

## Exploring crash-risk factors using Bayes' theorem and an optimization routine

**Dr Maria-Ioanna M. Imprialou**

Lecturer

School of Civil and Building Engineering

Loughborough University

Loughborough LE11 3TU

United Kingdom

Tel: +44(0)1509 228545

E-mail : [M.Imprialou@lboro.ac.uk](mailto:M.Imprialou@lboro.ac.uk)

**Professor Mike Maher**

Honorary Professor

Dept of Civil, Environmental & Geomatic Engineering

University College London

London WC1E 6BT

United Kingdom

E-mail: [m.maher@ucl.ac.uk](mailto:m.maher@ucl.ac.uk)

**Professor Mohammed Quddus\***

Professor of Intelligent Transport Systems

School of Civil and Building Engineering

Loughborough University

Loughborough LE11 3TU

United Kingdom

Tel: +44(0)1509 228545

E-mail : [M.A.Quddus@lboro.ac.uk](mailto:M.A.Quddus@lboro.ac.uk)

*\* Corresponding author*

Word Count: 5,544 + 5 Figures\*250+ 1 Table\*250 = 7,044

Paper submitted for presentation at the 95<sup>th</sup> Annual Meeting of the Transportation Research Board (TRB) in Washington D.C., USA

Submission date: *31 July 2015*

**ABSTRACT**

Regression models used to analyse crash counts are associated with some kinds of data aggregation (either spatial, or temporal or both) that may result in inconsistent or incorrect outcomes. This paper introduces a new non-regression approach for analysing risk factors affecting crash counts without aggregating crashes. The method is an application of the Bayes' Theorem that enables to compare the distribution of the prevailing traffic conditions on a road network (i.e. *a priori*) with the distribution of traffic conditions just before crashes (i.e. *a posteriori*). By making use of Bayes' Theorem, the probability densities of continuous explanatory variables are estimated using kernel density estimation and a posterior log likelihood is maximised by an optimisation routine (Maximum Likelihood Estimation). The method then estimates the parameters that define the crash risk that is associated with each of the examined crash contributory factors. Both simulated and real-world data were employed to demonstrate and validate the developed theory in which, for example, two explanatory traffic variables speed and volume were employed. Posterior kernel densities of speed and volume at the location and time of crashes have found to be different that prior kernel densities of the same variables. The findings are logical as higher traffic volumes increase the risk of all crashes independently of collision type, severity and time of occurrence. Higher speeds were found to decrease the risk of multiple-vehicle crashes at peak-times and not to affect significantly multiple-vehicle crash occurrences during off-peak times. However, the risk of single vehicle crashes always increases while speed increases.

*Keywords:* Bayes' theorem, maximum likelihood estimation, speed, volume, crashes.

## 1 INTRODUCTION

2 The development of successful road safety measures relies on the in-depth understanding of  
3 the factors and the mechanisms that are related to crash occurrences. There are numerous factors  
4 associated with driving attitudes, traffic or external conditions and often combinations of these factors  
5 that lead to crashes (1, 2). Since the first systematic crash analyses that emerged approximately seven  
6 decades ago, there have been constant and significant advances on our understanding of crash  
7 occurrences (3). The randomness and the complexity that characterise road crashes have not yet  
8 permitted a full explanation of these phenomena.

9 A large proportion of the current knowledge on crash contributory factors comes from crash  
10 count models that attempt to identify and explain crash precursors. Although the statistical approaches  
11 of count models are increasingly more accurate and integrated (4), crash models often provide  
12 inconsistent results even for the most significant traffic variables such as speed and traffic volume.  
13 The assumption that a relationship that applies at a group level (link-level) necessarily applies to an  
14 individual level (crash-level), termed as aggregation bias, is likely to be a source of inconsistent and  
15 erroneous results in crash analyses (5, 6). Count models incorporate by definition this characteristic:  
16 crashes are aggregated (typically with spatial criteria) and are represented in models by characteristic  
17 values of the examined variables that usually reflect the most typical situation on the roadway. In this  
18 way, the actual circumstances of when crashes occurred, which are likely to be unusual, are not  
19 necessarily represented in the models (7, 8).

20 To define accurately the role of potential crash contributory factors it is important to express  
21 the actual traffic conditions that are directly related with the crashes in the models. This paper aims to  
22 set out an approach to the modelling of collision data where the objective is to estimate the effect of  
23 traffic conditions at the time of the crash without aggregating crashes. The data then consist of the  
24 traffic conditions applying at or around the time and location of each crash. These conditions are  
25 therefore a biased sample of the conditions that generally prevail on the network, and this statistical  
26 conditionality needs to be accounted for in any modelling of the crash dataset. Hence the method to  
27 be employed is an application of Bayes' theorem for the case of a continuous variable aided by the  
28 maximum likelihood estimation algorithm. The prior is the distribution of a variable of interest (e.g.  
29 speed) on the network and the likelihood is the crash risk associated with the variable. The posterior  
30 density function describes the distribution of a variable of interest just before crashes. The approach is  
31 applicable for examination of the impact of one or more variables. This study focuses on the effect of  
32 speed and traffic volume on peak and off-peak by considering individual motorway crashes by  
33 severity and collision type.

## 34 LITERATURE REVIEW

35 Speed is considered as an important crash risk factor and is related with a large proportion of  
36 crashes (9–12). This is explainable considering that speed is a potential crash contributory factor that  
37 is always present on the network in contrast to many others that have random (e.g. rainfall) or periodic  
38 character (e.g. darkness) (13). It is justified that crashes that occur under higher speed conditions are  
39 more likely to have serious impact (9, 14–16). Higher travel speeds are also associated with many  
40 crash triggering factors such as lower reaction times, longer decisions, braking and stopping distances,  
41 reduction of manoeuvrability and increased possibilities of manipulation error, loss of control and  
42 exceeding the critical speed on a curve (9, 17–19). On the other hand, higher speeds are also related  
43 with more uniform distribution of speeds (i.e. lower speed variance) that is considered to be beneficial  
44 for road safety (19–22).

