construct a spatio-temporal correlation prediction model based on multiple time series and spatial relationships of railway station groups. Each freight station on the railway

line is treated as a node, with historical freight time series data and spatial characteristics among the stations studied. The STGNN model is used for short-term railway freight volume prediction. The experiment uses railway freight volume data from the Xiang-Yue section, and the results show that the model's prediction accuracy has improved compared to the benchmark models, demonstrating the effectiveness of the proposed prediction method.

## 2. Model Construction

### 2.1 Railway Network Spatial Relationships

In actual transportation organization, railway transportation differs from road transportation. Road transportation is a continuous network with fluid characteristics, where the correlation of traffic flow between adjacent roads is relatively high and can be represented by an adjacency matrix of physical roads. Railway transportation, on the other hand, is based on a physical network, with the transportation time and route determined by the train operation plan. Additionally, the interconnection between different stations in the railway network affects the flow of freight. Therefore, when analyzing the relationships in the railway network, it is necessary to consider both the connectivity structure between physical network nodes and more non-Euclidean relational matrices, fully reflecting the spatial correlations in railway freight. The Relational Graph Convolutional Network (R-GCN) method is used to extract spatial characteristics from the railway service network and station correlations. Three relationships in the railway network are analyzed: physical network relationship, service network relationship, and station relationship.

#### 2.1.1 Physical Network Relationships

Railway physical network $G$, set of freight stations $S$, freight station $i$、$j$, set of segments $e$, segment, the relationship of the railway physical network is as follows:

$$G = (S, E) \tag{1}$$

$$i、j \in S = \{0, \cdots, i, \cdots, |S|\} \tag{2}$$

$$e = (i, j) \in E \tag{3}$$

$$e_{i,j} = \begin{cases} 1 & (i,j) \in E \\ 0 & others \end{cases} \tag{4}$$

After normalization, the railway physical network is as follows:

$$R_{i,j}^p = \frac{e_{i,j}}{\sum\limits_{j \in S} e_{i,j}} \tag{5}$$

#### 2.1.2 Service Network Relationships

In the railway network, some freight stations are not physically connected directly, but they are directly linked in the service network due to train operation plans. The correlation of freight volumes among stations in the railway network does not solely depend on spatial distance but is also closely related to the operation plans, meaning that even freight stations that are physically distant may have strong correlations. Therefore, based on the physical relationships between freight stations, service relationships are comprehensively considered. The train operation plan is broken down into service arcs to construct a service network relationship between freight stations, fully exploring spatial correlations.

The set of service arcs between freight stations is $A$, and a service arc $a_{i,j}$ can be defined as:

$$a_{i,j} = \begin{cases} 1 & (i,j) \in A \\ 0 & others \end{cases} \tag{6}$$

After normalization, the railway service network is as follows:

$$R_{i,j}^s = \frac{a_{i,j}}{\sum\limits_{j \in S} a_{i,j}}$$

$$(7)$$

### 2.1.3 Station Relationships

To a certain extent, the cargo volume between railway freight stations can also reflect the correlation between them. The greater the cargo volume between stations, the stronger the correlation. The cargo volume between freight stations is denoted as $d_{i,j}$. If $d_{i,j}$ is within the top one-quarter of all inter-station cargo volumes, the value is retained; otherwise, it is set to 0.

$$d_{i,j} = \begin{cases} d_{i,j} .... & d_{i,j} \ is \ in \ the \ first \ quarter \ of \ the \ cargo \ flow \\ 0 & others \end{cases}$$

$$(8)$$

After normalization processing, the station relationship network matrix is:

$$R_{i,j}^{st} = \frac{d_{i,j}}{\sum\limits_{j \in S} d_{i,j}}$$

$$(9)$$

### 2.2 Model Structure

The STGNN model is used for railway freight volume prediction, aiming to extract spatial correlations between freight stations through a spatial convolution module. Since the railway network relationships are of multiple types, the Relational Graph Convolutional Network (R-GCN) is introduced in the spatial convolution layer to achieve multi-relation graph integration. In the time series module, a Recurrent Neural Network (RNN) is employed for freight volume prediction.

