

Research Article

Short-Term Traffic Volume Forecasting with Asymmetric Loss Based on Enhanced KNN Method

Zhiyuan Wang ¹, Shouwen Ji ¹, and Bowen Yu²

¹MOE Key Laboratory for Urban Transportation Complex Systems Theory and Technology, Beijing Jiaotong University, Beijing 100044, China

²China Railway Container Transport Co., Ltd Beijing Branch, Beijing 100055, China

Correspondence should be addressed to Shouwen Ji; shwji@bjtu.edu.cn

Received 27 December 2018; Revised 28 February 2019; Accepted 20 March 2019; Published 15 April 2019

Academic Editor: Mahmoud Mesbah

Copyright © 2019 Zhiyuan Wang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Short-term traffic volume forecasting is one of the most essential elements in Intelligent Transportation System (ITS) by providing prediction of traffic condition for traffic management and control applications. Among previous substantial forecasting approaches, K nearest neighbor (KNN) is a nonparametric and data-driven method popular for conciseness, interpretability, and real-time performance. However, in previous related researches, the limitations of Euclidean distance and forecasting with asymmetric loss have rarely been focused on. This research aims to fill up these gaps. This paper reconstructs Euclidean distance to overcome its limitation and proposes a KNN forecasting algorithm with asymmetric loss. Correspondingly, an asymmetric loss index, Imbalanced Mean Squared Error (IMSE), has also been proposed to test the effectiveness of newly designed algorithm. Moreover, the effect of Loess technique and suitable parameter value of dynamic KNN method have also been tested. In contrast to the traditional KNN algorithm, the proposed algorithm reduces the IMSE index by more than 10%, which shows its effectiveness when the cost of forecasting residual direction is notably different. This research expands the applicability of KNN method in short-term traffic volume forecasting and provides an available approach to forecast with asymmetric loss.

1. Introduction

With the booming development of detecting devices, mobile internet, and cloud computing technique, Intelligent Transportation System (ITS) is being implemented in real traffic management systems to improve the efficiency of traffic management. Short-term traffic volume forecasting, which can provide road information ahead of time, has been an essential part of ITS to support real-time traffic control and management [1]. Short-term traffic volume forecasting is the process of estimating directly the anticipated traffic conditions at a future time, given continuous short-term feedback of traffic information [2]. Distinguished from long-term traffic forecasting which serves for traffic planning, short-term traffic volume forecasting focuses on predicting traffic condition over time horizons ranging from few seconds to few hours. Most of traffic data used in short-time traffic forecasting models is collected by automatic detecting devices.

Besides abundant real-time traffic data, rapidly developing researches in this field using typical statistic models and machine learning algorithms also accelerate the application of short-term traffic volume forecasting. These approaches are reviewed in Section 2. However, most of these researches focus on forecasting with symmetric loss, which holds a hypothesis that the forecast value is larger or lower than the real value and it gets the same cost. Forecasting with symmetric loss is simple and mediocre but insufficient to satisfy the real need in the context of traffic volume, because if the forecasting value is larger than the real value, it only costs linear-increased management resource and guides travelers suboptimal routes. On the opposite, if forecasting traffic volume is lower than the real value, it may cause a traffic congestion and the whole traffic system will be more vulnerable and unpredictable.

This article tries to enhance normal KNN forecasting method to forecast with asymmetric loss. This paper is organized as follows. In Section 2, previous related researches

are reviewed. In Section 3, the basic concept of dynamic KNN forecasting is introduced firstly; then some detailed techniques are discussed including the enhanced Euclidean distance which takes the stability of difference between two traffic profiles into consideration, Loess smoothing, asymmetric loss index Imbalanced Mean Squared Error (IMSE), and corresponding algorithm. In Section 4, the basic data used in experiments is introduced firstly. Then profiles using Loess or not are both experimented in dynamic KNN forecasting model. The best ranges of three key parameters are discussed. The effectiveness of KNN with asymmetric loss is tested last. In Section 5, the conclusion of this paper is drawn and further issues in this direction are discussed.

2. Literature Review

Short-term traffic forecasting has been a classical research direction in ITS for nearly 40 years. After Box-Jenkins method was applied by Ahmed and Cook [3], enormous typical statistical approaches such as historical average algorithms [4], smoothing [4, 5], Kalman filtering [6], and ARIMA family models [7, 8] are widely used in this area. These well-founded mathematical approaches mostly are parametric models and act well in model specification; however they become insufficient when traffic pattern is complex and parameters of models are hard to adjust responsively [9, 10]. Corresponding entire spectrum literature before 2003 was critically reviewed by Vlahogianni E [2].

In recent 20 years, as automatic traffic detecting devices have been widely used and machine learning theories have made progress rapidly, data-driven empirical algorithms become prosperous in short-term traffic forecasting [11]. Such algorithms have advantages that they are free of any assumptions regarding the underlying model formulations and the uncertainty involved in estimating the model parameters. These algorithms include K nearest neighbor (KNN) [12, 13], Support Vector Machine (SVM) [14], Random Forest (RF) regression [15], and Artificial Neural Network (AI/NN) [16].

For KNN method, Smith and Demetsky tested the performance of KNN regression compared with neural networks, a historical average, and the ARIMA model and concluded that KNN was superior in the field of transferability and robustness [17]. Smith et al. used kernel neighborhoods and suggested that the method produced predictions with an accuracy comparable with that of the seasonal version of an ARIMA model [18]. Habtemichael concluded that previous researches using KNN method mostly used the simplest form of KNN [19]. Enhanced KNN method using weighted Euclidean distance and weight to the candidate value was proposed in this article.

The crucial step in KNN method is to define the similarity measurement between traffic profiles. It is also the subject in traffic volume clustering. Aghabozorgi et al. reviewed previous literature in time-series clustering and concluded Euclidean and Dynamic Time Warping (DTW) were useful similarity measurement for time series [20–23]. There is certain limitation in Euclidean distance and it is discussed and enhanced in Section 3.1. Other related researches in traffic volume clustering include Xia et al. [24] and Xia et

al. [25]. Lin et al. combined KNN with local linear wavelet neural network for short-term prediction of five-minute traffic volumes and get better performance compared with LLWNN and SVR [26].

