

# Derivation of the Empirical Bayesian method for the Negative Binomial-Lindley generalized linear model with application in traffic safety

Ali Khodadadi<sup>a,\*</sup>, Ioannis Tsapakis<sup>b</sup>, Mohammadali Shirazi<sup>c</sup>, Subasish Das<sup>d</sup> and Dominique Lord<sup>a</sup>

<sup>a</sup>Texas A&M University, 3136 TAMU, College Station, TX 77843-3136

<sup>b</sup>Texas A&M Transportation Institute 3500 NW Loop 410, Suite 315 San Antonio, TX 78229

<sup>c</sup>University of Maine, Orono, Maine, 04469

<sup>d</sup>Texas A&M Transportation Institute 3135 TAMU, College Station, TX 77843

## ARTICLE INFO

### Keywords:

Empirical Bayesian, Full Bayesian, Negative Binomial-Lindley, Hot spot identification, Expected crash values, Crash prediction model

## ABSTRACT

The expected crash frequency is the long-term average crash count for a specific site. It is extensively used to systematically evaluate the crash risk associated with roadway elements. To estimate the expected crashes, the Empirical Bayesian (EB) approach is typically employed. The EB method is a computationally convenient approximation to the Full Bayesian (FB) method, which gained popularity due to its simple interpretation, computational efficiency, and the ability to account for the regression to the mean bias. However, the common EB method used in traffic safety analysis is only applicable when the traditional Negative Binomial (NB) model is used. The NB model, however, is not appropriate when data is highly dispersed, skewed, or has a large number of zero observations. The Negative Binomial-Lindley (NB-L) model is a mixture of the NB and Lindley distributions and has shown superior fit compared to the NB model, especially when the dataset is characterized by excess zero observations. Even though several studies have used the NB-L in developing crash prediction models, the application of the NB-L in other safety-related tasks (e.g., hot spot identification) is largely neglected. This study proposed a framework to develop the EB method for the NB-L model and subsequently estimate the expected crash values. A comparison between the EB and FB estimates was performed to validate the approximation framework in general. The results indicated that the proposed EB framework is able to estimate expected crashes with comparable precision to the FB estimate, but with much less computational cost. In addition, a site ranking analysis using the EB estimates was conducted to validate the proposed approximation method in safety studies. However, it should be noted that any other type of safety analysis that requires access to the expected crashes can benefit from the proposed EB method. This study concluded that the proposed EB framework can properly approximate the underlying FB approach and can reasonably be considered as an alternative to the traditional EB formula derived from the NB model. The results of this study can help to extend the application of the advanced predictive models beyond predicting crashes to other safety-related tasks, with no additional computational efforts.

## 1. Introduction

The roadway safety management process involves multiple steps that are designed to monitor and reduce crash frequencies on existing roadways (Part, 2010). Of these steps, hot spot identification and safety effectiveness evaluation are two key approaches in safety evaluation and analysis. Hot spot identification identifies sites that can benefit the most from safety treatments. Safety effectiveness evaluation (e.g., before-after analysis, cross-sectional analysis) evaluates how safety has changed because of one or more specific treatments implemented to reduce the crashes. Both analyses

✉ a.khodadadi1994@tamu.edu (A. Khodadadi); i-tsapakis@tti.tamu.edu (I. Tsapakis); shirazi@maine.edu (M. Shirazi); s-das@tti.tamu.edu (S. Das); dlord@civil.tamu.edu (D. Lord)  
ORCID(s): 0000-0002-3413-8687 (A. Khodadadi)

46 require reliable and stable measures to quantitatively evaluate the crash risk associated with a roadway entity in a certain  
47 time period. There are three main steps associated with each of the aforementioned analyses. The first step involves  
48 developing a crash prediction model (also referred to as crash-frequency model). Crash prediction models are the main  
49 tool to predict crash frequencies and identify crash contributing factors. In the second step, the crash prediction models  
50 are used to assess the crash risk associated with each roadway element. The evaluated crash risk then can be used to  
51 determine the likelihood of crash occurrence for each specific site, as a function of site characteristics (*i.e.*, explanatory  
52 variables), in a certain time period. The final step involves ranking the sites in decreasing order based on the assessed  
53 crash risk (in case of hot spot identification), or determining the efficacy of the countermeasure(s) given the assessed  
54 crash risk before and after implementing the treatment.

55 The Negative Binomial (NB) is the most common statistical model to develop crash prediction models and estimate  
56 crash frequencies (Lord et al., 2021; Lord and Mannering, 2010; Mannering and Bhat, 2014). As opposed to the  
57 Poisson distribution which assumes the mean and variance of crash observations are equal, the NB distribution allows  
58 the variance of the response variable to be greater than the mean by using an additional parameter (referred to as over-  
59 dispersion parameter). Although research studies showed that the NB model addresses over dispersion commonly  
60 observed in crash data, this model does not necessarily account for issues related to other unique characteristics of  
61 crash data. In particular, crash datasets are often characterized by excess zero observations or low sample mean.  
62 The NB distribution is not flexible enough to deal with abundance of zero observations in the data (Geedipally et al.,  
63 2012). In addition, the NB model will output biased results as the sample mean goes lower (Lord, 2006). Different  
64 statistical models have been proposed by safety researchers to overcome limitations of the NB model. Poisson log-  
65 normal (Song et al., 2006; Park and Lord, 2007; Khazraee et al., 2018; Shirazi and Lord, 2019), Poisson-generalized  
66 inverse Gaussian (Zha et al., 2016; Zou et al., 2013), Conway-Maxwell-Poisson (Lord et al., 2010; Abdella et al.,  
67 2019), Semiparametric NB model (Shirazi et al., 2016), Poisson-Tweedie (Debrabant et al., 2018; Saha et al., 2020),  
68 Generalized Additive Models (Xie and Zhang, 2008), and Negative Binomial-Lindley (NB-L) (Zamani and Ismail,  
69 2010; Lord and Geedipally, 2011; Geedipally et al., 2012; Shirazi et al., 2017; Shaon et al., 2018; Khodadadi et al.,  
70 2021) are just a few examples of advanced count models developed over time to overcome or alleviate the limitations of  
71 the NB model. NB-L in particular is the subject of interest in this study. The NB-L model is a mixture of the negative  
72 binomial and Lindley distribution. This model was first proposed by Zamani and Ismail (2010), and then used in  
73 multiple research fields dealing with sparse count data modeling including crash data analysis (Lord and Geedipally,  
74 2011; Geedipally et al., 2012; Shaon et al., 2018; Khodadadi et al., 2021). The NB-L model offers extra flexibility using  
75 the Lindley distribution, resulting in a more powerful tool to fit to crash datasets (Shirazi et al., 2016). In particular,  
76 compared to the traditional NB models, the NB-L shows a better fit when a crash dataset contains many zero responses,  
77 or exhibits high dispersion, large skewness or long tail (Shirazi et al., 2017).

78 The three steps mentioned above (*i.e.*, developing crash prediction model, crash risk evaluation, and ranking/before-  
79 after analysis based on the evaluated crash risk) have been fully investigated for the well-known NB model. However,  
80 despite the superiority of the NB-L, no study has examined the application of the NB-L distribution or its general-  
81 ized linear model (GLM) beyond the first step (predicting crashes). This study fills this research gap by deriving the  
82 equations to estimate the expected crash frequency for the NB-L model based on Full Bayesian (FB) and Empirical  
83 Bayesian (EB) framework. Therefore, the primary objectives of this study are to (1) develop an EB framework to  
84 calculate the expected crash values for the NB-L models, (2) compare the EB and FB expected values to determine  
85 if the proposed EB framework properly approximates the underlying FB paradigm, and (3) test the application of the  
86 NB-L and its EB estimates of the expected crashes in other safety-related analyses (site ranking in this study) to ensure  
87 the applicability of the proposed framework.

## 88 **2. Background**

89 The NB-L distribution has been used by researchers in various fields, including safety analysis (crash prediction  
90 models). Lord and Geedipally (2011) examined the application of the NB-L distribution in highway safety. They  
91 applied both the NB and NB-L distributions to simulated and empirical sparse datasets. They found that the NB-L  
92 outperforms the traditional NB distribution. To extend the application of NB-L in safety analysis, Lord et al. (2012)  
93 introduced a generalized linear NB-L model (NB-L GLM) to link the crash frequencies to the site characteristics.  
94 The regression approach has been employed in numerous transportation-related studies to estimate the relationships  
95 between the response variable and influential factors (Safaei et al., 2021b; Darzian Rostami et al., 2020; Aman et al.,  
96 2021; Aman and Smith-Colin, 2020; Safaei et al., 2021a; ?; Asgharpour et al., 2021). Lord et al. (2012) observed  
97 that the NB-L GLM provides a better fit compared to the traditional NB GLM when analyzing a sparse or highly-  
98 dispersed dataset. Given the superiority of the NB-L over the traditional count models, different parameterizations  
99 of the NB-L model have been proposed, discussed, and applied in the literature. Two-parameters NB-L (Zamani and  
100 Ismail, 2010), three-parameters NB-L (Denthet et al., 2016), four-parameters NB-L (Tajuddin et al., 2020), Negative  
101 Binomial weighted-Lindley (NB-WLindley) (Khodadadi et al., 2022), and Negative Binomial-Lindley with different  
102 variance and dispersion structure (Khodadadi et al., 2021) are a few examples of the more advanced and more complex  
103 count models that are recently proposed to provide even greater flexibility to the original NB-L model.

