



Comparison of Confidence and Prediction Intervals for Different Mixed-Poisson Regression Models

Journal:	<i>Journal of Transportation Safety & Security</i>
Manuscript ID	UTSS-2018-0300.R1
Manuscript Type:	Original Article
Date Submitted by the Author:	22-Jun-2019
Complete List of Authors:	Ash, John E. ; University of Washington Zou, Yajie; Tongji University, Lord, Dominique; TAMU, Wang, Yinhai ; University of Washington
Keywords:	Mixed-Poisson, Confidence interval, Prediction interval, Highway safety

SCHOLARONE™
Manuscripts

Comparison of Confidence and Prediction Intervals for Different Mixed-Poisson Regression Models

Submitted by

John E. Ash

Research Assistant

Department of Civil and Environmental Engineering, University of Washington

Box 352700, Seattle, WA 98195-2700

Tel: (414) 467-9555 Email: jeash@uw.edu

Yajie Zou, Ph. D. (Corresponding Author)

Associate Professor

Key Laboratory of Road and Traffic Engineering Ministry of Education

Tongji University, Shanghai 201804, China

Tel: (86) 13681865023 Email: yajiezou@hotmail.com

Dominique Lord, Ph. D.

Professor

Zachry Department of Civil Engineering

Texas A&M University, 3136 TAMU

College Station, TX 77843-3136

Tel: (979) 458-3949, fax: (979) 845-6481 E-mail: d-lord@tamu.edu

Yinhai Wang, Ph.D.

Professor

Department of Civil and Environmental Engineering, University of Washington

Box 352700, Seattle, WA 98195-2700

Tel: (206) 616-2696, Fax: (206) 543-1543 Email: yinhai@uw.edu

ABSTRACT

A major focus for transportation safety analysts is the development of crash prediction models, a task for which an extremely wide selection of model types is available. Perhaps the most common crash prediction model is the negative binomial (NB) regression model. The NB model gained popularity due to its relative ease of implementation and its ability to handle overdispersion in crash data. Recently, many new models including the Poisson-Inverse-Gaussian, Sichel, Poisson-Lognormal, and Poisson-Weibull models have been introduced as they can also accommodate overdispersion and could potentially replace the NB model, since many have been found to perform better. All five of the aforementioned models, including the NB model, can be classified as mixed-Poisson models. A mixed-Poisson model arises when an error term, following a chosen mixture distribution, enters the functional form for the Poisson parameter. For the NB model, the mixture distribution is selected as gamma, hence the alternate model name of Poisson-Gamma model. In this paper, confidence intervals (CIs) for the Poisson mean (μ) as well as prediction intervals (PIs) for the Poisson parameter (m , alternately referred to as the safety), and the predicted number of crashes at a new site (y) are derived for each of the aforementioned types of mixed-Poisson models. After the derivations, the theory is put into practice when CIs and PIs are estimated for mixed-Poisson models developed from an animal-vehicle collision dataset. Ultimately, this study provides safety analysts with tools to express levels of uncertainty associated with estimates from safety-modeling efforts instead of simply providing point estimates.

Keywords: Mixed-Poisson; Confidence interval; Prediction interval; Highway safety

INTRODUCTION

Transportation safety analysts often develop statistical models to predict crash frequencies that take into account a variety of factors including geometric characteristics of facilities and traffic volumes among many others (Mannering & Bhat, 2014). With constant advances in statistical methodologies, a variety of potential model types are available to analyze crash frequency data (Lord & Mannering, 2010). Early efforts in crash frequency modeling typically focused around the use of a Poisson regression model to predict crash frequency (Gustavsson, 1969; Gustavsson & Svensson, 1976; Jovanis & Chang, 1986). Although the Poisson model is very straightforward to use, it is unable to handle overdispersion that is commonly observed in crash data (Lord & Mannering, 2010). Overdispersion is said to occur when the variance of the crash counts is found to be greater than the mean and is quite common in crash data (Lord & Mannering, 2010). Lord, Washington, and Ivan (2005) noted that overdispersion is a consequence of considering crash data as resulting from Poisson trials (i.e., Bernoulli trials where the probability of a crash in each trial is not constant). Common features of overdispersed crash datasets include high frequencies of zero-valued and/or large-valued crash counts that are not able to be modeled properly by a simple Poisson distribution (Lord et al., 2005). **While not a focus of this paper, is important to note that work has been done to address underdispersion in crash data as well (i.e., when the variance is less than the mean). For example, Lord, Geedipally, and Guikema (2010) and Giuffrè, Graná, Roberta, and Corriere (2011) applied the Conway-Maxwell Poisson model to examine underdispersed crash data.**

In an effort to accommodate overdispersion in crash data, a variety of models have been introduced. Perhaps the most popular is the negative binomial (NB) model which has been used by many researchers to model overdispersed crash data (Connors, Maher, Wood, Mountain, &

Ash et al.

