# Analysis of Traffic Congestion Using LSTM and Graph Theory

Boyang Han

*College of Transportation Engineering, Tongji University, Shanghai, 201804, China*

Keywords: Traffic Prediction, LSTM, Graph Theory, Medium-Sized City.

Abstract: As the urbanization process continues to advance, the per capita vehicle ownership in cities keeps increasing. However, with the growth in traffic volume on roads, the originally short commuting time is constantly lengthened. Meanwhile, as a large number of vehicles converge at intersections, congestion rates also rise. Despite the measures taken by traffic police and volunteers to guide traffic during peak hours, congestion often occurs randomly and exhibits complexity. Moreover, if congested roads are not addressed, congestion can spread. In this paper, the Long Short-Term Memory(LSTM) algorithm is employed to analyze and predict traffic volume based on traffic flow information at intersections. Considering that detection devices at intersections in some second and third-tier cities may experience aging or malfunctions, a gradient descent algorithm is utilized to calculate the turning intentions of vehicles at each intersection at different times. This information is then used to extrapolate the approximate traffic volume at neighboring intersections. This approach not only aids the work of traffic police but also allows drivers to choose routes based on current congestion conditions and future congestion predictions.

## 1 INTRODUCTION

With the continuous popularity of automobiles in second and third-tier cities, traffic congestion has gradually become a prominent issue. This is manifested by the fact that the road occupancy rate exceeds the originally designed values. At certain intersections, during peak hours, the number of vehicles allowed to pass within each green light cycle may not meet the actual demand, potentially leading to further congestion and traffic delays. A straightforward approach to address this issue is to dispatch traffic police to guide the dispersal of excess traffic. However, due to the relative lag in information, traffic police often cannot consider congestion conditions at surrounding intersections during enforcement, potentially resulting in inferior decisions.

On the one hand, with the integration of traffic and artificial intelligence technologies, the technology of predicting intersection traffic flow based on artificial intelligence algorithms and incorporating surrounding road factors has become increasingly mature. Many recent studies involve LSTM or other algorithms optimized based on RNN for predicting quantifiable indicators such as traffic flow (Luo et al

2019, Zhao et al 2020, He et al 2020 & Guo et al 2018).

Regarding the acquisition of traffic flow data, discrepancies exist in the timing and types of data obtained from various detectors. J. Guo et al. mainly integrates real-time discrete data to achieve real-time prediction of common indicators such as traffic flow (Ryu et al 2018). U. Ryu et al. and other research teams have more accurately predicted the values of different indicators, achieving good RMSE values (Zhang et al 2023, Fahs et al 2023 & Wang et al 2022). In terms of practicality, F. Zhao et al. used basic graph theory knowledge for the joint prediction of road traffic flow (Chen et al 2020). Tong Wang et al. further transformed traffic flow data into a direct basis for judging urban road congestion, enhancing its practical value (Cai et al 2020). However, in some second and third-tier cities in China and even in certain cities globally, issues such as deviations in measurements may arise due to factors such as road renovations and gravitational settling in circular loops, making it challenging to fundamentally resolve the source of data.

This paper addresses this often-overlooked problem. Given that various detectors on main roads are more complete than those on secondary roads, and their maintenance cycles and conditions are generally better, researchers can obtain accurate data from both

main and secondary roads in the previous month. This allows us to determine the probability of traffic flowing from main roads to secondary roads and vice versa at different time points. By obtaining the current traffic flow on the main road (or secondary road), researchers can predict the traffic flow on the current and surrounding roads in the subsequent period. This approach provides a certain direction for traffic police in managing congested intersections, facilitating timely or advanced prediction of congested intersections and manpower allocation by law enforcement agencies. Additionally, by comparing the predicted traffic flow on secondary roads based on main roads with the actual measured traffic flow on secondary roads, if significant differences are persisting for several days during peak hours and external factors such as road maintenance are ruled out, it may indicate that some detectors on secondary roads require maintenance. This paper, based on the LSTM algorithm for calculating traffic flow and using the gradient descent algorithm to analyze the correlation between adjacent roads in the road network, provides a reference value for improving road safety in small and medium-sized cities.