In the review of Vlahogianni [11], 10 future directions in short-term forecasting were proposed. In model selection, most previous researches just selected the model that provided the most accurate predictions regardless of whether certain modeling assumptions were violated or unrealistic. Forecasting with asymmetric loss has been researched in statistic field [27] and was widely used in economic issues [28–30]. Zhang et al. [31] used GJR-GARCH model with a conditional variance formulation to capture asymmetric response in the conditional variance in short-term traffic forecasting. GJR-GARCH allows the conditional variance to respond differently to the past negative and positive innovations, which is inspiring for this article. Lin et al. [32] used quantile regression to deal with the heteroscedasticity problem, which used asymmetric loss functions for prediction intervals calculation of short-term traffic volume.

In summary, the limitation of using Euclidean distance in similarity measurement of traffic profiles has rarely been discussed and using asymmetric loss forecasting of traffic volume is a relatively new issue and is practical in real ITS systems. This research aims to fill up these gaps.

3. Methodology

3.1. Basic Concepts of Dynamic KNN Forecasting Method. KNN is a nonparametric and data-driven method for classification and forecasting. The notion of KNN is “Whatever has happened before will happen again.” Similar pattern is extracted from historic data and compared with new data to determine the underlying classification label or value of new data. In traffic volume forecasting, KNN model needs a historic traffic profiles database. Given a certain subject traffic profile to make forecasting, KNN model compares the similarity between the subject profiles with profiles in database using predefined similarity measurement. Then K nearest neighbor profiles are chosen and aggregated in desired time horizon to make predictions of future traffic volume. To explain how to use KNN method in traffic volume forecasting, some terms are defined as follows. The use of terms generally keeps consistent with terms used by Habtemichael [19].

- (i) Subject profile: the traffic profile of one specific day to be forecasted. The data structure of subject profile can be $1 \times n$ (n timestamp) vector.
- (ii) Profiles database: a database contains historical traffic profile collected and preprocessed previously. The data structure of profiles database can be $m \times n$ (m is the size of historical profiles and n timestamp) two-dimensional table.
- (iii) Candidate profiles: the nearest neighbor selected from profiles database according to similarity between subject profiles with profiles in database. The number of candidate profiles is the parameter K in KNN method.

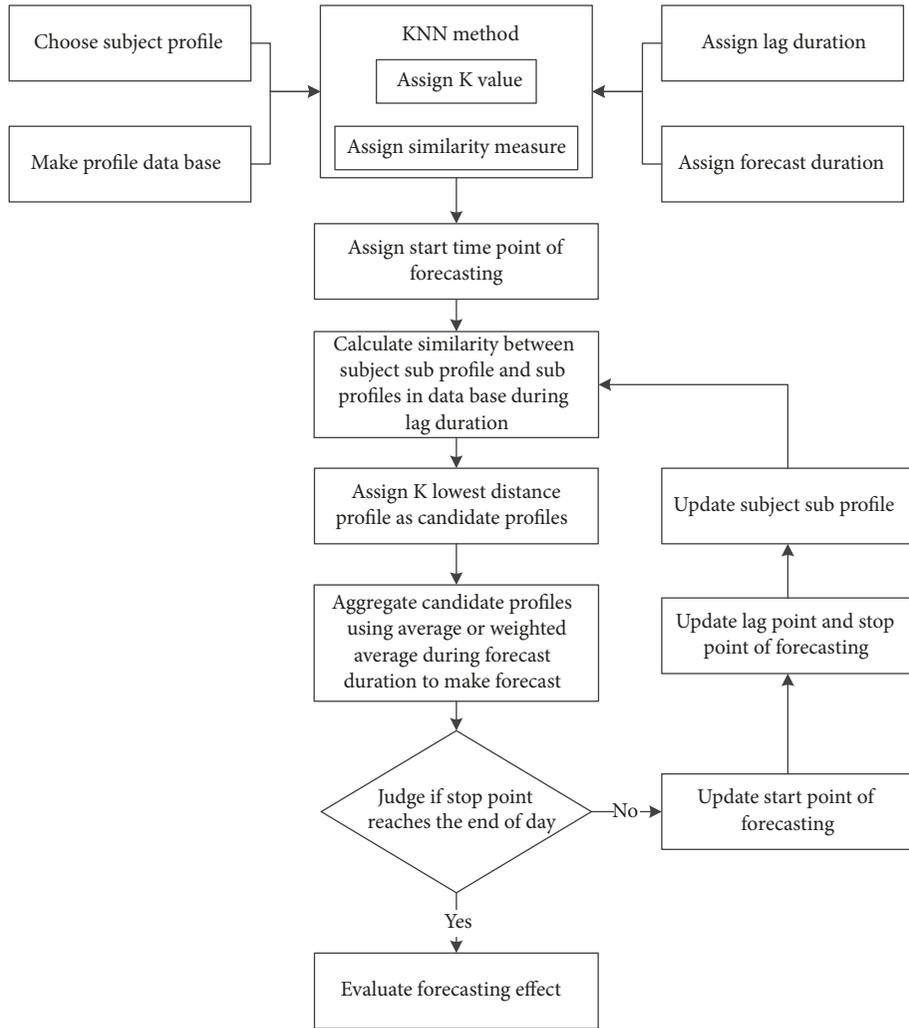


FIGURE 1: The flow chart of dynamic KNN traffic volume forecasting method.

