

HIGHWAY TOLLGATES TRAFFIC FLOW PREDICTION USING MACHINE LEARNING ALGORITHMS

Dai Yirui, Dong Meirong, Lim Chong Seng Hermann

Institute of Systems Science, National University of Singapore, Singapore 119615

ABSTRACT

Along with the fast development of interconnected economy, modern-day traffic networks have become larger and more complex. Ineffective management of these traffic networks have multiple negative impacts on the individual as well as the national economy.

The study aims to use machine learning algorithms such as XGBoost and LightGBM to predict the conditions at traffic bottlenecks such as highway tollgates to reduce congestion, enhance road users' experience and improve infrastructure management.

There are two main issues focused in this study, the first is to forecast the average traveling time through the tollgates and the second is to estimate the traffic volume at the tollgates.

Index Terms— One, two, three, four, five

1. INTRODUCTION

Traffic management is an increasingly difficult task for almost every government and traffic authorities due to rapid globalisation and greater population mobility. Poor traffic management can adversely impact quality of life and economic productivity due to rising fuel consumption, cost of movement and number of accidents. Although the consequences are alarming, traffic congestion remains a prevalent issue in many countries because they have inadequate resources to address the problem. Therefore, minimising traffic problems with optimized costs while satisfying the growing demand for more efficient traffic networks should be the core objective of today's traffic management studies.

To achieve this goal, many have applied advanced technologies including machines learning[1] and deep learning[2] in their works. In this study, the team aims to use machine learning algorithms such as XGBoost and LightGBM to predict traffic flow volume pattern and travel time patterns would occur in the near future(e.g.next 2 hours) on routes to help traffic authorities and commuters to make informed travel decisions, to divert traffic and to address the core objective.

The study uses Knowledge discovery in databases (KDD) 2017 dataset for analysis. The dataset is provided by “Hanzhou Jiaotong Amap”, the task is to predict highway tollgates traffic flow and the results are to be measured and ranked based

on the “Mean Absolute Percentage Error (MAPE)” score.

The dataset is chosen because manual tollgates are usually the bottlenecks of traffic network and can easily cause traffic congestion and overwhelm traffic authorities and commuters. One small improvement made to the bottlenecks could easily be one giant leap for traffic conditions across the route, by Pareto principle.

- The main report is provided in *.tex file
- The reference is provided in *.bib file
- The figures are provided as separate jpg/png files

2. PROBLEM OVERVIEW

The road network topology (as shown in Fig.1) is a directed graph formed by a sequence of interconnected road links. There are 6 routes in the network.

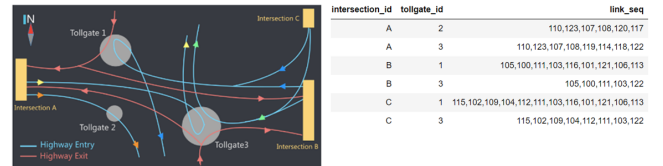


Fig. 1. Visualization of Road Network Topology

Problem 1 is to estimate for every 20-minute time window the average travel time of vehicles for a specific route:

- Routes from Intersection_id A to Tollgate_ids 2 3;
- Routes from Intersection_id B to Tollgate_ids 1 3;
- Routes from Intersection_id C to Tollgate_ids 1 3.

Problem 2 is to predict average tollgate traffic volume for every 20-minute time interval window where tollgates 1 and 3 has entry and exit traffic and 2 has only entry traffic.

There are 5 classes of data provided in the dataset which are links, routes, trajectories, volume and weather.

Mean Absolute Percentage Error (MAPE) is used to evaluate the accuracy.

Travel Time MAPE	Volume MAPE
$MAPE = \frac{1}{R} \sum_{r=1}^R \left(\frac{1}{T} \sum_{t=1}^T \left \frac{d_{rt} - p_{rt}}{d_{rt}} \right \right)$	$MAPE = \frac{1}{C} \sum_{c=1}^C \left(\frac{1}{T} \sum_{t=1}^T \left \frac{f_{ct} - p_{ct}}{f_{ct}} \right \right)$
<p>d_{rt}: actual average travel time for route r during time window t p_{rt}: predicted average travel time for route r during time window t R: the number of routes T: number of to-predict time windows</p>	<p>C: Number of tollgate and direction pairs. In this case $C = 5$. T: Number of time windows in test period. In this case $T = 12$. f_{ct}: actual volume for a particular tollgate & direction pair at a particular time window p_{ct}: predicted volume for a particular tollgate & direction pair at a particular time window</p>

Fig. 2. Visualization of Road Network Topology

3. RELATED WORK

The existing works on KDD 2017 have proposed many solutions for highway tollgates traffic prediction. Some of which are XGboost, Radial Basis Function (RBF) nets and deep neural networks. The most favourable model is XGboost. Therefore, the team uses tuned XGboost models (Chen, 2019)[3] as the main tool to generalize the regression model in this study.

