

# Predictions of Urban Flow Volumes and Incident Detection

T. Thomas University of Twente [t.thomas@utwente.nl](mailto:t.thomas@utwente.nl)

E.C. van Berkum University of Twente

## Abstract

Travel demand and supply increasingly are in a delicate balance in urban areas of the Netherlands. This has led to more awareness of the importance of accurate demand predictions and detections of incidents. In this paper we present a prediction scheme and detection algorithm which are based on an extensive study of volume patterns that were collected for about 20 urban intersections in the Dutch city of Almelo. Our scheme consists of: (1) base-line predictions for a given pre-selected day, (2) predictions with a 24 hours time horizon and (3) short term predictions with horizons smaller than 80 minutes. It appears that 24 hours predictions and short term predictions are significant more accurate than base-line predictions. In fact, in most cases prediction errors are negligible small for the short term predictions. We then use our knowledge about the random variations in volume measurements to construct a simple 3 plus 4-sigma clipping method to remove outlying measurements. Most of these measurements are caused by incidents or events. We briefly discuss how our methods can be used in practical applications.

**Keywords:** Travel demand, Predictions, Incident detection.

## 1. Introduction

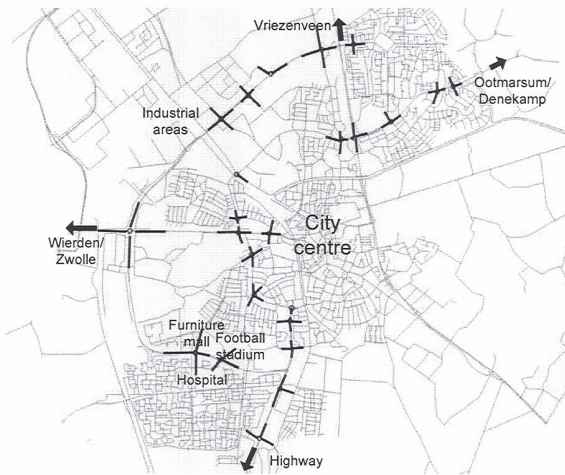
Congestion has increased significantly in the last few decades. The efficient use of existing infrastructure by dynamic traffic management (DTM) is one of the strategies to reduce congestion. An important requirement is the availability of detailed information about travel demand. In general demand cannot be measured directly, but must be estimated using information on volumes, i.e. traffic counts. In urban areas traffic information is scarce (i.e. compared to highways) and only since recently, traffic data are becoming available in traffic information centers (e.g. Hasberg and Serwill 2000, Kellerman and Schmid 2000, Leitsch 2000). However, an increasing number of volume measurements will lead to more reliable demand predictions and thus to better forecasts of the traffic circulation.

Different approaches exist for volume predictions. Unfortunately, few methods include random variations or noise in their predictions. Contrary to systematic variations (e.g. weekly, seasonal or weather related variations) these variations have no periodic character and they are unpredictable. Without knowing the nature and amount of noise it is impossible to evaluate prediction models properly. Besides, knowledge about noise levels can be used in the detection of outlying events.

In this paper we present a prediction scheme for travel demand and an incident detection method in which we have distinguished between systematic prediction errors and the noise. In section 2 we describe the data. In section 3 we pre-select urban volume patterns (daily profiles) into groups, and by using correlations within these groups we develop a prediction scheme for different prediction horizons (sections 4 and 5). In section 6 we present a fairly simple 3 plus 4-sigma clipping method which enables us to detect incidents (i.e. strong deviations from the predictions) in real time. In section 7 we briefly discuss possible applications.

## 2. Data

The study area for this research consists of the urban network of the Dutch city of Almelo. Traffic data were collected at about 20 intersections from September 2004 till 2005. Vehicles were detected by means of inductive loop detectors. These data were processed into volume measurements per link. In most cases measurements were provided in 5 minute intervals, which means that each daily time-series or volume profile contains 288 volumes. However, for about 30% off all links only 30 minute time-series were available. These profiles contain 48 volumes per profile. In Fig. 1 we show the study area. The thick lines in the figure correspond with intersection links for which traffic data were collected.



**Fig 1.** The city of Almelo. The thick lines correspond with intersection links for which traffic data were collected.

The volume measurements were inspected and invalid data were rejected (Weijermars et al. 2006). In Thomas et al. (2007a) we have shown that significant systematic short-term variations exist in travel demand. We found that recurrent patterns with typical periods of 30 minutes are quite common. An explanation for these patterns is that society is regulated by 30 minute intervals. The amplitude of these patterns can be as large as 20%

during rush hour. Because volumes generated by events also vary within short time intervals (minutes rather than hours), we argue that predictions should be made within 5 or 10 minute time intervals.

