

Theoretical and practical capabilities of probabilistic data fusion with Bayesian networks

Thorsten Neumann

*German Aerospace Center (DLR), Institute of Transportation Systems,
Rutherfordstr. 2, 12489 Berlin, Germany*

Marek Junghans

*German Aerospace Center (DLR), Institute of Transportation Systems,
Rutherfordstr. 2, 12489 Berlin, Germany*

Abstract

The increasing demand for high quality data in context of intelligent transportation systems more and more facilitates the use of data fusion methods in order to derive as much information as possible from existing sensors and sensor technologies. This paper discusses the application of simple Bayesian networks for this task with regard to fusing traffic state measurements, in particular. Theoretical aspects like model calibration and statistical quality of the results are discussed in a mathematically exact way. Moreover, a real-world example based on floating car data from two independent vehicle fleets is described in order to evaluate the Bayesian approach also from a practical perspective. Chances and restrictions of the presented model – also with regard to possible modifications – are critically discussed.

1. Introduction

Data fusion is an essential tool in the world of intelligent transportation systems (cf. [1]). Based on an increasing number of sensor technologies and data sources, it opens the field to new and better services concerning a wide range of traffic related applications (cf. [2]). That is, by cross-checking independent measurements of the same event and by combining complementary information from various sources, it helps to create a reliable and comprehensive impression of the real traffic situation of a considered road or transport network. With regard to that, data fusion is applied to all levels of data integration (cf. [3]), i.e. direct measurements from some given physical sensors can be fused as well as highly aggregated information about the traffic of a whole city, for instance.

Concerning the available data fusion methodologies, the existing literature (see [1]) distinguishes between statistical methods (e.g., least square weighted mean, data mining), probabilistic approaches (e.g., Kalman filtering, evidence theory), and other techniques based on neural networks or other kinds of artificial intelligence. Detailed implementations in context of traffic data fusion and intelligent transportation systems are described in [4, 5, 6, 7], for instance. In particular, many fusion models (e.g., [6, 8]) also take into account the underlying traffic dynamics in order to improve the results. The Bayesian approach (cf. [9]) – as it is discussed in the following with regard to traffic data fusion – gets along without modelling such dynamics in its simplest form. It belongs to the class of probabilistic methods and is characterized by a systematic utilization of the well-known Bayes theorem from probability theory.

The paper is structured as follows: Section 2 explains the layout and functional principles of the proposed generic Bayesian fusion approach which also includes the mathematical derivation of consistent estimates for the quality of the data fusion result. Then, the calibration task is considered from an analytical perspective in Section 3, in contrast to the common way of adaptive learning as used in [10]. Finally, Section 4 demonstrates and

practically evaluates the Bayesian approach based on real data from two independent probe vehicle fleets. Conclusions and ideas for future improvement are discussed in Section 5.

2. Basic concept

Assume there are n independent measurements (or estimates) X_i of the real traffic state Z for a given road section and a certain instant of time with Z and X_i for $i = 1, \dots, n$ being considered as random variables. Of course, X_i and Z should be correlated for all $i = 1, \dots, n$ because otherwise the X_i contained no information about the specific value of Z . Consequently, the X_i are usually correlated (via Z) as well given Z is unknown. Needless to say, if the value (or state) of Z is fixed, the X_i become independent again per definition. For illustration purposes, let $X_i := Z + E_i$ with stochastically independent error terms E_i . Obviously, X_i and X_j are correlated then for $i \neq j$ while they are independent if $Z = z$ is fixed because E_i and E_j are independent (cf. Figure 1).

