# A dynamic OD prediction approach for urban networks based on automatic number plate recognition data

Liu, Jing; Zheng, Fangfang; van Zuylen, Henk J.; Li, Jie

**Important note**
To cite this publication, please use the final published version (if applicable).
Please check the document version above.

22nd EURO Working Group on Transportation Meeting, EWGT 2019, 18-20 September 2019, Barcelona, Spain

# A dynamic OD prediction approach for urban networks based on automatic number plate recognition data

Jing Liu[a,d], Fangfang Zheng*[a,d], Henk J. van Zuylen [a,b,c,d], Jie Li[b]

[a] School of Transportation and Logistics, Southwest Jiaotong University, Western Hi-tech Zone Chengdu, Sichuan 611756, P.R.China
[b] Civil Engineering College, Hunan University, Lushan South Road, 410082 Changsha, Hunan Province, P.R. China
[c] Transport and Planning Department, Delft University of Technology, P. O. Box 5048, 2600 GA Delft, the Netherlands
[d] National Engineering Laboratory of Integrated Transportation Big Data Application Technology, Southwest Jiaotong University, Western Hi-tech Zone Chengdu, Sichuan 611756, P.R. China

## Abstract

OD flows provide important information for traffic management and planning. The prediction of dynamic OD matrices gives the possibility to apply anticipatory traffic management measures. In this paper, we propose an OD prediction approach based on the data obtained by Automated Number Plate Recognition (ANPR) cameras. The principal component analysis (PCA) is applied to reduce the dimension of the original OD matrices and to separate the main structure patterns from the noisier components. A state-space model is established for the main structure patterns and the structure deviations, and is incorporated in the Kalman filter framework to make predictions. We further propose three K-Nearest Neighbour (K-NN) based long-term pattern recognition approaches. The proposed approaches are validated with field ANPR data from Changsha city, P.R. China. The results show that the observed OD flows can be accurately predicted by our proposed approaches. Which prediction method performs best depends on the quality of the available data: for regular, periodic OD matrices the Kalman filter is better, for irregular OD matrices the pattern recognition that looks at different time periods in the historical data, gives better results.

*Keywords:* OD matrix prediction; principal component analysis; state-space kalman filter model; pattern recognition

## 1. Introduction

Traditionally, a dynamic origin-destination was not directly observable from traffic data. Several methods have been developed to derive origin-destination from traffic counts, historical data and partial actual data, such as origin

---

\* Corresponding author. Tel.: +86 18702828126
E-mail address: fzheng@swjtu.cn

destination patterns from probe vehicles. Antoniou et al. (2016) and Djukic (2014) give an extensive overview of the methods that have been developed in the past 35 years to estimate origin destination matrices from traffic counts.

New data acquisition methods make it possible to collect information about the origin destination flows in a road network. Data from mobile phones, probe vehicles (e.g. equipped with GPS devices), Automated Number Plate Recognition (ANPR) cameras (Antoniou, Ben-Akiva & Koutsopoulos 2004, Rao et al. 2018), Bluetooth scanners (Barcelo et al. 2010), data from tolling stations and video recordings from high altitude (e.g. drones) give information to derive dynamic Origin Destination (OD) matrices. For urban areas the methods to obtain dynamic OD matrices are more limited than for freeway networks. Mobile phone and Bluetooth data are less suitable in urban road networks (Li et al. 2011) and tolling station are seldomly present in urban areas. Especially ANPR cameras can give very accurate data about traffic flows through an urban network (Rao et al. 2018).

The step from the observation of the present OD matrix to a future one is methodologically challenging: There is a certain regularity in traffic patterns, like peak and off-peak flows, but for the prediction of the OD matrix for a short-term future a more sophisticated method is necessary.