Contrary to these recurrent patterns are random variations or noise. Random variations in successive traffic counts are uncorrelated and therefore unpredictable. On highways they are caused by variations in headways between cars. In urban areas traffic flows are interrupted, e.g. by traffic signals, which refract the random process. However, variable green times also add to quasi-random variations or noise. We therefore consider all processes which are locally constrained, which have short time-scales and which don't have a clear recurrent pattern as random. In Thomas et al. (2007a) we studied random variations and we have shown that for urban areas the noise in first order can be approximated by a Poisson distribution, although this probably is an under limit. Like e.g. Wild (1997), we decided to use 10 minute time lags for our predictions. These lags are short enough to follow significant systematic short-term variations, but are long enough to get the noise down to an acceptable level.

### 3. Base line predictions

Several authors (e.g., Wild 1997, Grol et al. 2000) have found that the shape of daily volume profiles depends on the day of the week. As a results they suggested that volume predictions can be improved when daily profiles are pre-selected into groups. For the Almelo data Weijermars et al. (2007) found significant differences between individual weekdays and between holidays and non-holidays, which they could explain by social-geographic factors. Based on their results we pre-selected our days in the following groups: Mondays, Tuesdays, Wednesdays, Thursdays, Fridays, Saturdays, Sundays and working days in the school holiday period (week 43 in 2004, and weeks 1, 7, 18 and 30 till 35 in 2005).

We then estimated the base-line prediction,  $q^{base'}$ , for day  $d$ , link  $l$  and time interval  $t$  as the historical mean of the pre-selected group to which day  $d$  belongs (note that we excluded event related traffic; see below):

$$q_{dl}^{base'} = \sum_{d' \in D} q_{d'l}^{obs} / N_D \quad (1)$$

in which  $q^{obs}$  are the historical measured profiles and in which the pre-selected group  $D$  consists of  $N_D$  days, denoted by  $d'$ . Note that we demanded that  $N_D \geq 10$ .

When large events take place traffic flows are influenced by the visitors of these events. At certain locations this will lead to a significant increase of traffic just before the event has started and after the event has finished. In Almelo home matches of the local professional football club can be counted among such events. In Thomas et al. (2007b) we isolated the event related volumes by subtracting the average profile without events from the average profile with events. Note that we corrected for possible time shifts. Per

link the event related peak was modelled by a Gaussian fit,  $q^{ev}$  (which in first order is a good approximation, although for large events the peaks sometimes have one-sided tails).

The final base-line prediction,  $q^{base}$ , for day  $d$ , link  $l$  and time interval  $t$ , is the sum of the non-event and event related volumes:

$$q_{dlt}^{base} = q_{dlt}^{base'} + q_{lt}^{ev} \quad (2)$$

The predictions can be validated by examining the distribution of the residuals (differences between measurements and predictions). In Thomas et al. (2007c) we show that, even when we take the noise into account, there still is a significant variation left in the residuals. In other words, the base-line predictions can en should be improved.

#### 4. 24 hour predictions

Travel demand may depend on the season. Seasonal variations can be quite substantial. The reason for these variations are not always clear. Sometimes they are caused by road works, and therefore they are not really 'seasonal'. Some variations however are periodic. In any case, volume residuals (with respect to the base-line predictions) are strongly correlated for successive days (e.g. Thomas et al. 2007a). In other words, actual volume measurements can be used to improve the base-line predictions for the next day. We call these 24 hour predictions.

We found that for our data the 24 hour predictions,  $q^{24}$ , were optimised, i.e. the rms (root-mean-square) of the residuals was minimised, when we used the following prediction model (for 10 minute time-series).

$$q_{dlt}^{24} = q_{dlt}^{base} \times \left[ \frac{\sum_{t'=t-9}^{t+9} q_{d-1lt'}^{obs}}{\sum_{t'=t-9}^{t+9} q_{d-1lt'}^{base}} \right]^{0.8} \quad (3)$$

In the formula  $q^{base}$  are the base-line predictions and  $q^{obs}$  are the measured volumes. According to the formula, the 24 hour prediction for day  $d$ , link  $l$  and time interval  $t$  is an update of the base-line prediction. The update is based on a power of 0.8 of the ratio between the averages of measurements and base-line predictions from the previous day. These averages are in fact central moving averages with a box width of 3 hours. Note that the power of 0.8 only is applied for successive weekdays. The 24 hour predictions of a Monday and a Saturday however can only be updated with the volumes of the previous Friday and Sunday respectively. For these non-successive weekdays the power of the ratio is 0.5, i.e. the update gets less weight. The reason for this is that volumes of non-successive weekdays are less strongly correlated.