- (iv) Lag duration: the time window considered to determine the similarity between subject and candidate profiles. For example, if time resolution is 5 minutes, the start point of forecasting is 6.00 AM (timestamp 72) and lag duration is 24 time intervals ($5 * 24 / 60 = 2$ hours); the subprofile from 4.00 AM to 6.00 AM of subject profile and candidate profiles will be extracted and used for similarity measure.
- (v) Forecast duration: the time window to make forecasting. For example, if time resolution is 5 minutes, the start point of forecasting is 6.00 AM and forecast duration is 12 time intervals ($5 * 12 / 60 = 1$ hour); the subprofile from 6.00 AM to 7.00 AM will be forecasted based on candidate subprofiles from 6.00 AM to 7.00 AM.

The KNN forecasting methodological approach proposed by Habtemichael can just achieve static forecast by one time [19]. The dynamic version of KNN forecasting is transformed as flow chart Figure 1 shows. This dynamic KNN approach

can achieve rolling traffic volume forecasting of a whole day by multisteps.

3.2. *Improved Similarity Measuring Method.* Many previous researches hold the view that Euclidean distance (or weighted Euclidean distance) is a proper choice of similarity measurement in traffic volume clustering and predicting. However, it is worth thinking about the limitation of Euclidean distance. For a sequence of points, x_i^n and y_i^n , Euclidean distance is calculated as

$$d_{ab} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \tag{1}$$

Take an example as Figure 2 shows. There are three imitating traffic volume time series y, y1, and y2. Use Euclidean distance to calculate the similarity between y and y1, y and y2. Because every point of y2 is closer to y compared with y1, the Euclidean distance between y2 and y (79.86) is lower than that between y1 and y (120).

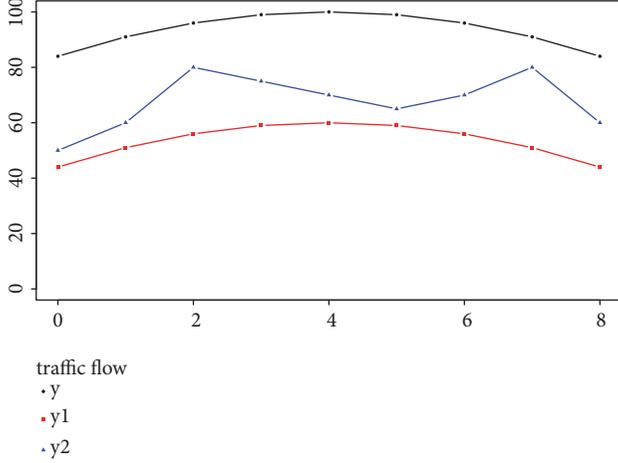


FIGURE 2: Three imitating traffic volume time series to explain the limitation of Euclidean distance in traffic volume similarity measurement.

However, in the context of traffic volume, profile y and $y1$ shows the traffic volume with peak in timestamp 4 and profile $y2$ with peak in timestamps 2 and 7. In this case, Euclidean distance can only measure similarity in absolute distance in total but ignores the difference in shape between profiles and cannot describe traffic volume more particularly.

To take the shape difference between profiles into consideration, adding the stability of difference measurement into Euclidean distance is a practical way as (2) shows:

$$d_{ab} = R_1 \sqrt{\sum_{i=1}^n z_i^2} + R_2 \sqrt{\sum_{i=1}^n (z_i - \bar{z}_i)^2} \quad (2)$$

where $z_i = x_i - y_i$ and \bar{z}_i is the average value of series z_i . $R1$ and $R2$ are the parameters to switch the balance of the weight between absolute distance and stability of difference, which can add flexibility in similarity measurement.

3.3. Loess Smoothing Technique to Reduce Noise. The traffic volume used in experiments is aggregated by 5 minutes and there is plenty of noise in profiles. Using raw data without any processing will make forecasting unstable and damp accuracy. This has been tested in Section 4.2.1. Technique of smoothing noise in time series should be used before forecasting.

Locally estimated scatterplot smoothing (Loess) is a mature nonparameter smoothing technique which has been widely used in previous related works. The theory of Loess technique is stated by Cleveland and Devlin [33]. Comparing experiments using raw profiles and using Loess smoothing profiles in KNN forecasting is conducted in Section 4.2. Span is the key parameter of Loess model, which can adjust the smoothness of smoothed profiles. Previous research [19] has proved that 0.2 is a proper value for span in traffic volume forecasting using KNN method.

3.4. Forecasting with Asymmetric Loss. In normal KNN forecasting model and most of other forecasting methods, there is a hypothesis: the cost of forecasting lower or higher is equal. Few researches focus on the asymmetric cost of forecasting. However, asymmetric cost is meaningful in traffic management. If forecasting traffic volume is higher than the real value, it only costs linear-increased management resource and makes travelers choose other possible routes. However, if forecasting value is lower than the real one, traffic chaos is more likely to occur and the whole traffic system is more unpredictable and vulnerable. It is sensible to make forecasting slightly higher according to the imbalance cost of direction of forecasting residual.

3.4.1. Enhanced Criterion Index IMSE. To measure the effect of asymmetric loss forecasting method, enhanced asymmetric criterion must be constructed firstly. The typical balanced criterion index MSE (Mean Squared Error) is calculated as (3), where O_t is observed value and F_t is forecast value.

$$MSE = \frac{1}{n} \sum_{t=1}^n (O_t - F_t)^2 \quad (3)$$

Because sum of squares does not take positive or negative of difference into consideration, the index MSE is balanced in forecasting evaluation. Think about the difference between the real value and forecasting value $Z_t = O_t - F_t$. When Z_t is positive, which means observed value is larger than forecasting value, the punishment should be greater. Oppositely, if Z_t is negative, the punishment should be less.