Remarkably, XGboost has received tremendous attention in recent years (Morde Setty, 2019)[4] because of its promising speed and performance (Brownlee, 2016)[5]. The top few winning teams of KDD cup 2017 competition have also recommended this learning algorithm for their predictions.

An alternative approach is to use a set of continuous distributions to approximate the regression via RBF. By using RBF, the idea is to combine (or sum) a set (kernels) of curves (Gaussian distributions) with various control parameters (std mean) to generalize the regression model (Deshpande, 2017)[6].

More complex approaches that can stack and vote from different model (LASSO, GBDT, ADABOOST and Random-Forest etc) layers and deep neural networks such as (CNN-LSTM-Attn) and Temporal-Spatial-LSTM (TSLSTM) are also examined. However, details including how the models are implemented and their parameters are not provided by the author. This complex approach managed to reach a MAPE value around 0.104 for volume prediction.

In the end, due to time constraints, the team decided to follow the proposal submitted by the 1st place team[7] in KDD cup 2017 competition to implement XGboost.

4. PROPOSED APPROACH

With reference from KDD cup 2017 winning team's proposal, the team has proposed following step by step approach:

Average Travel Time Prediction

Phase	Action taken	Action Category
Study-1	Aggregate raw trajectories data based on 20-mins time window	Coding
Study-2	Plot volume of all time window at a particular day (1st day)	Exploring
Feature-1	Process and select feature from link and route tables for data training and prediction (link_counts, route_length, double_in_link_count, double_out_link_count, min_width)	Exploring
Feature-2	Process the features from trajectories table and clean outliers (add time lag for 2 hours, because based on other research works, the average traveling time from previous 2 hours have more significant impacts on the data)	Exploring
Feature-3	Process and understand the data from weather table	Exploring
Feature-4	Combine all the tables	Exploring
Base-1	Implement baseline approach	Coding
Improve-1	Test the baseline with different features and parameters to improve accuracy	Experimenting

Average Volume Prediction

Phase	Action taken	Action Category
Study-1	Aggregate raw volume data based on time window (20 mins)	Coding
Study-2	Plot volume of all time window at a particular day (1st day)	Exploring
Study-3	Look for irregularities and try to explain the reason. Exclude irregularities if not helpful for prediction	Exploring
Study-4	Plot volume by day for the whole training period	Exploring
Study-5	Look for irregularities and try to explain the reason. Exclude irregularities if not helpful for prediction	Exploring
Study-5	Look for irregularities and try to explain the reason. Exclude irregularities if not helpful for prediction	Exploring
Study-5	Look for irregularities and try to explain the reason. Exclude irregularities if not helpful for prediction	Exploring
Study-6	Plot top (4) frequent ‘travellers’ and see if there are similar travelling patterns (Explain for peak and trough)	Exploring
Base-1	Implement baseline approach	Coding
Base-2	Experiment with denoised data (excluded holiday periods), calculate baseline MAPE as benchmark	Experimenting
Base-3	Parse and join weather data with volume data	Coding
Base-4	Further experiment with weather data, check if MAPE improve	Experimenting
Improve-1	Aggregate raw volume data based on time window (20 mins) but include ‘has _{etc} ’ as feature	Coding
Improve-2	Implement function to calculate mean volume values of each time window for past n (5) days	Coding
Improve-3	Exclude irregular high volumes at midnights	Experimenting
Improve-4	Experiment with volume vs. log(volume)	Experimenting
Improve-5	Experiment with different evaluation metrics: ‘mae’	Experimenting

5. EXPERIMENTAL RESULTS

For task 1, the travel time on a normal day is heavily influenced by route length. A longer route distance will usually result in longer travel time.

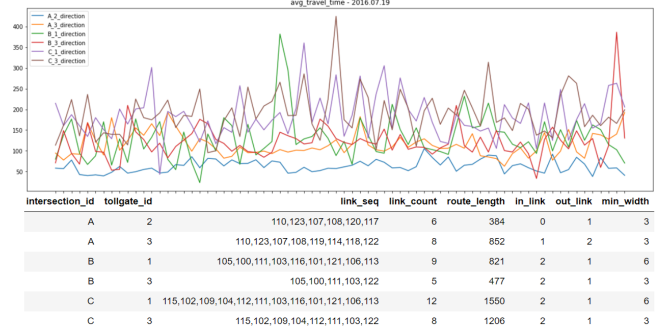


Fig. 3.