## 5. Short term predictions

Apart from weekly and seasonal variations there may be other variations that are more difficult to model. Weather for example can have an effect on travel demand. Because some of these variations have short time-scales (shorter than 24 hours), they are not included in 24 hours predictions. We can however combine actual measurements and 24 hour predictions to update predictions for the short term. This kind of short term predictions has been applied previously by e.g. Wild (1997).

The problem of such extrapolation methods is that the actual measurements contain noise which contaminates the short term prediction. A solution for this problem is to reduce the noise by using moving averages of actual measurements. Unfortunately, we were not able to improve the predictions, i.e. decrease the rms of the residuals, in this way. We therefore filtered the noise by a Kalman filter (Kalman 1960). The principle of this filter is that the noise in the measurements is smoothed by expected model values that are given by a state equation. We used the following state equation in which we estimated the volume of the next time step according to the expected increase (or decrease) in the 24 hour prediction.

$$q_t^{est} = q_{t-1}^{kal} + (q_t^{24} - q_{t-1}^{24}) \quad (4)$$

The filtered volume  $q_t^{kal}$  is then calculated by taking the linear combination of the state estimate  $q_t^{est}$  and the measured volume  $q_t^{obs}$ , in which the total variance due to model errors and measurement noise is minimised (for details see Kalman 1960). According to this recipe we need an estimate for the variance of the noise  $R_t$  and for the variance of the model error  $Q_t$ . We estimated these variances in the following way:

$$R_t = q_t^{24} \quad (5)$$

$$Q_t = (c' q_t^{24})^2 + (q_t^{24} + q_{t+1}^{24}) / N_D \quad (6)$$

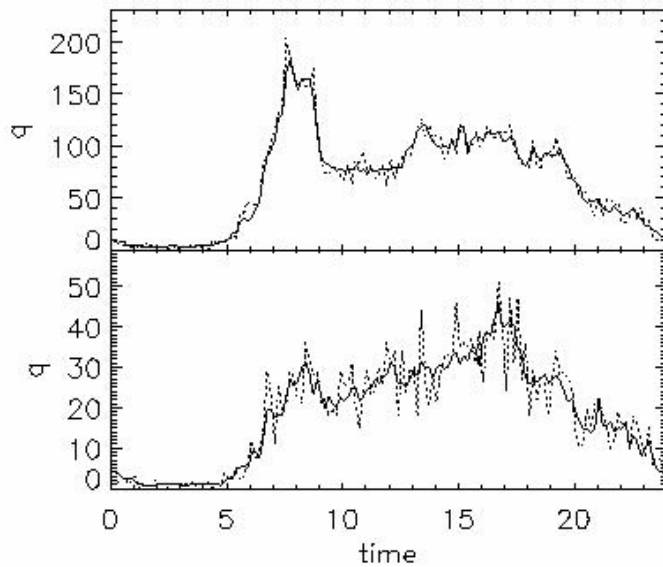
The noise in the measurements can be approximated by a Poisson distribution (Thomas et al. 2007a) with variance  $R$  equal to the predicted volume. The model error also contains some noise, which is described by the second term in (6). This noise is the result of the fact that the prediction model is based on a limited amount of historical data. If the number of days  $N_D$  in the group  $D$  is large this noise term becomes negligible small. The first term in (6) describes the systematic model error due to imperfections in the 24 hour predictions. In fact, these are the errors that we want to eliminate from the short term predictions. According to the state equation (3) we look for the error in the difference between two successive predictions. In Thomas et al. (2007c) we show that a good estimate for  $c'$  (the relative part of this error) is 0.03.

Given the filtered data we optimised the short term prediction  $q^{st}$  in the following way:

$$q_{t+T}^{st} = q_{t+T}^{24} \times \left[ \frac{\sum_{t'=t-5}^t q_{t'}^{kal}}{\sum_{t'=t-5}^t q_{t'}^{24}} \right]^{0.8-0.1T} \quad (7)$$

According to the formula we updated the 24 hour prediction with the filtered data of the previous hour. The maximum horizon  $T$  is 80 minutes (in 80 minutes the prediction will be equal to the 24 hour prediction).