## 2 METHOD

### 2.1 Data Source

When picking data, the need is to avoid excessive data in the minimal period, for example, big cities like Shanghai possess too many roads and traffic flows, which will be time-consuming. Bu too-limited data sources should also be avoided, such as small towns in rural areas of China. Thus, data in Xuancheng, Anhui province serves as a good case for further investigation. The open data in the paper collects the road locations, average speed, and traffic flows in Xuancheng during September 2022, with a minimum observation period of 5 minutes. The Surveillance camera and other sensors work from 0:00 to 23:55 and rest from 23:55 to 0:00 the next day. The approximated road network is shown in Fig. 1 below. Although there are many intersections in this graph,

the distribution is relatively scattered, the traffic flow is restrictive, and there are many broken roads a typical feature for non-large cities. Fig. 1 below shows the main road network in Xuancheng.



Figure 1: Road network of main roads in Xuancheng (Picture credit: Original).

Meanwhile, loop data captured by monitors in each intersection of Xuancheng is free and available online, which consists of data with 'ROAD_ID' in a combination of two intersections, for example, road '4589_4562'. The 'TURN' sections include four options, including right, straight, left, and unknown, corresponding to R|S|L|U (Table 1).

### 2.2 Data Preprocessing

Due to the primary focus of this study on predicting traffic volume, the actual geographical latitude and longitude information of the roads becomes less crucial. Additionally, in the experimental dataset, the "Unknown" category in the turning column includes instances of making U-turns, misidentification, and other problems like misclassifications, and as thus excluded. "Furthermore, as the open-source nature of the loop intersection data, the information on turns is provided solely based on direction rather than a specific mapping between roads. Consequently, it is necessary to integrate a series of main road network data, identify corresponding road segments, and complete the replacement to enhance the accuracy and completeness of turn-related information.

Table 1: Raw data of Traffic flow in Xuancheng.

| ROAD_ID | FTIME | TIME | COUNT | SPD | TURN |
|---------|-------|------|-------|-----|------|
| str | H: M:S | H: M:S | int | float | R|S|L|U |

Table 2: The modified road-network in Xuancheng, Anhui.

| ROAD_ID | FTIME | COUNT | HARM_SPD | END |
|---------|-------|-------|----------|-----|
| str | H:M:S | int | float | str |

By further integrating the data, sorting the data in the order of time, initial road, and end road is better, for the modified table as shown in Table 2.

## 2.3 Relevant Research Methods

The urban road network can be approximated as a graph, in which all roads can be equivalent to two directed edges of opposite directions and the same length, denoted by (U, V), where U is the set of all points u and V is the set of all edges v, each of which is connected by two different u.

In this paper, the smallest unit considered is the non-U-turn turning of a single vehicle at an intersection. Consequently, this behavior involves three nodes connected by edges with an additional intersection as the intermediate node. Based on the unidirectional nature of vehicle movement, the paper designates these three nodes as the predecessor (p), intermediate (m), and successor (s) nodes. Simultaneously, the paper uses the abbreviations p_m and m_s to denote the edges connecting the predecessor to the intermediate node and the intermediate to the successor node, respectively. Meanwhile, the Nomenclature and referred meanings used in this paper are all described in Table 3 below.

By obtaining the traffic flow data of Xuancheng, the paper constructs a comprehensive three-dimensional gradient descent model, with the special point that X and Y are exactly equal here, because in a minimum period, the traffic flow is from almost the same time to the outflow intersection, so the inflow traffic determinant coefficient for each intersection represented can be calculated.

The paper represents it in mathematical terms, which is，signifies the whole Road network's traffic flow

data, y is the current node calculated, with m means the dimension of X, step means the learning rate determined by artificial to be 0.01.

$$\text{Error} = X \cdot W - y \quad (1)$$
$$\text{gd} = X^T \cdot \text{error}/m \quad (2)$$
$$W' = W - \text{step} * \text{gd} \quad (3)$$

Thus, the final has the predicted weight of surrounded edges in different periods, with the sum equal to 1.

Simultaneously, the data are categorized based on the road to which the turning is directed, i.e., by grouping and swapping and then calculating the inverse matrix. Utilizing the weight distribution of the predecessor edge in the inverse matrix, this paper given the traffic volume on a particular road, determines the distribution of its traffic sources. By transposing the weight distribution of the predecessor edge on this road in the original matrix, and combining it with the previously obtained traffic source distribution, traffic flow data for adjacent edges will be able to calculate. LSTM was originally designed to address the prevalent issue of long-term dependencies in traditional recurrent neural networks. The use of LSTM enables effective transmission and representation of information over long sequences without the problem of neglecting (forgetting) useful information from distant past time steps. Simultaneously, LSTM can mitigate the problems of gradient vanishing/exploding commonly encountered in standard RNNs.