4

Ropkins, 2013; El-Basyouny & Sayed, 2006; Hauer, Ng, & Lovell, 1988; Lord & Mannering, 2010; Maycock & Hall, 1984; Park, Carlson, Porter, & Andersen, 2012; **Srinivasan, Baek, & Council, 2010**; Ye, Pendyala, Shankar, & Konduri, 2013). A key feature of the NB model is the assumption that the mean crash frequency (i.e., the Poisson parameter) for any site i , λ_i , follows a gamma distribution (Hauer, 1992). Thus, a formulation for the marginal mean and variance for a crash count, y_i , is obtained in which the variance can exceed the mean (Hauer, 1992; Lawless, 1987; Lord & Mannering, 2010). The NB model is alternately referred to as the Poisson-Gamma model as the crash count for site i , y_i , conditioned on the Poisson parameter λ_i (which follows the gamma distribution), is itself Poisson distributed. That being said, there is no reason to assume that λ_i must be gamma distributed (Hauer, 1997; Lord et al., 2005). In fact, researchers have investigated a variety of other distributions for the λ parameter, which result in several other types of mixed-Poisson regression models (i.e., models in which the crash count conditioned on the Poisson parameter, whose distribution is known as the mixing or mixture distribution, follows a Poisson distribution) (Cameron & Trivedi, 2013; Lawless, 1987). For example, one alternate choice of mixture distribution is the generalized inverse Gaussian (GIG) distribution which gives rise to the Sichel (SI) model for modeling overdispersed count data (Rigby, Stasinopoulos, and Akantziliotou, 2008). Zou, Wu, and Lord (2015) applied a Sichel model to analyze a highly-dispersed crash dataset from Texas and compared the results to those obtained from a traditional NB model. They found the SI model yielded lower values of both the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) (i.e., better statistical fit) than the NB model. Yet another choice of mixture distribution is the inverse Gaussian (IG) distribution, which gives rise to the Poisson-Inverse-Gaussian (PIG) model (Dean, Lawless, & Willmot, 1989). Zha, Lord, and Zou (2016) analyzed the aforementioned Texas crash dataset as well as a crash dataset from

1
2
3 Washington State and found that PIG regression models provided better fit (in terms of AIC and
4 BIC) than traditional NB models for their application. Other possible choices for the mixture
5 distribution include, but are not limited to, the Weibull and lognormal distributions, leading to the
6 Poisson-Weibull (PW) and Poisson-Lognormal (PLN) models, respectively. Both types of models
7 can accommodate overdispersion; the PW model has been applied to crash data by Cheng,
8 Geedipally, and Lord (2013) and the PLN model has been investigated by Lord and Miranda-
9 Moreno (2008), Aguero-Valverde and Jovanis (2008), **Lan and Persaud (2012), and Zhao, Liu,
10 Li, and Sharma (2018), among others.**

11
12 All of the aforementioned mixed-Poisson regression models provide only a point estimate
13 of the expected crash frequency at a given site. Although point estimates may provide some benefit
14 in prediction, in many cases a confidence interval for a given estimate is preferable (Casella &
15 Berger, 2001). Notably, confidence intervals are important for use in safety decision-making as
16 they express the uncertainty for a given point estimate (Lord, 2008; Lord, Kuo, & Geedipally,
17 2010). Wood (2005) derived formulae for the prediction intervals (PIs) for the predicted response
18 (i.e., crash frequency at a new site, y_i) and the gamma mean (m_i), as well as confidence intervals
19 (CIs) for the true mean crash frequency (alternately referred to as the mean response or Poisson
20 mean, μ_i), from the NB (Poisson-gamma) regression model. It is important to note the distinction
21 here between the Poisson parameter and the Poisson mean. In the case of standard Poisson
22 regression, the two values are in fact equal. However, in a mixed-Poisson model, introduction of
23 an error term into the Poisson parameter makes it such that the two terms are no longer equal.
24 Further details on this issue are provided later in the paper. Ultimately, Wood (2005) noted that
25 PIs and CIs may be especially useful to predict the number of crashes expected to occur at different
26 sites with similar features to sites considered in the model development. Lord (Lord, 2008)

Ash et al.