In general, dynamic OD prediction methods can be classified into parametric and non-parametric. Some widely used parametric methods include random walk models (Cremer & Keller, 1987), which considers only the OD matrix in the previous time step with a random correction, ARIMA models (Williams & Hoel, 2003), and State-Space Kalman Filter models (Okutani & Stephanedes, 1984, Zhou & Mahmassani 2007). In these prediction models, several features such as the average OD demand and the deviation between the OD demand and its historical average are used as state vectors. Okutani & Stephanedes (1984) applied Kalman filtering to obtain flow predictions from the weighted average of historical data and actual measurements. This can be applied also for the prediction of origin destination matrices, on the condition that history repeats itself: if traffic states at a present situation is similar to the traffic state on the same time period on previous days.

In a more sophisticated way, the relationship among these state vectors can be described linearly by q-order autoregressive model (Ashok & Ben-Akiva, 2000) or non-linearly by q-order recursive model (Zhou & Mahmassani, 2007).

Non-parameter methods include pattern recognition techniques such as Artificial Neural Networks (Nair, Liu, Rilett, & Gupta, 2001), K-nearest Neighbor (Clark, 2003, Zhang et al. 2013), Tensor Decomposition (Ren & Xie, 2017) and Bayesian Networks (e.g. Castillo et.al. 2008).

When dealing with large traffic networks, traffic demand prediction is rather computationally intensive. Therefore, a dimension reduction technique may be required. Principal component analysis (PCA) (Wold et al. 1987, Lakhina et al. 2004, Djukic 2014, Djukic et al. 2012) and Factor Analysis (FA) (Ma, Zhou & Antoniou, 2018 and Pragash et al. 2017) are widely used dimension reduction techniques. The difference between these two techniques is that the PCA approach extracts the linear combination of the original variables, while the FA method decomposes the original variables.

In this paper, we compare methodologies to predict OD demand based on historical observations. In most research papers it is assumed that the OD matrix can be predicted from historical data, apart from some error (e.g. Pragash et al. 2017). This so-called error might be due to a systematic deviation between OD matrices at time intervals on different days, for instance by events that do not occur every day at the same time. In order to deal with this possibility, we develop a method to identify similar traffic patterns on different times on the day.

We apply the proposed approaches to the ANPR data from Changsha city, P.R. China. The results show that even with limited history data sets, the proposed approaches can provide rather accurate prediction and have good transferability as well.

## 2. Methodology

In this paper we develop a methodology to make prediction of OD demand from traffic observations obtained from ANPR facilities. The OD matrix is directly observed and no estimation of the matrix from traffic volumes is needed. The focus of this paper is on prediction where the input is historical OD data. We apply these matrices as historical data in our prediction approach which is described briefly as follows:

- We use the PCA algorithm to reduce the dimension of the historical and predicted OD matrix and transform the OD demand into significant structure patterns, deviations from the structure and stochastic patterns;

- We establish a state-space model for significant structure patterns and structure deviations, and make a prediction basis from a historical data set using a Kalman filter predictor. In the meantime,
- We further develop K-Nearest Neighbours based pattern recognition methods to identify and predict structure patterns, and structure deviations.

In other words, we make a prediction under both shorter-term and longer-term, considering both the random trend as well as the latent pattern, in order to get a better prediction performance, also for the case that the historical observations are not applicable to the present and the future traffic state.

## 3. Data acquisition

The traffic data used in this study has been obtained from ANPR cameras in the city of Changsha, the capital of Hunan province in the P.R. China. Many intersections in Changsha are provided with ANPR cameras. Each camera can observe one lane and register the number plates of vehicles that pass the stop line of the intersection. The number plates of taxis can be separated from ordinary vehicles and OD matrices for taxis and other traffic can be separated (Sbaï et al. 2017). The moments of the passing number plates are registered in seconds. The Number plates of three days were collected for further analysis: 20, 21 and 22 April 2015. From the available ANPR data a selection was made of 22 intersections.