So the Imbalanced Mean Squared Error (IMSE) can be transformed as follows:

$$IMSE = \frac{1}{n} \sum_{i=1}^n P_t Z_t^2 \quad (4)$$

$$P_t = W_1 \quad \text{when } Z_t > 0$$

$$P_t = W_2 \quad \text{when } Z_t < 0$$

$$0 < W_2 \leq W_1 < 2$$

$$W_1 + W_2 = 2$$

3.4.2. Enhanced Asymmetric Algorithm of KNN Forecasting Method. In KNN forecasting model, the most important part is to calculate the similarity between subject profile and candidate profiles. It is natural to add asymmetric identification and operation into similarity calculation. A native way is to add asymmetric response into distance. If the profile in database is larger than the subject profile, it will not be brought into distance calculation and will be chosen more easily in nearest neighbor identification. On the opposite, if subject profile is larger, it will be brought into distance calculation. The distance of two profiles is calculated as follows:

$$D = R_1 \sqrt{\sum_{i=1}^n z_i^2} + R_2 \sqrt{\sum_{i=1}^n (z_i - \bar{z}_i)^2} \quad (6)$$

TABLE 1: Number of records each day in dataset.

Date	9/19	9/20	9/21	9/22	9/23	9/24	9/25	9/26
Record Number	287	284	282	277	276	271	275	282
Date	9/27	9/28	9/29	9/30	10/1	10/2	10/3	10/4
Record Number	274	0	197	288	288	288	287	288
Date	10/5	10/6	10/7	10/8	10/9	10/10	10/11	10/12
Record Number	288	288	288	283	278	276	278	194

$$\begin{aligned}
 d_t &= c_t - s_t \\
 z_t &= d_t \quad \text{when } d_t < 0 \\
 z_t &= 0 \quad \text{when } d_t > 0
 \end{aligned} \tag{7}$$

where c_t is candidate profile value and s_t is subject profile value at timestamp t .

The advantage of this asymmetric algorithm is that it is easy to be programmed. However, when profile is far larger than the subject value, this algorithm is vulnerable and can be affected by abnormal value in profiles. This algorithm is reasonable when profiles in database are all from one certain detecting device and contain few abnormal values in time series.

4. Result and Discussion

4.1. Data Use. The dataset used in this paper is from one of the traffic investigation stations in GuiZhou province in China. It was collected by an electromagnetic coil detecting device buried under ground of freeway. The dataset is supported by Transport Planning and Research Institute of China. The time span of experiment dataset is from September 19 to October 12, 2016. Traffic volume aggregates every 5 minutes. So there are 288 records ($60 / 5 * 24$) every day if there is no loss in detecting and aggregating.

However, missing value is unavoidable when detecting device is not completely reliable. The numbers of records every day in experiment dataset are shown in Table 1. There is no record on Sep 28, so it will not be used in later experiments. For other dates that contain missing value more or less, interpolation method is used to ensure the number of records every day is 288.

The date span of dataset contains a ‘‘National Day’’ holiday from September 30th to October 7th, which is useful when finding similar traffic patterns in later experiments. Traffic volume pattern between holiday and normal dates can be easily identified as Figure 3 shows.

4.2. Dynamic KNN Forecasting

4.2.1. Using Raw Data in Dynamic KNN Model. To test the feasibility of dynamic KNN forecasting method mentioned in Section 3, raw traffic volume data which is not smoothed by Loess technique is used in this section.

Firstly, as described in KNN workflow in Figure 1, the subject profile and profiles database should be appointed. For example, the series of 10/6 is chosen to be the subject profile,

and series of other dates including the date before and after 10/6 (because the scale of dataset is limited and data should be used as far as possible to maintain the scale of database) are chosen to be the element profile of database.

Secondly, key parameters of KNN model should be designated, including K the number of nearest neighbor, *lag duration* the time window to calculate similarity of profiles, and *forecasting duration* the time window to make forecasting. Using the best value of parameters in Section 4.3, K is 3, lag duration is 22 intervals (1 hour 50 minutes), and forecasting duration is 6 intervals (30 minutes).

Thirdly, the start time point of forecasting should be designated. Because the traffic volume in early of a day (like 0 AM to 6 AM) is quite low and not important in traffic management, it is assumed that the forecasting value is approximately equal to the real value before start time point of a day. In this experiment, the start time point is designated to timestamp 72 (6 AM).

After all these preparations, dynamic KNN forecasting method can be used to forecast traffic volume. Figure 4 shows that the forecast value fluctuates acutely, which may reduce the accuracy of forecasting. The criterion indexes MSE, MAE, and IMSE of this experiment lie in the first line of Table 5.

4.2.2. Using Loess Data in Dynamic KNN Model. As experiment in Section 4.2.1 shows, fluctuation of profiles increases the difficulty for KNN method to find nearest neighbor and meanwhile reduce the accuracy of forecasting result. As Section 3.3 states, Loess smoothing technique can reduce noise in profiles and improve the ability to identify similarity among profiles. In this section, experiment is conducted to verify the effect of Loess smoothing technique using dynamic KNN forecasting model.

The main steps of this experiment are similar to experiment in Section 4.2.1 except that every profile in database is smoothed by Loess technique with span of 0.2. The result of experiment is shown in Figure 5. In this figure, the forecasting value is relatively smooth after the start point (72 timestamps). The red line of forecasting value still contains some sharp breaks caused by the change of chosen nearest neighbor. The criterion indexes MSE, IMSE, and MAE of this forecasting experiment lie in the second line of Table 5. Three criterion indexes drop obviously compared with the first line. It means that the forecasting using smoothed profiles database gets better accuracy.