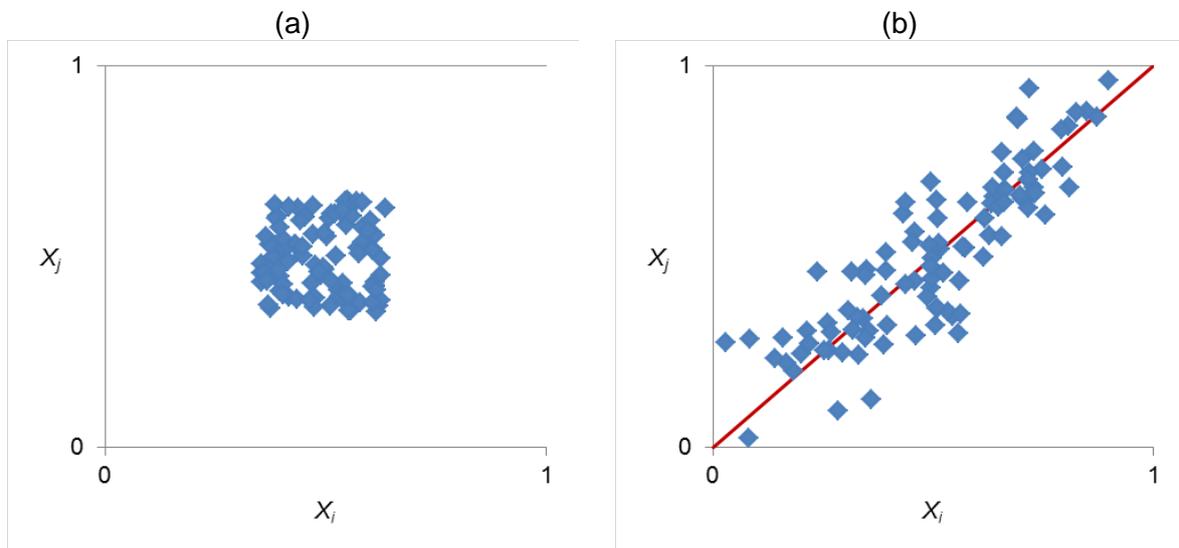


Figure 1: Random sample of (X_i, X_j) with $i \neq j$ where the realizations of Z (a) are fixed at 0.5 and (b) are varying randomly.

The same concept is now realized by using Bayesian networks (BN) which are a special class of so-called probabilistic graphical models (cf. [9, 11]). The corresponding model has the layout as depicted in Figure 2. It shows the direct influence (as represented by directed edges between nodes) of the real traffic state Z on what is measured at the nodes X_i . On the contrary, the X_i for $i = 1, \dots, n$ are not directly correlated but only via Z .

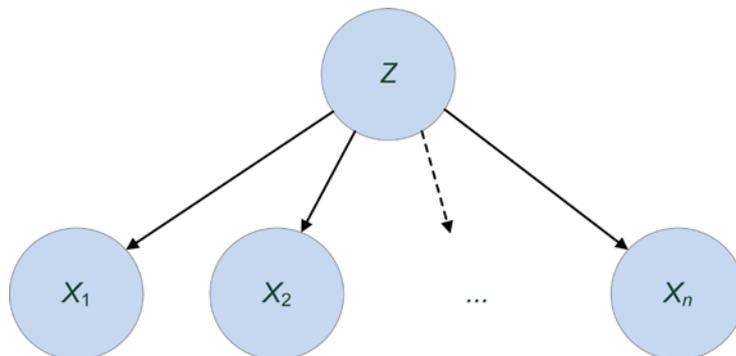


Figure 2: Generic Bayesian model for data fusion.

The main advantage of Bayesian networks is that they provide an efficient and compact, but also intuitive representation of the joint probability distribution $P(X_1, \dots, X_n, Z)$. Thus, all above-named correlations and independencies between arbitrary subsets of $\{X_1, \dots, X_n, Z\}$ are covered by this network. From the general theory of Bayesian networks (cf. [9, 11]) it follows that in fact $P(X_i, X_j | Z) = P(X_i | Z) \cdot P(X_j | Z)$ for the model in Figure 2 if $i \neq j$ while in general $P(X_i, X_j) \neq P(X_i) \cdot P(X_j)$.

The key observation in context of data fusion now is that – in addition to the direct relation between Z and each of the X_i – correlations among the observed measurements X_i may contain further information about the true state Z that is not available when the X_i remain separated from each other. As for the Bayesian model from Figure 2, that means that given sets of observations (comprising all or some X_i) can be used as so-called evidences that finally allow to compute new probabilities for the real state Z in a mathematically exact way by taking account of all statistical dependencies that has been discussed above (see [9] for algorithmic details). The result is an updated conditional probability distribution $P(Z | X_1, \dots, X_n)$ for Z from which the fusion result can be derived by various means.

Junghans and Jentschel (see [10]), for instance, used two common approaches for choosing the final fusion value Z^* which are the maximum a posteriori (MAP) estimator

$$Z^* = Z_{\text{MAP}}^* := \arg \max_z P(Z = z | X_1, \dots, X_n) \quad (1)$$

and the random wheel (RW) approach where $Z^* = Z_{\text{RW}}^*$ is just a random guess following the distribution $P(Z | X_1, \dots, X_n)$.