104 Sometimes crash risk is quantified by criteria such as short-term crash frequency, crash rate, crash severity, or  
105 crash cost (Miaou and Song, 2005; Huang et al., 2009; Guo et al., 2020); however, ignoring the influential crash fac-  
106 tors (*e.g.*, Annual Average Daily Traffic, roadway characteristics) could make these methods inefficient. In addition, the  
107 uncertainty associated with using the raw crash data could reduce the accuracy of the results, especially for long-term  
108 planning processes. The limitations associated with using the historical crash records alone led the researchers and

109 transportation agencies to develop statistical approaches to more accurately predict the crash risk (*i.e.*, expected crash  
110 risk); they then used these approaches to rank the sites by the magnitude that their estimated crash risk exceeded the  
111 normal crash risk, which is estimated using sites with similar characteristics (Huang et al., 2009). The expected crash  
112 frequency is the long-term average crash count for a specific site. Considering the expected number of crashes in hot  
113 spot identification can overcome or minimize issues such as the regression-to-the-mean (RTM) bias (Hauer, 1997) or  
114 limited sample size (Miaou and Lord, 2003). Furthermore, given that the expected crash frequency uses both observed  
115 crash data and the number of crashes estimated from a crash prediction model, it can also account for the fundamentally  
116 non-linear relationship between the crash frequency and explanatory variables, the unobserved heterogeneity among  
117 the sites (Lord and Mannering, 2010), and the uncertainty associated with parameters of the underlying regression  
118 model (Miaou and Lord, 2003). The FB and EB are the two methods that are applied to estimate the expected crash  
119 frequencies. The FB method requires access to the hierarchical representation of the underlying predictive model in  
120 order to draw random samples from the posterior distribution of the parameters of interest. The hierarchical repre-  
121 sentation of Bayesian models makes the FB approach more flexible than other methods since it eliminates the need  
122 for the closed form representation of the model. Hierarchical models are frequently used in crash data analysis. One  
123 of the main advantages of the hierarchical models is the ability to incorporate relevant prior knowledge and common  
124 beliefs about the parameters into the modeling process in a natural probabilistic way. The FB method has broadly been  
125 used in various safety-related analyses such as estimating crash prediction models, before-and-after studies, and hot  
126 spot identification (Guo et al., 2019; Farid et al., 2017; Aguero-Valverde and Jovanis, 2009; Miaou and Song, 2005;  
127 Miranda-Moreno et al., 2013; Shirazi et al., 2017; Lan and Persaud, 2011; Persaud et al., 2010; Pu et al., 2020).

128 Despite the broad applications of the FB approach, this method is often computationally intensive. In particular,  
129 for complex models involving a large number of observations and many variables, the FB method can be a time-  
130 consuming task due to the integration over the distribution of many parameters. Furthermore, the FB approach requires  
131 consideration of a prior distribution on all the unknown parameters. However, finding a well-reasoned and well-defined  
132 prior distribution for the problem in hand could be quite challenging. The EB method is a promising alternative to the  
133 standard FB paradigm. The EB approach is a special case of the general FB framework when some assumptions are  
134 simplified. Unlike the FB method where each parameter is defined as a random variable, the EB approach assumes  
135 that the parameters in the highest level of hierarchy are known without any uncertainty (Huang et al., 2009). The EB  
136 method could be thought of a computationally convenient approximation to the FB method, and has gained popularity  
137 among the safety analysts and transportation agencies due to its simple interpretation, computational efficiency, and  
138 the ability to account for the RTM bias (Miaou and Lord, 2003; Huang et al., 2009; Persaud et al., 2010; Khattak  
139 et al., 2018; Das et al., 2019). Despite the fact that the EB method is a reliable method to estimate the expected crash  
140 risks, it is just an approximation to a more general FB paradigm. First of all, the EB method does not account for the

141 uncertainty embedded in the parameters. Parameters of the crash prediction model are estimated from the observed  
142 crash data which are naturally subjected to uncertainty. Ignoring these uncertainties might lead to overestimating the  
143 precision and/or less accurate estimates (Miaou and Lord, 2003). Secondly, the EB method might be criticized for  
144 a double usage of the data (Huang et al., 2009; Hauer, 1997). In an ideal EB procedure, two sources of data should  
145 be used. One source is to develop the crash prediction model and get the predicted crash values, and another is to  
146 independently enrich the model with prior knowledge. However, in practice, the safety performance functions (SPF)  
147 are obtained from the recorded crash frequency and thus both predicted crash frequency and observed crash frequency  
148 are derived from the same source of information. Despite these limitations, the EB expected crash frequency is a good  
149 approximation for the expected values derived from the FB method as it still accounts for the RTM, can refine the  
150 predicted mean of an entity (Zou et al., 2013), and yields similar estimates as FB estimates with comparable precision  
151 but less computational cost. All these confirm that the EB approach is a proper, yet less expensive alternative compared  
152 to the FB method.

153 The EB method has been used for the NB model where the expected crash frequency is defined as a linear combi-  
154 nation of the predicted crash frequency (derived from the SPF model) and the observed crash frequency (Hauer et al.,  
155 2002). Similarly, the EB procedure proposed in the highway safety manual (HSM) is only applicable when the tra-  
156 ditional NB model is being used. As mentioned earlier, different extensions of the NB model have been introduced  
157 to deal with problematic characteristics of crash data. As these extensions get more complex and go deeper in the  
158 hierarchy, the EB procedure becomes harder to implement as the closed form of such distributions are unavailable or  
159 hard to compute analytically. Consequently, there is a clear need to examine the application of the EB method when  
160 more advanced models are being used. To this end, some studies attempted to translate the EB framework for more  
161 complex models such as Sichel (Zou et al., 2013) or finite mixture NB (Zou et al., 2018).

162 In terms of ranking procedures, generally, there are two types of ranking approaches, naive ranking and model-  
163 based ranking (Huang et al., 2009). The naive ranking method uses the raw crash data to make an ordered list of  
164 sites for hot spot identifications. The model-based ranking approaches, however, take the expected crash values as an  
165 indication of the crash risk (Huang et al., 2009). The expected crash values are extensively used to sort a list of roadway  
166 entities. However, the obtained sorted list could potentially differ among the ranking criteria since they are based on  
167 different measures and logics. Several model-based ranking criteria for hot spot identification have been investigated  
168 in the literature to better represent the stochastic nature of the crash data. Many studies have recommended ranking  
169 sites based on the FB or EB expected value of the Poisson mean, and they concluded that the use of posterior Poisson  
170 mean would result in a more reliable and more accurate order compared to the naive ranking criteria (Guo et al., 2019;  
171 Lee et al., 2019; Meng et al., 2020; Lan and Persaud, 2011). In the same token, Shen and Louis (1998) discussed that  
172 the posterior Poisson mean is an optimal choice when inferences about the expected crashes are of interest. However,

173 it might perform poorly if the rank of the expected crashes is the subject of interest. Consequently, some studies have  
 174 attempted to directly take uncertainties in rankings into considerations and employed a Bayesian framework in ranking  
 175 criteria as well (Laird and Louis, 1989; Miaou and Song, 2005; Liu and Sharma, 2018; Shen and Louis, 1998). The  
 176 posterior ranking criteria (*e.g.*, posterior expected, mode, or median rank) takes all posterior simulations into account  
 177 for each site's crash risk (not only the posterior mean), and then outputs a ranked list of sites for each simulation run,  
 178 accordingly. In a comparison study between the EB and FB approach for hot spot identification, Lan and Persaud  
 179 (2011) explored eight different ranking criteria including posterior expected, posterior mode, and posterior median  
 180 ranking. The authors concluded that, in general, the posterior rank criteria would perform better than other model-  
 181 based and naive ranking methods. Similar results were observed in the study done by Laird and Louis (1989). They  
 182 also concluded that the posterior distribution of a parameter's rank typically carries more information about the true  
 183 ranking in comparison to the integer rank of that parameter.

### 184 3. Methodology

185 The NB-L distribution is a mixture of the NB and the one-parameter Lindley distribution. The NB-L distribution  
 186 offers a more flexible structure with more degrees of freedom compared to the traditional NB distribution. Different  
 187 hierarchical variations of the NB-L distribution have been introduced and analyzed (Geedipally et al., 2012; Zamani  
 188 and Ismail, 2010; Gomez-Deniz and Calderin-Ojeda, 2017). This study used the original representation developed by  
 189 Zamani and Ismail (2010) and the generalized linear model proposed by Geedipally et al. (2012) to derive FB and EB  
 190 procedures. Let  $Y_i$  denote the crash frequency following an NB distribution with shape parameter,  $p_i$ , and rate (or over  
 191 dispersion) parameter,  $\phi$ . The hierarchical expression for the NB-L generalized linear model (NB-L GLM) is defined  
 192 as follows (Geedipally et al., 2012):

$$\begin{aligned}
 Y_i &\sim NB(p_i, \phi); \quad \phi > 0, \quad 0 < p_i < 1 \\
 p_i &= e^{-\eta_i} \\
 \eta_i &\sim Lindley(\theta_i) \\
 \theta_i &= \mu_i = e^{\beta X_i} \\
 \phi &\sim \pi_\phi \\
 \beta &\sim \pi_\beta
 \end{aligned} \tag{1}$$

193 where,  $\theta$  is the Lindley parameter,  $\mathbf{X}_i = (1, X_1, X_2, \dots, X_m)$  is the vector including the contributing variables for site  $i$ ,  
 194  $\beta = (\beta_0, \beta_1, \dots, \beta_m)$  is the vector of regression coefficient to be estimated, and  $\pi_\phi$  and  $\pi_\beta$  are the prior distribution for

195  $\phi$  and  $\beta$ , respectively.

