### SHORT-TERM CONGESTION FORECASTING: TIME SERIES VERSUS FUZZY SETS

# G Huisken

# Department of Civil Engineering & Management, University of Twente, P O Box 217, 7500 AE, Enschede, The Netherlands

# 1. INTRODUCTION

Congestion is a source of negative effects on the economy, the environment, and the human state of mind. Up to the 1980's the traditional salvation to avoid congestion was to simply expand the infrastructure, especially when roads were concerned. After a rise of ecological and environmental awareness and an increase of the costs of constructing additional lanes or roads combined with smaller budgets other solutions were sought in order to increase capacity of the existing infrastructure. Together with efforts to convince people to diminish their mobility or use other modes of transportation in stead of automobiles many transportation governments are taking Dynamic Traffic Management (DTM) measures.

DTM measures intend to increase infrastructure capacity by making more efficient use of the infrastructure resulting in spatial and/or temporal traffic flow effects. Unfortunately, it is widespread policy to activate these measures *after* congestion has already occurred and the traffic flow has broken down resulting in a time consuming traffic flow regeneration. It is therefore important to have some knowledge about traffic conditions in advance so anticipating DTM measures could be taken to keep the traffic flow going. The objective of the paper is to give some clues regarding the best choice of methodology, i.e. time series analysis or fuzzy sets, in the particular situation to forecast congestion on freeways within the 5 - 15 minutes period (short-term period).

# 2. METHODS

The methods that have been chosen to forecast congestion on freeways are time series analysis and fuzzy sets. Both methods are data driven, meaning that the model development is heavily dependent on the data set. In order to compare the performance of the methods they should provide results of the same quantity, in this case a binary one: 'congestion' or 'no congestion'.

# 2.1 Time Series Analysis

The Auto Regression Moving Average (ARMA) time series analysis method [1] is widely used as a prediction method [2]. Given F observations  $X_0, X_1, ..., X_{F-1}$ , ARMA(f, g) time series processes can be written in the form:

$$X_{t} = \mu + a_{1}X_{t-1} + a_{2}X_{t-2} + \dots + a_{f}X_{t-f} + b_{0}\varepsilon_{t} + b_{1}\varepsilon_{t-1} + \dots + b_{g}\varepsilon_{t-g}$$
(1)

for t = f, f + 1, ... and with  $\varepsilon$  is white noise with mean zero and variance  $\sigma^2$ . The parameters  $\mu$ ,  $a_1$ ,  $a_2$ , ...,  $a_f$ ,  $b_0$ ,  $b_1$ , ...,  $b_g$ , and  $\sigma^2$  are to be estimated by deriving least squares and maximum likelihood estimators. If we rewrite (1) as

$$RX = M\varepsilon + PX^{0} + \mu I \quad , \quad X = \begin{bmatrix} X_{f} \\ X_{f+1} \\ \dots \\ X_{F-1} \end{bmatrix} \quad X^{0} = \begin{bmatrix} X_{0} \\ X_{1} \\ \dots \\ X_{f-1} \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_{f-g} \\ \varepsilon_{f-g+1} \\ \dots \\ \varepsilon_{F-1} \end{bmatrix}$$
(2)

and R, M, and P representing matrices containing the a and b parameters and I being the unity vector we can minimise (2) to obtain least squares estimation and maximise (2) to obtain maximum likelihood estimation.

#### 2.2 Fuzzy Sets

Fuzzy Set theory [3] is an expansion of the classic set theory using uncertainties and probabilities and is in practice synonymous with Fuzzy Logic (FL). In FL an element has a probability *between* 0 and 1 of belonging to a certain set whereas in classic set theory it is either a member (probability 1) or not (probability 0). In other words: FL is a theory that relates to classes of objects with unsharp boundaries in which membership is a matter of degree.

In order to apply FL one has to map the input variables onto sets of membership functions and this process of converting a precise input value to a fuzzy value is called "fuzzification". These fuzzy values are subsequently processed by a set of fuzzy logic rules, e.g.