In the construction of a model, the train data is set to be 14 days, while the test data is set to be 5 days.

The LSTM model and the basic nomenclature can be interpreted as the model following(as Fig. 2).

Table 3: The Nomenclature and referred meanings.

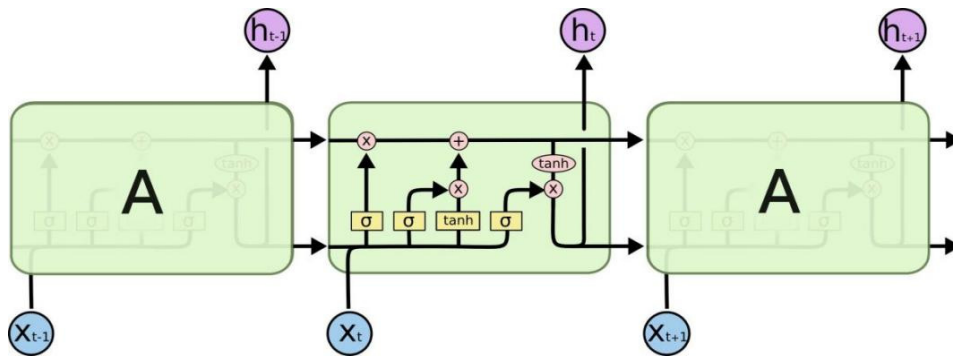| Nomenclature and referred meanings | |
|---|---|
| p | head of the minimal unit, the first node calculated in each Trajectory; |
| m | Middle of the unit, or the node signifying the intersection; |
| s | Back of the unit, the last node calculated in each; Trajectory; |
| T $(u, u')_{(p,m,s)}$ | The minimal unit in the research of this paper, with (p, m) be the road before a turn, and (m, s) as $u'$, the road after a turn |
| $f_0, f_5, \dots$ | Forecast of the future 0,5,10 based on LSTM ... The trend of minute traffic flow converges on the model at a macro level |
| $\Delta_{(u,v)}$ | The axis of congestion rate difference between different roads, calculated by the division of average speed and scheduled speed in each road by Chinese national standard |
| $w_{(T)}, w_{(T)}^{-1}$ | Weight/Inverse weight matrix of how many cars turn from u to v in duration of T, which can help predict the traffic flow of edge (u, v) after time T |

Figure 2: The model representing LSTM methods (Picture credit: Original).

# 3 RESULT AND DISCUSSION

## 3.1 Accuracy Prediction

By fitting the model with train data, the predicted data is available by calculating the last data, and by crisscrossing and, in time can be used to predict in time T. The accuracy of LSTM in the predicted 5 days and then the prediction of traffic flows in the following time after the test data, which is 2020/9/27,0:00:00, is listed as shown in Fig. 3, with the best RMSE reaching 1.76. This Fig.3 shows that the fitting process is successful and the prediction will be accurate at least in normal conditions.
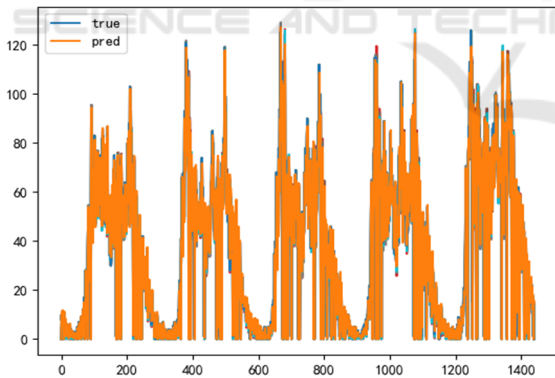


Figure 3: The accuracy of LSTM in prediction (Picture credit: Original).

## 3.2 Intersection Prediction

By using the original LSTM model, the data of the last time step obtain the function to predict traffic volume for 5 minutes ahead. Furthermore, by using the predicted value for y as input for another round of prediction, the data for 10 minutes ahead is then available (Table 4).

