

1 **Approximating incident occurrence time with a**
2 **change-point latent variable framework**

3 Francisco Câmara Pereira

4 Singapore-MIT Alliance for Research and Technology, Future Urban Mobility
5 1 CREATE Way, #09-02 CREATE Tower, Singapore 138602
6 Tel: 65-6601 1547, Fax: 65-6778 5654
7 Email address: camara@smart.mit.edu

8 Oren Lederman

9 Singapore-MIT Alliance for Research and Technology, Future Urban Mobility
10 1 CREATE Way, #09-02 CREATE Tower, Singapore 138602
11 Tel: 65-6601 1547, Fax: 65-6778 5654
12 Email address: orenled@mit.edu

13 Moshe Ben-Akiva

14 Massachusetts Institute of Technology
15 Room 1-181, 77 Massachusetts Avenue, Cambridge, MA, 02139
16 Telephone: 617.253.5324
17 Email address: mba@mit.edu

18 4944 words + 7 figures + 2 tables= 7194 (limit=7500 words)

19 October 28, 2013

1 **ABSTRACT**

2 We propose a methodology to approximate actual incident occurrence time by analyzing down-
3 stream volume sensor data. We model the time difference between actual occurrence time and
4 reported time (or *delay*) as a latent variable that becomes a parameter in a change-point time series
5 model. We then apply a *maximum a posteriori* (MAP) framework to infer the most probable delay.
6 This MAP framework uses the time series model as the likelihood function and a bayesian prior
7 based on field knowledge.

8 We applied our model on 5 months of traffic sensor data and accident reports from 3 Singa-
9 pore expressways and corrected the accident start times for 1086 accidents in total. We compared
10 the results with a manually constructed baseline and obtained a mean absolute error (MAE) be-
11 tween 5.7 and 7.4 minutes and a root mean squared error (RMSE) between 10 and 12.

1 INTRODUCTION

2 By their nature, accurate traffic incident occurrence is difficult to detect both spatially and tem-
3 porally. Plenty of intelligent transportation systems (ITS) research exist that focus on incident
4 detection (e.g. Weil et al. (1), Karim and Adeli (2) and Tang and Gao (3)) but the focus is mostly
5 on detecting that an incident *has* occurred, rather than precisely *when* it happened. The problem is
6 particularly complex whenever the incident occurs away from sensors or when the overall scenario
7 is highly prone to sensor noise (e.g. due to harsh weather, sensor quality, traffic conditions).

8 The detection and consequent reporting of incidents to the traffic managers can happen in
9 different ways, namely via the above mentioned ITS incident detection systems as well as driver
10 reporting (e.g. people involved in the accident), personnel reporting (e.g. traffic police) or through
11 camera monitoring in the traffic management center. The actual incident occurrence is very often
12 not observed at all.

13 Accurate information of incident occurrence time can be relevant in a few situations. Traffic
14 prediction systems can improve their performance by incorporating incident details at the right time
15 and place. Without such information, they may rely on wrong assumptions thus generating faulty
16 results. Even with an efficient self-adaptive mechanism, they need time to correct the parameters,
17 particularly under complex scenarios. For example, dynamic traffic assignment (DTA) models
18 need to update capacity parameters accordingly for affected areas. Data-driven algorithms (e.g.
19 Neural Networks, ARIMA) need to adapt to different *regimes* (e.g. Antoniou et al. (4)). Of course,
20 this incident information will arrive itself with some delay, and this needs to be considered by the
21 traffic prediction algorithm itself (e.g. by a “roll-back” mechanism).

22 The ability of traffic prediction engines to roll-back and re-generate their calculations is
23 fundamental because they rely on spatial and temporal correlations, i.e., the error due to lack of
24 incident information will be propagated unless the system integrates it properly in the right time.

25 The second general motivation regards to post-hoc incident analysis. From the point of
26 view of traffic emergency management, it is important to assess the performance of incident re-
27 sponse systems, and more accurate information will lead to better informed decisions in a context
28 where timing is crucial.

29 A third reason, yet very specific to our case, is to correct incident duration information
30 in our automated incident analysis framework where we estimate capacity reduction and incident
31 clearance duration sequentially in time Pereira et al. (5). These two variables alone, capacity
32 reduction and incident duration, sufficiently specify the role of an incident for traffic prediction
33 systems, particularly when based on simulation models (e.g. DynaMIT, from Ben-Akiva et al.
34 (6)).

35 In this paper, we will focus on estimating actual incident occurrence time, t_{occ} , at some
36 reporting time t_0 or later, where obviously $t_0 > t_{occ}$. Our incident start detection model will rely
37 on a signal feed that consists of volumes aggregated by 5 minutes intervals, as observed by the
38 closest sensor downstream to the incident. This specific setting is determined by our case study,
39 the Singapore expressways, from which we have volumes for a period of 5 months.

40 Our task can be described as follows. At time t_0 , a report is received about an incident that
41 occurred at a certain location. We have a feed of volume information from the sensor downstream
42 to that location, aggregated by 5 minute intervals. We want to determine the most likely period
43 when this incident has occurred. For the purposes of this paper, we assume complete availability
44 of data (before and after time t_0), leaving for future work a real-time sequential version, where this
45 process is run at time t_0 and subsequently as new information arrives.

1 More often than not, the incident will happen far from the downstream sensor, so its impact
2 on the volume signal will be itself delayed. In our case, due to the short distance between sensors
3 in the Singapore expressways, it has been observed that this delay is negligible, particularly taking
4 into account the 5 minutes aggregation Mak (7).

5 Another note relates to the reporting times. Even though traffic entities give their best to
6 efficiently streamline the process, the heterogeneous characteristics of the incidents and their re-
7 sponse sometimes lead to considerable delay in the reporting itself. Besides the detection delay,
8 there may be other operational sources of delay that are determined by the incident characteris-
9 tics (e.g. the traffic police may detect it soon but first try to unblock the road or provide local
10 assistance; in case of multiple incidents, some reports may be keyed in the system with more de-
11 lay). As a result, we may observe significant differences between reported time and actual incident
12 occurrence.