196 The GLM representation in Eq.(1) was suggested by Geedipally et al. (2012) as an alternative to NB-L GLM, but it  
197 has not yet been used for modeling due to its complexity. Note that the above NB-L GLM representation is available in  
198 closed-form. Therefore, this representation works well for addressing the objectives of this paper since the EB analysis  
199 usually requires the maximum likelihood estimates (MLE) of the parameters, which require access to the closed form  
200 formulation of the probability mass function (pmf). In the next section, the derivation of the expected crash values is  
201 discussed in detail for both FB and EB methods.

### 202 3.1. Full Bayesian expected values

203 In the FB paradigm, both parameters and hyper-parameters are assumed to follow pre-defined distributions (*i.e.*,  
204 prior distribution) which represent the underlying uncertainty in the parameters. The FB method treats all the pa-  
205 rameters as unknown random variables and takes all their uncertainties into account by integrating over the prior  
206 distributions. In the Bayesian context, either for EB or FB methods, the posterior predictive distribution is used to  
207 estimate the expected crash values. The posterior predictive distribution represents the distribution of the expected  
208 data given the observed data and predictive model. The posterior predictive distribution for an expected data point,  
209  $y_{exp}$ , given the observed value,  $y_{obs}$ , could be written as follows:

$$p(y_{exp}|y_{obs}) = \int_{\gamma} p(y_{exp}|\gamma, y_{obs}) p(\gamma|y_{obs}) d\gamma \quad (2)$$

210 where,  $p(y_{exp}|\gamma, y_{obs})$  is the likelihood of the expected data given the observed data and model parameters ( $\gamma$ ), and  
211  $p(\gamma|y_{obs})$  is the posterior distribution of the parameters give the observed data. In this section, first we document  
212 the derivation of the posterior predictive distribution and the FB expected values for the NB model; then, the same  
213 procedure is extended to derive the FB expected crash frequencies for the NB-L model.

214 The NB distribution itself could be re-parameterized as a continuous mixture of the Poisson and Gamma distribu-  
215 tions, where the Poisson mean follows a Gamma distribution. The hierarchical representation of the NB GLM with  
216 mean,  $\mu$ , and over-dispersion parameter,  $\phi$ , could be written as follows:

$$\begin{aligned} Y_i &\sim \text{Poisson}(\lambda_i) \\ \lambda_i &\sim \text{Gamma}(\phi, \phi/\mu_i) \\ \mu_i &= e^{\beta X_i} \\ \phi &\sim \pi_{\phi} \\ \beta &\sim \pi_{\beta} \end{aligned} \quad (3)$$

217 Using the definition of the posterior predictive distribution in Eq.(2), the probability of the expected crashes,  $y_{exp}$ ,  
 218 given the observed crash data is as follows:

$$p(y_{exp}|y_{obs}) = \int_{\lambda} p(y_{exp}|\lambda, y_{obs}) p(\lambda|y_{obs}) d\lambda \quad (4)$$

219 where,  $p(y_{exp}|\lambda) \sim Poisson(\lambda)$  and  $p(\lambda|y_{obs})$  is the posterior distribution of the Poisson mean,  $\lambda$ , which could be  
 220 written as follows by definition:

$$p(\lambda|y_{obs}) = \int_{\phi, \beta} p(\lambda|y_{obs}, \phi, \beta) \pi_{\phi, \beta} d(\phi, \beta) \quad (5)$$

221 Given Eq.(4) and Eq.(5), the Full Bayesian posterior predictive distribution of the NB GLM could be written as  
 222 follows:

$$p(y_{exp}|y_{obs}) = \int_{\lambda} p(y_{exp}|\lambda, y_{obs}) \left( \int_{\phi, \beta} p(\lambda|y_{obs}, \phi, \beta) \pi_{\phi, \beta} d(\phi, \beta) \right) d\lambda \quad (6)$$

223 Let  $\mathcal{P}$  and  $\mathcal{G}$  denote the Poisson and Gamma distributions, respectively. Given the definition of the posterior  
 224 distribution, we know that:

$$p(\lambda|y_{obs}, \phi, \beta) \propto \mathcal{P}(y_{obs}|\lambda) \mathcal{G}(\lambda|\phi, \beta) \quad (7)$$

225 Therefore, given the fact that the Gamma distribution is a conjugate prior for the Poisson distribution, Eq.(6) could  
 226 be further simplified as follows:

227

$$p(y_{exp}|y_{obs}) = \int_{\lambda} \mathcal{P}(\lambda) \left( \int_{\phi, \beta} \mathcal{G}(y_{obs} + \phi, 1 + \phi/\mu) \pi_{\phi, \beta} d(\phi, \beta) \right) d\lambda \quad (8)$$

228 Using the FB approach, instead of solving the integral or calculating the closed form representation, we can take  
 229 the Monte Carlo Markov Chain (MCMC) approach to draw random samples from the posterior predictive distribution.  
 230 The following steps describe the procedure to draw a random sample from the posterior predictive distribution at each  
 231 site  $i$ :

- 232 • Draw a random sample from the prior distributions,  $\pi_{\beta}$  and  $\pi_{\phi}$ ; then, calculate  $\mu_i$ ;
- 233 • Plug in the samples from the previous step in  $\mathcal{G}(y + \phi, 1 + \phi/\mu_i)$ , and then draw a random sample from the  
 234 distribution. It gives us a random sample from the  $\lambda$ 's posterior distribution;



235 • Plug in the posterior  $\lambda$  sample from the previous step in  $\mathcal{P}(\lambda)$ , and then draw a random sample from the dis-  
 236 tribution. It gives us a random sample from the posterior predictive distribution of the crash frequency at site  
 237  $i$ .

238 By repeating the hierarchical procedure described above, we can have the necessary samples from the posterior  
 239 predictive distribution to estimate the expected crash frequency. For this purpose, we can use any measure of centrality  
 240 (*i.e.*, mode, mean, median) to average out the predictive distribution and achieve the expected value. Note that in the  
 241 last step,  $(y_{exp}|\lambda, y_{obs})$  follows a Poisson distribution. The parameter of the Poisson distribution shows its mean value.  
 242 Therefore, by drawing random samples and then taking the average, we can find the conditional expectation of  $\lambda$  given  
 243 the observed data,  $E(\lambda|y_{obs})$ . This means that if we parameterize the crash frequency as a Poisson mixture model with  
 244 parameter  $\lambda$ , the posterior predictive distribution would be the same as the posterior distribution of  $\lambda$ . This concept  
 245 will be useful when developing EB estimates for the NB-L GLM.

246 Even though the NB-L has been discussed and documented in the literature, no study has yet outlined the derivation  
 247 of expected crash values for the NB-L GLM. A similar procedure as that used in developing FB for NB GLM is also  
 248 applicable in the case of NB-L model. Using the NB-L GLM formulation written in Eq.(1), the posterior predictive  
 249 distribution could be expressed as follows:

$$p(y_{exp}|y_{obs}) = \int_{\eta, \phi} p(y_{exp}|\eta, \phi, y_{obs}) p(\eta, \phi|y_{obs}) d(\eta, \phi) \quad (9)$$

250 Putting the full posterior expression of  $\eta$  parameter in Eq.(9), the above formulation could be re-written as follows  
 251 ( $\mathcal{NB}$  denotes the NB distribution):

$$p(y_{exp}|y_{obs}) = \int_{\eta, \phi} \mathcal{NB}(e^{-\eta}, \phi) \left( \int_{\beta} p(\eta|\beta, y_{obs}) \pi_{\beta} d(\beta) \right) \pi_{\phi} d(\eta, \phi) \quad (10)$$

252 The Lindley distribution does not have any conjugate prior; hence, the integral above cannot be further simplified.  
 253 The following procedure should be followed to draw random samples:

- 254 • Draw a random sample from each prior distribution,  $\pi_{\beta}$  and  $\pi_{\phi}$ .
- 255 • Plug in the sample  $\beta$  from the previous step in  $p(\eta|\beta y_i)$ , and draw a random sample from the distribution. It  
 256 gives us a random sample from posterior distribution of  $\eta$ .
- 257 • Plug in the posterior  $\eta$  sample from the previous step and  $\phi$  from the first step in  $\mathcal{NB}(e^{-\eta}, \phi)$ , and draw a random  
 258 sample from the distribution. It gives us a random sample from the posterior predictive of the crash frequency  
 259 at site  $i$ .

260 This procedure could be easily formulated and summarized in statistical software developed for MCMC analysis  
 261 such as WinBUGS (Lunn et al., 2000), or JAGS (Plummer et al., 2016).

### 262 3.2. Empirical Bayesian expected values

263 As mentioned in the previous section, the main motivation behind the EB method is simplifying the computationally  
 264 intensive steps of the FB procedure. Unlike the FB method where all the parameters are random variables specified  
 265 by prior distributions, the EB method does not consider the uncertainty associated with the parameters; instead, the  
 266 point estimate of the parameters, either maximum likelihood (MLE) or method of moment (MOM) estimates, is used  
 267 in the highest levels of hierarchy. In the following, the EB approach for the NB model is reviewed, and then the EB  
 268 approximation for the NB-L model is developed.

269 The three-step procedure explained for the FB analysis is simplified by some approximations to obtain the expected  
 270 crash values in the EB paradigm as follows:

271 Step one - Estimate the parameters of the highest level of hierarchy.

272 To obtain parameter estimates  $\hat{\beta}$  and  $\hat{\phi}$ , closed form representation of the NB GLM is essential. The closed form  
 273 expression of NB GLM is available and could be easily achieved by marginalizing  $\lambda$  variable out:

$$p(y|\phi, \beta) = \int_{\lambda} p(y|\lambda) p(\lambda|\phi, \beta) d\lambda = \frac{\Gamma(y+\phi)}{\Gamma(y+1)\Gamma(\phi)} \left(\frac{\phi}{\mu+\phi}\right)^{\phi} \left(\frac{\mu}{\mu+\phi}\right)^y \quad (11)$$

274 The above expression is the pmf of the NB distribution.  $\hat{\beta}$  and  $\hat{\phi}$  are called the marginal maximum likelihood  
 275 estimates (MMLE) and could be simply calculated through MLE or MOM approaches.