For both variants, the quality of the fusion result in terms of the statistical chance of its correctness can be derived exactly from $P(Z | X_1, \dots, X_n)$. Namely, let z_1, \dots, z_m be the possible values of Z and $p_k := P(Z = z_k | X_1, \dots, X_n)$ their probabilities for $k = 1, \dots, m$ given the set of observed measurements X_1, \dots, X_n . Then,

$$P(Z_{\text{MAP}}^* \text{ is correct} | X_1, \dots, X_n) = P(Z_{\text{MAP}}^* = Z | X_1, \dots, X_n) = \max\{p_1, \dots, p_m\} \quad (2)$$

and

$$P(Z_{\text{RW}}^* \text{ is correct} | X_1, \dots, X_n) = P(Z_{\text{RW}}^* = Z | X_1, \dots, X_n) = \sum_{k=1, \dots, m} p_k^2 \quad (3)$$

It follows that the MAP estimator always performs better in this statistical sense as can be shown easily. For,

$$P(Z_{\text{RW}}^* = Z | X_1, \dots, X_n) \leq \max\{p_1, \dots, p_m\} \cdot \sum_{k=1, \dots, m} p_k = P(Z_{\text{MAP}}^* = Z | X_1, \dots, X_n). \quad (4)$$

On the other hand, the MAP approach has the drawback that it can never yield correct estimates for unlikely events, of course. The largest difference with regard to the quality of Z_{MAP}^* and Z_{RW}^* is obtained if one of the probabilities p_k is significantly larger than the others, but is not too large at the same time. Both qualities are identical instead if $P(Z | X_1, \dots, X_n)$ is a Dirac or Laplace distribution.

Note further that the Bayesian approach as described above takes into account the statistical correlations between Z and X_i for $i = 1, \dots, n$ only. Consequently, other advanced fusion methods as discussed in [6, 8], for instance, make use of explicit models of relevant traffic flow dynamics. By that, they are able to specify the stochastic (in-)dependencies from the joint probability distribution $P(X_1, \dots, X_n, Z)$ more precisely in terms of a physical (and not only probabilistic) model. Clearly, this opens the chance for even better fusion results as will be further discussed in Section 5. However, adaption and calibration of such models are often more complex so that they may not be applicable to all situations. Moreover, there are ways to integrate traffic dynamics or other influencing factors into the Bayesian network approach as well (cf. [10, 12]).

3. Calibration

With regard to the network from Figure 2, the Bayesian theory implies that the joint probability distribution $P(X_1, \dots, X_n, Z)$ can be factorized as

$$P(X_1, \dots, X_n, Z) = P(Z) \cdot \prod_{i=1, \dots, n} P(X_i | Z). \quad (5)$$

This is also called the chain rule for Bayesian networks (cf. [9]). According to that, the probabilities on the right hand side of Equation (5) need to be calibrated in order to fully define the Bayesian model. Based on reference measurements this could be done via adaptive learning as in [10], for instance. Instead of that, this contribution proposes an analytical approach based on the following system of equations that is obtained from Equation (5) via marginalization over Z . Namely,

$$P(X_1, \dots, X_n) = \sum_{k=1, \dots, m} P(Z = z_k) \cdot \prod_{i=1, \dots, n} P(X_i | Z = z_k). \quad (6)$$

Given a sufficiently large set of real observations $\{X_1 = x_1, \dots, X_n = x_n\}$, the distribution on the left hand side is known as it is directly estimated in terms of the relative shares of all possible realizations of (X_1, \dots, X_n) within this “reference” set. Note that this does not incorporate any knowledge about true traffic states or the like so far.

Let now n_i be the number of possible states for X_i where $i = 1, \dots, n$, and m the number of possible states for Z as above. The number of parameters in Equation (6) is then given by $m + \sum_{i=1, \dots, n} n_i \cdot m$ of which $1 + n \cdot m$ are redundant due to trivial normalization constraints. Table 1 exemplarily lists the parameters in case of 3 sources X_i and $m = n_i = 2$ for all i where $x_{i,l}$ denotes the l -th possible state of X_i . At this, parameters that are not trivially redundant, are marked in blue color and are abbreviated as α_r for $r = 1, \dots, 7$.