13 We propose to approach this problem in the following way:

- 14 • each accident has its time series of volume data from the downstream sensor, we trans-
15 lated this into flow. The actual accident occurrence should originate a change in the flow
16 time series parameters. This corresponds to the concept of "change point" in time series
17 analysis literature;
- 18 • the true accident occurrence time is, in general, unobserved or *latent*;
- 19 • we will determine t_{occ} through a latent variable whose best approximation should support
20 the maximum likelihood of the observed time series of flow for the downstream signal;
- 21 • since we do not have the ground truth values for t_{occ} , we will build a manual (visual)
22 baseline of the start times from the time series signals.

23 The next section will be dedicated to the literature review, followed by the description of
24 the methodology (Section 4). In Section 5, we describe the experiments. The paper will end with
25 a discussion (Section 6) and the conclusions (Section 7).

26 LITERATURE REVIEW

27 Literature exclusively about inference of incident start times is not abundant, but, on the other
28 hand, there exist many works on the related topic of incident detection. In this review, we will
29 approach these two major topics.

30 Within an analysis of incident detection methods in a Singapore expressway, Mak (7) ver-
31 ified that report times do not reflect the actual start time of accidents. And, due to the delay/time
32 lag between reported and actual start time, simply using the reported start time would lead to mis-
33 classified traffic patterns and affect the accuracy for incident detection models. On the other hand,
34 he concluded that, based on the distances between the accident and the adjacent detector, the time
35 it takes for the disturbance to reach the sensor falls between 0.4 and 1.39 minutes.

36 Finding the incident occurrence time in the sensor signal is however a difficult challenge.
37 Dia (8) proposes that a 20% disturbance in traffic parameters (speed, occupancy and volume at the
38 upstream station) could be used to indicate either the start or end of an incident. According to (7),
39 compared to a visual inspection, this method is more efficient and requires less experience from
40 the user. However, in his experiments the two methods gave different start and end time values.

41 Besides, a strict threshold value may be too strong a constraint if we also consider other aspects,
1 such as sensor quality, general traffic conditions, weather status or time of day.

2 Rather than taken by itself, the task of incident start time inference has often been inte-
3 grated within automatic incident detection and incident analysis works (e.g. (9)). In these cases,
4 the manual analysis is chosen for practical reasons: incident start times are only needed to help cal-
5 ibrate a classification algorithm (of incident detection), not being a necessary input of the model.
6 In fact, it is arguable that, in some applications, the exact start time is of secondary importance.
7 On an operational setting, the need is to detect the occurrence and characterization of incidents as
8 soon as possible, regardless of its precise on-set time. Although the start times are not secondary
9 in incident analysis, the traditional practice is still to manually analyze the signal (e.g. using 5D
10 stacked bar charts (Lee et al. (10))) or use simple heuristics such as mentioned earlier (Dia (8)).

11 However, there are situations where approximating the start time is important. For example,
12 Jeong (11) explains that recent studies to validate AID algorithms had to rely on simulated data
13 since reported start time normally maintained by freeway patrols and incident management systems
14 is not precise. The difference might be a couple of minutes or even more, and creates an undesirable
15 shift in the incident data.

16 In traffic prediction too, an accurate characterization of incidents is important. For exam-
17 ple, DynaMIT (6) uses such information to roll-back the network state estimation process with
18 revised assumptions about capacity reduction on the affected links. By doing that, the system will
19 better understand the flow changes and generate more accurate predictions. Notice that, in such a
20 context, one needs a fully-automated system for incident start time inference rather than one base
21 on manual/visual analysis.

22 A final motivating argument for incident start time inference relates to incident analysis,
23 particularly related to clearance duration and response times. Traffic agencies regularly need to as-
24 sess and revise their incident management procedures and wrong incident start times will obviously
25 affect accuracy of such evaluations.

26 It is thus remarkable that little more has been done on this specific topic despite its rel-
27 evance. This may be explained by the almost total lack of observability of incident occurrence
28 times and for the higher focus on incident detection systems (which do not necessarily need to
29 know *when* it happen, rather they are focused on that it *had* happened) rather than for traffic pre-
30 diction. Current growing emphasis on traffic prediction services may possibly reverse this trend.

31 **METHODOLOGY**

32 Our goal is to estimate the start time of an incident, given a time series signal with traffic flow data.
33 It is known that an accident *has* occurred as well as its reported location. We are also given the
34 sensor data from the vicinities of the accident.

35 From the point of view of the signal analysis, we can think of two general methodologies:
36 heuristics application; anomaly detection. The heuristics solution follows the line proposed by
37 Dia (8), that assumes a 20% disturbance in traffic parameters. As with others (e.g. *largest flow*
38 *drop*), this rule is obviously too rigid, and this threshold would probably vary from dataset to
39 dataset. Even with the best choice of threshold, one would face problems in applying the rule. For
40 example, what would the time window be? Incident disturbance may be gradual through time, or
41 be abrupt and recover quickly.

42 The anomaly detection approach defines the model of *expected* behavior and captures de-
43 viations as being anomalies. It lends itself to a more flexible approach than using heuristics in the

44 sense that it does not expect a specific threshold or a deterministic rule.

1 In our case, we know that an accident occurred and expect it to somehow impact on the
 2 traffic signal. We prefer not to make strong heuristic assumptions but still need to assume that there
 3 was in fact an observable change in the signal. A concept that fits nicely with this description is that
 4 of *time series change-point*: when different subsequences of a data series follow different statistical
 5 distributions, commonly of the same functional form but having different parameters (12). In other
 6 words, an incident should originate a change in the traffic time series signal parameters that lasts
 7 during a period of time. It is also plausible that this change ranges from a simple shift in the mean
 8 to general parameter redefinition.

9 Since the actual incident occurrence time is not observed, we will determine it through a
 10 latent variable. Our task is thus to obtain the most likely value for this variable given the change-
 11 point model assumption that the signal before the incident occurrence time will have a certain set
 12 of parameters, and after it will have another set of parameters.