276 Step two - Derive the expected value of posterior Poisson mean.

277 As mentioned before, Gamma is a conjugate prior for Poisson distribution. As a result, the posterior distribution  
 278 of the Poisson mean,  $\lambda$ , given the data as well as its expected value, is available in closed form (Zou et al., 2018):

$$p(\lambda|\beta, \phi, y) = \text{Gamma}(y+\phi, 1+\phi/\mu) \quad (12)$$

$$E(\lambda|\beta, \phi, y) = \frac{y+\phi}{1+\phi/\mu} = \left(\frac{\mu}{\mu+\phi}\right)y + \left(\frac{\phi}{\mu+\phi}\right)\mu \quad (13)$$

279 The above formula (Eq.(13)) is the known EB formula for the expected crash value, which is extensively used in  
 280 safety analysis (Hauer et al., 2002). Finally, by plugging in  $\hat{\beta}$  and  $\hat{\phi}$  from the previous step in Eq.(13), the EB expected  
 281 crash frequency for each site is calculated.

282 The outlined procedure for derivation of the EB estimates for the NB GLM is also applicable in the NB-L GLM.  
 283 Each step is thoroughly discussed in the following:

284 Step one - Estimate the parameters of the highest level of hierarchy.

285 Estimating  $\phi$  and  $\beta$  requires the closed form representation of the NB-L GLM. The hierarchical representation of  
 286 the NB-L model outlined in Eq.(1) can be expressed in closed form by marginalizing  $\eta$  parameter out:

$$p(y|\beta, \phi) = \int_{\eta} p(y|\phi, \eta) p(\eta|\beta) d\eta \quad (14)$$

287 Solving for the above integral would result in the pmf of the NB-L GLM. It follows from the pmf of the NB-L  
 288 distribution which was developed by (Zamani and Ismail, 2010):

$$p(y_i|\phi, \beta) = \frac{e^{2\beta X_i}}{1 + e^{\beta X_i}} \binom{\phi + y_i - 1}{y_i} \sum_{j=0}^{y_i} (-1)^j \binom{y_i}{j} \frac{e^{\beta X_i + \phi + j + 1}}{(e^{\beta X_i + \phi + j})^2} \quad (15)$$

289 The MLEs,  $\hat{\beta}$  and  $\hat{\phi}$ , could then be calculated by maximizing the likelihood (or log-likelihood) function. The  
 290 log-likelihood function of NB-L GLM is given as follows:

$$ll = \sum_{i=1}^n \left[ \log \binom{\phi + y_i - 1}{y_i} + 2(\beta X_i) - \log(1 + e^{\beta X_i}) \right. \\ \left. + \log \left( \sum_{j=0}^{y_i} (-1)^j \binom{y_i}{j} \frac{e^{\beta X_i + \phi + j + 1}}{(e^{\beta X_i + \phi + j})^2} \right) \right] \quad (16)$$

291 The first partial derivative with respect to the unknown parameters could be written as follows:

$$\frac{\partial ll}{\partial \beta} = \sum_{i=1}^n \left( 2X_i - \frac{X_i}{1 + e^{\beta X_i}} \right) + \frac{\sum_{j=0}^{y_i} (-1)^{j+1} \binom{y_i}{j} \frac{X_i (e^{\beta X_i + \phi + j + 2})}{(e^{\beta X_i + \phi + j})^3}}{\sum_{j=0}^{y_i} (-1)^j \binom{y_i}{j} \frac{e^{\beta X_i + \phi + j + 1}}{(e^{\beta X_i + \phi + j})^2}} \quad (17)$$

292

$$\frac{\partial ll}{\partial \phi} = \frac{\partial}{\partial \phi} \left[ \sum_{i=1}^n \log \binom{\phi + y_i - 1}{y_i} \right] + \frac{\sum_{j=0}^{y_i} (-1)^{j+1} \binom{y_i}{j} \frac{e^{\beta X_i + \phi + j + 2}}{(e^{\beta X_i + \phi + j})^3}}{\sum_{j=0}^{y_i} (-1)^j \binom{y_i}{j} \frac{e^{\beta X_i + \phi + j + 1}}{(e^{\beta X_i + \phi + j})^2}} \quad (18)$$

293 The first part in Eq.(18) could be re-written as follows (Klugman et al., 2012; Tajuddin et al., 2020):

$$\frac{\partial}{\partial \phi} \sum_{i=1}^n \log \binom{\phi + y_i - 1}{y_i} = \sum_{i=1}^n \sum_{m=0}^{y_i-1} \frac{1}{\phi + m} \quad (19)$$

294 As a result, the partial derivative of the log-likelihood with respect to  $\phi$  is given as:

$$\frac{\partial ll}{\partial \phi} = \sum_{i=1}^n \left( \sum_{m=0}^{y_i-1} \frac{1}{\phi + m} \right) + \frac{\sum_{j=0}^{y_i} (-1)^{j+1} \binom{y_i}{j} \frac{e^{\beta X_i + \phi + j + 2}}{(e^{\beta X_i + \phi + j})^3}}{\sum_{j=0}^{y_i} (-1)^j \binom{y_i}{j} \frac{e^{\beta X_i + \phi + j + 1}}{(e^{\beta X_i + \phi + j})^2}} \quad (20)$$

295 The above derivative equations could be simultaneously solved using numeric methods (gradient descend, Newton-  
296 raphson, etc.) in order to estimate the unknown parameters.

297 Step two - Derive the expected value of posterior Poisson mean.

298 The hierarchical representation of the NB-L GLM defined in Eq.(1) does not involve the Poisson mean,  $\lambda$ , parameter  
299 since  $\lambda$  has already been marginalized out in the definition of the NB distribution. However, we can formulate the NB-L  
300 GLM as a functional of  $\lambda$  by breaking down the NB distribution to a mixture of Poisson and Gamma distribution:

$$\begin{aligned} Y_i &\sim \text{Poisson}(\lambda) \\ \lambda &\sim \text{Gamma}\left(\phi, \frac{e^{-\eta_i}}{1 - e^{-\eta_i}}\right) \\ \eta_i &\sim \text{Lindley}(\theta_i) \\ \theta_i &= \mu_i = e^{\beta X_i} \end{aligned} \quad (21)$$

301 Following from the above hierarchical representation, the pmf of the NB-L GLM could be re-written as a function  
302 of the Poisson mean,  $\lambda$ :

$$p(y|\phi, \beta) = \int_{\lambda} p(y|\lambda) \left( \int_{\eta} p(\lambda|\phi, \eta) p(\eta|\beta) d(\eta) \right) d\lambda \quad (22)$$

303 Clearly, the above expression is the pmf of a Poisson mixture distribution with mixing distribution as follows:

$$p(\lambda|\phi, \beta) = \int_{\eta} p(\lambda|\phi, \eta) p(\eta|\beta) d(\eta) \quad (23)$$

304 Neither the mixing distribution itself,  $p(\lambda|\phi, \beta)$ , nor its posterior distribution,  $p(\lambda|\phi, \beta, y)$ , can be parameterized  
305 in closed-form. Hence, we can't directly calculate the posterior expectation of the Poisson mean,  $E(\lambda|\phi, \beta, y)$ , in the  
306 same way we did in the case of the NB model. Instead, we can take advantage of a useful property of the Poisson  
307 mixture distributions documented by Karlis and Xekalaki (2005) and Willmot (1986). Suppose Y follows a mixture  
308 Poisson distribution with pmf  $p(x)$ . Then, the posterior moments of any order of the Poisson mean,  $E(\lambda^r|X = x)$ ,

309 could be calculated as follows:

$$E(\lambda^r | X = x) = \frac{p(x+r)}{p(x)}(x+1)\dots(x+r) \quad (24)$$

310 where,  $p(y)$  is the mixed Poisson pmf.

311 We can use this property to calculate the posterior expectation of the Poisson mean when the NB-L model is being  
 312 used. In the same way as before, let  $Y$  follow the NB-L distribution derived in Eq.(15). Then, the posterior expectation  
 313 of the Poisson mean could be written as follows:

$$E(\lambda | y, \phi, \beta) = \frac{p(y+1 | \phi, \beta)}{p(y | \phi, \beta)}(y+1) \quad (25)$$

$$= \frac{\frac{e^{2\beta X}}{1+e^{\beta X}} \binom{\phi+y+1-1}{y+1} \sum_{j=0}^{y+1} (-1)^j \binom{y+1}{j} \frac{e^{\beta X + \phi + j + 1}}{(e^{\beta X + \phi + j})^2}}{\frac{e^{2\beta X}}{1+e^{\beta X}} \binom{\phi+y-1}{y} \sum_{j=0}^y (-1)^j \binom{y}{j} \frac{e^{\beta X + \phi + j + 1}}{(e^{\beta X + \phi + j})^2}}(y+1) \quad (26)$$

314 The formula in Eq.(26) can be further summarized as follows:

$$E(\lambda | y, \phi, \beta) = \frac{A(y+1)}{A(y)}(y+1) \quad (27)$$

315 where,  $A(y) = \sum_{j=0}^y (-1)^j \binom{y}{j} \frac{e^{\beta X + \phi + j + 1}}{(e^{\beta X + \phi + j})^2} = \sum_{j=0}^y (-1)^j \binom{y}{j} \frac{\mu + \phi + j + 1}{(\mu + \phi + j)^2}$

316 Eq.(27) can be used to estimate the posterior mean of  $\lambda$  without knowing any information about the mixing distri-  
 317 bution or to solve for the expectation definition itself. The proposed EB formula is comparable with the famous EB  
 318 formula derived for the NB model indicated in Eq.(13). Finally, we need to incorporate the MLEs,  $\hat{\phi}$  and  $\hat{\beta}$ , in Eq.(27).  
 319 The resulting value is the desired EB expected crash value for the NB-L model. As observed, it does not involve any  
 320 intensive computation (like the FB method) or solving any complex integral. The next sections describe the dataset  
 321 used for empirical evaluation of the proposed EB framework as well as the implementation details.