The MAPE result is 0.242141.

The raw data is presented in a way that timestamp is taken when a vehicle passed through a particular toll gate at a certain time:

	time	tollgate_id	direction	vehicle_model	has_etc	vehicle_type
0	2016-09-19 23:09:25	2	0	1	0	NaN
1	2016-09-19 23:11:53	2	0	1	0	NaN
2	2016-09-19 23:13:54	2	0	1	0	NaN
3	2016-09-19 23:17:48	1	0	1	1	NaN
4	2016-09-19 23:16:07	2	0	1	0	NaN

Fig. 4.

To estimate the volume of target time window, firstly, individual vehicle count should be summed up to every 20 minutes time interval. The team has used our reference code to perform the aggregation.

	tollgate_id	time_window	direction	volume
0	3	[2016-09-19 00:00:00,2016-09-19 00:20:00)	0	17
1	3	[2016-09-19 00:00:00,2016-09-19 00:20:00)	1	181
2	1	[2016-09-19 00:00:00,2016-09-19 00:20:00)	0	13
3	1	[2016-09-19 00:00:00,2016-09-19 00:20:00)	1	140
4	2	[2016-09-19 00:00:00,2016-09-19 00:20:00)	0	2

Fig. 5.

5.1. A dive into the data:

The volume pattern of each tollgate direction pair on the first day is shown in below chart. Notice that tollgate 2 has only 1 direction thus leave us total 5 tollgate direction pairs. From

the chart, it is clear to see that there are some peaks at rush hour periods.

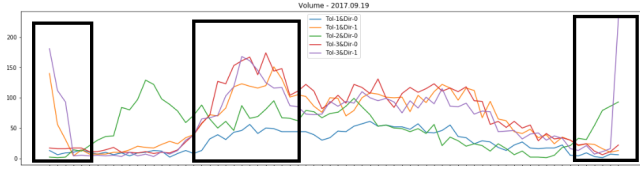


Fig. 6.

The other observation is that there is high volume for certain tollgates at midnight. However, by comparing with the overall data at all times, it is difficult to conclude meaningful features for those passing at midnight. In this work, the team decided to exclude the data from midnight while training since they are not relevant to the rush hour data but add much noise.

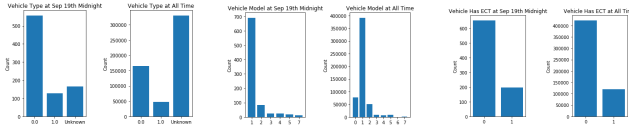


Fig. 7.

To study the commuting behavior of individual vehicle, the team has selected the top 4 vehicles by counting the number of appearances made at any tollgate. It is noticeable that each of these vehicles has its own unique travelling patterns. It is not feasible to rationalize peak periods by only looking at these patterns.

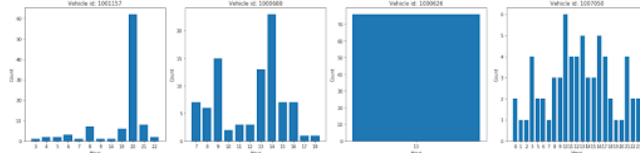


Fig. 8.

The volume of all tollgate direction pairs on a daily basis is shown in below chart. It is clear to see that the volume pattern between Sep 30th and Oct 7th is very different from the rest. This behavior is aligned with the fact that this time period is a long holiday in China, and thus many people will not travel in the same way as in working days. In this case, the target time period that requires our prediction is not in holiday period. Therefore, in order to give a more precise prediction, it is important to only look at those non-holiday data.

5.2. Incorporating with weather data as features:

Weather conditions can significantly influence traffic operations (Hani et al., 2009)[8]. To better predict traffic flow, the

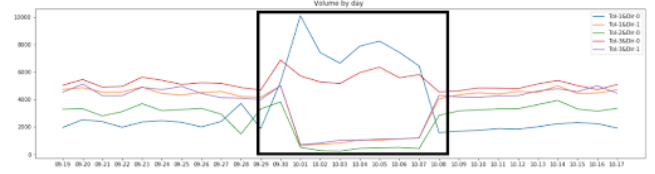


Fig. 9.

team has done some experiments incorporating with weather data. With reference from github code, the team has done some data preprocessing to join weather data with volume data:

- Align time window: weather - 3 hours vs. volume - 20 mins
- Convert 'wind_direction' to numeric (360 degrees)
- Fill in 'NaN' fields with overall mean of corresponding feature

Below table shows the comparison between MAPE with and without weather data:

	Overall MAPE	Rush Hours MAPE
Weather (N)	0.268922	0.191247
Weather (Y)	0.429293	0.187615