IF 
$$(f_1 = k_1)$$
 AND  $(g_1 = l_3)$  THEN  $(h_1 = m_2)$  (3)

where f and g are fuzzificated input variables, k, l, and m are membership function values, and h is a fuzzy output variable. All fuzzy output values are combined in order to find a precise output value: this process is - not surprisingly - named 'defuzzification'. FL has emerged as a tool to translate linguistic into mathematical information (computing with words) and is primarily used as a control and/or decision tool [4].

#### **3. DATA GATHERING**

The field data set contains information gathered from a part of the outer western roadway section (southbound) of the A10 - the beltway around Amsterdam (figure 1). The data were collected from 35 traffic lane induction loops during the period from 13.38 January 7<sup>th</sup> until 23.59 February 5<sup>th</sup> 1999 through the MONICA - MONItoring CAsco - data management system into one minute aggregated time bins. The data consist of information about volume, mean speed, standard deviation of speed, occupancy, and an indication of congestion. These parameters with the exception of the latter two were given for each of three categories: vehicles up to 5 meters, vehicles with a length between 5 meters and 12.5 meters, and vehicles over 12.5 meters.



Figure 1: Data gathering road section

Due to technical reasons there were four gaps in the data collection: the period of 02.49 January 10<sup>th</sup> until 09.04 January 11<sup>th</sup>, the period of 02.51 until 05.58 of January 16<sup>th</sup>, the entire day of January 19<sup>th</sup> and the period of 03.17 until 08.23 of January 23<sup>rd</sup>. In the remaining data set 9.9% of data was excluded because of unreliability.

# 4. MODEL DEVELOPMENT

In order to compare the performance of both methods the data set was divided into four equal sized subsets. Each method's performance was estimated by testing a subset on a model of which the parameters were described and trained using the remaining three subsets. This procedure was carried out on every subset.

# 4.1 Input features

The input features are the same for both methods: they include volume, mean speed and occupancy. All features are gathered into 1-minute time bins. An additional feature is the standard deviation of speed within this time bin and this feature can be regarded as an indicator of the chaos or turbulence of the traffic flow.

The data that was used to describe, train and test the models consisted of four input features for both methods: the above mentioned features which belong to the first category, i.e. vehicles with a length up to 5 meters (see chapter 3). The data was gathered from two points (target detectors; see figure 1) and consisted of temporal information of the input features. These two target detectors were chosen because they are situated on the spot where congestion usually originates. The reason to exclude other data, e.g. data from other induction loops, is twofold: firstly this way the models that are developed according to both methods use exactly the same input data so no method has advantages of extra information. The second reason was due to computer memory capacity: an expansion of input features was not possible in the case of the fuzzy sets method (the number of parameters that have to be established to form a complete set of rules that describe the fuzzy sets is exploding when the amount of input variables exceeds four or five).

# 4.2 **Output features**

The output features or targets consisted of binary congestion indicators of the target detectors and shifted in time over 5, 10, and 15 minutes in order to estimate the predictive performance. Since there are two target detectors there were 2 (methods) \* 2 (detectors) \* 3 (indicators) \* 4 (subsets) = 48 models developed.

# 4.3 **Performance measures**

The performance was measured by testing each subset on a - according to the method and remaining subsets - calibrated model. The generated outputs were rounded and compared with the actual congestion indicator. If errors occurred they were categorised as *false alarm* (falsely forecasting 'congestion') or just *error* (falsely forecasting 'no congestion'). Both methods were evaluated using the same criteria.



# 5. ANALYSIS

Graph 1a-f: Error percentages of the *bottom* target detector models (figure 1)

Graphs 1a-f contain the results of both method's performance through the models developed according to that particular model and the data from the bottom target detector as indicated in figure 1. As clearly can be seen the performance on subset 2 and 3 is poorer than that on subset 1 and 4. This can be explained by the percentage of congestion time that occurred in every subset (graph 3): subset 2 and 3 have a much higher percentage of congested time than subsets 1 and 4. The same comments with respect to the top target detector's method performances can be made viewing graphs 2a-f.