The forecasting residual error after start time point (from timestamps 73 to 288) is shown in Figure 6. It shows the residual error distributes nearby 0 and contains no apparent

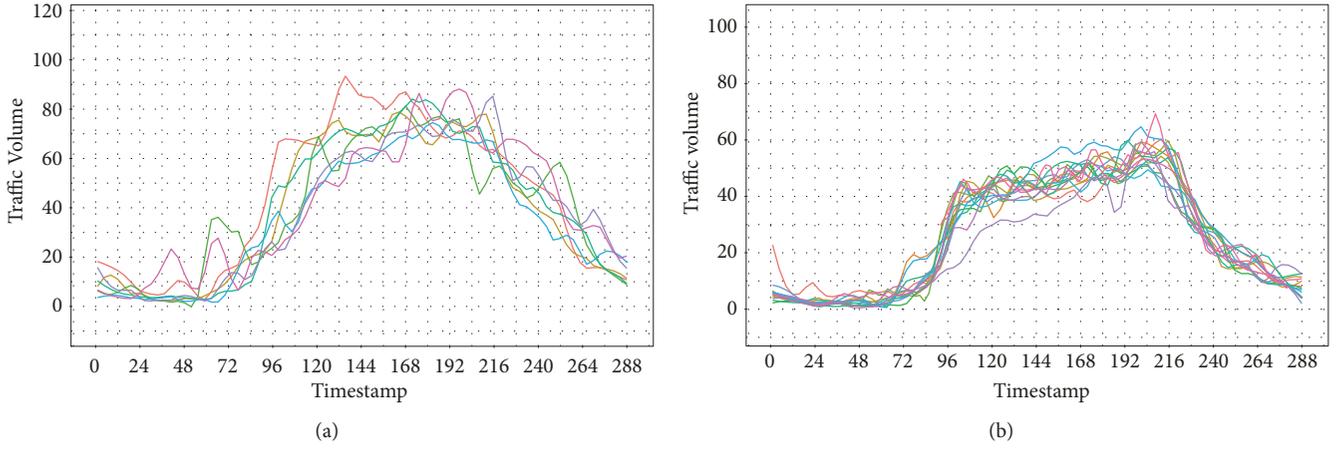


FIGURE 3: (a) Traffic pattern in holiday and (b) traffic pattern in normal date.

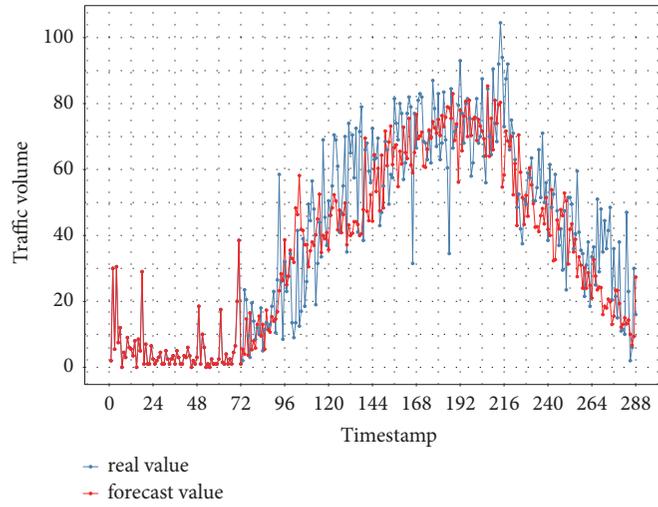


FIGURE 4: The result of dynamic KNN forecasting using raw traffic data. The blue line represents the real traffic volume and the red line represents forecasting value.

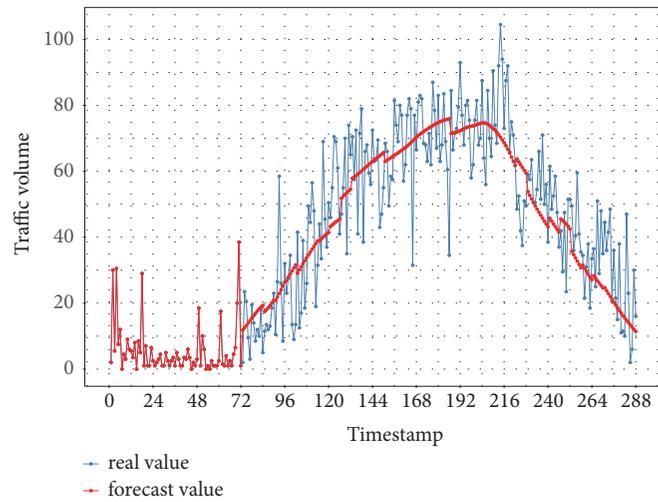


FIGURE 5: The result of dynamic KNN forecasting using smoothed profile database. The blue line represents the real traffic volume and the red line represents forecasting value.

TABLE 2: Criterion indexes in experiment of different K value.

K Value	MSE in symmetric loss model	IMSE in asymmetric loss model
2	128.30	145.44
3	126.13	131.39
4	128.50	131.98
5	130.47	126.46
6	134.16	130.22
7	136.69	137.23
8	145.34	149.40
9	153.35	164.97
10	161.30	180.90

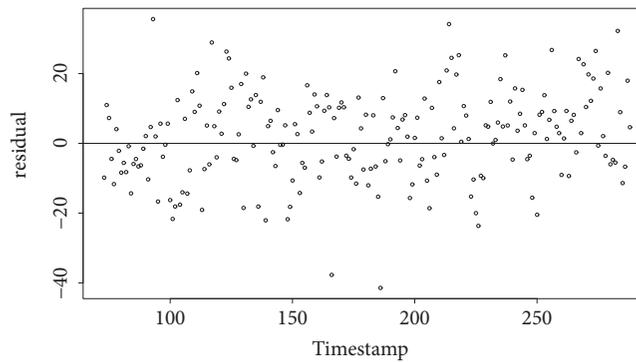


FIGURE 6: Residual error with time.

pattern with time. However, residual error above 0 is slightly more than error below 0. It means that in normal KNN forecasting model, forecasting value tends to be lower than real value. In traffic volume forecasting, this inclination may cause bad effect in traffic management which has been discussed in Section 3.4. To reduce this effect, asymmetric loss and asymmetric KNN algorithm have been introduced in and related experiment is conducted in Section 4.4.

Figure 7(a) shows Auto Correlation Function (ACF) of residual error and shows there is a little self-correlation in residual error. Figure 7(b) tests normality of residual error by Quantile-Quantile plot (Q-Q plot) and shows the normality of residual is acceptable.