5.3. Further Denoising Data Adding 'has_etc' feature:

One of the findings from previous studies is that irregular high volumes happen often at midnight. These irregularities are hard to explain and not helpful for prediction at target time period. In addition, as proposed by the winning team, features including 'has_etc' may help to improve prediction accuracy. This makes sense because vehicles that have ETC need not wait for fee collection at the tollgate. Below table shows the comparison between MAPE before and after:

	Overall MAPE	Rush Hours MAPE
MidN(Y) & ETC(N)	0.429293	0.187615
MidN(N) & ETC(Y)	0.304405	0.156712

5.4. Adding in Statistical Features:

As proposed by KDD cup 2017 winning team, statistical features (mean, min, max etc) of corresponding time window of last n days, overall statistics of last n days, and rush hours statistics may help improve prediction accuracy. To experiment the impacts of statistical feature, the team has calculated the mean of corresponding time window of last 5 days for all training/testing instances. The calculation process can

consume significant time and resources, so caution should be taken when input large data. Below table shows the comparison between MAPE before and after adding the statistical feature:

	Overall MAPE	Rush Hours MAPE
Stats (N)	0.304405	0.156712
Stats (Y)	0.202389	0.137198

Visualization of final predicted volume vs actual volume in target time periods:

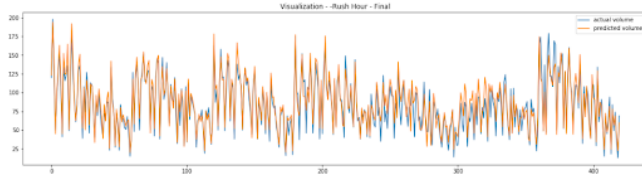


Fig. 10.

6. CONCLUSIONS

By integrating above findings and techniques, MAPE is improved from 0.191247 to 0.137198 for the target prediction period. The accuracy itself ranks top 10 in KDD cup 2017 volume competition. More importantly, this work can serve as a prototype that can incorporate forecasting weather data and historical statistics for real time traffic estimation. However, many other proposed approaches that could potentially improve the accuracy are not included in this work, such as using more statistics, step by step modelling (use of latest available data to update model), ensembled and weighted systems (XG-Boost, NN, and LightGBM). Some approaches includes use of logarithm for labelling, and M-estimator as objective function are explored but not applied due to reduction of accuracy and lack of implementation details.

7. REFERENCES

- [1] Jaros law Rzeszotko and Sinh Hoa Nguyen, *Machine Learning for Traffic Prediction*, 2011.
- [2] Kaishun Wu Zhidan Liu, Zhenjiang Li and Mo Li, *Urban Traffic Prediction from Mobility Data Using Deep Learning*, 2018.
- [3] Zeyu Chen, “KDD CUP 2017 code,” <https://github.com/chenzeyuczy/KDD2017>, [Online; accessed October 30, 2019].
- [4] Venkat Anurag Setty Vishal Morde, “XG-Boost Algorithm: Long May She Reign!,” <https://towardsdatascience.com/https-medium-com-vishalmorde-xgboost-algorithm-long-she-may-rein-edd9f99be63d>, [Online; accessed October 30, 2019].
- [5] Jason Brownlee, “A Gentle Introduction to XGBoost for Applied Machine Learning,” <https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning/>, [Online; accessed October 30, 2019].
- [6] Mohit Deshpande, “Using Neural Networks for Regression: Radial Basis Function Networks,” <https://pythonmachinelearning.pro/using-neural-networks-for-regression-radial-basis-function-networks/>, [Online; accessed October 30, 2019].
- [7] Huan Chen Ke Hu, Pan Huang and Peng Yan, *KDD CUP 2017 Travel Time Prediction Predicting Travel Time – The Winning Solution of KDD CUP 2017*, 2017.
- [8] Jiwon Kim Hani S. Mahmassani, Jing Dong, Roger B. Chen, and Byungkyu (Brian) Park, *Incorporating Weather Impacts in Traffic Estimation and Prediction Systems*, 2009.

$$B_{r,c} = \sum \{f(i,j) | (i,j) \in \Omega_{r,c}\}. \quad (1)$$

$$\sum_x = a + b + \hat{c}, \quad (2)$$

An inline equation is $a + b = c$. An example of two-column figure is provided in Figure 11, and the single-column figures is provided in Figure 12.

Table 1. The performance comparison.

Approach	Ref. [?]	Ref. [?]	Proposed approach
Metric A	0.8181	0.9171	0.9616
Metric B	0.8236	0.7654	0.8615



INSTITUTE OF SYSTEMS SCIENCE

Fig. 11. Test figure (two-column).



Fig. 12. Test figure (single column).