13 Formally, our goal is to find the incident occurrence time t_{occ} , which is determined by
 14 $t_{occ} = t_0 - z$, where t_0 is the reported time and z is the latent *delay* observed from the signal. This
 15 delay in practice corresponds to the reporting delay, and it is in fact our latent variable. We want
 16 to maximize the probability $p(z|\mathbf{y}, \theta)$, such that

$$p(z|\mathbf{y}, \theta) \propto p(\mathbf{y}|z, \theta) * p(z)$$

17 with \mathbf{y} being the vector of time-series signal and θ the vector of parameters. $p(z)$ is the
 18 bayesian prior for the latent variable z , which will be discussed later. Notice that we can obtain the
 19 exact posterior probability for certain z_i by normalization:

$$p(z_i|\mathbf{y}, \theta) = \frac{p(\mathbf{y}|z_i, \theta) * p(z_i)}{\sum_j p(\mathbf{y}|z_j, \theta) * p(z_j)}$$

20 In general, we are looking for maximizing the probability of z (a process also known as
 21 *maximum a posteriori*, or MAP)¹. Thus, we will vary choices for θ parameters as well as for the
 22 value of z . Formally, we want:

$$\operatorname{argmax}_{z, \theta} p(\mathbf{y}, z|\theta) = \operatorname{argmax}_z \left[p(\mathbf{y}|z, \operatorname{argmax}_{\theta} p(\mathbf{y}|z, \theta)) * p(z) \right]$$

23 The likelihood, $p(\mathbf{y}|z, \theta)$, is expanded as

$$p(\mathbf{y}|z, \theta) = p(y_{t_0-l}, y_{t_0-l+1}, \dots, y_{t_0}, \dots, y_{t_0+r-1}, y_{t_0+r}|z, \theta) \quad (1)$$

24 where l and r correspond respectively to the left and right time window boundaries to
 25 inspect. Notice that, if the algorithm is intended to run in real-time, r will be constrained to the
 26 most recent data. The likelihood function is determined by the typical multivariate gaussian, such
 27 that

$$\ln(p(\mathbf{y}|z, \theta)) = -\frac{1}{2} \ln((2\pi)^n |\Sigma|) - \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^T \Sigma (\mathbf{y} - \hat{\mathbf{y}}) \quad (2)$$

¹In fact, the posterior distributions of z should not be ignored since they carrier potentially relevant information about uncertainty of the start time approximation. However, for simplicity in this paper, we will only focus on the MAP.

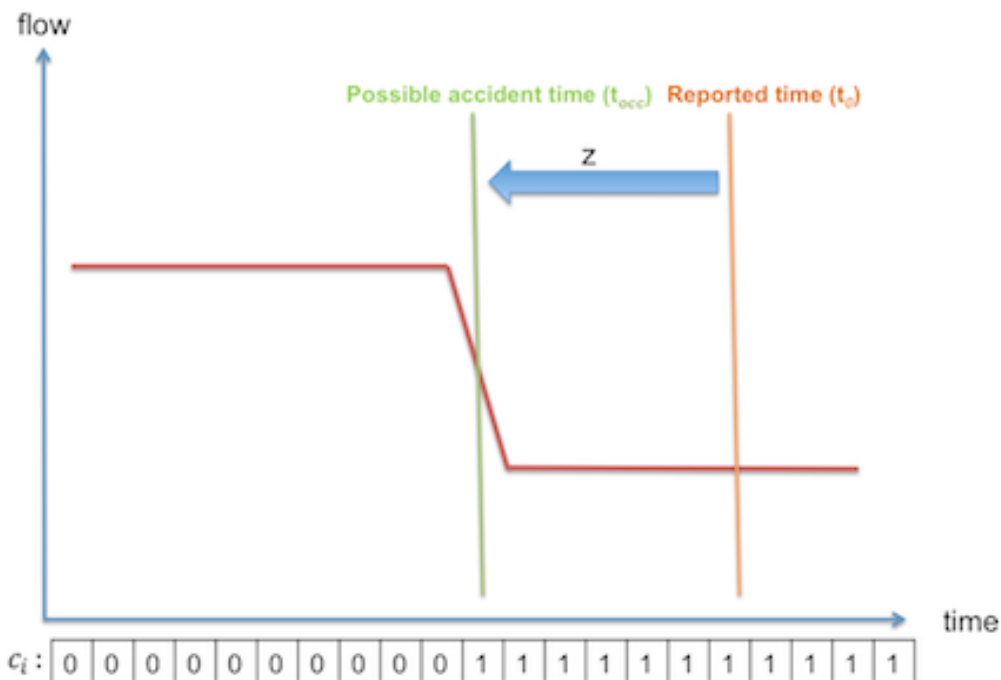


FIGURE 1 The role of z in determining the overlay and the occurrence time.

28 for $n = r + l + 1$ (all points in the time window). Σ is the covariance matrix created with $\Sigma_{i,j}$
 1 containing the auto covariance for lag $|i - j|$. \hat{y} corresponds to the prediction for y according to a
 2 time series model that considers the change point specified in z . In practice, we use the time series
 3 library from the Weka package (Hall et al. (13)). It formulates the time series prediction task as a
 4 (potentially non-linear) regression problem and transforms each time series point as a single input
 5 vector that carries temporal relations (e.g. lags) as well as other features, called “overlays”. In
 6 this way, we can run any Weka regression method over our data, such as support vector machines
 7 (SVM), neural networks and so on. The overlays correspond to the indicator functions in time
 8 series literature and are crucial to our context.

9 The change-point defined by z is turned into the following overlay:

$$c_i = \begin{cases} 1 & t_i + z > t_0 \\ 0 & \text{otherwise} \end{cases}$$

10 where c_i is the value of the overlay at time t_i . The visual intuition for this effect is in Figure
 11 1.

12 In this way, the regression algorithm will be able to distinguish between the two regimes
 13 while still keeping the whole time series in a single model. With adaptable regression models,
 14 this should both allow for superficial, mean-shift changes, as well as to entire regime changing
 15 situations, maximizing the global coherence, as opposed to explicitly breaking the series into two
 16 separate parts and re-estimating the models separately.