Graph 2a-f: Error percentages of the *top* target detector models (figure 1)

Another feature that easily can be concluded is that the shorter the forecasting period the better the performance. This is - not surprising - also true for both target detector models and for both methods.

A closer look at the graphs displays an interesting difference between both methods: apparently the time series method has a better *false alarm* record than the fuzzy sets method and the fuzzy sets method has a better *error* record than the time series method. This implies that in case of a 'conservative' strategy (i.e. giving a higher priority to not falsely forecasting 'no congestion') one would prefer the fuzzy sets method and in case of a 'progressive' strategy (i.e. giving a higher priority to not falsely forecasting 'congestion') one would prefer the fuzzy forecasting 'congestion' one would prefer the time series method.



Graph 3: Congestion percentage (in time) per subset

What is the overall performance of both methods? To come to this number all error minutes were added together and divided by the total amount of 1-minute time bins. The results are displayed in graph 4. As can be seen both methods perform almost equally well; the time series method has a slight advantage in the overall performance.



Graph 4 : Performance of the methods

# 6. CONCLUSIONS

As can be seen in the analysis chapter, both methods perform almost equally well. The time series method has a slightly better overall performance. The method which one should use to model short-term congestion forecasting is depending on one's strategy: if *false alarm* has to be omitted as much as possible the time series method would be the better method. If, on the other hand, a higher priority is given to *error* prevention the fuzzy sets method would be the appropriate method.

A few remarks must be made, however. It is very likely that both methods will give a poorer result in case of an incident. This because the model was never calibrated with the according data. The data that was used to calibrate the model originates from one point, however. This suggests that the results in case of an incident would not be dramatically poor.

Another remark concerning this research is that it is part of a bigger project in which attention will be given to method's performances when other (spatial) information is used, so further research is needed.

# 7. ACKNOWLEDGEMENTS

The author wishes to acknowledge the University of Twente for providing financial funds and the Transport Research Centre for providing the field data.

# 8. **REFERENCES**

- [1] Box, G.E.P. and Jenkins, G.M. (1970). '*Time Series Analysis, Forecasting and Control*'. Holden-Day, San Francisco.
- [2] Williams, B.M., Durvasula, P.K. and Brown, D.E. (1998). 'Urban Freeway Traffic Flow Prediction Application of Seasonal Autoregressive Integrated Moving Average and Exponential Smoothing Models'. *Transportation Research Record*, **1644**, pp 132-141.
- [3] Zadeh, L.A. (1965). 'Fuzzy Sets'. *Information and Control*, **8**(3), pp 338 353.
- [4] Sasaki. (1988). 'Traffic Control Process of Expressway by Fuzzy Logic', *Fuzzy Sets and Systems*, **26**(6), pp 165-178.

#### SHORT-TERM CONGESTION FORECASTING: TIME SERIES VERSUS FUZZY SETS

#### G Huisken

#### Department of Civil Engineering & Management, University of Twente, P O Box 217, 7500 AE, Enschede, The Netherlands

#### Short Curriculum Vitae of Giovanni Huisken

Giovanni Huisken is currently enrolled as a Ph.D. candidate at the Department of Transportation Engineering & Management, University of Twente, The Netherlands. His main interest is in Dynamic Traffic Management, more specific: comparing several techniques (e.g. neural networks, fuzzy logic, time series analysis, genetic algorithms) that can be used to predict congestion. During the last quarter of 1999 he was visiting researcher at the Institute of Transportation Studies at the University of California, Irvine, U.S.A. Other interests of his are Dynamic Traffic Information and Trip Generation Modelling. Before joining the University of Twente, Giovanni earned his M.Sc. in applied physics at the University of Groningen, The Netherlands, where he also worked with neural networks (nuclear medicine field).