4.3. Identifying Suitable Parameter Value in KNN Methods

4.3.1. Suitable Number of Nearest Neighbors K. Number of nearest neighbors K is one of the most important parameters, which determines the candidate scale used to make forecasting. When K is too small, the profiles used to forecast are insufficient and forecasting may fluctuate sharply by extreme value of candidate profiles. When K is too large, the candidate profiles are more likely to contain dissimilar traffic profiles, which may reduce the accuracy of forecasting. So choose an appropriate K value that is critical for this model.

In symmetric loss model, different K value is experimented and MSE is chosen to be the criterion index. In asymmetric loss model, IMSE is chosen to be the criterion

index. The result is shown in Table 2. MSE is lowest when K is 3 in symmetric loss model and IMSE is lowest when K is 5 in asymmetric loss. It means asymmetric model needs little more nearest neighbors to make the best forecasting.

4.3.2. Suitable Value of Lag Duration. Lag duration determines the length of profile to calculate similarity among the subject profile and candidate profiles. When lag duration is too short, the selection of nearest neighbor profiles changes notably, which makes forecasting unstable. When lag duration is too long, the forecasting is more stable; however the flexibility of model is worse, which makes accuracy decline.

Lag duration from 4 intervals (20 minutes) to 48 intervals (4 hours) is tested in both symmetric model and asymmetric model. The corresponding forecasting criterion indexes are shown in Table 3. MSE is lowest when lag duration is 22 in symmetric loss model and IMSE is lowest when lag duration is 30 in asymmetric model.

4.3.3. Suitable Value of Forecasting Duration. Forecasting duration determines the time window of one step. Longer forecasting duration can provide more traffic information in future, which is more helpful. However, if forecasting duration is too long, the predicting accuracy may decline and may cause unnecessary traffic chaos. This section tries to test the ultimate limit of forecasting duration in dynamic KNN traffic forecasting model.

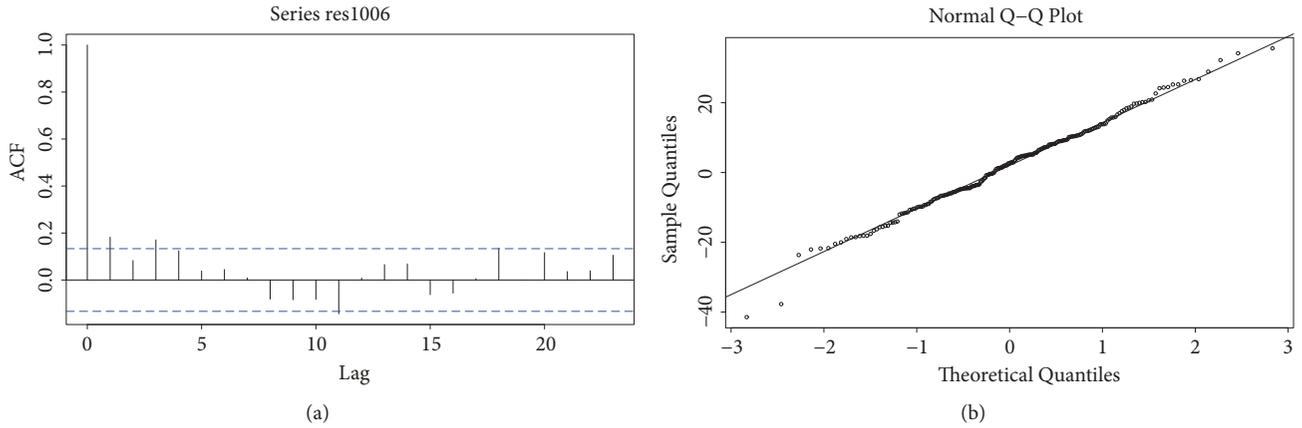


FIGURE 7: Residual error diagnosis. (a) Self-correlation and (b) normality of residual error.

TABLE 3: Criterion indexes in experiment of different lag duration.

Lag Duration	MSE in symmetric model	IMSE in asymmetric model	Lag Duration	MSE in symmetric model	IMSE in asymmetric model
---	---	---	26	128.32	125.11
4	150.56	126.38	28	126.37	125.02
6	134.94	128.41	30	127.73	124.71
8	136.62	128.41	32	131.88	125.94
10	136.62	128.36	34	134.61	125.94
12	132.63	128.36	36	136.02	126.32
14	129.55	128.36	38	133.34	128.65
16	127.54	128.40	40	132.28	128.65
18	125.86	127.46	42	134.48	129.08
20	128.48	127.46	44	136.77	129.08
22	125.35	126.59	46	137.62	129.08
24	126.13	126.46	48	141.94	129.08

TABLE 4: Criterion indexes in experiment of different forecasting duration.

Forecasting duration	MSE in symmetric model	IMSE in asymmetric model
2	128.82	145.44
4	125.23	131.39
6	124.02	131.9
8	126.13	126.46
12	126.13	130.22
18	131.83	137.23
24	134.55	149.40
36	149.59	164.97

TABLE 5: Criterion indexes in experiment of normal algorithms and asymmetric algorithms.

Criterion Index	MSE	MAE	IMSE
Normal model using raw data	168.44	8.83	200.90
Normal model using smoothed data	124.30	7.74	138.77
Asymmetric loss model using smoothed data	143.05	7.97	124.71

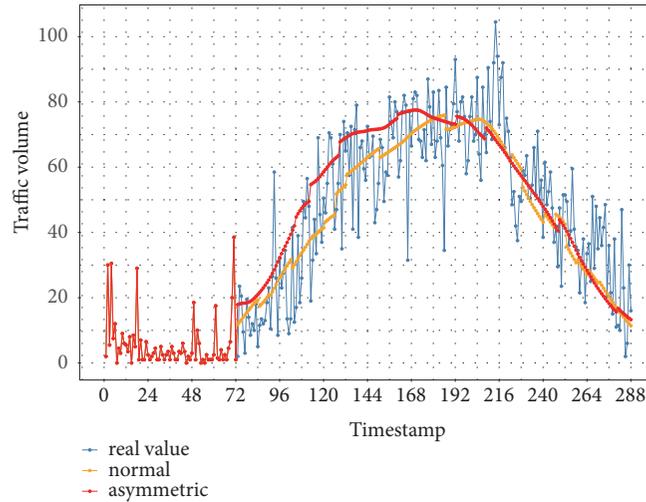


FIGURE 8: The result of dynamic KNN forecasting using asymmetric algorithm. The blue line represents the real traffic volume. The orange line represents the forecasting using normal algorithm. The red line represents the forecasting using asymmetric algorithm.