(1) Spatial Convolution Module

This module considers different types of connections between freight stations, including physical network relationships, service network relationships, and station relationships. The Relational Graph Convolutional Network (R-GCN) processes the three network relationship matrices to extract and integrate spatial features, which are then fed into the time series module.

(2) Time Series Module

In traffic flow prediction research, the commonly used time series prediction models are Recurrent Neural Networks (RNN), which can learn implicit sequential relationships in historical data. However, they may suffer from gradient explosion or vanishing problems during computation. Therefore, improved time series prediction models with higher accuracy have been developed based on the RNN structure. The time series module in this paper adopts the Gated Recurrent Unit (GRU) to predict the freight volume of railway station groups.

(3) Overall Model Structure

The STGNN model consists of spatial convolution and time series modules. In the spatial convolution phase, the Relational Graph Convolution (R-GCN) processes various spatial relationships in the railway network at each time step, extracting spatial correlation features. The outputs of the graph convolution at each time step are concatenated into time series data and input into the time series model to capture the temporal characteristics of freight volume at each station, ultimately predicting the freight volume at each station through a fully connected layer. The STGNN model structure is shown in Figure 1.
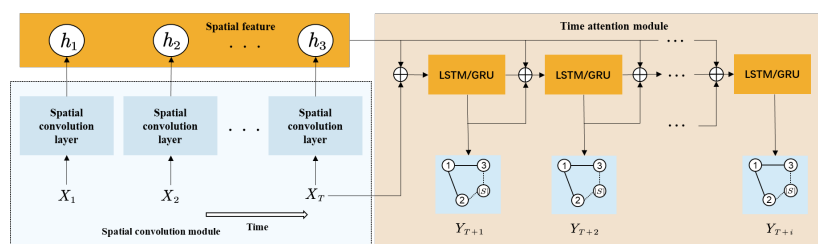


*Figure 1: STGNN framework structure.*

## 3. Case Study and Results Analysis

### 3.1 Data Preprocessing

The original dataset used in this paper comes from the Guangzhou Railway Bureau's Xiang-Yue section, including freight data from several railway stations between 2011 and 2017. The data fields include shipment date, billing weight, freight rate, etc.

(1) Removal of Duplicates and Invalid Tickets

Each freight ticket record has a unique ticket ID and an invalid ticket identifier. During the cleaning of the original data, a large number of duplicate ticket records were found. Therefore, only one record was retained for records with the same ID and attribute values, while invalid tickets were removed based on the identifier.

(2) Freight Station Selection

When analyzing the original data of freight stations, it was found that the density of recorded stations was high, with many small stations having low freight volume and discontinuous dates. To facilitate subsequent spatiotemporal analysis and training of the STGNN model, the freight volume of all stations within a prefecture-level city was consolidated at the station with the highest volume.

(3) Extraction of Valid Freight Information

The original data contains many attributes, but including all attributes in model training would result in high computational cost and potentially affect prediction accuracy. Therefore, only attributes strongly related to freight volume prediction, such as date and billing weight, were selected and aggregated by date to obtain the daily total freight volume at each station.

(4) Missing Value Imputation

Since the original freight data contains a small number of missing values, it is necessary to impute the missing data. Given that freight volume changes are relatively stable in the short term, linear interpolation is used for missing value imputation. The formula of linear interpolation is as follows:

$$x_i = \begin{cases} x_{i+1} - (x_{i+2} - x_{i+1}), & i = 0; \\ x_{i-1} + (x_{i-1} - x_{i-2}), & i = n; \\ \dfrac{x_{i-1} + x_{i+1}}{2}, & else. \end{cases}$$

(10)