45 The studies that examine the role of speed on crashes using statistical modelling have  
46 provided different crash-speed relationships. Most of the studies that applied linear regression models  
47 found driving speeds to be linearly or exponentially related with crashes (23–26). However, a few  
48 studies contradicted the common belief proposing that speed is inversely related to crash frequency  
49 (20) or that the crash-speed relationship is statistically insignificant (22). Some more recent papers  
50 that explored the impact of speed on crashes using advanced crash models have also reported a  
51 statistically insignificant relationship between speed and crashes (27–29) but others suggested that  
52 this relationship is negative (30). However, in models that represent aggregated pre-crash conditions,  
53 speed was found to trigger crashes (7). Pei et al. (16) attempted to explain the results' inconsistencies  
54 suggesting that the crash-speed relationship strongly depends on the selected measure of exposure; the  
55

1 relationship was shown to be negative when the exposure variable was vehicle miles travelled but  
2 positive when the exposure was expressed as vehicle hours travelled.

3 The relationship of speed with crashes cannot be accurately defined without controlling for  
4 the simultaneous effect of other traffic and road characteristics (9). Traffic flow is one of the most  
5 important and well-studied crash precursors in the literature. Naturally, the number of vehicles on the  
6 roadway is directly proportional with the number of vehicle interactions that can potentially lead to  
7 collisions (19). AADT has been found by a large number of researchers to have an increasing  
8 relationship with crash frequency (31–35), that in other words means that the busier a road is, the  
9 more crashes are expected to occur on it. However, not all crash types are equally related with high  
10 traffic; single vehicle crashes have been found to be related with lower traffic than multiple vehicle  
11 crashes (36–40).

12 Inconsistencies or errors in the results might be related with a variety of methodological and  
13 data limitations such as the use of less appropriate statistical models, omission bias and data  
14 aggregation bias (4). The latter is considered to have significant impact on the validity of the models'  
15 results (5, 6), however it is relatively less studied and addressed. Count models that examine crash  
16 frequency encompass some kind of crash data aggregation. The dominant crash aggregation method is  
17 based on topological and temporal criteria. Crash counts that occurred on pre-defined road links (i.e.  
18 link-based models) or areas (i.e. area-wide models) during a certain time period are aggregated and  
19 modelled against selected explanatory variables under the assumption that the crash frequency on a  
20 particular location can be explained by its most frequently observed, or average, conditions. Research  
21 shows however that crashes are related with sudden and unexpected circumstances that may be not  
22 possible to be represented by characteristic values such as annual averages (8).

23 The recently proposed condition-based crash data aggregation approach is an attempt to  
24 address aggregation bias issues by developing crash datasets that represent the actual pre-crash traffic  
25 and geometric conditions (7). In condition-based modelling crashes are allocated to one out of a  
26 number of equally likely pre-crash scenarios that resembled more accurately the conditions on the  
27 roadway just before the crash occurrence. This approach enabled an improved representation of the  
28 conditions that are related with crashes, but it still involves a level of aggregation which might have  
29 impact on the results. This paper proposes an alternative method for defining the role of traffic related  
30 potential risk factors that do not require any form of crash aggregation. The method will be  
31 demonstrated firstly using synthetic data and then it will be applied to real-world motorway crashes in  
32 England with the aim of developing crash-speed-volume relationships

## 33 34 THEORETICAL AND EMPIRICAL DERIVATIONS

35 The purpose is to develop a non-regression based approach for the modelling of crash data  
36 that does not require crash aggregation with the aim of estimating the effect of explanatory variables  
37 at the time of the crash. This is achieved by making use of Bayes' Theorem for the case of a  
38 continuous variable to derive the conditional probability density of the explanatory variables (e.g. the  
39 speed and flow values), given that a crash has occurred. We consider the occurrence of a crash in a  
40 location / time interval combination, which will be referred to here as *segments*. First, the method is  
41 formulated by considering a single variable (i.e. speed). Following, the method is extended to the case  
42 of multiple explanatory variables related to crashes.

### 43 44 One explanatory variable

45 Consider having available  $N$  traffic measurements on all the segments of a network over a  
46 year and  $n$  traffic measurements that represent the traffic conditions at or around the time and location  
47 of each individual crash that occurred on the network during the study period. The probability of a  
48 crash occurring in any of such traffic measurements is very small (i.e.  $N \gg n$ ). This probability is  
49 assumed to be dependent upon the conditions that apply at that point and time (e.g. traffic flow,  
50 speed). For simplicity of presentation it will be firstly assumed that this probability is a function only  
51 of speed ( $v$ ) and so is denoted by  $g(c|v)$  which might typically be of the form  $g(c|v) = bv^\beta$  in  
52 which  $b$  is a constant and  $\beta$  is the speed parameter.

53 Assuming that  $f(v)$  denotes the prior probability density function of speed values over all  
54 segments, and  $g(c|v)$  is the risk of a crash in a segment with speed  $v$ , then Bayes' Theorem for the

1 case of a continuous variable can be applied to derive the posterior density function  $p(v|c)$  describing  
 2 the distribution of speeds applying amongst the crash set:  
 3

$$4 \quad p(v|c) = \frac{f(v)g(c|v)}{\int f(v)g(c|v)dv} = \frac{f(v) v^\beta}{\int f(v) v^\beta dv} \quad (1)$$