Figure 1 The road network of the CBD of Changsha. The intersections with bold numbers have ANPR

## 4. Prediction models for the OD matrix

### 4.1. Data reduction

In order to reduce the data of the OD matrix a principal component analysis was executed. That is shown in the upper right part of Figure 2. First of all, the original OD matrix is centralized by column to derive the average value for each column and the covariance matrix of the centralized matrix. Secondly, we calculate the eigenvalue and the eigenvector of the covariance matrix and determine the number of principal component k. We select the first k column of the eigenvector as the principle components. In this way, we have transformed the high-dimensional OD matrix to the low-dimensional coordinate. Finally, we can calculate the score of each principal component by multiplying the centralized OD matrix with k principal components (first k column of the eigenvector).

Each whole OD matrix contains many elements but there are regularities in the data: correlations exist and it is possible to reduce the whole matrix to a few components that can represent the whole matrix in such a way that the matrix can be reproduced from a limited number of components. It appeared in our OD matrices that only 3 principal components were needed to explain 76% of the variation in the OD matrix over the whole day, 9 components explain 85%. In this paper we use 5 components for the prediction procedure, explaining 83%. The temporal behaviour of the score values of the principal components is rather regular for the first 3 and noisier for the higher components (see Figure 3).
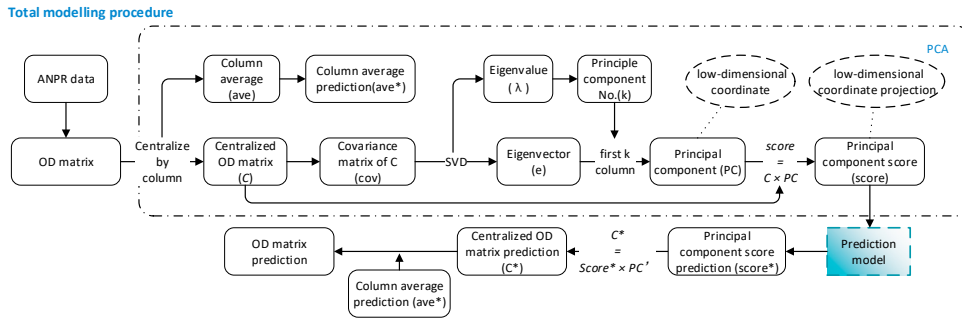
**Total modelling procedure**



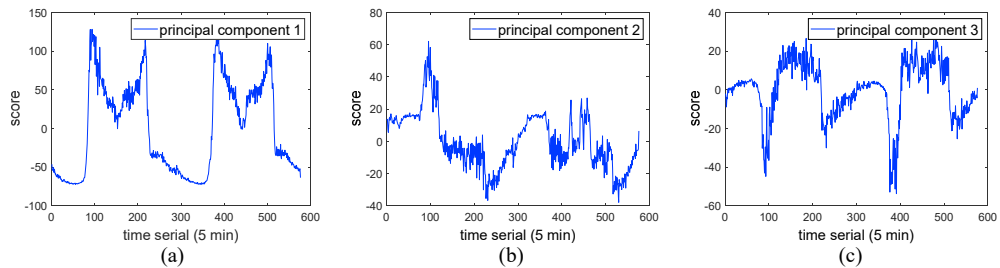Figure 2 The generic flow chart for the prediction procedure



Figure 3 The magnitude of the first 3 principal components as a function of the time of the day for 20 and 21 April 2015: (a) component 1; (b) component 2; (c) component 3