## 322 4. Data description

323 In the previous section, the derivation of the expected crash value using FB and EB frameworks was discussed for  
 324 the NB-L GLM. In order to examine the developed FB and EB frameworks, this study used two datasets. Both Vir-  
 325 ginia (2014-2019) and Texas (2014-2019) datasets represent the crash statistics of the non-federal aid system (NFAS)  
 326 roadways discussed in Khodadadi et al. (2021) and Das et al. (2021). NFAS roadways are typically characterized by  
 327 lower volumes and lower crash frequency in comparison with other roadway functional classifications (Khodadadi  
 328 et al., 2021). Consequently, many NFAS segments experience zero crashes. Further, there are a lot of missing data in  
 329 many roadway characteristics (e.g., shoulder width) that could potentially be used as predictors. Therefore, a limited

**Table 1**  
Summary Statistics of datasets

Dataset	Variables	Min	Max	Average	Standard Deviation
Texas	Number of crashes	0	15	0.86	1.65
	Average 5-years AADT (vpd)	43	1166	313.8	253
	Segment length (miles)	0.10	4.41	0.96	0.93
Virginia	Number of crashes	0	8	2.01	2.09
	Average 5-years AADT (vpd)	163	5180	694	625
	Segment length (miles)	0.13	5.67	1.35	1.08

330 number of variables were available to use in the generalized linear modeling framework. However, as the emphasis  
 331 of this study is to assess and compare the developed framework, using fewer variables is not an issue. The summary  
 332 statistics of both datasets are provided in Table 1.

## 333 5. Modeling results

334 In this section, the modeling results for both NB and NB-L GLMs are presented, then the proposed EB procedure  
 335 for NB-L GLM is examined. First, the EB and FB estimates of the expected crashes were compared to generally  
 336 validate the proposed EB method and show how well the EB estimates mimic the FB estimates. Then, the EB and FB  
 337 estimates were used in site ranking analysis to determine how similar the produced ranks and identified hot spots were.

### 338 5.1. Crash prediction models

339 The NB and NB-L GLMs were developed for each dataset. For each model, only the AADT and segment length  
 340 were included as contributing covariates. It should be noted that as the models are developed using the same functional  
 341 form and compared using the same dataset, therefore, as noted earlier, the omitted variable bias would not be an  
 342 issue. Also, as discussed before, this study aims to develop EB estimates for the expected crash values and explore  
 343 whether they approximate the FB estimate properly. Including more variables might enhance the predictive models'  
 344 performance; however, it will not affect the underlying theoretical framework of deriving EB estimates.

345 For each GLM, Full Bayesian and maximum likelihood estimates were calculated. We employed the MCMC  
 346 method using an open-sourced R package, called "RJAGS" (Plummer et al., 2003), to estimate the posterior of param-  
 347 eters. This study assumed a non-informative gamma, and a non-informative normal distribution for the prior distribu-  
 348 tion of  $\beta$ 's and  $\phi$  parameters, respectively. In total, three chains and 60,000 iterations were set up to ensure the MCMC  
 349 convergence. The first 4,000 samples of each chain were discarded. Also, to reduce the potential auto-correlation  
 350 among the random draws, every third sample of the rest was used for estimations of unknown parameters.

351 The maximum likelihood estimates are needed in order to develop the EB estimates. Unlike the NB distribution,

**Table 2**  
Modeling results for Virginia dataset

Variables	NB GLM		NB-L GLM	
	FB estimate (s.d.)	MLE	FB estimate (s.d.)	MLE
Intercept ( $\beta_0$ )	-3.48 (0.22)	-3.47	7.86 (0.36)*	8.27*
Log (AADT) ( $\beta_1$ )	0.51 (0.03)	0.51	-0.51(0.04)*	-0.53*
Length ( $\beta_2$ )	0.57 (0.02)	0.58	-0.60 (0.03)*	-0.59*
$\phi$	3.94 (0.42)	3.87	75.07 (14.74)	96.63
DIC		8301		7135
WAIC		6859		6630
MAD		0.89		0.86
Log-likelihood		-3119		-2888

\* The estimates for the NB-L GLM do not carry information regarding the causal association of covariates and crashes (see discussion below)

352 the NB-L is not a part of natural exponential family distributions. Hence, its log-likelihood function is not strictly  
353 concave. Numerical approaches equipped with proper initial values are needed to approach the global maximum  
354 point. Due to the significant sensitivity observed among the NB-L partial derivative equations and the initial values,  
355 a meta-heuristic genetic algorithm, together with a gradient descent approach, was utilized to ensure convergence to  
356 the global maximum point. For this purpose, the "GA" package in R (Scrucca et al., 2013) was used to solve the  
357 optimization problem. A total of 200 iterations with 500 initial populations were considered to maximize the objective  
358 function. A gradient descant approach was also employed in each iteration of the genetic algorithm to locally search  
359 for better estimates and further enhance the maximization process.

360 The FB estimates and MLEs for the Texas and Virginia datasets for both Nb and NB-L GLMs are summarized in  
361 Table 2 and Table 3, respectively. Three performance measures, namely Deviance Information Criteria (DIC), Mean  
362 Absolute Deviance (MAD), and Widely Applicable Information Criteria (WAIC) were used for model comparisons  
363 (Lord et al., 2021). WAIC was developed by Vehtari et al. (2017) and it appeared to be a robust alternative for DIC in  
364 the Bayesian framework (Watanabe and Oppen, 2010; Khodadadi et al., 2021). All the performance measures showed  
365 that the NB-L models fit the data better than the NB models. These results were expected since both datasets were  
366 characterized by a large number of zeros and high skewness (the domain under which the NB-L model performs better  
367 than the NB).

368 As indicated in Table 2 and Table 3, the signs and magnitudes of estimates are different across models. This issue  
369 could be accredited to the particular representation of the NB-L model used in this study. As opposed to the NB model  
370 where the mean function has a log-linear association with covariates, the mean of the NB-L model has a non-linear  
371 relationship with covariates. Consequently, the magnitude and sign of the estimates do not necessarily represent the  
372 causal relationship between covariates and crash frequencies. This issue and its associated limitations are covered in  
373 the Discussion section further below.

**Table 3**  
Modeling results for Texas dataset

Variables	NB GLM		NB-L GLM	
	FB estimate (s.d.)	MLE	FB estimate (s.d.)	MLE
Intercept ( $\beta_0$ )	-5.49 (0.12)	-5.49	8.00 (0.28)*	7.21*
Log (AADT) ( $\beta_1$ )	0.80 (0.02)	0.79	-0.73 (0.02)*	-0.68*
Length ( $\beta_2$ )	0.66 (0.02)	0.66	-0.60 (0.02)*	-0.61*
$\phi$	1.07 (0.04)	1.07	17.83 (3.16)	10.71
DIC		26155		22471
WAIC		21031		20761
MAD		0.45		0.44
Log-likelihood		-9023		-8847

\* The estimates for the NB-L GLM do not carry information regarding the causal association of covariates and crashes (see discussion below)

## 374 5.2. Expected crash values

375 Using the above modeling results, we attempted to calculate both FB and EB expected crash values. Estimating  
 376 the FB expected values involves multiple random sampling steps, which are doable using any software developed for  
 377 MCMC analysis. However, the MCMC analysis requires the model tree defined using the standard distributions. The  
 378 Lindley distribution is not a standard distribution but can be re-parameterized as a two-component gamma mixture  
 379 (Zamani and Ismail, 2010):

$$\epsilon \sim Lindley(\theta) \equiv \frac{1}{1+\theta} Gamma(2, \theta) + \frac{\theta}{\theta+1} Gamma(1, \theta) \quad (28)$$

380 To calculate the EB expected crashes, we merely plugged in the MLEs (*i.e.*,  $\hat{\phi}$ ,  $\hat{\beta}$ ) from Table 2 and Table 3 in the  
 381 EB formula developed in Eq.(26). The EB expected crash value for site  $i$  would be as follows:

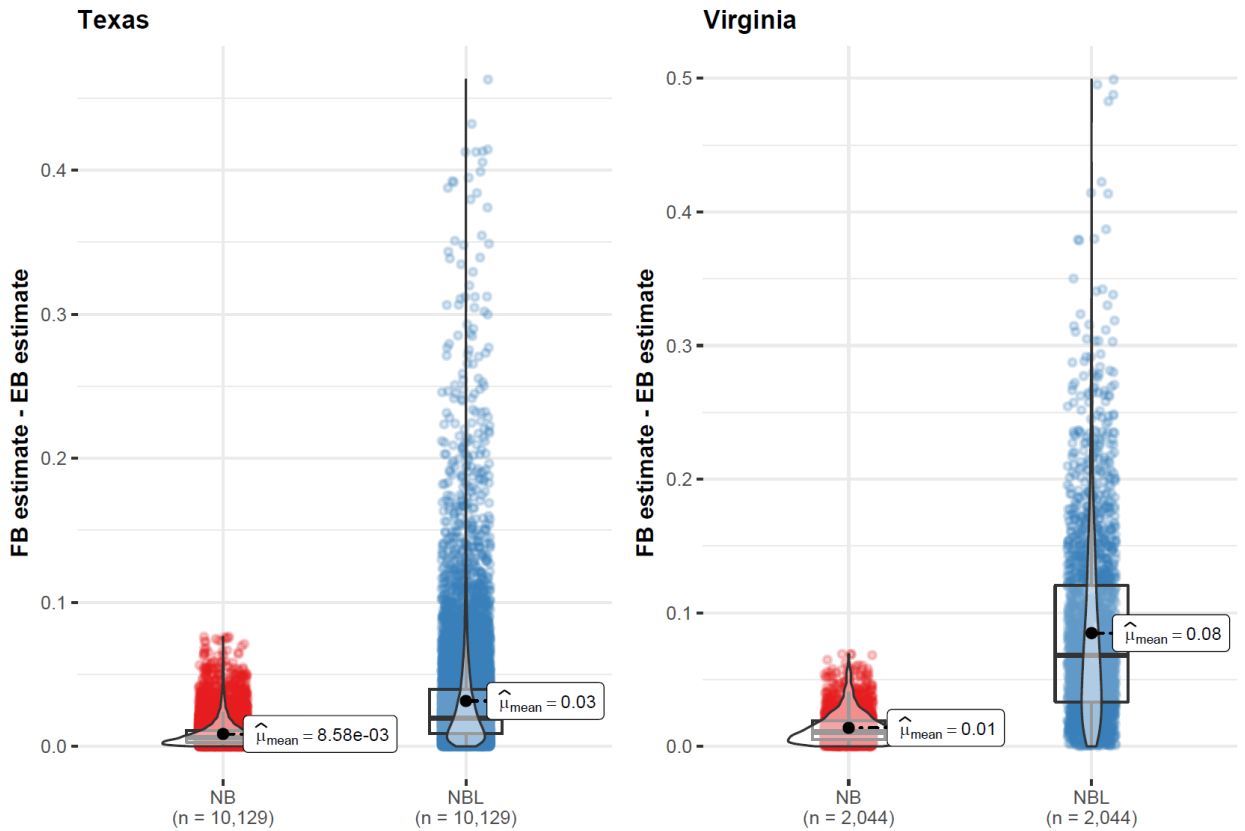
$$E(\lambda_i | y_i) = \frac{A(y_i + 1)}{A(y_i)} (y_i + 1) \quad (29)$$

382 where,  $A(y_i) = \sum_{j=0}^{y_i} (-1)^j \binom{y_i}{j} \frac{e^{\hat{\beta}X + \hat{\phi} + j + 1}}{(e^{\hat{\beta}X_i + \hat{\phi} + j})^2}$

383 Results from both methods were observed and compared to see whether EB estimates properly approximate the FB  
 384 estimates. The absolute difference between the estimates for the NB and NB-L models are plotted in Figure 1. These  
 385 violin plots show the extent to which the FB expected values and their approximated EB counterparts are different. As  
 386 seen, the difference between the mean of the expected values are quite small for both models; this indicates that the  
 387 proposed EB formula for the NB-L GLM can accurately approximate the FB procedure, thus avoiding the demanding  
 388 MCMC analysis.

389 However, the difference between the EB and FB expected values were larger for the NB-L in comparison to those





**Figure 1:** The absolute difference between the EB and FB estimates of the expected crashes

390 of the NB model. This issue is fully covered in the discussion section. Even though the difference between the EB and  
 391 FB estimates seem relatively larger for the NB-L models, the authors observed that the relative differences will not  
 392 exceed 25% for sites with crash experiences. This indicates that the EB and FB estimates are close, and it is anticipated  
 393 that the EB estimates will be adequate for safety applications.

### 394 5.3. Application in site ranking

395 The empirical results in the previous section showed that the EB estimates of the expected crash values are appro-  
 396 priate alternatives to the FB estimates. Both the absolute and relative differences between the EB and FB estimates  
 397 of expected crashes were small, indicating that the EB estimates well approximate their FB counterparts. However,  
 398 sometimes the order of the expected crash values is of interest rather than the magnitude itself. In site ranking studies,  
 399 the aim is to sort the study sites in decreasing order of their evaluated crash risk (expected crash value). As a result,  
 400 this study examined the application of the proposed EB framework for the NB-L GLM in site ranking.

401 In order to assess how a ranking criterion performs, a reference ranking is needed as the basis for the comparison.  
 402 Several studies found Posterior ranking criteria to be a better alternative to integer ranking when the model is imple-

**Table 4**  
Site ranking results

Dataset	Risk	Unordered Ranked Groups					
	Evaluation Criteria	1-10	1-20	1-50	1-100	1-200	1-500
Texas	NB-L FB	10/10 *	20/20	50/50	100/100	200/200	500/500
	NB-L EB	9/10	20/20	46/50	97/100	190/200	486/500
Virginia	NB-L FB	10/10	19/20	49/50	99/100	196/200	498/500
	NB-L EB	8/10	15/20	49/50	95/100	196/200	492/500

\*Posterior expected ranking is the reference ranking criteria

403 mented as a Bayesian framework (Laird and Louis, 1989; Miaou and Song, 2005; Liu and Sharma, 2018; Shen and  
 404 Louis, 1998). Laird and Louis (1989) observed that the posterior distribution of ranks carries more information than  
 405 the integer rank that is usually assigned to the parameter mean. The Posterior ranking criterion takes into account all  
 406 posterior simulations for each site's crash risk (not only the posterior mean), and results in a ranking for each simulation  
 407 run. Eventually, by taking the average of the simulated ranks for each site, the posterior expected rank is achieved. This  
 408 study assumed the posterior expected ranks as the reference ranks in order to compare the ranking criteria produced  
 409 by FB and EB expected crash values for the NB-L model. The first 10, 20, 50, 100, 200, and 500 top ranked sites  
 410 were identified for each ranking criterion. Table 4 shows the number of sites that appeared in both the ranking criteria  
 411 being evaluated and the reference ranking (posterior expected ranks). As seen in both datasets, the hot spots identified  
 412 by the EB approach are quite similar to those identified by the FB approach. Similarity of the ranks indicate that the  
 413 proposed EB approach can be a proper alternative to the FB approach not only in crash prediction, but also in hotspot  
 414 identification. Site ranking is the only safety analysis evaluated in this study; however, any other type of safety studies  
 415 that use the long-term crash mean can benefit from the proposed EB method.

## 416 6. Discussion

417 Modeling results from previous works and the current study indicated that compared to the traditional NB model,  
 418 the NB-L model provides a superior fit when analyzing crash datasets with excess zero observations. Consequently,  
 419 the NB-L has a better performance in evaluating the crash risk associated with each site. However, the NB-L expected  
 420 values were only available using the FB approach, which could be difficult to compute for large datasets. This study  
 421 introduced the EB framework to approximate the FB estimates of the NB-L model. Results from the previous section  
 422 showed that the proposed EB framework for the NB-L can properly approximate the underlying FB procedure, so it  
 423 can be used for analyses that require expected crash values (*e.g.*, site ranking, before-after analysis). Some interesting

424 findings, results, and limitations are discussed below.

425 As indicated in Table 2 and Table 3, neither the signs nor the magnitudes of the estimated coefficients are com-  
426 parable between the NB and NB-L models. This issue can be attributed to the way the mean function is structured.  
427 The NB model in this study is parameterized by its mean,  $\mu$ , and overdispersion parameter,  $\phi$ . The mean function is  
428 assumed to have a log-linear relationship with the site characteristics, (*i.e.*,  $\mathbf{X}$ ) through the regression coefficients (*i.e.*,  
429  $\beta$ ). As a result, the coefficients are directly related to the mean crashes so, their magnitude and sign carry information  
430 about how and to what extent each covariate affects the crash frequencies. However, the NB-L model is parameterized  
431 differently. Geedipally et al. (2012) introduced two different parameterizations for the NB-L GLM. The first one uses  
432 the NB formulated by mean and overdispersion parameter where each site-specific mean value is multiplied by an  
433 adjustment factor (or as indicated in the original paper, frailty term),  $\epsilon$ :

$$\begin{aligned} Y &\sim NB(y; \epsilon\mu, \phi) \\ \epsilon &\sim Lindley(\theta) \\ \mu &= e^{\beta X} \end{aligned} \tag{30}$$

434 This parameterization is easy to interpret, and given that  $E(Y) = \mu$ , the regression coefficients are directly related to  
435 the mean response. This representation, however, is not available in closed form and hence, cannot be used in the EB  
436 framework proposed in this study. Instead, we used the second parameterization, which is available in closed form:

$$\begin{aligned} Y &\sim NB(y; p, \phi) \\ -\ln(p) &\sim Lindley(\theta) \\ \theta &= e^{\beta X} \end{aligned} \tag{31}$$

437 This parameterization, which follows the original NB-L parameterization discussed in Zamani and Ismail (2010),  
438 links the site characteristics to the Lindley parameter,  $\theta$ , not the mean. The mean of this parameterization can be  
439 written as follows (Zamani and Ismail, 2010):

$$E(Y) = \phi \left( \frac{\theta^3}{(\theta + 1)(\theta - 1)^2 - 1} \right) \tag{32}$$

440 As seen in the above definition, the mean response is a non-linear and non-invertible function of regression co-  
441 efficients and overdispersion parameter. Consequently, the signs and magnitudes of the regression coefficients in the  
442 second parameterization do not necessarily show the causal relationship between the covariates and crash frequency.