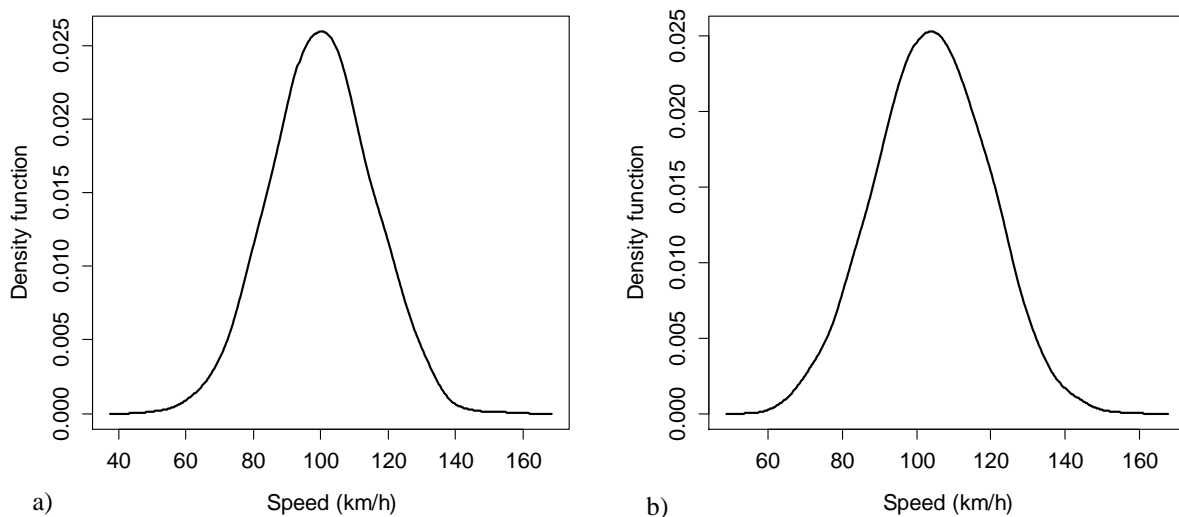
5  
 6 Note that the denominator of equation (1) is (apart from the constant  $b$ ) the probability of a collision  
 7 on any randomly chosen segment, and is a function of the parameter  $\beta$ . Then, if we have a set of  $n$   
 8 crashes with associated speed values  $v_1 \dots v_n$  the log likelihood function can be written as follows:  
 9

$$10 \quad \log L(\beta) = \sum_{j=1}^n \log(p(v_j|c_j)) \quad (2)$$

11 Equation (2) is maximised with respect to  $\beta$  using the *maxLik* function in R (41).

12 To illustrate this approach, we generate a large sample of segments' speeds (say  $N_h$ ) from a normal  
 13 distribution with mean 100 km/h and standard deviation 13 km/h i.e.  $N(100,13^2)$ . For each one, we  
 14 generate a random uniform variate to determine whether or not a crash will occur (with probability  
 15  $bv^\beta$ ) using  $\beta = 2$ , and the values of speed samples ( $N_h$ ) and  $b$  such that about 1,600 crashes are  
 16 generated.

17 Based on the value that was selected for  $\beta$  it is assumed that higher speeds are more  
 18 represented than lower speeds amongst the crash speeds (e.g. the risk of a collision at 120 km/h is  
 19 nine times greater than at 40 km/h, because  $\beta = 2$ ). Figure 1 presents the plots of the density functions  
 20 of a) the prior and b) the posterior distributions of speeds. It can be seen that the latter is slightly  
 21 shifted to the right, as higher speeds are relatively more represented in the crash set than in the  
 segments as a whole.



22 **Figure 1: a) Prior and b) posterior speed distributions**

23 In applying equations (1) and (2), the prior distribution  $f(v)$  must be given – or an estimate of  
 24 it should be employed. Since the distribution of  $f(v)$  is unknown, the kernel density estimate of a  
 25 sample of all speeds is used. The integration in the denominator of Equation (1) is obtained by  
 26 summing over the points used in the kernel density estimation. Using code in R that generates a  
 27 sample of artificial data as described above and then maximises the log likelihood, the estimate of  $\hat{\beta}$   
 28 is found to be 1.92 with a standard error of 0.19, confirming the good behaviour of the method.

## 1 Extension to more than one explanatory variable

2 Crash risk cannot be determined solely by the speed levels on a segment. The model outlined  
 3 above can be extended to include two or more explanatory variables. Assuming that the risk  $g(c|q, v)$   
 4 of a collision on a segment is now a function of traffic volume ( $q$ ) and speed ( $v$ ) so that  $g(c|q, v) =$   
 5  $bq^\alpha v^\beta$ , it is possible to estimate the parameters  $\alpha$  and  $\beta$  from the observed values of the traffic  
 6 volume and speed that prevailed at the time and location of each of a set of collisions. As before, in  
 7 addition to the observations on the crash set, the values of volume and speed for a randomly chosen  
 8 sample of segments are required, so as to provide an estimate of the prior density function  $f(q, v)$ .  
 9 Then Bayes' Theorem (see equation (3)) can again be applied to derive the posterior density function  
 10  $p(q, v|c)$  describing the distribution of flows and speeds applying amongst the crash set:  
 11

$$p(q, v|c) = \frac{f(q, v)g(c|q, v)}{\int \int f(q, v)g(c|q, v)dqdv} = \frac{f(q, v) q^\alpha v^\beta}{\int \int f(q, v) q^\alpha v^\beta dqdv} \quad (3)$$

12

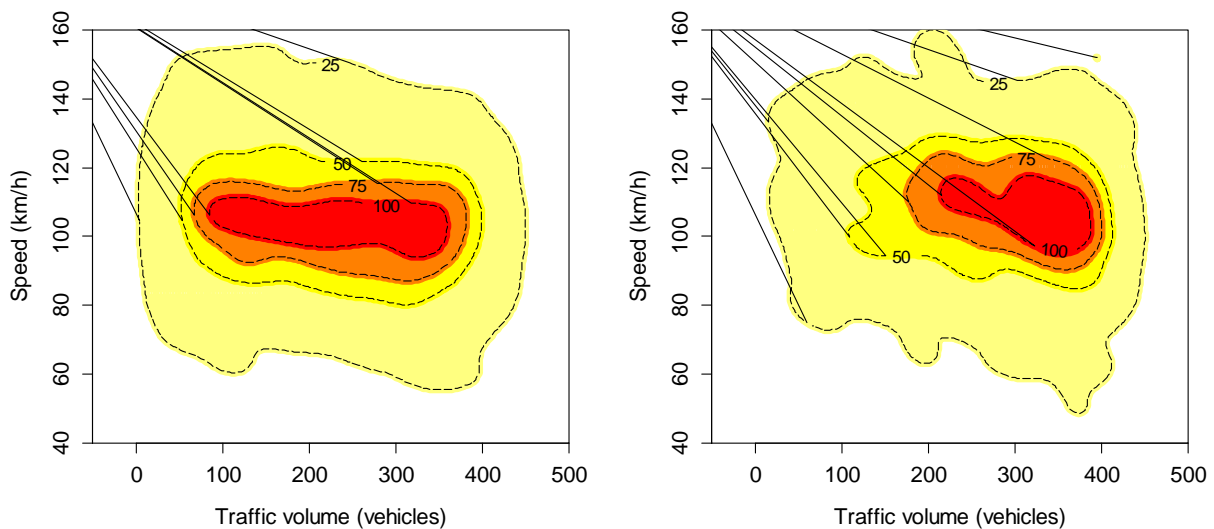
13

The log likelihood is then:

$$\log L(\alpha, \beta) = \sum_{j=1}^n \log (p(q_j, v_j|c_j)) \quad (4)$$

14 The log likelihood was again maximised with respect to the parameters  $\alpha$  and  $\beta$  using the  
 15 *maxLik* function in R (41).