In order to make the temporal behaviour of the first five principal components more regular and eliminate the irregular outliers, we approximate the principal components by third order polynomials. We divided the total number of time steps (each time step is 5min) into 6 segments within which a good fit can be obtained between the polynomial and the historical values of the principal components. The segments are time step 0 to 75 (early period 00:00 to 06:15), 75 – 90 (first part of the morning peak, 06:15 to 07:30), 90 – 155 (second half morning peak and morning off-peak, 07:30 to 12:55), 155 – 220 (off-peak period and early evening peak, 12:55 to 18:20), 220 – 230 (evening peak, 18:20 to 19:10), 230 – 288 (evening period, 19:10 to 24:00). Figure 4 shows the polynomial fit of the scores of the first three principal components.
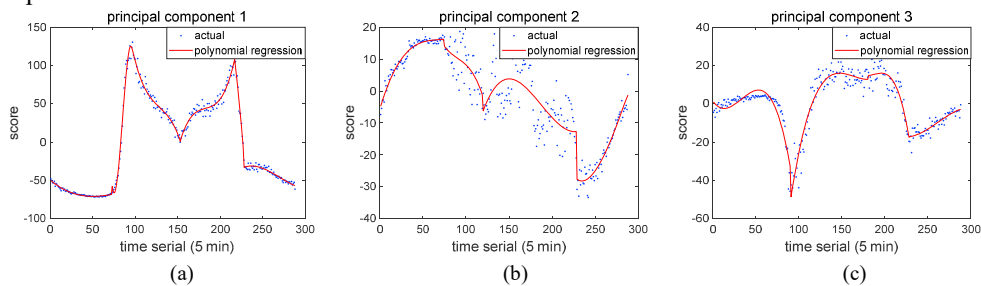


Figure 4 The first 3 Principal Components fitted with polynomials: (a) component 1; (b) component 2; (c) component 3

### 4.2. Kalman filter

The first prediction model uses a Kalman filter. The basic assumption is that OD matrices show every day similar dependencies on the time of the day with some structural deviations. These structural deviations are estimated using a Kalman filter (e.g. Okutani & Stephanedes, 1984, Zhou & Mahmassani 2007). The prediction procedure is shown

in Figure 5 where we apply the polynomial fitting approach to the historical data (as reduced by PCA) to derive the feature pattern for each main component.
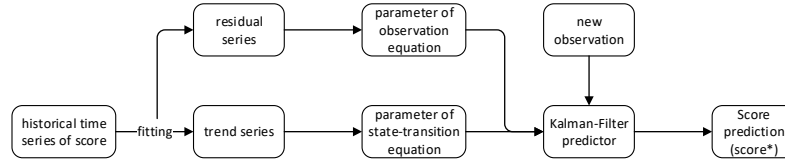


Figure 5 Scheme of the Kalman filtering method

The principal components $C_{t+1}$ are estimated from the historical data model and the observed scores. The historical data model is the polynomial approximation of the principal component scores.

$$C_{t+1} = \Phi_t C_t + \omega_t \tag{1}$$

$$y_t = H_t C_t + \upsilon_t \tag{2}$$

The random variables $\omega_t$ and $\upsilon_t$ represent the error in the prediction model and the observation respectively. The matrix $\Phi_t$ relates the future value of the score value to the value in the previous time-step. The matrix $H_t$ is the transformation of the observed OD matrix $y_t$ to the scores of the principal components. Using this state-space model and a Kalman filter predictor (Zhou and Mahmassani 2007, Djukic 2014 Pragash et al. 2017), we made the prediction for the first component as shown in Figure. It can be seen that the prediction results are quite accurate for the scores of the first component. Similarly, we apply the proposed approach to the second to the fifth components and combine them for the final prediction, as described in the next chapter.
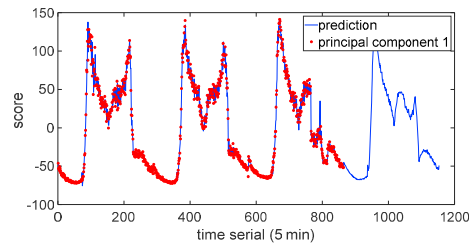