Table 4: Traffic prediction of "4589_4562" after the test data.

| 2022/9/27,16:00:00 | |
| --- | --- |
| 0-5 minute delayed traffic forecast: | 13.4028 |
| 5-10 minute delayed traffic forecast: | 13.8148 |
| 10-15 minute delayed traffic forecast: | 12.0810 |
| 15-20 minute delayed traffic forecast: | 12.5041 |
| 20-25 minute delayed traffic forecast: | 14.6358 |
| 25-30 minute delayed traffic forecast: | 13.2260 |

Furthermore, by examining the weights and the inverse matrix for the roads, for example, the data of "4589_4562" road at 16:05 can calculate predictions for the traffic volume on surrounding roads "4562_4651", "4605_4589" and "4660_4589"(Table 5).

Table 5: weight index and inverse weight index of test data in 16:05.

| | In | In_axis | Out | Out_axis |
| --- | --- | --- | --- | --- |
| $W_{4589\_4562}$ | 4562_4651 | 0.538 | 4605_4589 | 0.866 |
| $W_{4589\_4562}$ | 4562_4651 | 0.538 | 4660_4589 | 0.133 |
| $W^{-1}_{4589\_4562}$ | 4605_4589 | 1 | 4562_4651 | 1 |
| $W^{-1}_{4589\_4562}$ | 4660_4589 | 1 | 4562_4651 | 1 |

The traffic volume on the road "4589_4562" is entirely provided by the road "4562_4651," with the traffic flow entering this lane from "4562_4651" accounting for 53.8% of the traffic on that road. Therefore, the paper has:

$$f_{4562\_4651(5)} = f_{4589\_4562(5)} \times \frac{w_{out}^{-1}}{w_{in}} \qquad (4)$$

$$f_{4605\_4589/4660\_4589(5)} = f_{4589\_4562(5)} \times \frac{w_{in}^{-1}}{w_{out}} \qquad (5)$$

As a result, it can be seen that the traffic volume of "4562_4651" in 5 minutes is 13.4 * 1 / 0.538 = 24.9, and the traffic volume of "4605_4589" in 5 minutes is 13.4 / 1 * 0.866 = 11.6, with "4660_4589" be at 1.8. Whereas steps (4) and (5) can be further extended to address traffic congestion relief issues:

In a situation, someone has to go out in three hours, but the baggage to bring is delayed in large transport vehicles in rush hour, and calculate the average length of the cars as follows.

$$q = N / T \qquad (6)$$

For q is the traffic flow, N is the number of cars calculated in minimal time split T.

$$K = N / L \qquad (7)$$

K is the occupation rate, L is the harm length, and is accessible given the location.

$$T * v * K/q = L_{car} \qquad (8)$$

For average circumstances, Lcar =4.5m, which is not useful in this question, and suppose Lsmall_car =3m, Lbig_car =3m, it can predict how many trucks are delivering the baggage by solving a coefficient-to-be-determined-equation.

When directing traffic in the aftermath of a traffic accident, traffic police should also take into account the planned traffic flow in adjacent lanes, rather than allowing all vehicles to proceed freely. Otherwise, it may lead to congestion on the road ahead. In this scenario, let's denote the additional traffic flow that the traffic police can clear within a unit area and time as X. Therefore, when there is a change in X, the increase in traffic flow for the subsequent road segment is given by:

$$\Delta y = w_{out} \cdot \Delta x \qquad (9)$$

Here, $\Delta y$ represents the increase in traffic flow for the subsequent road segment, $\Delta x$ represents the change in the additional traffic flow cleared by the traffic police within a unit area and time, and k is a coefficient that signifies the impact of the traffic flow cleared by the traffic police on the subsequent road segment's traffic flow. The specific value of this coefficient is influenced by various factors such as traffic flow characteristics, road structure, and vehicle speed.

By plotting the variation in average vehicle length throughout the day, it is easy to identify the time periods when the proportion of large vehicles is highest. This information can serve as a basis for

further research and analysis focused on understanding the dynamics and implications of heavy traffic during specific times of the day.