(5) Noise Reduction with Isolation Forest

Due to significant noise in the original data, abnormal samples were detected using the Isolation Forest algorithm, and interfering data was removed. For missing values resulting from the removal of abnormal samples, linear interpolation was still applied. The part of the freight volume data of Hengyang Station after noise reduction is shown in Figure 2.
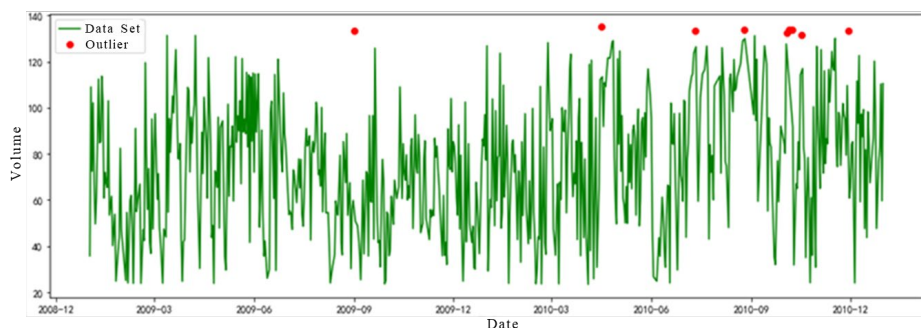


*Figure 2: Some freight data of Hengyang station after noise reduction.*

After preprocessing, over 3,000 records of railway freight data from the Xiang-Yue section are retained as the experimental dataset, which is then split into training and testing sets in an 8:2 ratio.

### 3.2 Model Parameters and Evaluation Metrics

The model is built using the deep learning framework PyTorch, with additional use of libraries such as Numpy, Pandas, and Sklearn. The railway freight volume data is standardized using Z-score normalization to smooth the training process and improve convergence speed. In the time series module, the number of hidden layers in the LSTM is set to 128. The model's optimizer is Adam, the loss function is mean squared error (MSE), the activation function is Tanh, the learning rate is set to 0.001, the batch size is 64, and the total training epochs are set to 100.

To provide a more intuitive display of the model's prediction results, root mean square error (RMSE) and mean absolute error (MAE) are used as comprehensive evaluation metrics, calculated as follows:

$$I_{MAE} = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}|_i$$

(11)

$$I_{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$$

(12)

Where $y_i$ and $\hat{y}_i$ represent the predicted and actual freight volumes at time $i$, represent the predicted and actual freight volumes at time $n$ is the total number of data points.

### 3.3 Comparison with Other Models

To validate the performance of the proposed model, common forecasting models are selected for comparison. The baseline models are:

(1) Historical Average (HA), Uses the historical average as the predicted freight volume.

(2) Support Vector Regression (SVR)[8], Fits the freight volume data to a curve and predicts the next time step.

(3) Random Forest (RF)[9], Constructs multiple decision trees and averages the results for regression prediction. The number of decision trees and maximum depth are determined using grid search to find the optimal parameters.

(4) Long Short-Term Memory (LSTM), A variant of recurrent neural networks that better retains long-term information. The number of hidden neurons is set to 64, the loss function is MSE, the optimizer is Adam, and the number of iterations is set to 100.

(5) Gated Recurrent Unit (GRU), Another variant of RNNs with a simpler structure than LSTM. The parameters are set similarly to those of LSTM.

(6) Diffusion Convolutional Recurrent Neural Network (DCRNN)[10], Views the convolution operation as a diffusion process, using diffusion graph convolution to improve the recurrent unit's capture of spatio-temporal features. The parameters are set similarly to the STGNN model.

### 3.4 Experimental Results and Evaluation

Considering the autocorrelation of freight flow sequences at different freight stations and model training time, three forecasting steps are set: 1, 5, and 10 days, to compare model performance at different time granularities. The experimental results are shown in Table 1. Overall, the STGNN model outperforms traditional models across various time steps and metrics. Among the models, the performance of the traditional statistical algorithm HA is the worst, as it only averages historical data and fails to learn changes in spatial and temporal features for regression prediction. Random Forest and SVR, which are traditional machine learning methods, perform better than LSTM and GRU-based sequence prediction methods. However, the prediction accuracy of DCRNN, which incorporates spatial correlation, is significantly improved in all forecasting steps, demonstrating the importance of spatial correlation. The STGNN model's graph neural network further integrates multiple types of spatial correlations, achieving better performance in one-step and five-step predictions, with MAE reduced by 8.49% and 4.59%, respectively, compared to DCRNN. Although the MAE of the STGNN model is slightly lower than that of DCRNN in ten-step prediction, its RMSE remains superior, indicating that the

STGNN model achieves the best overall performance.