Z		X ₁		X ₂		X ₃	
		x _{1,1}	x _{1,2}	x _{2,1}	x _{2,2}	x _{3,1}	x _{3,2}
z ₁	α_1 [P(Z = z ₁)]	α_2 [P(X ₁ = x _{1,1} Z = z ₁)]	(1 - α_2)	α_4 [P(X ₂ = x _{2,1} Z = z ₁)]	(1 - α_4)	α_6 [P(X ₃ = x _{3,1} Z = z ₁)]	(1 - α_6)
	(1 - α_1)	α_3 [P(X ₁ = x _{1,1} Z = z ₂)]	(1 - α_3)	α_5 [P(X ₂ = x _{2,1} Z = z ₂)]	(1 - α_5)	α_7 [P(X ₃ = x _{3,1} Z = z ₂)]	(1 - α_7)

Table 1: Parameters of the Bayesian network model ($n = 3 \mid m = n_i = 2$).

For illustration purposes, consider the following numerical example where the states $x_{i,1}$ and z_1 reflect “free-flow traffic” while $x_{i,2}$ and z_2 are representing “congested traffic”. The probabilities $\beta_s := P(X_1 = x_1, \dots, X_n = x_n)$ for $s = 1, \dots, 8$ are given in Table 2.

X ₁	X ₂	X ₃	P(X ₁ = x _{1,1} , ..., X _n = x _n)
x _{1,1} (free-flow)	x _{2,1} (free-flow)	x _{3,1} (free-flow)	0.65445 =: β_1
x _{1,1} (free-flow)	x _{2,1} (free-flow)	x _{3,2} (congested)	0.11655 =: β_2
x _{1,1} (free-flow)	x _{2,2} (congested)	x _{3,1} (free-flow)	0.0378 =: β_3
x _{1,1} (free-flow)	x _{2,2} (congested)	x _{3,2} (congested)	0.0162 =: β_4
x _{1,2} (congested)	x _{2,1} (free-flow)	x _{3,1} (free-flow)	0.0748 =: β_5
x _{1,2} (congested)	x _{2,1} (free-flow)	x _{3,2} (congested)	0.0192 =: β_6
x _{1,2} (congested)	x _{2,2} (congested)	x _{3,1} (free-flow)	0.02295 =: β_7
x _{1,2} (congested)	x _{2,2} (congested)	x _{3,2} (congested)	0.05805 =: β_8

Table 2: Numerical example.

This all leads to a set of eight equations with seven variables ($\alpha_1, \dots, \alpha_7$). Due to its specific structure, the resulting system of equations however is still under-constrained. In fact, only three variables can be mathematically derived given the others, i.e. four of the α_r must be estimated externally. This, for instance, could be done by reference measurements in context of validating the data sources under some (or all) possible traffic conditions. In general, the statistical quality in terms of the relevant α_r needs to be known for one sensor completely while for the others the quality with regard to one of the true traffic states z_k may be missing without affecting the unique solvability of the system of equations in Equation (6).

Hence, assume that $\alpha_2 = 0.9$ and $\alpha_3 = 0.15$ as well as $\alpha_4 = 0.95$ and $\alpha_6 = 0.85$. From Equation (6), it follows that

$$\beta_1 + \beta_2 + \beta_3 + \beta_4 = \alpha_1 \cdot \alpha_2 + (1 - \alpha_1) \cdot \alpha_3, \tag{7a}$$

$$\beta_1 + \beta_2 = \alpha_1 \cdot \alpha_2 \cdot \alpha_4 + (1 - \alpha_1) \cdot \alpha_3 \cdot \alpha_5, \tag{7b}$$

$$\beta_1 = \alpha_1 \cdot \alpha_2 \cdot \alpha_4 \cdot \alpha_6 + (1 - \alpha_1) \cdot \alpha_3 \cdot \alpha_5 \cdot \alpha_7. \tag{7c}$$

Consequently, Equation (7a) yields $\alpha_1 = 0.9$. Based on that, Equation (7b) shows that $\alpha_5 = 0.1$. And finally, $\alpha_7 = 0.25$ because of Equation (7c). Table 3 summarizes the complete calibration results for the considered numerical example. The four input values are marked in blue color. Cells with computed values are left white. In particular, note that the statistical quality of the measurements X_2 and X_3 in case of congested traffic ($Z = z_2$) as well as the a priori distribution of the real traffic state Z are not input for the calibration but are derived from the data instead.