16 For a similar illustration as for the univariate model, artificial data are generated with  $\alpha = 0.85$   
 17 and  $\beta = 2$ , and with some partial relationship between the flow values and the speed values (with  
 18 mean speed reducing as flow increases), but with random (normally distributed) error terms added to  
 19 the mean speed value. The size ( $n_h$ ) of the collision set is around 1,600, as is that of the previous  
 20 example. The estimates obtained from one such dataset are:  $\hat{\alpha} = 0.93$  (standard error = 0.07), and  $\hat{\beta} =$   
 21  $2.04$  (standard error = 0.22), which are reasonably close to the true values of 0.85 and 2.0 respectively.  
 22 The contour plots of kernel density estimates of the prior (left) and posterior (right) are shown in  
 23 Figure 2. The lower 25%, 50% and 75% contours of the density regions are indicated. The underlying  
 24 linear relationship between speed (vertical axis) and flow (horizontal axis) can be seen, as can be seen  
 25 the shift towards relatively higher flow and higher speed as we move from the prior to the posterior.  
 26  
 27



1 **Figure 2: Contour plot of the a) prior and the b) posterior kernel density functions for speed**  
 2 **and volume.**

3  
 4 **DATA DESCRIPTION**

5 As shown in the previous section, the developed models have provided acceptable parameter  
 6 estimates for the synthetic datasets. It is therefore our expectation that the model would exhibit the  
 7 same behaviour for the case of real world data. As has been outlined in the previous section, in order  
 8 to apply the model shown in equation (3) the distributions of the traffic conditions on the road  
 9 network and just before crashes must be known.

10 Network traffic data of the year 2012 were available for the entire Strategic Road Network of  
 11 England (SRN). The SRN that consists of all the motorways and the major A-roads, is the busiest and  
 12 the most significant network of the country (41). Traffic data were extracted from the UK Highways  
 13 Agency Journey Time Database (JTDB) (42) which includes link-level traffic measurements obtained  
 14 by inductive loop detectors which are installed all over the road network. The measurement interval is  
 15 15 minutes and the total number of (junction to junction) links is 2,505, resulting in a dataset of  
 16 approximately 88 million observations. The variables used for this analysis are: 15-minute average  
 17 speed (km/h) and 15-minute volume (vehicles).

18 Crash data were obtained from the National Crash Database, the so-called STATS19 reports  
 19 (43). During 2012 there were overall 10,520 reported crashes of which 1.9% were fatal, 11.7% serious  
 20 and 86.4% slight. Along with a unique reference code there are a number of elements of information  
 21 that are reported for each crash such as:

- 22 • Crash date
- 23 • Time of the crash
- 24 • Crash location (coordinates)
- 25 • Road name
- 26 • Road type
- 27 • Crash Severity
- 28 • Number of involved vehicles
- 29 • Vehicle(s)' direction

30 The traffic conditions on the roadway that are related with crashes are not directly available  
 31 from the STATS19 reports. To obtain this information, attributes of the network and the crash datasets  
 32 were merged. The objective of this merger was to identify as precisely as possible the prevailing  
 33 traffic pre-crash conditions. The reported crash locations were employed to indicate the road links  
 34 where each of the crashes occurred. Crash locations in the form of coordinates are likely to contain  
 35 errors. To improve the accuracy of the matching between links and crashes, crash locations were  
 36 refined using a Fuzzy-logic based crash mapping algorithm that provides almost 99% crash allocation  
 37 (44). In this way, after identifying a road link for each crash, the reported crash date and time were  
 38 used to identify the traffic conditions prior the collision. Initially, the set of measurements on the date  
 39 of the crash that includes the crash time was identified (e.g. for a crash "A" that occurred at 07:03 the  
 40 set of traffic measurements is the 07:00-07:15). The measurement interval of network data is quite  
 41 long (15 minutes) thus for crashes that occurred on the beginning of the interval, such as the one in  
 42 the example, these measurements might be more representative of the traffic conditions after the  
 43 collision rather than before. To avoid this undesirable effect that may lead in misrepresentation of the  
 44 traffic conditions, the measurement set prior to the time of the crash was also considered (e.g. for  
 45 crash "A" this is the 06:45-07:00 traffic measurements set). The final traffic conditions for each crash  
 46 were estimated using a weighted average of the values of these two measurement sets with the  
 47 weights determining the proportion of the 15-minute interval before the crash that corresponds to each  
 48 set. (e.g. for crash "A" the weight for the interval 06:45-07:00 is  $\frac{12}{15} = 0.8$  and for the interval 07:00-  
 49 07:15 is  $\frac{3}{15} = 0.2$ ).