In Fig. 2 we show two examples of predictions with a 10 minute horizon. The predictions appear to follow the measurement quite well. In Thomas et al. (2007c) we show that these predictions are a significant improvement compared to base-line predictions, and that in fact prediction errors are often negligible small for working days. Note that travel demand predictions are less accurate for weekends and school holiday periods.



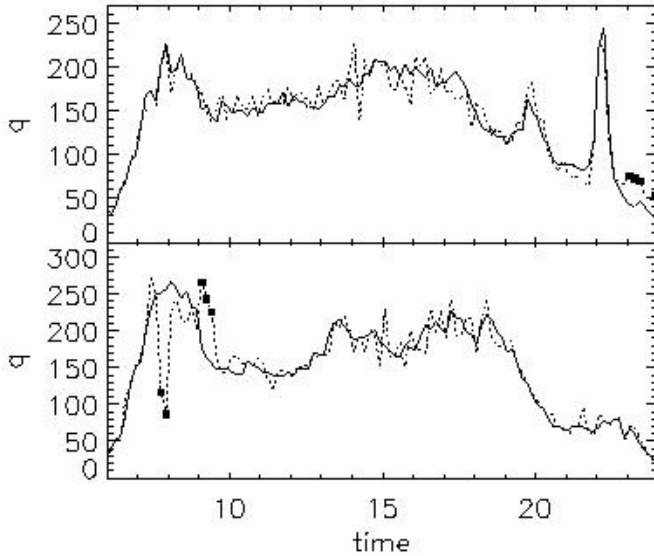
**Fig 2.** Two examples of predictions (solid line) of volume time-series (dotted line)

## 6. Identification of incidents

When traffic accidents, road works, or other unique events occur, traffic volumes can differ significantly from the average. Most of these events are unexpected and can have a large impact on traffic circulation. Therefore it is important to identify these so called outliers as quickly as possible. We used a 3 plus 4-sigma clipping method which can detect outlying events as they occur. When the absolute value of  $q^{obs} - q^{st}$  ( $q^{st}$  is the short term prediction with a 10 minute horizon) exceeds  $4\sqrt{q^{st}}$  ( $\sqrt{q^{st}}$  is the standard deviation of the noise) then the measurement is identified as an outlier. The chance that this happens for a measurement that is not an outlier is  $6.3 \cdot 10^{-5}$ . The number of detected 4 sigma outliers is much higher (about 0.2% of all measurements, which is still a marginal

fraction). Thus, the chance that a detected outlier is in fact not an outlier (false alarm) is about 3% (equal to  $6.3 \cdot 10^{-5} / 0.2 \cdot 10^{-2}$ ). For a detection limit of 3 sigma the fraction of 'false alarms' is much higher (30% or more), which makes it less suited as detection limit. However, the chance that two successive 'false alarms' occur, is very small for the 3 sigma detection limit ( $< 1\%$ ). With the latter criterion we can identify outlying events that are not very extreme, but that have relatively long timescales (20 minutes or more).

In Fig. 3 we show two examples of outlying events (flagged by symbols). The dotted lines are time-series of two observed daily profiles. The solid lines are the predictions with a 10 minute time horizon. In the top panel the outlying event occurs late in the evening. In fact it occurs on the tail of a peak which was generated by visitors of a football match. Note that in this case, the peak itself actually has been predicted quite well by the Gaussian model. In the bottom panel the outlying event is caused by an accident that occurred during the morning rush hour.



**Fig 3.** Two examples of volume time-series (dotted line) with outliers (solid symbols). The predictions are the solid lines.

Many algorithms use threshold values for detecting incidents. Some authors have pointed out that the choice of threshold values often is rather arbitrary (e.g. Ihler et al. 2006). Although we also choose thresholds, our algorithm is very successful for two reasons. Our expected volumes are, contrary to some algorithms, robust predictions based on historical data of many days. More importantly however, the threshold is not arbitrarily chosen, but depends on the random variation of the volumes, which can be described in a uniform way.

## **7. Applications**

In a congested free area like Almelo traffic counts are a measure for travel demand. Travel demand predictions can be useful at a local level. It can be applied for travel-time predictions or for managing of traffic control systems (e.g. Wang et al. 2005, Smith et al. 2001). In the latter case demand predictions can be used to optimise intersection traffic light split times. In fact, some authors already have developed traffic control systems which can adapt to a changing travel demand (e.g. Yang et al. 2005, Yang 2004). In these cases artificial neural networks were used.