Z		X ₁		X ₂		X ₃	
		X _{1,1}	X _{1,2}	X _{2,1}	X _{2,2}	X _{3,1}	X _{3,2}
z ₁	90%	90%	10%	95%	5%	85%	15%
z ₂	10%	15%	85%	10%	90%	25%	75%

Table 3: Calibration results for the numerical example (cf. Table 1).

4. Demonstration and evaluation

In order to demonstrate the Bayesian data fusion approach, floating car data (FCD) as provided by two independent vehicle fleets operating in Athens (Greece) were analyzed. That is, for both systems, mean travel times Δt were derived separately via common FCD methods (cf. [13]) based on the observed positions of the vehicles. Based on that, mean delays d for each road section were computed by subtracting the free-flow travel times that were determined by the length L of the road section and the specific speed limit v_{\max} , i.e. $d := \Delta t - L / v_{\max}$. Given sufficient data availability, this allowed to reconstruct the cumulative delays for every possible route in the considered road network. Most of all, it was possible to compare these “system-based” delays to the exact delays that were observed by the individual vehicles of both fleets along their driven trajectories. Note that data of the particular vehicle were excluded from the computation of the “system-based” delays, of course, in order to avoid any circular reasoning. The whole approach is called “self-evaluation” of FCD and has been described comprehensively by Kuhns *et al.* in [14].

As the result of this preprocessing, there is a set of trajectories with known true delays ($\sim Z$) and two independent estimates in terms of “system-based” delays ($\sim X_i$). For simplicity, all delays were then transformed into three levels of service (A, B, C) that were defined more or less randomly such that in each case low delays up to the 50th percentile are considered as “A”, delays up to the 75th percentile as “B” and the rest as “C”. Figure 3 shows the resulting Bayesian network (without evidences) that is used in the following for data fusion. The full

calibration can be found in Table 4 and was done based on the available data from December 2012 and 2013.

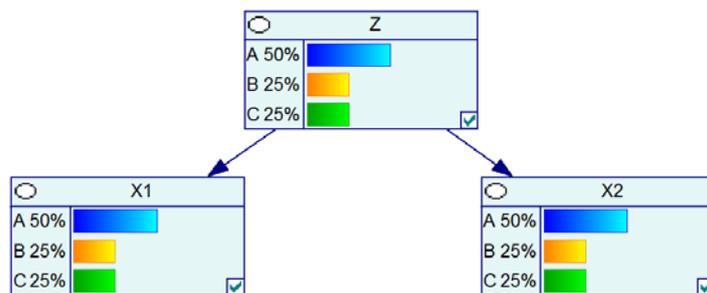


Figure 3: Bayesian data fusion model for the FCD example (generated with GeNIe 2.0, [15]).

	Z	X ₁			X ₂		
		A	B	C	A	B	C
A	50%	79%	17%	4%	78%	18%	4%
B	25%	28%	42%	30%	29%	41%	30%
C	25%	14%	24%	62%	15%	24%	61%

Table 4: Calibration results for the FCD example (cf. Table 1).

The general probability that X_i is correct in the sense of agreeing with Z , is hence given by

$$\begin{aligned}
 P(X_i = Z) &= P(X_i = A \mid Z = A) \cdot P(Z = A) \\
 &\quad + P(X_i = B \mid Z = B) \cdot P(Z = B) \\
 &\quad + P(X_i = C \mid Z = C) \cdot P(Z = C).
 \end{aligned}
 \tag{8}$$

Based on the numbers from Table 4, one obtains the qualities $P(X_1 = Z) = 65.5\%$ and $P(X_2 = Z) = 64.5\%$ with regard to the simplified definition of the levels of service.

Obviously, the same formula holds for the fusion result when X_i is replaced with Z_{MAP}^* or Z_{RW}^* (cf. Section 2). It followed that $P(Z_{MAP}^* = Z) = 66.2\%$ and $P(Z_{RW}^* = Z) = 60.8\%$ when the calibrated model from Figure 3 was applied to the data set from December 2012 and 2013 (cf. Figure 4). The full list of obtained probabilities $P(Z_{MAP}^* = \cdot \mid Z = \cdot)$ as appearing in Equation (8) is given in Table 5 which also includes the results when validating the above model for the fusion of corresponding data from January 2013/2014 and February 2013/2014 respectively.