Figure 6 The first component prediction by Kalman Filter predictor

### 4.3. K-NN prediction models

The models 2 to 4 are developed are based on the method of the K-Nearest Neighbours (Altman 1992): we try to find patterns in the temporal behaviour in the historic files which are not necessarily at the same time of the day as in the present OD matrix. The basic idea is that the historical data sets of features are identified. Each member of the set represents a certain historic OD matrix. We distinguish state vectors, traffic states preceding a certain moment, and label vectors, states after that moment. In order to better capture the variation trend in the whole historical data set, both state vectors and label vectors are not exact score values but derived as the difference of the PCA scores between the current time step and the previous time step. We call these vectors as state trend vectors and label trend vectors, respectively. In model 2, the simple K-nearest neighbour algorithm finds the state trend vectors that has the shortest distance to the present traffic state and uses that to predict the next OD matrix. In our calculation the Euclidian distance is used. The new prediction is made by distance-weighted $k$ label trend vectors corresponding to $k$ selected state trend vectors. Figure 7(a) illustrates the detailed prediction procedure using model 2.

In order to better capture the trend pattern in the training data sets, we apply the polynomial fitting approach as described in Model 1 to the state trend vectors and the label trend vectors. The feature residual sets are calculated as the difference between the original trend vectors (state and label) and the fitted pattern vectors. Figure 7(b) shows the prediction procedure of Model 3 and 4. The main difference between Model 3 and Model 4 is the prediction of the

variance of residual (VoR). In model 3, we assume the residual is Gaussian white noise with zero mean, and the residual vector at time step $i+1$ is subject to the same distribution as that at the previous time step $i$. The optimal estimation of the VoR in the historical training data is derived by applying the polynomial fitting approach. We simply consider the VoR at the future time step $i+1$ is the same as that at the previous time step $i$. While in model 4, we again apply the K-NN method as described in Model 2 to predict the VoR.

We distinguish two data sets: the 'feature set' consisting of $d$ past values of the scores and a 'label' data set of $D$ future values of the scores. It can be compared with the 'tail' and the 'whiskers'. For every time step $i$, the tail consists of values of the time periods $i-d, i-d+1, \ldots, i$. The 'whiskers' are the values of the $D$ time steps after the moment $i$. In Figure 8(a) the 'whiskers' are shown for 6 time-steps ahead. For 1 time-step ahead the predictions are much more accurate at all times (Figure 8(b)).
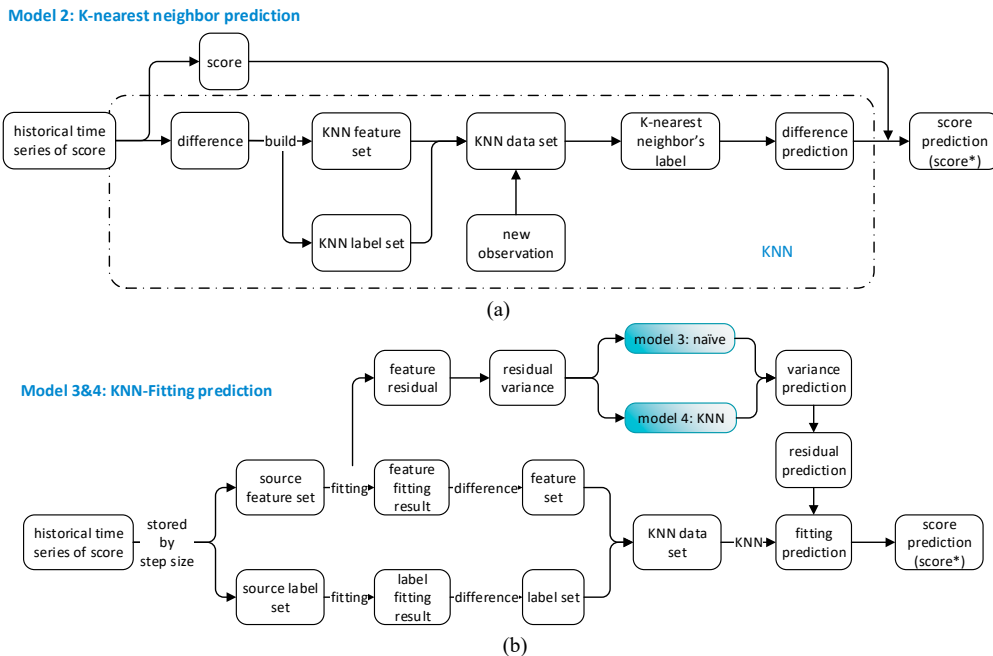