50  
 51  
 52  
 53

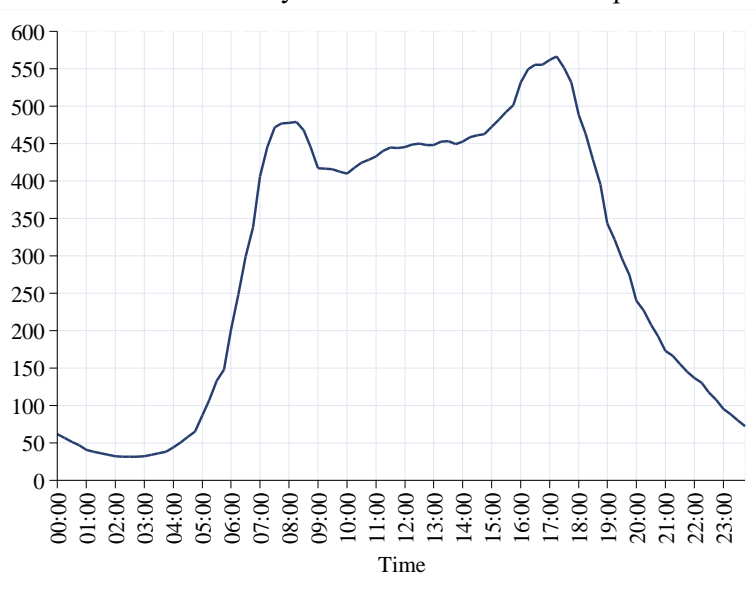
1

2 **RESULTS**

3 Analysing all crashes together might not be the most appropriate way to evaluate crash  
 4 precursors because crashes of different severity levels and collision types might be related with  
 5 dissimilar contributory factors (36, 39). Therefore, main carriageway crashes of the database were  
 6 divided into five categories:

- 7 • All
- 8 • Fatal and serious injury
- 9 • Slight injury
- 10 • Single vehicle
- 11 • Multiple vehicle

12 In addition, to investigate the potential differences between the triggering factors of crashes  
 13 that occur at peak and off-peak hours the aforementioned five crash categories were examined  
 14 separately for peak and off-peak hours. The peak and off-peak time intervals were defined from the  
 15 hourly annual average traffic volume measurements that are presented in Figure 3. Observing the  
 16 volume levels it can be seen that the volume is higher from 7am to 7pm and so this time interval was  
 17 set to define “peak times”. The rest of the day was considered to be “off-peak”



18

19

**Figure 3: Annual average traffic volume over a day.**

20 For each crash category the impact of speed and volume was examined first separately  
 21 (Equations (1), (2)) and then combined (Equations (3), (4)). The speed only models will be henceforth  
 22 referred as SO and the volume only as VO. Models that take into account the effect of speed and  
 23 volume simultaneously will be abbreviated SV.

24 From Table 1 it can be seen that the results between crashes that occurred at peak times are  
 25 significantly different than those for off-peak crashes. The risk of motorway crashes, apart from the  
 26 single vehicle ones, during peak times was found to increase at lower speeds and higher volumes.  
 27 This finding implies that crash frequency at peak times is more related with congestion than speed.  
 28 This is consistent with studies that examined the impact of congestion on crash frequency and severity  
 29 (45, 46). Multiple vehicle crashes were found to have the strongest relationship with traffic volume a  
 30 finding similar to those of other researchers (36, 46). The exponent of the risk function for traffic  
 31 volume was estimated to be 1.1 which indicates, for instance, that the risk of crashes when the traffic  
 32 volume is 400 is 2.14 (i.e. 21.1) times higher than when it is 200, *ceteris paribus*. The risk of single  
 33 vehicle crashes that occurred at peak times is proportional to both the traffic speed and volume, which  
 34 is also in line with literature that suggests that single vehicle crashes are one of the most speed-related  
 35 crash types (47). Figure 4 shows that the posterior density is shifted to the right compared to the prior



1 density. This implies that high speeds are more represented in the crash dataset than the network  
2 dataset.

3 Table 1 summarises the estimated parameters of each model along with their standard errors.  
4 The estimates represent the exponents ( $\alpha$  and  $\beta$ ) of the risk functions (i.e.  $g(c|v)$  for the SO model,  
5  $g(c|q)$  for the VO model and  $g(c|v, q)$  for the SV model). Negative exponents show that as the value  
6 of the examined variable decreases the probability of a crash increases and positive exponents imply  
7 the opposite. Figures 4 and 5 show the contour plots of the prior and the posterior speed-volume  
8 distributions for all peak and off-peak crashes respectively.

9 From Table 1 it can be seen that the results between crashes that occurred at peak times are  
10 significantly different than those for off-peak crashes. The risk of motorway crashes, apart from the  
11 single vehicle ones, during peak times was found to increase at lower speeds and higher volumes.  
12 This finding implies that crash frequency at peak times is more related with congestion than speed.  
13 This is consistent with studies that examined the impact of congestion on crash frequency and severity  
14 (45, 46). Multiple vehicle crashes were found to have the strongest relationship with traffic volume a  
15 finding similar to those of other researchers (36, 46). The exponent of the risk function for traffic  
16 volume was estimated to be 1.1 which indicates, for instance, that the risk of crashes when the traffic  
17 volume is 400 is 2.14 (i.e.  $2^{1.1}$ ) times higher than when it is 200, *ceteris paribus*. The risk of single  
18 vehicle crashes that occurred at peak times is proportional to both the traffic speed and volume, which  
19 is also in line with literature that suggests that single vehicle crashes are one of the most speed-related  
20 crash types (47). Figure 4 shows that the posterior density is shifted to the right compared to the prior  
21 density. This implies that high speeds are more represented in the crash dataset than the network  
22 dataset.