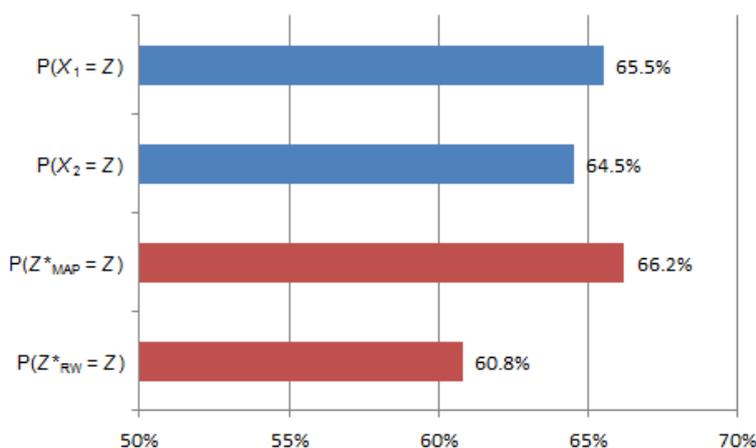


Figure 4: Quality of the X_i and the fusion results in the FCD example (December 2012/2013).

Z		Z_{MAP}^* (December 2012/2013)			Z_{MAP}^* (January 2013/2014)			Z_{MAP}^* (February 2013/2014)		
		A	B	C	A	B	C	A	B	C
A	50%	83%	11%	6%	82%	11%	7%	82%	11%	7%
B	25%	33%	30%	37%	33%	28%	39%	31%	29%	40%
C	25%	16%	15%	68%	17%	15%	68%	16%	15%	69%
$P(Z_{MAP}^* = Z)$		66.2%			65.1%			65.7%		

Table 5: Detailed quality of the fusion result (MAP) in the FCD example.

As can be seen, Figure 4 and Table 5 imply that the approach of Z_{MAP}^* yields small improvements compared to the quality of the X_i in most (but not all) cases. In fact, the MAP estimator seems to prefer the extreme values (i.e., A and C instead of B) for some reason what finally narrows its quality at the same time. Probably, one of the problems here is the very simple way of defining the levels of service that do not reflect any knowledge about traffic dynamics or the like as well as possible systematic differences in the data of both FCD fleets. The random wheel approach (Z_{RW}^*), by the way, provides results that are consistently worse than the original measurements. Clearly, this is because of the additional randomization when choosing the final fusion value as has already loomed in the discussion from Section 2.

5. Conclusion

All in all, the theoretical and practical results showed that the Bayesian approach for data fusion works. However, at least in the considered example (cf. Section 4), the quality lags behind the expectations from the theoretical analysis. That is, the overall benefits in the demonstration example are relatively small (in case of the MAP estimator). Moreover, even a reduction of quality has been observed for some regimes. All this implies that more complex fusion models that take into account not only statistical correlations but also physical relations between the measurements might be more successful in general.

Despite that, the generic Bayesian network from Figure 2 is valid for a huge number of applications and combinations of data sources. Specific sensor properties, for instance, do not need to be known for this simple approach. More detailed knowledge should, of course, be used for improving the model via additional nodes (cf. [10]). So far, just the statistical numbers as in Table 1 (and even that not for all possible states) need to be calibrated for situations as in the practical example that was studied in this paper. Interestingly, some information about the sensor qualities is even directly derived from the measurements via solving the system of equations from Equation (6). In context of a rough traffic state estimation with only two possible states (e.g., “free-flow” and “congested” traffic), that means that new sensors need to be validated under free-flow conditions only, for instance. This is an interesting aspect because validating sensors often is an expensive and time-consuming task. Moreover, the limitation to free-flow traffic avoids the problem of reproducing congested traffic conditions for the reference measurements.

With regard to the limited benefits in the considered real data example, however, the generic model from Figure 2 cannot be recommended unrestrictedly for practical purposes at the moment although it is theoretically interesting. Perhaps, other forms of discretizing the traffic states in terms of levels of service (cf. Section 4) might help to obtain better results. Furthermore, Bayesian models with higher complexity than discussed in this paper might overcome at least some of the encountered problems. To this end, it should be considered whether so-called dynamic Bayesian networks can help to integrate explicit traffic flow dynamics into the described model (cf. [12]), for instance. Finally, also the potential of

additional nodes concerning external factors such as weather conditions or the like might be interesting to be analyzed (cf. [10]) within the above theoretical framework.