17 Regarding the bayesian prior model for z , its role is to introduce local knowledge about
 18 the delays. Although the precise occurrence time is generally impossible to capture, for many

19 accidents traffic operators are able to provide rough estimates. In our case, traffic operators shared
1 the intuition that reporting delays should average around 5 minutes or less. This obviously provides
2 only the mean for the prior, so its form is left to the modeler. Care must be taken to avoid negative
3 delays (i.e. an accident cannot occur after reported time), therefore symmetric forms such as the
4 normal distribution are not an option. Moreover, it is reasonable to expect high mass in lower
5 delays (e.g. 5 minutes and below) and a long tail for the high delays. A log-normal distribution
6 can provide such a behavior.

7 After determining the most likely occurrence time in the signal, there is one last step to
8 consider. This analysis only gave us the moment in time where the incident disturbance reached
9 the sensor. We also need to consider the time it takes such disturbance to go from the actual incident
10 location to that sensor's location. We can consider two situations, depending on the sensor position
11 relative to the accident, namely whether it is upstream or downstream.

12 In the upstream case, the disturbance propagation should be dependent on the queue forma-
13 tion rate. For such case, one needs to model the queue shockwave (also using information on speed
14 and capacity) to determine the time it takes for the disturbance to reach the sensor. A few solutions
15 exist for this model (e.g. cell transmission model (Daganzo (14)), shockwave speed model (Kuhne
16 and Michalopoulos (15))).

17 When the sensor is downstream to incident, the time difference of the disturbance should
18 be much smaller. Mak (7) showed that, for the case of Singapore expressways, the time it takes for
19 the disturbance to reach the sensor is between 0.4 and 1.39 minutes. If we're working with larger
20 time intervals (e.g. 5 minutes), this time difference becomes negligible. Our own observations on
21 the same dataset reinforce this conclusion.

22 As mentioned before, this paper deals with the downstream case, leaving a solution that
23 simultaneously considers upstream and downstream sensors to further work.

24 Summarizing, our method approximates the incident occurrence time by maximizing the
25 likelihood of a change-point model. This model is parameterized by the latent variable, z , that
26 represents the delay and directly defines the change-point. The other parameters are determined
27 by a general training procedure. This method does not demand any particular functional form
28 for the change-point time series model and uses a bayesian prior for the variable z that builds on
29 operational field knowledge.

30 As in earlier works, since we don't have a set of ground-truth observations, we will man-
31 ually define a baseline for comparison. Of course, such baseline will be affected by our own
32 perception biases. For example, one intuition is that a "very large drop" in flow should indicate
33 the incident occurrence. While this may be a reasonable assumption, this drop may itself follow
34 a trend (e.g. peak hour effects) or be a second consequence of the incident (e.g. police arrival).
35 To allow a clearer analysis, we also assign a confidence value (between 1-low and 5-high) to each
36 case. We have high confidence when the signal is stable enough and the disturbance is clear enough
37 (as in Figure 1) or supported by the text report. We will evaluate our model with the root mean
38 squared error (RMSE) as well as the Mean Absolute Error (MAE).

39 EXPERIMENTS

40 Data

41 The dataset comprises 5 months of traffic flow data from 268 sensors located in 3 Singapore ex-
42 pressways. During this period a total of 1086 traffic incidents were recorded. The distance between
43 incidents and downstream sensors is distributed according to Figure 2.

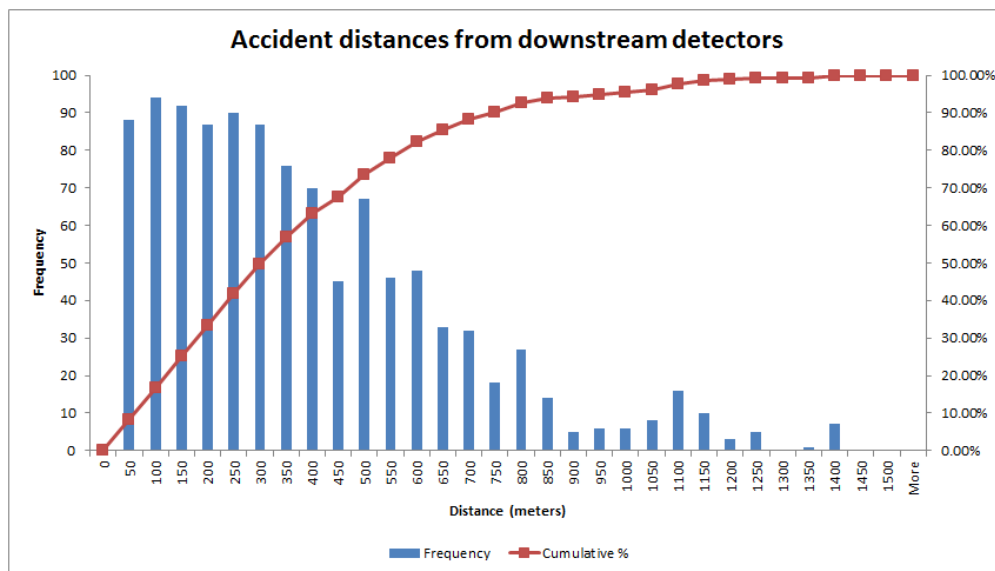


FIGURE 2 Distance between incidents and downstream sensors.

44 For each incident, we have the report start time as well as the traffic volumes for the down-
 1 stream sensor, aggregated on 5 minute windows, the observed incident duration and a small text
 2 report. In Figure 3, we show two incident cases where we consider that the disturbance is clear
 3 (top) and unclear (bottom). We also show the results according to our algorithm (" t_{occ} ").