For large urban networks macroscopic models are often used to estimate traffic circulation. These models include an estimate of an origin-destination (OD) matrix and a dynamical assignment of OD relations to the network. Reliable estimates of OD matrices are essential in macroscopic models. Camus et al. (1997) argued that volume predictions can be used to improve predictions of OD matrices. It is our intention to use travel demand predictions in macroscopic models.

In the United Kingdom, automatic incident detection systems are being integrated into adaptive traffic signal systems in urban areas (e.g. Ash 1997, Bowers et al. 1996). In the Netherlands, authorities like to include incident detection and event predictions into the management system of Dutch motorways (Taale et al. 2004). We suggest that our detection algorithm might be included into traffic management systems of Dutch cities.

## **Acknowledgements**

We want to thank Vialis and the municipality of Almelo for providing us the volume data in Almelo. This research is funded by Transumo.

## **References**

- Ash A. (1997) Incident detection in urban areas controlled by SCOOT. IEE Colloquium on Incident Detection and Management.
- Bowers, D.J., Bretherton, R.D., Bowen, G.T., Wall G.T. (1996) Traffic congestion incident detection. TRL Report 217, Transport Research Laboratory, UK.
- Camus, R. Cantarella, G.E., Domenico I. (1997) Real-time estimation and prediction of origin-destination matrices per time slice. *International Journal of Forecasting* 13, 13-19.
- Grol R. van, Inaudi D., Kroes E. (2000) On-line traffic condition forecasting using on-line measurements and a historical database. *Proceedings of 7<sup>th</sup> World Congress on Intelligent Transport Systems*, Turin.



Hasberg P., Serwill D. (2000). Stadtinfoköln – a global mobility information system for the Cologne area. In: 7<sup>th</sup> World Congress on Intelligent Transport Systems, Turin, Italy.

Ihler A., Hutchins J., Smyth P. (2006) Adaptive event detection with time-varying Poisson processes. Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, Philadelphia, USA.

Kalman R.E. (1960) A new approach to linear filtering and predictions problems. Transactions of the ASME, Ser. D, Journal of Basic Engineering, 82, 34-45.

Kellerman A., Schmid A. (2000) Mobinet: Intermodal traffic management in Munich – control centre development. In: 7<sup>th</sup> World Congress on Intelligent Transport Systems, Turin, Italy.

Leitsch B. (2002) A Public-privat partnership for mobility – Traffic management Center Berlin. In: 9<sup>th</sup> World Congress on Intelligent Transport Systems, Chicago, USA.

Smith, B. L., Scherer, W. T., Hauser, T. A. (2001) Data-mining Tools for the Support of Signal-Timing Plan Development. Transportation Research Record 1768, 141-147.

Taale H., Westerman M., Stoelhorst H., Van Amelsfort D. (2004) Regional and sustainable traffic management in The Netherlands: methodology and applications. European Transport Conference 2004, Strasbourg, France.

Thomas T., Weijermars W.A.M., Van Berkum E.C. (2007a) Periodic variations in urban traffic flows. European Journal of Transport and Infrastructure Research, submitted.

Thomas T., Van Berkum E.C. (2007b) Events and incidents in urban flow volumes. Proceedings of 3<sup>rd</sup> International Symposium on Transportation Network Reliability, Delft, The Netherlands

Thomas T., Weijermars W.A.M, Van Berkum E.C. (2007c) Predictions of urban flow volumes. Transportation Research –C, submitted.

Wang X., Cottrell W., Mu S. (2005) Using K-Means Clustering to Identify Time-of Day Break Points for Traffic Signal Timing Plans. In: Proceedings of the 8<sup>th</sup> IEEE Conference of Intelligent Transportation System Conference, Vienna, Austria.

Weijermars W.A.M., Van Berkum E.C. (2006) Detection of invalid loop detector data in urban areas. Transportation Research Record, forthcoming.

Weijermars W.A.M., Thomas T., Van Berkum E.C. (2007) Clustering of urban traffic patterns. Transactions on Intelligent Transportation Systems, submitted.

Wild D. (1997) Short-term forecasting based on a transformation and classification of traffic volume time series. International Journal of Forecasting 13, pp 63-72.

Yang, J. S. (2004) Traffic Signal Timing Control for a Small-scale Road Network. *Control and Applications* 441, 48-49.

Yang Z. S., Chen X., Tang Y.S., Sun J.P. (2005) Intelligent Cooperation Control of Urban Traffic Networks. In: *Proceedings of the Fourth International Conference on Machine Learning and Cybernetics, Guangzhou, China.*