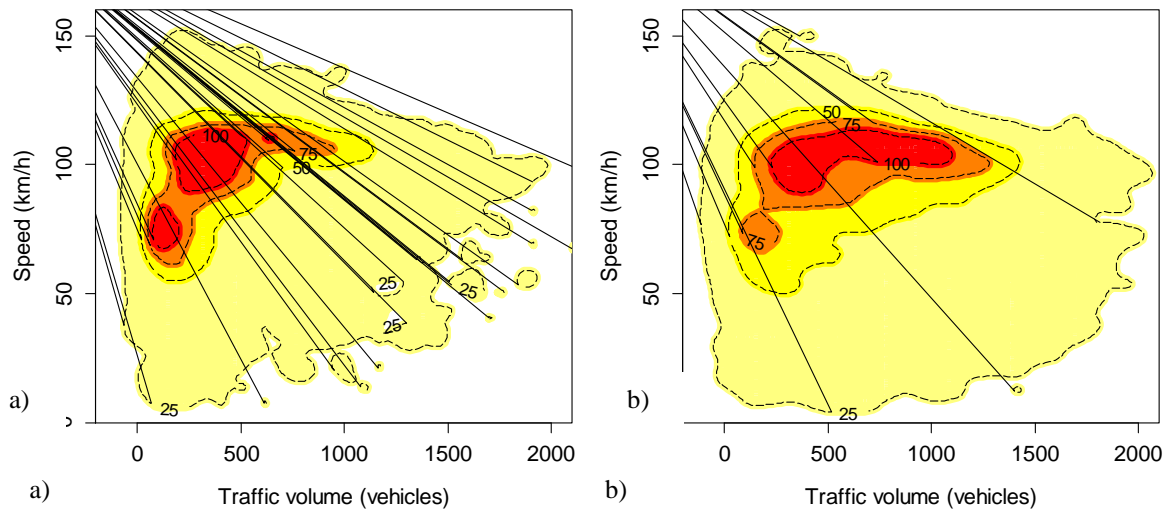
23 **Table 1: Models' parameter estimations for peak and off-peak crashes.**

Peak Times										
Variable	All crashes		Fatal & serious-injury crashes		Slight-injury crashes		Single vehicle crashes		Multiple vehicle crashes	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
Speed (SO)	<b>-0.655</b>	0.034	<b>-0.329</b>	0.114	<b>-0.691</b>	0.035	<b>0.565</b>	0.108	<b>-0.865</b>	0.034
Volume (VO)	<b>0.774</b>	0.019	<b>0.338</b>	0.050	<b>0.836</b>	0.020	<b>0.139</b>	0.033	<b>1.004</b>	0.023
Speed (SV)	<b>-0.989</b>	0.030	<b>-0.545</b>	0.109	<b>-1.035</b>	0.031	<b>0.461</b>	0.107	<b>-1.213</b>	0.031
Volume (SV)	<b>0.859</b>	0.019	<b>0.387</b>	0.051	<b>0.924</b>	0.020	<b>0.104</b>	0.033	<b>1.102</b>	0.022
Off-Peak Times										
Variable	All crashes		Fatal & serious-injury crashes		Slight-injury crashes		Single vehicle crashes		Multiple vehicle crashes	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
Speed (SO)	<b>0.761</b>	0.119	<b>0.650</b>	0.288	<b>0.786</b>	0.141	<b>0.848</b>	0.183	<b>0.682</b>	0.175
Volume (VO)	<b>0.525</b>	0.022	<b>0.266</b>	0.048	<b>0.587</b>	0.025	<b>0.216</b>	0.029	<b>0.858</b>	0.033
Speed (SV)	0.092	0.120	0.148	0.279	0.077	0.137	<b>0.448</b>	0.191	-0.235	0.158
Volume (SV)	<b>0.472</b>	0.023	<b>0.241</b>	0.046	<b>0.532</b>	0.026	<b>0.190</b>	0.027	<b>0.815</b>	0.034

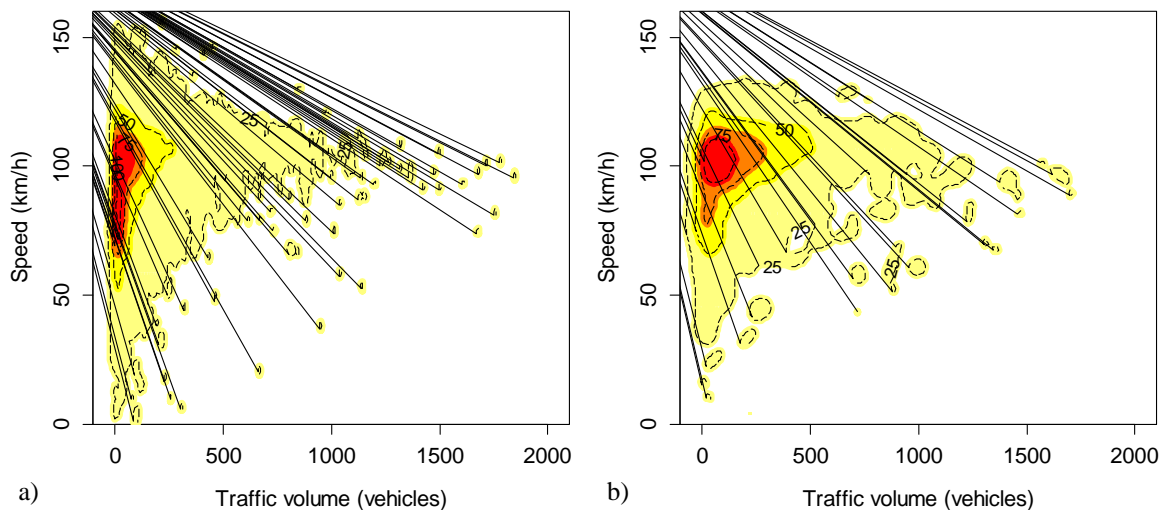
\*Statistically significant estimations are marked with bold font.

24 Off-peak crashes were found to be more related with high speed conditions, a result rather  
25 expected as driving with inappropriate speed tends to be more frequent at night time (48). When speed  
26 was examined individually (SO models), the exponent of the risk function was found to be positive  
27 for all crash types and severities. However, when the models controlled for traffic volume the effect  
28 of speed on the risk function was eliminated. In SV models the parameter of speed for all crashes  
29 excluding the single vehicle ones, was found to be statistically insignificant. For single vehicle  
30 crashes the exponent for speed is 0.448 which means that the risk of single vehicle crash at off-peak  
31 times when speed is 120 km/h is approximately 1.2 (i.e.  $1.5^{0.448}$ ) times higher than when speed is 80  
32 km/h. The effect of traffic volume on off-peak crash frequency is positive for all off-peak crashes but  
33

1 relatively lower compared to those of peak time crashes. From Figure 5 it can be seen that the  
 2 posterior density has a small rightwards shift towards the axis of volume. There is also a slight shift  
 3 up in the lower section of the contours, which indicates that low speeds were less observed in the off-  
 4 peak crash dataset than in the network dataset.  
 5