4 For each incident, we obtained the downstream sensor data from 180 minutes before and 60
 5 minutes after the reported time (i.e. $l = 180$ and $r = 60$ in equation 1). We standardized the signal
 6 having as reference the same time of day and type of day (weekend/weekday) throughout the entire
 7 dataset. The total number of cases is 1086. From this set, we manually inspected 401 cases for
 8 later comparison and validation. This manual inspection was essentially visual, very occasionally
 9 using the report text for further verification.

10 Experimental design

11 For each incident, we independently ran our model. After a trial period, where we tested with
 12 the complete regression portfolio from Weka (Hall et al. (13)), we decided to use a support vector
 13 machine algorithm. We also defined the number of lags to be 12 and kept the remaining parameters
 14 at default values after some exploratory experiments.

15 We defined the bayesian prior for z to be the lognormal with mean 1.87 and shape (σ^2) of
 16 0.26.

17 Since our dataset is aggregated on 5 minute intervals, our values for z will also vary dis-
 18 cretely in the same fashion. For each incident and each possible value of z within the range $[-l, r]$
 19 (with z being a multiple of 5), we estimate a time series model as above described. We calculate
 20 the likelihood using equation 2 and multiply it with the prior probability. The highest value will
 21 indicate the maximum a posteriori, i.e. the best estimate for the incident occurrence time.

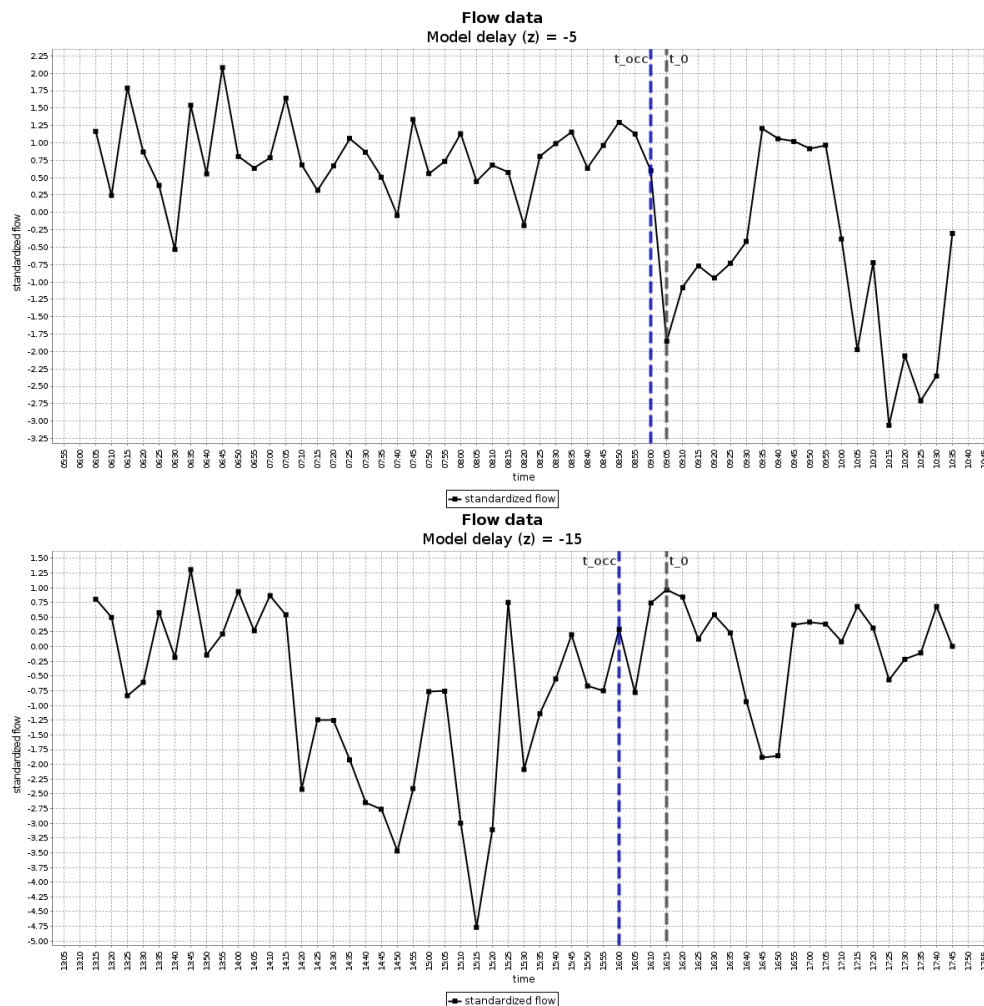


FIGURE 3 Two incident signals.

22 Results

23 The left plot in Figure 4 shows the distribution of the general results through the dataset. The
 24 values are negative with respect to reported time (i.e. -5 corresponds to "5 minutes before the
 1 reported time"). Expectably, we see a high frequency of 0 and 5 minute delays, but also that there
 2 is a relatively long tail up to 45 minutes. One can argue that such concentration in short delays is
 3 artificially induced by the prior. To understand its influence, we depict the results without the prior
 4 effect on the right.

5 It turns out that the bayesian prior does have a strong influence. Its role is essentially to
 6 bias the maximum a posteriori towards the 5 minutes delay other things being approximately equal.
 7 Notice that, for each incident, the dataset has a reasonably large number of points (50 points, from
 8 -180 to 60 minutes) and therefore the prior is only relevant when multiple MAP candidates exist
 9 with competitive probabilities. This behavior is expectable for a prior. Whether it is desirable
 10 or correct is subject to the context. In our case, we will compare these results with the manual
 11 baseline mentioned above.

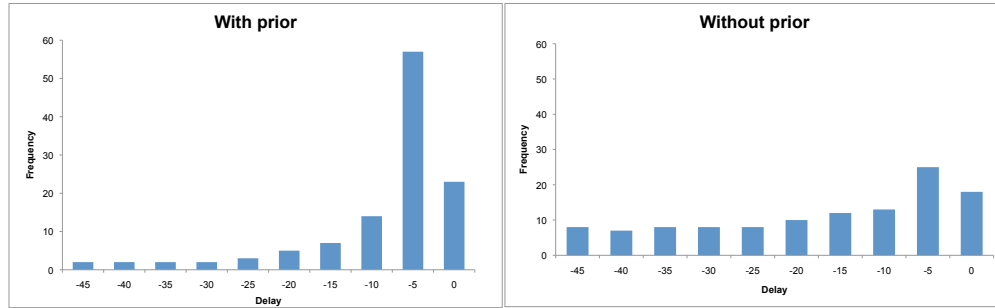


FIGURE 4 Distribution of results through the dataset.