As Table 4 shows, the criterion indexes of different forecasting duration fluctuate slightly from 4 intervals (20 minutes) to 12 intervals (1 hour). When forecasting duration is larger than 1 hour, indexes increase obviously, which means forecasting is becoming unreliable. The best forecasting duration is 6 (30 minutes) in symmetric model and 8 (40 minutes) in asymmetric model.

4.4. Using Asymmetric Loss and Asymmetric Algorithm in KNN. As analysis in Section 3.4 shows, traffic volume forecasting with asymmetric loss is more practical in traffic management and traffic information service. This section uses asymmetric loss criterion index IMSE and asymmetric algorithm with the best parameters as Section 4.3 to test whether new-designed asymmetric algorithm can achieve asymmetric forecasting. The forecasting result is shown in Figure 8.

In Figure 8, the red line is between the orange line and the highest edge of blue line, which means forecasting using asymmetric algorithm is inclined to forecast traffic volume a little higher in a reasonable extent.

However, after timestamp 192 (16 PM), the red line is close to the orange line. By examining the chosen nearest neighbors, it is found that the chosen neighbors are the same to neighbors of normal model. It is inferred that this result is related to the scale of profiles database. Larger scale of profile database may act better in asymmetric forecasting. This inference will be examined in further research. The forecasting criterion indexes lie in the last line of Table 5.

From Table 5, asymmetric algorithm performs worse in MSE/MAE but does better in IMSE. Using the newly designed algorithm, IMSE index drops more than 10%. It means that asymmetric algorithms are more useful when loss of prediction direction is different, which is reasonable in traffic management. So the asymmetric algorithm achieves the aim of design.

5. Conclusions

Using KNN method in short-term traffic forecasting has a history of nearly 20 years. However some detailed operations in KNN still have potential to be enhanced to satisfy the growing need of real ITS systems. In this paper, several limitations of previous researches in this direction are discussed and corresponding methods are proposed and tested.

First, the limitation of Euclidean distance is discussed using a counterexample. Only the absolute distance of every pair of points of two traffic profiles will be calculated and the shape difference is not taken into consideration. To solve this problem, Euclidean distance is reconstructed to contain the stability of the difference between two profiles and use parameters to adjust the balance between absolute difference and shape difference.

Second, Loess technique has been widely used in previous related researches. However the comparison of using raw profiles and smoothed profiles in dynamic KNN method in traffic volume forecasting was seldom made. This research provides strong evidence suggesting that the criterion indexes MSE and MAE drop sharply when Loess is used in KNN method.

Third, asymmetric loss was seldom discussed in previous short-term traffic forecasting researches and it has realistic meaning in real traffic management and service. The asymmetric loss index IMSE is proposed and the asymmetric loss version of KNN algorithm is constructed and tested in later experiment. The results show that IMSE index drops more than 10% and the forecasting value is closer to the upper edge of real traffic volume, which means that the newly designed algorithm can achieve better performance when the cost of forecasting direction has significant difference.

There are still some limitations in this study. The concept of asymmetric loss and newly designed method is still immature and relatively rough, which will be easily affected

if profiles contain extreme values. More refined traffic forecasting methods and algorithms with asymmetric loss can be proposed in this direction in future researches. And more corresponding simulation experiments can be conducted to test the usefulness of asymmetric loss in traffic management and traffic guidance.

Data Availability

The data used to support the findings of this study have been deposited in the <https://github.com/ahorawzy/TFTSA/tree/master/data-raw> repository.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Acknowledgments

The study is funded by scientific and technological support program of the Ministry of Science and Technology of People's Republic of China (2014BAH23F01). WenPeng Zhao (zhaowp@cahs.com.cn) gave useful advice about preprocessing of raw data and R language programming. Chuantao Wang (wangchuantao@bucea.edu.cn) gave useful advice about major and minor revise and undertook some revision work.