443 To put it concisely, the proposed EB framework and the NB-L representation we utilized in this study (indicated in  
444 Eq.31) are only applicable when the expected crash values are of interest. The expected crashes can then be used in  
445 various safety-related studies such as hot spot identification or before-after analysis. However, if the goal is to deter-  
446 mine the underlying relationship between the crash frequencies and contributing factors, the other parameterization  
447 of the NB-L model indicated in Eq.(30) should be employed (see Khodadadi et al. (2021); Geedipally et al. (2012);  
448 Shirazi and Lord (2019); Shirazi et al. (2016)).

449 In addition, we observed that the log-likelihood function of the NB-L model behaves unpredictably when large-  
450 valued parameters or large inputs are involved. This issue can be attributed to the summation term existing in the  
451 NB-L closed-form expression,  $\sum_{j=0}^y \binom{y}{j} (-1)^j \frac{e^{\beta X + \phi + j + 1}}{(e^{\beta X + \phi + j})^2}$ . This summation part results from the following substitution,  
452  $(1 - e^{-\lambda})^y = \sum_{j=0}^y \binom{y}{j} (-1)^j e^{-\lambda j}$ , which was used in Zamani and Ismail (2010), and Khodadadi et al. (2022) to derive  
453 the pmf of the NB-L distribution. This part outputs small negative values when large  $y$ 's or large-valued parameters  
454 are input. Consequently, proper initial values are required when maximizing the likelihood function to avoid large  
455 estimates and negative likelihoods, and ensure valid estimates and inferences.

456 Furthermore, a larger difference between the FB and EB expected values was observed in the NB-L compared  
457 to the NB model. This issue can be attributed to two reasons. First, the NB-L likelihood is not strictly concave and  
458 hence, the global optimization is not possible or very difficult to get. Numerical approaches are needed to approach  
459 to the global maximum point as much as possible which eventually lead to a range of local maxima and a range of  
460 estimates. Unlike in the case of the NB model where a single set of MLEs achieve the global maximum, in the NB-L  
461 model, a range of local estimates would be achieved whose accuracy depend of the initial values. Therefore, the MLEs  
462 and the FB estimates are quite different for the NB-L model (see Table2). This difference between the MLEs and FB  
463 estimate will result in the different EB and FB estimates of the expected values. Another potential reason could be  
464 the bias-variance trade-off. Due to the flexible structure of the NB-L model, it tends to output less-biased and hence,  
465 high-variance results. High variability of the NB-L model is mirrored in high variance of the expected values.

466 Aside from the high-variability of the NB-L model, the absolute and relative differences showed negligible values.  
467 Similarity of the ranking from the EB and FB estimates also confirmed that the EB estimates will be adequate for safety  
468 applications. The proposed framework will be specially useful in situations where the traditional NB models are not  
469 flexible enough (*e.g.*, abundance of zeros in the data or high skewness), or output biased results (*e.g.*, data with low  
470 sample mean).

## 471 7. Summary and conclusions

472 Even though there is rich literature on the advanced predictive models in traffic safety, little has been done to extend  
473 their application to other roadway safety tasks. The NB-L model has been proposed for sparse count data modeling

474 and, as indicated in several studies, provides superior performance compared to the common NB model used in traffic  
475 safety. However, its application has not been examined in other roadway safety tasks.

476 Expected crash values estimated from the crash prediction models are the main evaluation tool in safety analysis and  
477 represent the long-term risk associated with a roadway entity. This study proposed an EB framework to approximate  
478 the underlying FB method and derive the expected crash values for the NB-L model. The derived expected crashes  
479 can be used in various safety-related studies (e.g., hot spot identification, before-after analysis). The results showed  
480 that the proposed EB framework is able to estimate expected crashes with comparable precision to the FB estimate but  
481 with much lower computational costs. The proposed framework was further examined in site ranking analysis. We  
482 observed that ranks produced by the EB estimates were similar to those of FB estimates, indicating that the proposed  
483 framework can be safely employed in other highway safety tasks such as hot spot identification analysis.

484 The EB approach introduced in this study can be utilized in any type of analysis that requires access to the expected  
485 crash values. The resulting EB expected crashes take advantage of the probabilistic structure of the FB paradigm while  
486 avoiding its time-consuming computational efforts. For future studies, a similar framework as introduced in this study  
487 can be used to develop an EB method for other advanced predictive models in traffic safety. Also, further work should  
488 be performed to validate the application of the framework in other safety-related analyses such as before-after analysis.

## 489 **Author Contribution Statement**

490 The authors confirm contribution to the paper as follows: study conception and design: Ali Khodadadi and Do-  
491 minique Lord; data collection: Ali Khodadadi, Ioannis Tsapakis, and Subasish Das; analysis and interpretation of  
492 results: Ali Khodadadi, Mohammadali Shirazi; draft manuscript preparation: Ali Khodadadi, Mohammadali Shirazi,  
493 Ioannis Tsapakis, Subasish Das, and Dominique Lord. All authors reviewed the results and approved the final version  
494 of the manuscript.

## 495 **Declaration of Competing Interest**

496 The authors declare that they have no known competing financial interests or personal relationships that could have  
497 influenced the work reported in this paper.

## 498 **Funding Sources**

499 The work was completed as part of SAFE-D University Transportation Center project Use of Disruptive Technolo-  
500 gies to Support Safety Analysis and Meet New Federal Requirements. The study was partly funded by the A.P. and  
501 Florence Wiley Faculty Fellow.

## 502 **References**

- 503 Abdella, G.M., Kim, J., Al-Khalifa, K.N., Hamouda, A.M., 2019. Penalized conway-maxwell-poisson regression for modelling dispersed discrete  
504 data: The case study of motor vehicle crash frequency. *Safety Science* 120, 157–163.
- 505 Aguero-Valverde, J., Jovanis, P.P., 2009. Bayesian multivariate poisson lognormal models for crash severity modeling and site ranking. *Transporta-  
506 tion research record* 2136, 82–91.
- 507 Aman, J.J., Smith-Colin, J., Zhang, W., 2021. Listen to e-scooter riders: Mining rider satisfaction factors from app store reviews. *Transportation  
508 Research Part D: Transport and Environment* 95, 102856.
- 509 Aman, J.J.C., Smith-Colin, J., 2020. Transit deserts: Equity analysis of public transit accessibility. *Journal of Transport Geography* 89, 102869.
- 510 Asgharpour, S., Javadinasr, M., Bayati, Z., et al., 2021. Investigating severity of motorcycle-involved crashes in a developing country. *arXiv preprint  
511 arXiv:2110.00381* .
- 512 Darzian Rostami, A., Katthe, A., Sohrabi, A., Jahangiri, A., 2020. Predicting critical bicycle-vehicle conflicts at signalized intersections. *Journal  
513 of advanced transportation* 2020.
- 514 Das, S., Bibeka, A., Sun, X., Zhou, H., Jalayer, M., 2019. Elderly pedestrian fatal crash-related contributing factors: applying empirical bayes  
515 geometric mean method. *Transportation research record* 2673, 254–263.
- 516 Das, S., Tsapakis, I., Khodadadi, A., 2021. Safety performance functions for low-volume rural minor collector two-lane roadways. *IATSS Research  
517* .
- 518 Debrabant, B., Halekoh, U., Bonat, W.H., Hansen, D.L., Hjelmborg, J., Lauritsen, J., 2018. Identifying traffic accident black spots with poisson-  
519 tweedie models. *Accident Analysis & Prevention* 111, 147–154.
- 520 Denthet, S., Thongteeraparp, A., Bodhisuwan, W., 2016. Mixed distribution of negative binomial and two-parameter lindley distributions, in: 2016  
521 12th International Conference on Mathematics, Statistics, and Their Applications (ICMSA), IEEE. pp. 104–107.
- 522 Farid, A., Abdel-Aty, M., Lee, J., Eluru, N., 2017. Application of bayesian informative priors to enhance the transferability of safety performance  
523 functions. *Journal of safety research* 62, 155–161.
- 524 Geedipally, S.R., Lord, D., Dhavala, S.S., 2012. The negative binomial-lindley generalized linear model: Characteristics and application using  
525 crash data. *Accident Analysis & Prevention* 45, 258–265.
- 526 Gomez-Deniz, E., Calderin-Ojeda, E., 2017. An alternative representation of the negative binomial-lindley distribution. new results and applications.  
527 *arXiv preprint arXiv:1703.04812* .
- 528 Guo, X., Wu, L., Lord, D., 2020. Generalized criteria for evaluating hotspot identification methods. *Accident Analysis & Prevention* 145, 105684.
- 529 Guo, X., Wu, L., Zou, Y., Fawcett, L., 2019. Comparative analysis of empirical bayes and bayesian hierarchical models in hotspot identification.  
530 *Transportation research record* 2673, 111–121.
- 531 Hauer, E., 1997. Observational before/after studies in road safety. estimating the effect of highway and traffic engineering measures on road safety.
- 532 Hauer, E., Harwood, D.W., Council, F.M., Griffith, M.S., 2002. Estimating safety by the empirical bayes method: a tutorial. *Transportation Research  
533 Record* 1784, 126–131.
- 534 Huang, H., Chin, H.C., Haque, M.M., 2009. Empirical evaluation of alternative approaches in identifying crash hot spots: Naive ranking, empirical  
535 bayes, full bayes methods. *Transportation Research Record* 2103, 32–41.
- 536 Karlis, D., Xekalaki, E., 2005. Mixed poisson distributions. *International Statistical Review/Revue Internationale de Statistique* , 35–58.
- 537 Khattak, Z.H., Magalotti, M.J., Fontaine, M.D., 2018. Estimating safety effects of adaptive signal control technology using the empirical bayes  
538 method. *Journal of safety research* 64, 121–128.
- 539 Khazraee, S.H., Johnson, V., Lord, D., 2018. Bayesian poisson hierarchical models for crash data analysis: Investigating the impact of model choice

540 on site-specific predictions. *Accident Analysis & Prevention* 117, 181–195.