(a)

(b)

Figure 7 The prediction process of K-NN models



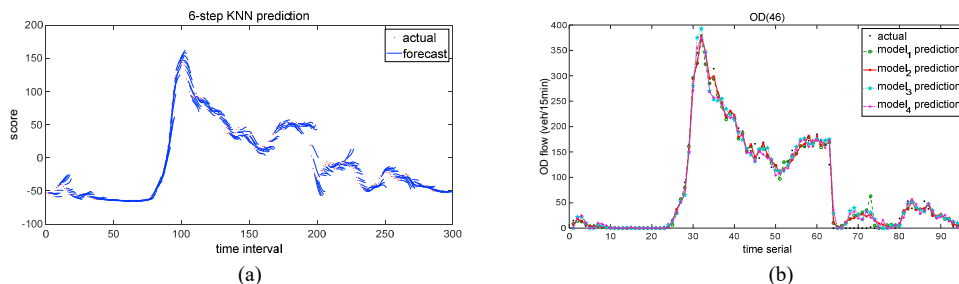(a)                                             (b)

Figure 8 (a) The prediction of 6 steps ahead for the scores of the first Principal Component; (b) predictions according to the 4 methods for one particular OD pair (15 minutes time steps).

## 5. Evaluation

We use three performance metrics, namely, Mean Absolute Error (MAE), Root Mean Square Error (RMSE) and Mean Absolute Percentage Error (MAPE), to measure the prediction accuracy. As shown in Figure 9 and Table 1, model 1 gives more accurate prediction results compared with models 3 and 4 for the morning peak, but in the

afternoon peak, when the traffic pattern is less regular, all K-NN methods perform better. The relatively simple method 2 gives better results than methods 3 and 4.



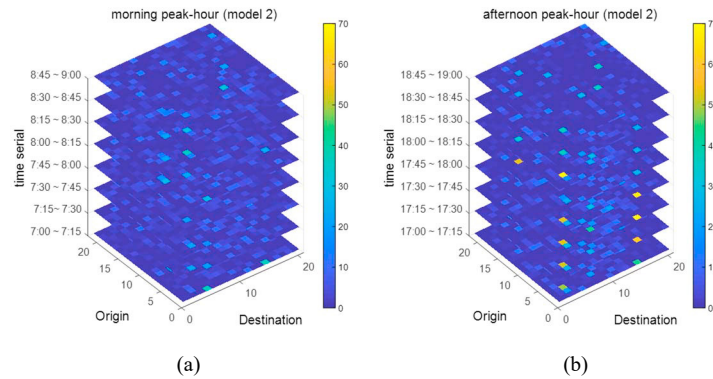(a)                                                    (b)

Figure 9 Prediction performance of Model 2 in terms of Mean average error (MAE): (a) morning peak-hour; (b) afternoon peak-hour

Table 1 Errors of the estimation of flow 46 in the morning and evening peak (compute by mean average error, root mean squared error and mean absolute percent error)

| OD pair 46 | Morning peak | | | | Evening peak | | | |
|---|---|---|---|---|---|---|---|---|
| | Model 1 | Model 2 | Model 3 | Model 4 | Model 1 | Model 2 | Model 3 | Model 4 |
| MAE | 12.17 | 11.96 | 19.49 | 14.01 | 23.74 | 18.98 | 20.81 | 20.60 |
| RMSE | 14.87 | 13.59 | 25.92 | 22.01 | 30.32 | 20.22 | 24.64 | 21.92 |
| MAPE | 5.87 | 6.84 | 10.59 | 8.04 | 25.61 | 20.41 | 22.38 | 22.12 |