6 **Figure 4: Contour plots of the a) prior and b) posterior kernel density of speed and volume for**  
 7 **all crashes at peak times**  
 8



9 **Figure 5: Contour plots of the a) prior and b) posterior kernel density of speed and volume for**  
 10 **all crashes at off-peak times**  
 11

## 12 CONCLUSIONS

13 The purpose of this paper was to demonstrate and apply a new method for modelling crash  
 14 data in order to define the impact of the levels of traffic speed and volume on motorway crashes. In

1 contrast to count regression models, that are typically employed for examining crash frequency, the  
 2 proposed method does not require crash data aggregation but deals explicitly with the traffic  
 3 conditions prevailing at the time and location of each recorded crash. In this way, aggregation bias  
 4 problems that are expected to decrease the reliability of modelling outcomes are avoided. As it can be  
 5 expected that traffic conditions influence the likelihood of crashes, it follows that the distribution of  
 6 traffic conditions amongst the crash set (the posterior distribution) will be different from that  
 7 prevailing in the population of traffic conditions generally (the prior distribution). The method is  
 8 theoretically founded on the estimation of conditional probability as expressed by the Bayes' Theorem  
 9 and, by effectively comparing the prior and posterior distributions of traffic conditions, it aims to  
 10 estimate the parameters that define the risk that is associated with each of the examined crash  
 11 contributory factors. The application of the proposed model on synthetic datasets validated its  
 12 capability of providing reliable estimates before being applied to the real-world datasets.

13 The results provide useful insight on the relationship of traffic conditions with crashes. In  
 14 summary, peak-time crashes, apart from the single vehicle ones but including those that had serious  
 15 outcomes, were found to be mainly related with congested traffic; they increase while speed drops and  
 16 volume rises. The relationship of off-peak crashes with speed was found to be positive when it was  
 17 examined individually but in models that control for the effect of traffic volume, its effect was  
 18 negligible. However, the risk of single vehicle crashes independently of the time of their occurrence  
 19 was found to be proportional with speed even when volume was included in the model. From the  
 20 modelling outcomes it is confirmed that high traffic volumes are linked with more crashes of all types  
 21 and severities although the maximum effect is on multiple vehicle collisions. The relationship of  
 22 speed with crashes varies according to the general traffic conditions and therefore it is not  
 23 straightforward to reply to the question of whether speed has negative impact of road safety. However,  
 24 considering that speed was found to contribute to crashes that are related with low density conditions  
 25 (i.e. single vehicle, off-peak) when drivers tend to be free to choose their own speeds, speeding can be  
 26 regarded as crash triggering action.

27 In order to better understand crash risk, more factors should be taken into consideration.  
 28 Traffic density, speed variance and v/c ratio are some of them. The absence of these factors in the  
 29 developed models may have affected the estimations. Another limitation of this study is that the  
 30 measurement interval of the traffic measurements (i.e. 15 minutes) was quite long to provide precise  
 31 pre-crash conditions estimations. In the future, similar analyses should include a more complete  
 32 specification of risk that will lead to a broader view of the impact of traffic characteristics on crashes  
 33 and also higher resolution measurements (e.g. 15 seconds). An interesting extension of the model  
 34 would also include the development risk functions that will incorporate spatial characteristics that are  
 35 possibly related with crash occurrences such as road geometry.

## 36 REFERENCES

- 38 1. Brown ID (1982) Exposure and experience are a confounded nuisance in research on driver  
 39 behaviour. *Accid Anal Prev* 14:345–352.
- 40 2. Montella A (2011) Identifying crash contributory factors at urban roundabouts and using  
 41 association rules to explore their relationships to different crash types. *Accid Anal Prev*  
 42 43:1451–63.
- 43 3. Hagenzieker MP, Commandeur JJF, Bijleveld FD (2014) The history of road safety research:  
 44 A quantitative approach. *Transp Res Part F Traffic Psychol Behav* 25:150–162.
- 45 4. Lord D, Mannering F (2010) The statistical analysis of crash-frequency data: A review and  
 46 assessment of methodological alternatives. *Transp Res Part A Policy Pract* 44:291–305.
- 47 5. Davis GA (2004) Possible aggregation biases in road safety research and a mechanism  
 48 approach to accident modeling. *Accid Anal Prev* 36:1119–1127.
- 49 6. Clark W, Avery K (1976) The effects of data aggregation in statistical analysis. *Geogr Anal*  
 50 8:428–438.
- 51 7. Imprialou M-IM, Quddus M, Pitfield DE, Lord D (2016) Re-visiting crash–speed relationships:  
 52 A new perspective in crash modelling. *Accid Anal Prev* 86:173–185.