12 Tables tables 1 and 2 show the MAE and RMSE results, respectively, in comparison with
 13 our manual baseline, with and without using the bayesian prior for z . We show the results for the
 14 cases with higher confidence (ranked either 4 or 5) and for all cases evaluated.

TABLE 1 Comparison with manual baseline (MAE)

Prior	No Prior	
5.7692	9.9573	High confidence
7.3852	12.615	all cases

TABLE 2 Comparison with manual baseline

Prior	No Prior	
10.096	15.832	High confidence
11.874	17.739	all cases

1 Regarding the bayesian prior, it has a non-negligible role in both MAE and RMSE. We also
 2 tested with different scale parameters and the chosen mean and scale yielded the best results. It
 3 is arguable that our own baseline is itself biased, but the intuition that the prior helps avoid over-
 4 fitting is relevant. Since we estimate an individual change-point model for each incident, without
 5 the prior, the only knowledge considered would be the flow time series for that specific window,
 6 which would make the model too sensitive to local context (e.g. sensor noise, secondary accidents).
 7 Of course, the parameters of the prior themselves should be realistic as much as possible.

8 On a less positive note, the MAE and RMSE errors in comparison with the baseline are con-
 9 siderably high, even when using the bayesian prior. There may be several general explanations for
 10 this fact. Upon visual inspection, there are incidents with more than one plausible occurrence time

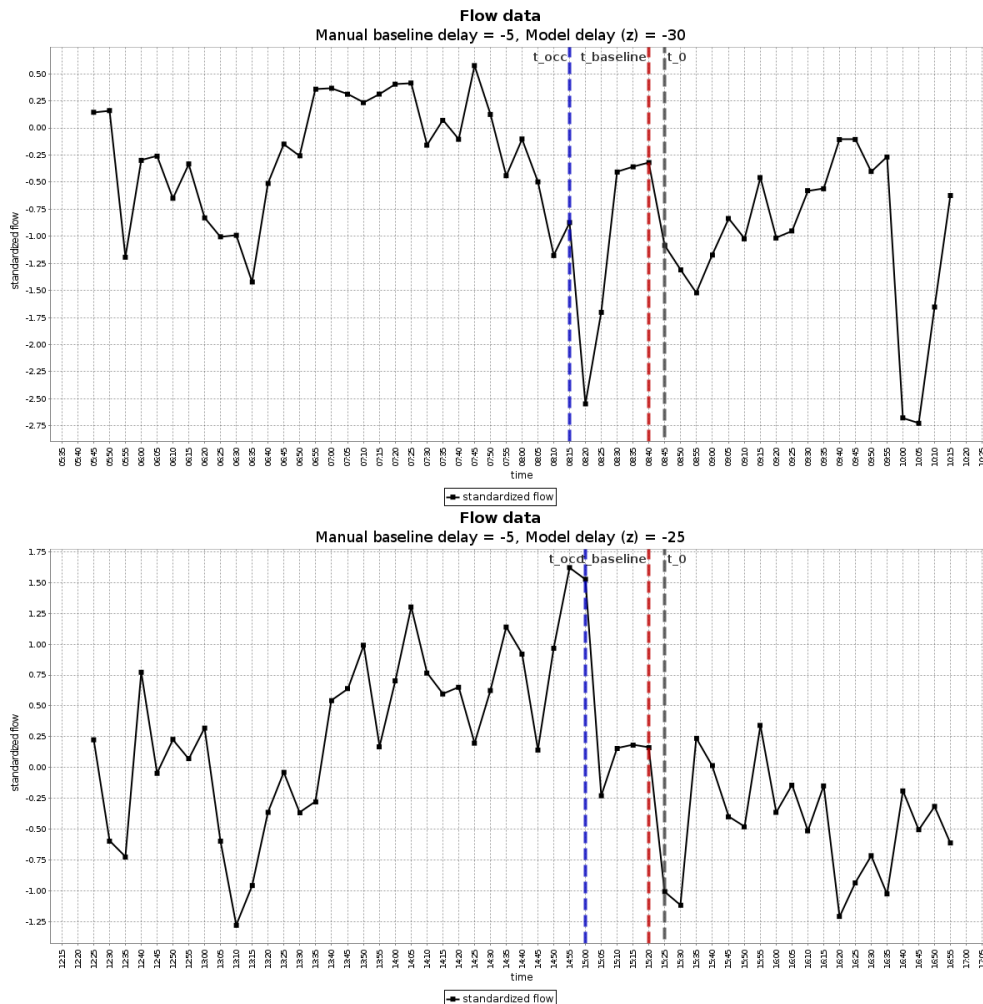


FIGURE 5 Both baseline and adjusted are plausible incident occurrence times.

11 and, in these cases, intuition provides contradictory answers: delay time shouldn't be unreasonably
 12 high (i.e. it should be as close as possible to report time, t_0); the incident may generate several flow
 13 disturbances (i.e. occurrence should be the earliest possible, followed by other episodes, maybe
 14 more intense). Figure 5 gives two examples.

15 Another explanation for the differences to the baseline is that our algorithm captures time
 16 series model changes rather than single signal drops, therefore it may “see” patterns that the human
 17 eye cannot. In Figure 6, we show an example where our algorithm apparently caught the beginning
 18 of a new time series pattern rather than the earlier big flow drop.

19 Another problem is ambiguity in the change-point determination due to the discrete nature
 20 of the signal. In some baseline cases, we chose the beginning of a large drop to correspond to the
 21 incident occurrence time, while the algorithm may also choose the end or middle of this drop (see
 22 Figure 7 for an example) depending on the maximum likelihood of the change point model.