References

- [1] M. Cheslow, S. Hatcher G, and M. Patel V, "An initial evaluation of alternative intelligent vehicle highway systems architectures," *System Architecture*, Article ID 92W0000063, 1992.
- [2] E. I. Vlahogianni, J. C. Golias, and M. G. Karlaftis, "Short-term traffic forecasting: overview of objectives and methods," *Transport Reviews*, vol. 24, no. 5, pp. 533–557, 2004.
- [3] M. S. Ahmed and A. R. Cook, "Analysis of freeway traffic time-series data by using box–jenkins techniques," *Transportation Research Record*, no. 722, pp. 1–9, 1979.
- [4] B. L. Smith and M. J. Demetsky, "Multiple-interval freeway traffic flow forecasting," *Transportation Research Record*, no. 1554, pp. 136–141, 1996.
- [5] B. M. Williams, P. K. Durvasula, and D. E. Brown, "Urban traffic flow prediction: application of seasonal autoregressive integrated moving average and exponential smoothing models," *Transportation Research Record*, no. 1644, pp. 132–141, 1998.
- [6] I. Okutani and Y. J. Stephanedes, "Dynamic prediction of traffic volume through Kalman filtering theory," *Transportation Research Part B: Methodological*, vol. 18, no. 1, pp. 1–11, 1984.
- [7] M. Levin and Y. D. Tsao, "On forecasting freeway occupancies and volumes," *Transportation Research Record*, vol. 773, pp. 47–49, 1980.
- [8] G. A. Davis, N. L. Niham, M. M. Hamed, and L. N. Jacobson, "Adaptive forecasting of freeway traffic congestion," *Transportation Research Record*, no. 1287, pp. 29–33, 1991.
- [9] B. L. Smith, B. M. Williams, and R. K. Oswald, "Comparison of parametric and non-parametric models for traffic flow forecasting," *Transportation Research Part C: Emerging Technologies*, vol. 10, no. 4, pp. 303–321, 2002.
- [10] S. Clark, "Traffic prediction using multivariate nonparametric regression," *Journal of Transportation Engineering*, vol. 129, no. 2, pp. 161–168, 2003.
- [11] E. I. Vlahogianni, M. G. Karlaftis, and J. C. Golias, "Short-term traffic forecasting: where we are and where we're going," *Transportation Research Part C: Emerging Technologies*, vol. 43, pp. 3–19, 2014.
- [12] S. Cheng, F. Lu, P. Peng et al., "Short-term traffic forecasting: an adaptive ST-KNN model that considers spatial heterogeneity," *Computers Environment & Urban Systems*, vol. 71, pp. 186–198, 2018.
- [13] F. Guo, J. W. Polak, and R. Krishnan, "Predictor fusion for short-term traffic forecasting," *Transportation Research Part C: Emerging Technologies*, vol. 92, pp. 90–100, 2018.
- [14] Z. Sun and G. Fox, "Traffic flow forecasting based on combination of multidimensional scaling and SVM," *International Journal of Intelligent Transportation Systems Research*, vol. 12, no. 1, pp. 20–25, 2014.
- [15] B. Hamner, "Predicting travel times with context-dependent random forests by modeling local and aggregate traffic flow," in *Proceedings of the 10th IEEE International Conference on Data Mining Workshops (ICDMW '10)*, pp. 1357–1359, Sydney, Australia, December 2010.
- [16] H. Van Lint and C. Van Hinsbergen, "Short-term traffic and travel time prediction models," *Transportation Research E-Circular*, vol. 22, no. 1, pp. 22–41, 2012.
- [17] B. L. Smith and M. J. Demetsky, "Traffic flow forecasting: comparison of modelling approaches," *Journal of Transportation Engineering*, vol. 123, no. 4, pp. 261–266, 1997.
- [18] B. L. Smith, B. M. Williams, and K. R. Oswald, "Parametric and nonparametric traffic volume forecasting," *Transportation Research Board 79th Annual Meeting*, p. 29, 2000.
- [19] F. G. Habtemichael and M. Cetin, "Short-term traffic flow rate forecasting based on identifying similar traffic patterns," *Transportation Research Part C: Emerging Technologies*, vol. 66, pp. 61–78, 2016.
- [20] S. Aghabozorgi, A. S. Shirkhorshidi, and T. Y. Wah, "Time-series clustering - a decade review," *Information Systems*, vol. 53, pp. 16–38, 2015.
- [21] C. Faloutsos, M. Ranganathan, and Y. Manolopoulos, "Fast subsequence matching in time-series databases," in *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 419–429, 1994.
- [22] J. Aach and G. M. Church, "Aligning gene expression time series with time warping algorithms," *Bioinformatics*, vol. 17, no. 6, pp. 495–508, 2001.
- [23] S. Chu, E. Keogh, D. Hart, M. Pazzani et al., "Iterative deepening dynamic time warping for time series," in *Proceedings of the 2nd SIAM International Conference on Data Mining*, pp. 195–212, 2002.
- [24] J. Xia and M. Chen, "A nested clustering technique for freeway operating condition classification," *Computer-Aided Civil and Infrastructure Engineering*, vol. 22, no. 6, pp. 430–437, 2007.
- [25] J. Xia, W. Huang, and J. Guo, "A clustering approach to online freeway traffic state identification using ITS data," *KSCCE Journal of Civil Engineering*, vol. 16, no. 3, pp. 426–432, 2012.
- [26] L. Lin, Y. Li, and A. Sadek, "A k nearest neighbor based local linear wavelet neural network model for on-line short-term traffic volume prediction," *Procedia - Social and Behavioral Sciences*, vol. 96, pp. 2066–2077, 2013.

- [27] C. W. J. Granger, "Prediction with a generalized cost of error function," *Journal of the Operational Research Society*, vol. 20, no. 2, pp. 199–207, 1969.
- [28] C. Pierdzioch, J.-C. Rülke, and G. Stadtmann, "A note on forecasting the prices of gold and silver: asymmetric loss and forecast rationality," *The Quarterly Review of Economics and Finance*, vol. 53, no. 3, pp. 294–301, 2013.
- [29] Y. B. Ahn and Y. Tsuchiya, "Asymmetric loss and the rationality of inflation forecasts: evidence from South Korea," *Pacific Economic Review*, 2017.
- [30] Y. Tsuchiya, "Assessing macroeconomic forecasts for Japan under an asymmetric loss function," *International Journal of Forecasting*, vol. 32, no. 2, pp. 233–242, 2016.
- [31] Y. Zhang and A. Haghani, "A hybrid short-term traffic flow forecasting method based on spectral analysis and statistical volatility model," *Transportation Research Part C: Emerging Technologies*, vol. 43, no. 1, pp. 65–78, 2013.
- [32] L. Lin, J. C. Handley, Y. Gu, L. Zhu, X. Wen, and A. W. Sadek, "Quantifying uncertainty in short-term traffic prediction and its application to optimal staffing plan development," *Transportation Research Part C: Emerging Technologies*, vol. 92, pp. 323–348, 2018.
- [33] W. S. Cleveland and S. J. Devlin, "Locally weighted regression: an approach to regression analysis by local fitting," *Journal of the American Statistical Association*, vol. 83, no. 403, pp. 596–610, 1988.



Hindawi

Submit your manuscripts at
www.hindawi.com