## 3.3 Limitations and Future Outlook

Through the methods employed in this study, the traffic volume of adjacent lanes with greater accuracy based on the traffic flow of a specific road will be available, especially main arteries such as thoroughfares. This capability aids traffic police in efficiently managing congestion and allows for prompt assignment of the next duty location, facilitating smoother transitions between tasks. However, this paper acknowledges certain limitations. In terms of data preprocessing, the outright removal of lanes with 'Unknown' turning information is inappropriate, as it may encompass instances of left turns, right turns, or straight movements that were not identified. Moreover, this action results in the deletion of approximately one-third of the data, introducing a certain level of bias into the results. Addressing this issue requires more comprehensive support from road infrastructure or conducting on-site investigations to gather firsthand information for comparison with the extensive dataset, thereby optimizing the results. It's also important to note that due to the reliance on predicted values, the accuracy of verified predictions in LSTM may gradually decrease over subsequent time steps.

While the emphasis of this paper lies in prediction, the incorporation of direct vehicle speed data allows for theoretical integration of road length, congestion prediction, and congestion coefficient calculation (i.e., the ratio of a road's design speed to theoretical speed). This, in conjunction with spatiotemporal mapping, could yield a more sophisticated road planning approach. Unlike most current navigation software that guides based on historically optimal routes or real-time congestion predictions, the theoretical framework holds significant research potential for offering more refined road guidance.

## 4 CONCLUSION

This paper aims to address the pressing issue of daily commuting congestion in mid-sized cities by considering the perspectives of both drivers and traffic police. The primary objective was to leverage real-time traffic flow data from specific road segments for estimating traffic conditions near congested intersections. By employing the LSTM (Long Short-Term Memory) algorithm, accurate

predictions of traffic flow data at individual intersections were achieved.

The incorporation of graph theory, coupled with the calculation of edge-weighted correlation matrices, facilitated the determination of traffic flow proportions between roads, encompassing both forward and reverse matrices. This methodology allowed for the computation of traffic flow on adjacent edges and could be extended to the broader road network after considering multiple layers of correlation. However, the paper acknowledges the potential compromise in accuracy as unaccounted vehicles, especially those entering or exiting specific locations, may influence the results.

A noteworthy contribution of this study lies in addressing the application gap of LSTM algorithms for predicting traffic flow and enhancing traffic safety management in mid-sized cities with the lack of the data accuracy. Despite the valuable insights provided, challenges arose due to insufficient precision in data collection, leading to the discarding of significant data during initialization and potentially resulting in an overall underestimation of predictive values.

Future research directions were proposed, emphasizing the importance of obtaining more precise information about vehicle speed and incorporating it into a spatiotemporal map to overlay temporal changes in traffic flow data. The paper also suggested employing congestion calculations based on national standards(China), to assess congestion levels on different roads.

These proposed enhancements not only contribute to improving the accuracy of real-time planning in navigation systems but also offer valuable insights into the challenges and opportunities within the domains of civilian traffic forecasting and road traffic safety management. In light of these findings, the paper concludes by highlighting the significant research value in extending the application of LSTM algorithms in this context, prompting further investigation and exploration in this crucial field.

## REFERENCES

X. Luo, D. Li, Y. Yang, and S. Zhang, J. Adv. Transp., vol. 2019, pp. 1–10, Feb. (2019).

F. Zhao, G.Q. Zeng,and K.D. Lu, IEEE Trans. Veh. Technol., 60(1), 101–113, Jan. (2020).

P. He, G. Jiang, S.-K. Lam, and Y. Sun, Inf. Sci., 512, 1394–1406, (2020).

J. Guo, Z. Liu, W. Huang, Y. Wei, and J. Cao, IET Intell. Transp. Syst.,12(2), 143–150, Mar. (2018).

U. Ryu, J. Wang,T. Kim, S. Kwak, and U. Juhyok, Res. C, Emerg. Technol., 96,.55–71, (2018).

T.Zhang, J. Xu, S. Cong, C. Qu, and W. Zhao, IEEE Access, 11, 36471-36491, (2023).

W. Fahs, F. Chbib, A. Rammal, R. Khatoun, Procedia Computer Science 220 , 202-209,(2023).

T. Wang, A.r Hussain, Q. Sun, IEEE Intellegent Transportation Systems Magazine, 102-120, July/August (2022).

C. Chen, Z. Liu, S. Wan. IEEE Transactions on Intelligent Transportation Systems, (2020).

L. Cai, Y. Yu, S. Zhang, Y. Song, Z. Xiong and T. Zhou, in IEEE Access, 8, 22686-22696, (2020).