Acknowledgements

The authors would like to thank Günter Kuhns for his support in data preprocessing. The original FCD has been provided by two Greek companies (BK Telematics and Zelitron) via the EU-project “Simple Fleet” (Grant Agreement No. FP7-ICT-2011-SME-DCL-296423).

References

- [1] N.-E. El Faouzi, H. Leung, A. Kurian: “Data fusion in intelligent transportation systems: Progress and challenges – a survey”, *Information Fusion*, 12 (1), pp. 4-10, (2011).
- [2] J. Hu, I. Kaspriaris, M. G. H. Bell: “Current state and future outlook of traffic data fusion in London”, 15th International IEEE Conference on Intelligent Transportation Systems (ITSC), September 16-19, 2012, Anchorage, Alaska, USA, (2012).
- [3] M. E. Liggins, D. L. Hall, J. Llinas: “Handbook of Multisensor Data Fusion: Theory and Practice”, 2nd edition, Taylor & Francis, Boca Raton, Florida, USA, (2008).
- [4] N.-E. El Faouzi: “Data-driven aggregative schemes for multisource estimation fusion: a road travel time application”, in: *Multisensor, multisource information fusion: architectures, algorithms, and applications 2004*, ser. SPIE Proceedings, pp. 351-359, (2004).
- [5] C. Nanthawichit, T. Nakatsuji, H. Suzuki: “Application of probe vehicle data for real-time traffic state estimation and short-term travel time prediction on a freeway”, TRB 82nd Annual Meeting, January 12-16, 2003, Washington, D.C., USA, (2003).
- [6] E. Cipriani, S. Gori, L. Mannini: “Traffic state estimation based on data fusion techniques”, 15th International IEEE Conference on Intelligent Transportation Systems (ITSC), September 16-19, 2012, Anchorage, Alaska, USA, (2012).
- [7] Q.-J. Kong, Z. Li, Y. Chen, Y. Liu: “An approach to urban traffic state estimation by fusing multisource information”, *IEEE Transactions on Intelligent Transportation Systems*, 10 (3), pp. 499-511, (2009).
- [8] T. Z. Qiu, X.-Y. Lu, A. H. F. Chow, S. Shladover: “Estimation of Freeway Traffic Density with Loop Detector and Probe Vehicle Data”, *Transportation Research Record*, 2178, pp. 21-29, (2010).
- [9] D. Koller, N. Friedman: “Probabilistic graphical models – Principles and Techniques”, The MIT Press, Cambridge, Massachusetts, (2009).
- [10] M. Junghans, H.-J. Jentschel: “Correction of Selection Bias in Traffic Data by Bayesian Network Data Fusion”, *Journal of Advances in Information Fusion*, 3 (1), pp. 50-62, (2008).
- [11] E. Charniak: “Bayesian Networks without Tears”, *AI Magazine*, Winter 1991, pp. 50-63, (1991).
- [12] T. Neumann, P. L. Böhnke, L. C. Touko Tcheumadjeu: “Dynamic representation of the fundamental diagram via Bayesian networks for estimating traffic flows from probe vehicle data”, 16th International IEEE Conference on Intelligent Transportation Systems (ITSC), October 6-9, 2013, The Hague, The Netherlands, (2013).
- [13] J. F. Ehmke, S. Meisel, D. C. Mattfeld: “Floating Car Data Based Analysis of Urban Travel Times for the Provision of Traffic Quality”, in: J. Barceló, M. Kuwahara (eds.): “Traffic Data Collection and its Standardization”, *International Series in Operations Research & Management Science*, 144, Springer, New York, pp. 129-149, (2010).
- [14] G. Kuhns, R. Ebendt, P. Wagner, A. Sohr, E. Brockfeld: “Self evaluation of floating car data based on travel times from actual vehicle trajectories”, *IEEE Forum on Integrated and Sustainable Transportation Systems (FISTS)*, Vienna, Austria, (2011).
- [15] GeNle 2.0, <https://dslpitt.org/genie/>, last access: 27/04/2015.