*Table 1: Freight volume forecast results*

| Step Size | Evaluation Index | Prediction Model | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | HA | SVR | RF | LSTM | GRU | DCRNN | STGNN |
| Single step Prediction | MAE | 34.37 | 29.59 | 30.19 | 31.68 | 30.21 | 24.67 | 16.18 |
| | RMSE | 55.63 | 49.21 | 49.26 | 52.08 | 51.81 | 43.14 | 35.59 |
| Five-step Prediction | MAE | 37.00 | 36.57 | 35.99 | 36.92 | 36.23 | 32.56 | 27.97 |
| | RMSE | 64.10 | 64.46 | 62.56 | 64.34 | 63.33 | 53.79 | 48.95 |
| Ten-step Prediction | MAE | 42.78 | 39.29 | 39.13 | 41.52 | 40.96 | 34.17 | 34.52 |
| | RMSE | 71.92 | 66.84 | 66.75 | 70.39 | 67.57 | 55.60 | 54.35 |

## 4. Conclusion

This paper constructs an STGNN model for railway freight volume prediction. In the spatial convolution module, multi-relational graph convolution (R-GCN) is introduced to extract and integrate spatial features, while the time module uses multi-layer Gated Recurrent Units to capture temporal characteristics for time series prediction. The comparison with other models demonstrates the effectiveness of spatial feature analysis, significantly improving prediction accuracy. Future research can incorporate more dynamic features to further enhance prediction performance and better serve railway transportation organization.

## References

*[1] Wu Wei, Fu Zhuo, Wang Xiao. Passenger Volume Prediction Method for Transport Corridors [J]. Journal of Railway Science and Engineering, 2012, 9(05): 96-102. DOI: 10.19713/j.cnki.43-1423/u. 2012. 05.018.*

*[2] Li Jianguo, Xiang Wanli, Wang Jiugeng. Railway Freight Volume Prediction Based on Grey Genetic BP Neural Network [J]. Science & Technology and Industry, 2022, 22(01): 119-124.*

*[3] Guo Dongdong, Wang Xifu. Railway Freight Volume Prediction Based on BP Network [J]. Railway Freight, 2006(02): 21-23.*

*[4] Zhang Guandong, Yang Chen. Freight Volume Prediction Based on Multi-Dimensional Long Short-Term Memory Network [J]. Statistics and Decision, 2022, 38(12): 180-183. DOI: 10.13546/j.cnki.tjyjc. 2022. 12.036.*

*[5] Liu Qidong, Liu Chaoyue, Qiu Zixin et al. Time-Aware Transformer-Based Traffic Flow Prediction Method [J]. Computer Science, 2023, 50(11): 88-96.*

*[6] Xu Li, Fu Xiangyuan, Li Haoran. Spatio-Temporal Traffic Flow Prediction Model Based on Gated Convolution [J]. Computer Applications, 2023, 43(09): 2760-2765.*

*[7] Li, Yaguang, Yu, Rose, Shahabi, Cyrus, Liu, Yan. Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting [J]. arXiv, 2017, 20(12): 38.*

*[8] He Bisheng. Research on Coordination Optimization Theory and Methods of High-Speed Railway Train Operation Plans and Timetables [D]. Beijing: Beijing Jiaotong University, 2014: 101-110.*

*[9] Chen Rong, Liang Changyong, Lu Wenxing et al. Forecasting Tourism Flow Based on Seasonal PSO-SVR Model [J]. Systems Engineering-Theory & Practice, 2014, 34(5): 1290-1296.*

*[10] Chen Zhonghui, Ling Xianyao, Feng Xinxin et al. Short-Term Traffic State Prediction Approach Based on FCM and Random Forest [J]. Journal of Electronics & Information Technology, 2018, 40(8): 1879-1886.*