- 1 8. Hossain M, Muromachi Y (2013) Understanding crash mechanism on urban expressways  
2 using high-resolution traffic data. *Accid Anal Prev* 57:17–29.
- 3 9. Aarts L, Van Schagen I (2006) Driving speed and the risk of road crashes: A review. *Accid*  
4 *Anal Prev* 38:215–224.
- 5 10. Clarke DD, Ward P, Bartle C, Truman W (2010) Killer crashes: fatal road traffic accidents in  
6 the UK. *Accid Anal Prev* 42:764–70.
- 7 11. Shibata A, Fukuda K (1994) Risk factors of fatality in motor vehicle traffic accidents. *Accid*  
8 *Anal Prev* 26:391–397.
- 9 12. Assum T (1997) Attitudes and road accident risk. *Accid Anal Prev* 29:153–159.
- 10 13. Elvik R, Christensen P, Amundsen A (2004) Speed and road accidents: An evaluation of the  
11 Power Model.
- 12 14. Joksch HC (1993) Velocity change and fatality risk in a crash—A rule of thumb. *Accid Anal*  
13 *Prev* 25:103–104.
- 14 15. Kloeden CN, Mclean AJ, Moore VM, Ponte G (1997) Travelling Speed and the Risk of Crash  
15 Involvement Volume 1 - Findings. South Australia
- 16 16. Pei X, Wong SC, Sze NN (2012) The roles of exposure and speed in road safety analysis.  
17 *Accid Anal Prev* 48:464–71.
- 18 17. Fildes B, Lee S (1993) The Speed review: Road Environment, Behaviour, Speed Limits,  
19 Enforcement and Crashes. Victoria
- 20 18. Hale AE (1990) Safety and speed. A systems view of determinants and control measures.  
21 *IATSS Res.* 14:
- 22 19. Navon D (2003) The paradox of driving speed: two adverse effects on highway accident rate.  
23 *Accid Anal Prev* 35:361–7.
- 24 20. Lave C (1985) Speeding , Coordination , and the 55 MPH Limit. *Am Econ Assoc* 75:1159–  
25 1164.
- 26 21. Graves P, Lee D, Sexton R (1993) Speed variance, eforcement and the optimal speed limit.  
27 *Econ Lett* 42:237–243.
- 28 22. Garber NJ, Gadiraju R (1989) Factors Affecting Speed Variance and Its Influence on  
29 Accidents. *Transp Res Rec J Transp Res Board* 1213:64–71.
- 30 23. Fildes BN, Rumbold G, Leening A (1991) Speed Behaviour and Drivers' Attitude to Speeding.  
31 Victoria,Austalia
- 32 24. Baruya A, Finch DJ (1994) Investigation of traffic speeds and accidents on urban roads.  
33 *Traffic Manag. Road Safety. Proc. Semin. J held PTRC Eur. Transp. Forum,September 12-16,*  
34 *Vol. P381*
- 35 25. Quimby A, Maycock G, Palmer C, Buttress S (1999) The factors that influence a driver ' s  
36 choice of speed — a questionnaire study. Transport Research Laboratory, Crowthorne,  
37 England
- 38 26. Taylor MC, Lynam DA, Baruya A (2000) The effects of drivers' speed on the frequency of  
39 road accidents. Transport Research Laboratory, Crowthorne, England
- 40 27. Kweon Y-J, Kockelman K (2005) Safety effects of speed limit changes: Use of panel models,  
41 including speed, use, and design variables. *Transp Res Rec J Transp Res Board* 148–158.
- 42 28. Quddus M (2013) Exploring the Relationship Between Average Speed, Speed Variation, and  
43 Accident Rates Using Spatial Statistical Models and GIS. *J Transp Saf Secur* 5:27–45.
- 44 29. Kockelman KM, Ma J (2007) Freeway Speeds and Speed Variations Preceding Crashes ,  
45 Within and Across Lanes. *Transp Res Forum* 46:43–61.
- 46 30. Baruya A (1998) Speed-accident relationships on European roads. 9th Int. Conf. Road Saf. Eur.  
47 1:
- 48 31. Milton J, Mannering F (1998) The relationship among highway geometrics , traffic-related  
49 elements and motor-vehicle accident frequencies. *Transportation (Amst)* 25:395–413.
- 50 32. Abdel-Aty M, Radwan AE (2000) Modeling traffic accident occurrence and involvement.  
51 *Accid Anal Prev* 32:633–42.
- 52 33. Miaou SP, Lum H (1993) Modeling vehicle accidents and highway geometric design  
53 relationships. *Accid Anal Prev* 25:689–709.
- 54 34. Anastasopoulos PC, Mannering FL (2009) A note on modeling vehicle accident frequencies  
55 with random-parameters count models. *Accid Anal Prev* 41:153–9.

- 1 35. Chang L-Y (2005) Analysis of freeway accident frequencies: Negative binomial regression  
2 versus artificial neural network. *Saf Sci* 43:541–557.
- 3 36. Qin X, Ivan JN, Ravishanker N (2004) Selecting exposure measures in crash rate prediction  
4 for two-lane highway segments. *Accid Anal Prev* 36:183–191.
- 5 37. Martin JL (2002) Relationship between crash rate and hourly traffic flow on interurban  
6 motorways. *Accid Anal Prev* 34:619–629.
- 7 38. Lord D, Manar A, Vizioli A (2005) Modeling crash-flow-density and crash-flow-V/C ratio  
8 relationships for rural and urban freeway segments. *Accid Anal Prev* 37:185–99.
- 9 39. Kim DG, Washington S, Oh J (2006) Modeling Crash Types: New Insights into the Effects of  
10 Covariates on Crashes at Rural Intersections. *J Transp Eng* 132:282–292.
- 11 40. Bham GH, Javvadi BS, Manepalli URR (2012) Multinomial Logistic Regression Model for  
12 Single-Vehicle and Multivehicle Collisions on Urban U.S. Highways in Arkansas. *J Transp*  
13 *Eng* 138:786–797.
- 14 41. Henningsen A, Toomet O (2011) maxLik: A package for maximum likelihood estimation in R.  
15 *Comput Stat* 26:443–458.
- 16 42. Department for Transport (2011) Road Network Policy Consultation. London, England
- 17 43. Highways Agency (2011) HATRIS JTDB Reference Manual.
- 18 44. Department for Transport (2011) STATS19 road accident injury statistics – report form.  
19 [https://www.gov.uk/government/uploads/system/uploads/attachment\\_data/file/230590/stats19](https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/230590/stats19)  
20 .pdf.
- 21 45. Imprialou M-IM, Quddus M, Pitfield DE (2014) High accuracy crash mapping using fuzzy  
22 logic. *Transp Res Part C Emerg Technol* 42:107–120.
- 23 46. Wang C, Quddus M a., Ison SG (2011) Predicting accident frequency at their severity levels  
24 and its application in site ranking using a two-stage mixed multivariate model. *Accid Anal*  
25 *Prev* 43:1979–1990.
- 26 47. Yan X, Radwan E, Abdel-Aty M (2005) Characteristics of rear-end accidents at signalized  
27 intersections using multiple logistic regression model. *Accid Anal Prev* 37:983–95.
- 28 48. Lang SW, Waller PF, Shope JT (1996) Adolescent driving: Characteristics associated with  
29 single-vehicle and injury crashes. *J Safety Res* 27:241–257.
- 30 49. Ivan JN, Pasupathy RK, Ossenbruggen PJ (1999) Differences in causality factors for single  
31 and multi-vehicle crashes on two-lane roads. *Accid Anal Prev* 31:695–704.
- 32