6

1
2
3 developed a methodology for calculating the predicted confidence intervals for the multiplication
4 of NB regression models with crash modification factors (CMFs). Geedipally and Lord (2008)
5 compared PIs for the predicted response (y) and the gamma mean (m) as well as CIs for the mean
6 response/Poisson mean (μ) as computed from NB models with fixed and varying dispersion
7 parameters, as well as for univariate and bivariate NB models (Geedipally & Lord, 2010). Lord,
8 Kuo, and Geedipally (2010) estimated PIs for the number of crashes (y) as obtained from an NB
9 model with several covariates and compared them to those estimated from a “baseline model” (i.e.,
10 flow-only model) that was adjusted with crash-modification factors (CMFs). Connors et al. (2013)
11 plotted CIs and PIs for predicted values of μ_i and m for variable flows and segment lengths for the
12 NB and PLN cases; they did not, however, provide explicit formulae for the CIs and PIs associated
13 with the PLN model.
14
15
16
17
18
19
20
21
22
23
24
25
26
27

28
29 The goal of this study is to extend the work of Wood (2005) and develop the associated
30 CIs and PIs for a broader range of mixed-Poisson regression models commonly used by safety
31 analysts today. In order to ensure this work aligns with Wood (2005), the use of his notation of PI
32 in reference to y and m (dependent variables) and the use of CI in reference to μ (a model
33 parameter) will continue henceforth in this paper. Besides reviewing the derivation of the PIs for
34 y and m and the CI for μ , respectively, for the NB model per Wood (2005), derivations of CIs
35 and PIs for the aforementioned three values (where m is now generalized to be the Poisson
36 parameter which follows the given mixture distribution, alternately referred to as the safety) will
37 be provided for the Poisson-Inverse-Gaussian, Sichel, Poisson-Weibull, and Poisson-Lognormal
38 regression models. Then, a case study making use of an animal-vehicle collision dataset will be
39 conducted, in which the five mixed-Poisson models in consideration will be estimated. Once the
40 models have been estimated, PIs and CIs will be estimated and plotted for y , m , and μ , as
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

determined from each model. These PIs and CIs will then be compared and discussed with regards to the case study and in general terms. Ultimately, this study provides safety analysts with tools to express levels of uncertainty associated with estimates from safety-modeling efforts instead of simply providing point estimates.

DERIVATION OF CONFIDENCE AND PREDICTION INTERVALS

This section provides the derivations for the confidence and prediction intervals for each type of mixed-Poisson model considered in this study. First, however, brief background information on mixed-Poisson models is provided. In general, there are three-levels in the hierarchy of a mixed-Poisson model. At the lowest level in the hierarchy is the mean response (μ_i), also known as the Poisson mean, which itself follows a normal distribution, $N(\mu_0, \sigma_0^2)$. One level up is the Poisson parameter (m_i), alternately known as the safety, which when conditioned on the Poisson mean follows the mixture distribution in consideration. Finally, there is the predicted response (y_i), i.e., the crash frequency at site i , which when conditioned on the Poisson parameter (m_i), follows a Poisson distribution.

Mixed-Poisson Models and Formulation

A mixed-Poisson model is defined by two primary criteria. First, the count in consideration (i.e., number of crashes λ_i) follows a Poisson distribution conditional on the Poisson parameter λ_i (Cameron & Trivedi, 2013). Following the terminology and notation of Wood (2005), the Poisson parameter λ_i will be referred to as the “safety” and denoted m_i :

$$f(y_i | m_i) = \frac{\exp(-m_i) * m_i^{y_i}}{y_i!}, y_i = 0, 1, 2, \dots \quad (1)$$