541 Khodadadi, A., Shirazi, M., Geedipaly, S., Lord, D., 2022. Evaluating alternative variations of negative binomial-lindley distribution for modeling  
542 crash data. *Transportmetrica A: Transport Science* .

543 Khodadadi, A., Tsapakis, I., Das, S., Lord, D., Li, Y., 2021. Application of different negative binomial parameterizations to develop safety perfor-  
544 mance functions for non-federal aid system roads. *Accident Analysis & Prevention* 156, 106103.

545 Klugman, S.A., Panjer, H.H., Willmot, G.E., 2012. *Loss models: from data to decisions*. volume 715. John Wiley & Sons.

546 Laird, N.M., Louis, T.A., 1989. Empirical bayes ranking methods. *Journal of Educational Statistics* 14, 29–46.

547 Lan, B., Persaud, B., 2011. Fully bayesian approach to investigate and evaluate ranking criteria for black spot identification. *Transportation research*  
548 *record* 2237, 117–125.

549 Lee, A.S., Lin, W.H., Gill, G.S., Cheng, W., 2019. An enhanced empirical bayesian method for identifying road hotspots and predicting number of  
550 crashes. *Journal of Transportation Safety & Security* 11, 562–578.

551 Liu, C., Sharma, A., 2018. Using the multivariate spatio-temporal bayesian model to analyze traffic crashes by severity. *Analytic methods in accident*  
552 *research* 17, 14–31.

553 Lord, D., 2006. Modeling motor vehicle crashes using poisson-gamma models: Examining the effects of low sample mean values and small sample  
554 size on the estimation of the fixed dispersion parameter. *Accident Analysis & Prevention* 38, 751–766.

555 Lord, D., Geedipally, S.R., 2011. The negative binomial–lindley distribution as a tool for analyzing crash data characterized by a large amount of  
556 zeros. *Accident Analysis & Prevention* 43, 1738–1742.

557 Lord, D., Geedipally, S.R., Guikema, S.D., 2010. Extension of the application of conway-maxwell-poisson models: analyzing traffic crash data  
558 exhibiting underdispersion. *Risk Analysis: An International Journal* 30, 1268–1276.

559 Lord, D., Mannering, F., 2010. The statistical analysis of crash-frequency data: a review and assessment of methodological alternatives. *Trans-*  
560 *portation research part A: policy and practice* 44, 291–305.

561 Lord, D., Park, B.J., Model, P.G., 2012. Negative binomial regression models and estimation methods. *Probability Density and Likelihood Functions*.  
562 Texas A&M University, Korea Transport Institute , 1–15.

563 Lord, D., Qin, X., Geedipally, S.R., 2021. *Highway Safety Analytics and Modeling*. Elsevier, B.V., Amsterdam, The Netherlands.

564 Lunn, D.J., Thomas, A., Best, N., Spiegelhalter, D., 2000. Winbugs-a bayesian modelling framework: concepts, structure, and extensibility. *Statistics*  
565 *and computing* 10, 325–337.

566 Mannering, F.L., Bhat, C.R., 2014. *Analytic methods in accident research: Methodological frontier and future directions*. *Analytic methods in*  
567 *accident research* 1, 1–22.

568 Meng, Y., Wu, L., Ma, C., Guo, X., Wang, X., 2020. A comparative analysis of intersection hotspot identification: Fixed vs. varying dispersion  
569 parameters in negative binomial models. *Journal of Transportation Safety & Security* , 1–18.

570 Miaou, S.P., Lord, D., 2003. Modeling traffic crash-flow relationships for intersections: dispersion parameter, functional form, and bayes versus  
571 empirical bayes methods. *Transportation Research Record* 1840, 31–40.

572 Miaou, S.P., Song, J.J., 2005. Bayesian ranking of sites for engineering safety improvements: decision parameter, treatability concept, statistical  
573 criterion, and spatial dependence. *Accident Analysis & Prevention* 37, 699–720.

574 Miranda-Moreno, L.F., Heydari, S., Lord, D., Fu, L., 2013. Bayesian road safety analysis: Incorporation of past evidence and effect of hyper-prior  
575 choice. *Journal of safety research* 46, 31–40.

576 Park, E.S., Lord, D., 2007. Multivariate poisson-lognormal models for jointly modeling crash frequency by severity. *Transportation Research*  
577 *Record* 2019, 1–6.

578 Part, D., 2010. Highway safety manual. American Association of State Highway and Transportation Officials: Washington, DC, USA .

579 Persaud, B., Lan, B., Lyon, C., Bhim, R., 2010. Comparison of empirical bayes and full bayes approaches for before–after road safety evaluations.  
580 Accident Analysis & Prevention 42, 38–43.

581 Plummer, M., et al., 2003. Jags: A program for analysis of bayesian graphical models using gibbs sampling, in: Proceedings of the 3rd international  
582 workshop on distributed statistical computing, Vienna, Austria.. pp. 1–10.

583 Plummer, M., et al., 2016. rjags: Bayesian graphical models using mcmc. R package version 4.

584 Pu, Z., Li, Z., Jiang, Y., Wang, Y., 2020. Full bayesian before-after analysis of safety effects of variable speed limit system. IEEE transactions on  
585 intelligent transportation systems .

586 Safaei, B., Safaei, N., Masoud, A., Seyedekrami, S., 2021a. Weighing criteria and prioritizing strategies to reduce motorcycle-related injuries using  
587 combination of fuzzy topsis and ahp methods. Advances in transportation studies 54.

588 Safaei, N., Zhou, C., Safaei, B., Masoud, A., 2021b. Gasoline prices and their relationship to the number of fatal crashes on us roads. Transportation  
589 Engineering 4, 100053.

590 Saha, D., Alluri, P., Dumbaugh, E., Gan, A., 2020. Application of the poisson-tweedie distribution in analyzing crash frequency data. Accident  
591 Analysis & Prevention 137, 105456.

592 Scrucca, L., et al., 2013. Ga: a package for genetic algorithms in r. Journal of Statistical Software 53, 1–37.

593 Shaon, M.R.R., Qin, X., Shirazi, M., Lord, D., Geedipally, S.R., 2018. Developing a random parameters negative binomial-lindley model to analyze  
594 highly over-dispersed crash count data. Analytic methods in accident research 18, 33–44.

595 Shen, W., Louis, T.A., 1998. Triple-goal estimates in two-stage hierarchical models. Journal of the Royal Statistical Society: Series B (Statistical  
596 Methodology) 60, 455–471.

597 Shirazi, M., Dhavala, S.S., Lord, D., Geedipally, S.R., 2017. A methodology to design heuristics for model selection based on the characteristics  
598 of data: Application to investigate when the negative binomial lindley (nb-l) is preferred over the negative binomial (nb). Accident Analysis &  
599 Prevention 107, 186–194.

600 Shirazi, M., Lord, D., 2019. Characteristics-based heuristics to select a logical distribution between the poisson-gamma and the poisson-lognormal  
601 for crash data modelling. Transportmetrica A: Transport Science 15, 1791–1803.

602 Shirazi, M., Lord, D., Dhavala, S.S., Geedipally, S.R., 2016. A semiparametric negative binomial generalized linear model for modeling over-  
603 dispersed count data with a heavy tail: Characteristics and applications to crash data. Accident Analysis & Prevention 91, 10–18.

604 Song, J.J., Ghosh, M., Miaou, S., Mallick, B., 2006. Bayesian multivariate spatial models for roadway traffic crash mapping. Journal of multivariate  
605 analysis 97, 246–273.

606 Tajuddin, R.R.M., Ismail, N., Ibrahim, K., Bakar, S.A.A., 2020. A four-parameter negative binomial-lindley distribution for modeling over and  
607 underdispersed count data with excess zeros. Communications in Statistics-Theory and Methods , 1–13.

608 Vehtari, A., Gelman, A., Gabry, J., 2017. Practical bayesian model evaluation using leave-one-out cross-validation and waic. Statistics and computing  
609 27, 1413–1432.

610 Watanabe, S., Opper, M., 2010. Asymptotic equivalence of bayes cross validation and widely applicable information criterion in singular learning  
611 theory. Journal of machine learning research 11.

612 Willmot, G., 1986. Mixed compound poisson distributions. ASTIN Bulletin: The Journal of the IAA 16, S59–S79.

613 Xie, Y., Zhang, Y., 2008. Crash frequency analysis with generalized additive models. Transportation Research Record 2061, 39–45.

614 Zamani, H., Ismail, N., 2010. Negative binomial-lindley distribution and its application. Journal of Mathematics and Statistics 6, 4–9.

615 Zha, L., Lord, D., Zou, Y., 2016. The poisson inverse gaussian (pig) generalized linear regression model for analyzing motor vehicle crash data.



- 616 Journal of Transportation Safety & Security 8, 18–35.
- 617 Zou, Y., Ash, J.E., Park, B.J., Lord, D., Wu, L., 2018. Empirical bayes estimates of finite mixture of negative binomial regression models and its  
618 application to highway safety. *Journal of Applied Statistics* 45, 1652–1669.
- 619 Zou, Y., Lord, D., Zhang, Y., Peng, Y., 2013. Comparison of sichel and negative binomial models in estimating empirical bayes estimates. *Trans-*  
620 *portation research record* 2392, 11–21.