Due to the different characteristics of the morning and evening peak data, the flow pattern of the morning peak is more regular, while the pattern of the evening peak on the third day (22 April 2015) is significantly different from the historical data (20 and 21 April 2015). Therefore, the prediction accuracy of the evening peak is not as good as the morning peak. Time series decomposition combined with the Kalman filtering prediction method can well explore the structural time variation of historical data, establish a state transition model based on the historical data trend, and correct the model with observation. This method has a rather high prediction accuracy under the condition that the future development trend of the data is not much different from the state transition relationship present in the historical situation. Otherwise, the prediction for a certain time step will rely heavily on the observation information of the previous step, so that the result shows a significant one-step lag.

The data-driven pattern recognition method (K-NN models) is in general slightly less accurate than the Kalman filter model in the case of data showing obvious periodic variation patterns, but it is more adaptable when dealing with abnormal conditions. By defining the identification object, searching for the historical trend that is close to the current and predicted matrix according to the subsequent development of the history, the recognition range is not limited to the historical period, but can jump out of the limitation of periodicity, globally identify the possible trend, and thus generate a good prediction result.

## 6. Conclusions

Under normal conditions the OD flow shows a regular pattern over the day. In this situation, the Kalman filtering method based on the time series analysis theory will give more accurate prediction results. However, for the anomalies that do not conform to the historical pattern, the prediction accuracy of the Kalman filtering method is significantly reduced. At this time, the K-NN method based on pattern recognition can respond to the abnormal changes more quickly, and the recognition area is not limited to the same period of history. It is a global search for similar trends at other time periods and gives more reliable predictions.

Some considerations for the possible improvement of the methodology are as follows:

- The current algorithm uses Euclidean distance as an indicator for finding neighbouring states. In the future, different distance measures can be considered such as e.g. the information measure.
- An anomaly detection module can be added to the prediction method to distinguish the normal mode from the abnormal condition and use different models for the prediction.
- Using the K-NN method to predict the search process when the amount of data is large may be time consuming, and the historical state library can be clustered to reduce unnecessary searches.

## Acknowledgements

## References

Altman, N. S. 1992. An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician*. 46 (3): 175–185.

Antoniou, C., Barceló, J., Breen, M., 2016 Towards a generic benchmarking platform for origin–destination flows estimation/updating algorithms: Design, demonstration and validation, *Transportation Research Part C: Emerging Technologies, vol. 66*, May 2016, pp. 79-98

Antoniou, C., M. Ben-Akiva, & H. N. Koutsopoulos. 2004 Incorporating Automated Vehicle Identification Data into Origin–Destination Estimation. Transportation Research Record: Journal of the Transportation Research Board, No. 1882, 2004, pp. 37–44.

Ashok, K., & Ben-Akiva, M. E. 2000. Alternative approaches for real-time estimation and prediction of time-dependent Origin-Destination flows, *Transportation Science, 34*(1), 21-36. doi:10.1287/trsc.34.1.21.12282.

Barceló, J., Montero, L., Marqués, & L., Carmona, C. 2010 Travel Time Forecasting and Dynamic Origin-Destination Estimation for Freeways Based on Bluetooth Traffic Monitoring, *Transportation Research Record: Journal of the Transportation Research Board, vol. 2175,* pp 19-27, doi: 10.3141/2175-03

Castillo, E., Menéndez, J. M., & Sánchez-Cambronero, S. 2008. Predicting traffic flow using Bayesian networks, *Transportation Research Part B: Methodological, 42*(5), 482-509. doi:https://doi.org/10.1016/j.trb.2007.10.003

Clark, S. 2003. Traffic prediction using multivariate nonparametric regression. *Journal of Transportation Engineering, 129*(2), 161-168.

Cremer, M. & Keller. H. 1987. A new class of dynamic methods for the identification of origin-destination flows. *Transportation Research B: Methodological 21*(2), 117-132.