23 Finally, an important limitation relates to having only one change point. The sensor signal
 24 may have more than two regime changes due to the incident. Indeed, if we want to include the

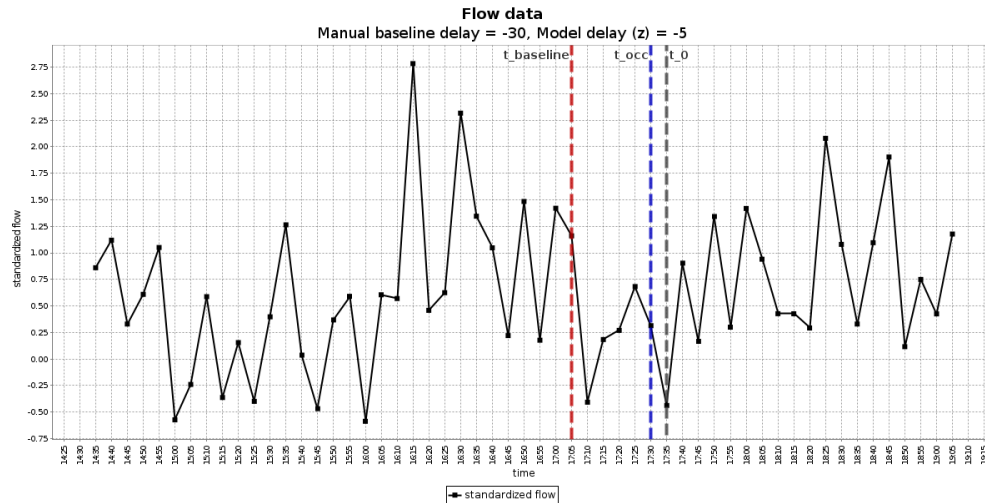


FIGURE 6 Chosen time is the beginning of a new time series.

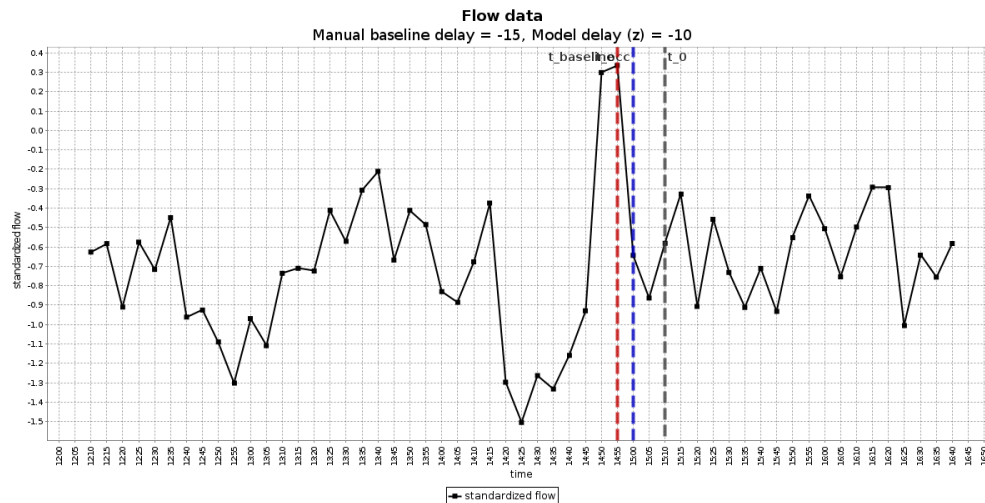


FIGURE 7 Ambiguity due to signal discretization.

1 recovery phase, we'd need at least two change points.

2 **DISCUSSION**

3 As with any other model designed to solve a concrete real-world problem, our methodological
 4 decisions were sensitive to the available dataset, which had 5 min aggregations of flow data and
 5 its own characteristics in terms of sensor quality and spatial distribution. With speed information
 1 for the same sensor locations, we could extend the model in a number of ways. Ideally we would
 2 integrate the speeds into the joint distribution of equation 1, thus Σ in equation 2 would also have
 3 to consider speed-flow signal correlations, for example using the fundamental traffic flow diagram.
 4 Alternatively we could use a quasi-likelihood solution, by assuming independence between speed
 5 and flow.

6 Another extension would be to include upstream sensor data, in a dual-sensor model. Fol-
7 lowing the same principles from this paper, we could obtain two independent approximations, but
8 we should take advantage of their spatial/temporal relationships. Depending on distance to both
9 sensors as well as upstream queue buildup, the disturbance should arrive at potentially distinct, and
10 physically plausible, times to both locations. Furthermore, unless there is an intersection between
11 the sensors, their time series models should be somehow correlated, particularly before the inci-
12 dent occurs. These considerations imply a joint distribution model with both signals, in equation
13 1, and again particular care with covariance matrix Σ to reflect the cross-correlations between the
14 two signals.

15 An extension of this model that considers speed data as well as both upstream and down-
16 stream sensors will be presented in a subsequent article.

17 The empirical evaluation of the results shows that, in general, our model proposes plausible
18 incident occurrence times. However, the lack of observability makes this evaluation essentially
19 subjective. An alternative validation methodology is to use a microsimulation traffic model that
20 is able to simulate incidents (e.g. MITSIMLab (Ben-Akiva et al. (16))). Having noise and spatial
21 models for sensors we can understand the sensitivity of our proposal with respect to these aspects.

22 From the point of view of the modeler, the role of the bayesian prior needs particular
23 attention. While one should not “tweak” it to influence the results towards some subjective goal, it
24 may be a mistake to ignore field knowledge and intuition. The results showed that the knowledge
25 introduced by the field operators was relevant to the quality of the model, as compared to the
26 baseline. More objective solutions could have been explored, such as conditioning the prior on
27 some general heuristic, such as the *largest drop* or the 20% rule suggested by Dia (8).

28 Finally, this model intends ultimately to be applied on a real-time basis. This implies
29 that the data available will be limited, particularly the right bound, r , of the time window will be
30 increasing sequentially in time. An obvious next step is thus to simulate such a sequential model
31 and observe how its predictions evolve accordingly.