Djukic, T*., Dynamic OD demand estimation and prediction for dynamic traffic management, Doctoral Dissertation TUDelft, 2014, ISBN: 9789055841790,* https://repository.tudelft.nl/islandora/object/uuid:ab12d7a7-e77b-424d-b478-d58657f94dd1

Djukic, T., Flötteröd, G., van Lint, H., & Hoogendoorn, S., 2012, Efficient real time OD matrix estimation based on Principal Component Analysis, *Proceedings of the 15th International IEEE Conference on Intelligent Transportation Systems (ITSC2012)*, pp. 115-121

Lakhina, A., Papagiannaki, K., Crovella, M., Diot, C., Kolaczyk, E. D., & Taft, N. 2004. Structural analysis of network traffic flows. ACM SIGMETRICS Performance Evaluation Review *32*(1), pp. 61-72. doi:10.1145/1012888.1005697

Li, J., Van Zuylen, H., Liu, C., & Lu, S. ,2011, Monitoring travel times in an urban network using video, GPS and Bluetooth Procedia - Social and Behavioral Sciences, 20, pp. 630-637.

Ma, T., Zhou, Z., & Antoniou, C. 2018. Dynamic factor model for network traffic state forecast. *Transportation Research Part B: Methodological, 118*, 281-317.

Nair, A. S., Liu, J.-C., Rilett, L., & Gupta, S. 2001. Non-linear analysis of traffic flow.  ITSC 2001 IEEE Intelligent Transportation Systems. Proceedings (Cat. No.01TH8585)681-685. doi:10.1109/ITSC.2001.948742

Okutani, I., & Stephanedes, Y.J. 1984. Dynamic prediction of traffic volume through Kalman filtering theory. *Transportation Research Part B: Methodological, 18*(1), 1-11. doi:https://doi.org/10.1016/0191-2615(84)90002-X

Pragash, A.A., Seshadri, R., Antoniou, C., Pareira, F.C., & Ben-Akiva, M.-E. 2017. Reducing the Dimension of Online Calibration in Dynamic Traffic Assignment Systems*, Transportation Research Record: Journal of the Transportation Research Board, No. 2667, pp. 96–107*

Rao, W., Wu, Y.-J.,Xia.J., Ou, J. & Kluger, R. 2018, Origin-destination pattern estimation based on trajectory reconstruction using automatic license plate recognition data. Transportation Research Part C 95  29–46

Ren, J., & Xie, Q. ,2017, Efficient OD Trip Matrix Prediction Based on Tensor Decomposition*. Paper presented at the 2017 18th IEEE International Conference on Mobile Data Management (MDM).*

Sbaï, A., van Zuylen, H. J., Li, J., Zheng, F. 2017. Estimation of an Urban OD Matrix Using Different Information Sources. In: Gervasi O. et al. (eds) *Computational Science and Its Application*s – ICCSA 2017. Lecture Notes in Computer Science, vol 10405. Springer, Cham

Williams, B. M., & Hoel, L. A. 2003. Modeling and Forecasting Vehicular Traffic Flow as a Seasonal ARIMA Process: Theoretical Basis and Empirical Results. *Journal of Transportation Engineering  Volume 129 Issue 6129*(6), 664-672.

Wold, S., Esbensen, K., & Geladi, P. 1987. Principal component analysis. Chemometrics and Intelligent Laboratory Systems, 2(1), 37-52.

Zhang, L. , Liu, Q. , Yang, W. , Wei, N. , & Dong, D. 2013. An improved k-nearest neighbor model for short-term traffic flow prediction. Procedia - Social and Behavioral Sciences, 96, 653-662.

Zhou, X., & Mahmassani, H. S. 2007. A structural state space model for real-time traffic origin-destination demand estimation and prediction in a day-to-day learning framework. *Transportation Research Part B: Methodological, 41*(8), 823-840.