32 CONCLUSION

33 We proposed a methodology that fully automates the approximation of incident occurrence time
34 by analysis of downstream sensor flow data. We apply a latent variable framework that uses a
35 change-point time series model as the likelihood function. The actual incident occurrence time
36 has been generally neglected in literature, either being reduced to a calibration parameter during
37 training of incident detection models; or in post-hoc incident analysis works. In both cases, the
38 typical approach is to manually analyze the signal (e.g. (Mak (7))) or apply simple heuristics such
39 as disturbance thresholds in traffic parameters (e.g. (Dia (8))).

40 We built the model over a few principled statements: unless in exceptional cases, the re-
41 ported time should have a delay with respect to actual incident occurrence time; unless under very
42 low flow/high capacity, incidents should affect the traffic flow signal during a period of time right
43 after its occurrence; unless there is an intersection or the distance to the sensor is too high, such
44 disturbance should reach the downstream signal shortly after the incident.

45 We tested our model on a dataset with flows from Singapore expressways for a period of
46 5 months, together with an incident records database. In order to evaluate our model, we manu-
47 ally built a baseline on a sub-set of this dataset. Results show that our model generally proposes
48 plausible incident occurrence time approximations, even when disagreeing with our baseline.

49 One of the biggest challenges of such a framework is effectively the lack of ground-truth

39 and consequently an objective validation. A next step is to use traffic microsimulator model to
40 generate incidents as well as sensor data, and study the quality of our model with respect to sensor
41 data quality and distance to the incidents.

42 We will also extend the model to consider other types of data (e.g. speed information) as
43 well as both the downstream and upstream sensors in a single formulation.

1 ACKNOWLEDGEMENTS

2 The authors gratefully acknowledge Land Transport Authority of Singapore for providing the
3 dataset and helping with local expertise. This research was supported by the National Research
4 Foundation Singapore through the Singapore MIT Alliance for Research and Technology's FM
5 IRG research programme.

6 REFERENCES

- 7 [1] Weil, R., J. Wootton, and A. GarcŠa-Ortiz, Traffic incident detection: Sensors and algorithms.
8 *Mathematical and Computer Modelling*, Vol. 27, No. 9D11, 1998, pp. 257 – 291.
- 9 [2] Karim, A. and H. Adeli, Incident Detection Algorithm using Wavelet Energy Representation
10 of Traffic Patterns. *J. Transp. Eng.*, Vol. 128, No. 3, 2002, p. 232D242.
- 11 [3] Tang, S. and H. Gao, Traffic-incident detection-algorithm based on nonparametric regression.
12 *Intelligent Transportation Systems, IEEE Transactions on*, Vol. 6, No. 1, 2005, pp. 38–42.
- 13 [4] Antoniou, C., H. N. Koutsopoulos, and G. Yannis, Dynamic Data-Driven Local Traffic State
14 Estimation and Prediction. *Transport Research - Part C*, forthcoming.
- 15 [5] Pereira, F., F. Rodrigues, and M. Ben-Akiva, Text analysis in incident duration prediction.
16 *submitted, in review*, 2013.
- 17 [6] Ben-Akiva, M., H. N. Koutsopoulos, C. Antoniou, and R. Balakrishna, Traffic Simulation
18 with DynaMIT. In *Fundamentals of Traffic Simulation* (J. Barceló, ed.), Springer, New York,
19 NY, USA, 2010, pp. 363–398.
- 20 [7] Mak, C. L., *New dual-variable algorithms for detecting lane-blocking incidents on express-*
21 *ways*. Ph.D. thesis, Nanyang Technological University, Singapore, 2002.
- 22 [8] Dia, H., *Artificial Neural Network Models for Automated Freeway Incident Detection*. Ph.D.
23 thesis, Monash University, Australia, 1996.
- 24 [9] Mak, C. and H. Fan, Heavy SSow-based incident detection algorithm using information from
25 two adjacent detector stations. *Journal of Intelligent Transportation Systems*, Vol. 10, No. 1,
26 2006, pp. 23–31.
- 27 [10] Lee, D.-H., S.-T. Jeng, and P. Chandrasekar, Applying data mining techniques for traffic
28 incident analysis. *Journal of The Institution of Engineers, Singapore*, Vol. 44, No. 2, 2004.
- 29 [11] Jeong, Y., M. Castro-Neto, M. Jeong, and L. Han, A Wavelet-Based Freeway Incident De-
30 tection Algorithm with Adapting Threshold Parameters. *Transportation Research Part C D*
31 *Emerging Technologies*, Vol. 19C, No. 1, 2011, pp. 1–19.

- 32 [12] Hawkins, D. M., Fitting multiple change-point models to data. *Computational Statistics &*
33 *Data Analysis*, Vol. 37, No. 3, 2001, pp. 323 – 341.
- 34 [13] Hall, M., E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, The WEKA
35 Data Mining Software: An Update. *SIGKDD Explorations*, Vol. 11, No. 1, 2009.
- 1 [14] Daganzo, C. F., The cell transmission model: A dynamic representation of highway traffic
2 consistent with the hydrodynamic theory. *Transportation Research Part B: Methodological*,
3 Vol. 28, No. 4, 1994, pp. 269–287.
- 4 [15] Kuhne, R. and P. Michalopoulos, Continuum flow models. In *Traffic flow theory: A state of*
5 *the art report - revised monograph on traffic flow theory* (N. H. Gartner, C. Messer, and A. K.
6 Rathi, eds.), Oak Ridge National Laboratory, 1997.
- 7 [16] Ben-Akiva, M., H. Koutsopoulos, T. Toledo, Q. Yang, C. Choudhury, C. Antoniou, and
8 R. Balakrishna, Traffic Simulation with MITSIMLab. In *Fundamentals of Traffic Simula-*
9 *tion* (J. BarcelÓ, ed.), Springer New York, Vol. 145 of *International Series in Operations*
10 *Research & Management Science*, 2010, pp. 233–268.