Second, the Poisson parameter, m_i , has a multiplicative error term following a chosen mixture distribution (e.g., gamma, inverse Gaussian etc.) that is expressed in the conditional mean (i.e., the safety conditioned on the mean response μ_i) (Cameron & Trivedi, 2013). Note without the error term, the expression reduces to that of the mean response (μ_i), alternately referred to as the Poisson mean.

$$\begin{aligned}
 m_i &= \exp(\beta_0 + \sum_{j=1}^K x'_{ij} * \beta_j + \varepsilon_i) \\
 &= \exp(\beta_0 + \varepsilon_i) * \exp\left(\sum_{j=1}^K x_{ij} * \beta_j\right) \\
 &= \exp\left(\beta_0 + \sum_{j=1}^K x_{ij} * \beta_j\right) * \exp(\varepsilon_i) \\
 &= \mu_i * \nu_i
 \end{aligned} \tag{2}$$

Where,

i = site index;

β_j = j^{th} regression coefficient;

x_{ij} = j^{th} covariate for site i ; and

ε_i = error term such that $\exp(\varepsilon_i)$, itself referred to as ν_i , follows the chosen mixture distribution.

As aforementioned, $Y|m_i \sim \text{Poisson}(m_i)$. The marginal distribution for Y is derived by integrating out the error term ν_i as follows, note $h(\nu_i)$ is the mixture distribution (Cameron & Trivedi, 2013):

$$\begin{aligned}
 f(y_i|\mu_i) &= \int_0^\infty g(y_i|\mu_i, v_i) * h(v_i) dv_i \\
 &= E_v[g(y_i|\mu_i, v_i)]
 \end{aligned}
 \tag{3}$$

It can then be shown, via application of the equality $m_i = \mu_i v_i$, that the Poisson parameter m_i does indeed follow the mixture distribution (just like v_i) (Cameron & Trivedi, 2013).

Parametrizations of Mixture Distributions

A total of five mixed-Poisson models are considered in this study. The models, corresponding mixture distributions for m_i and v_i , and parameterizations for the mixture distributions are presented in the following where the section header notes the model type and the distribution in parentheses is the mixture distribution. The subscript i is left out without loss of generality.

Negative Binomial [NB] Model (Gamma)

The NB model arises when the choice of mixture distribution for the Poisson-mixture model is chosen to be the gamma distribution. Specifically, $v \sim \text{Gamma}(\delta, \phi)$; however, in order to properly identify the intercept in the regression equation, $E[v] = 1$. This result is obtained by setting $\delta = \phi$, and thus leading to a one-parameter gamma distribution. It then follows that $\text{Var}[v] = 1/\phi$, alternately stated $\text{Var}[v] = \alpha$, and further that $m|\phi, \mu \sim \text{Gamma}\left(\phi, \frac{\phi}{\mu}\right)$ (Cameron & Trivedi, 2013).

Poisson-Inverse-Gaussian [PIG] Model (Inverse Gaussian)

The PIG model arises when the inverse Gaussian (IG) distribution is selected as the mixture distribution. Specifically, $v \sim \text{IG}(\mu_{IG}, \lambda)$, where the subscript ‘‘IG’’ is used to distinguish μ_{IG} (the

Ash et al.

10

mean of the IG distribution) from the Poisson mean (μ). As was the case with the NB model, the intercept identification condition calls for $E[v] = 1$. If $\mu_{IG} = 1$, then $E[v] = 1$, and further $Var[v] = 1/\lambda$ (Rigby and Stasinopolous, 2009).

Sichel [SI] (Generalized Inverse Gaussian)

If the mixture distribution is selected to be the generalized inverse Gaussian (GIG) distribution, the Sichel model is obtained. Here, $v \sim GIG(\mu_{GIG}, \sigma_{GIG}, \nu_{GIG})$. If $E[v] = 1$ (intercept identification condition), the variance of v is expressed as follows (Rigby et al., 2008; Rigby and Stasinopolous, 2009):

$$Var[v] = \frac{2\sigma_{GIG}(\nu_{GIG} + 1)}{c} + \frac{1}{c^2} - 1 \quad (4)$$

Where,

$$c = R_{\nu_{GIG}}\left(\frac{1}{\sigma_{GIG}}\right);$$

$$R_{\lambda}(t) = K_{\lambda+1}(t)/K_{\lambda}(t); \text{ and}$$

$$K_{\lambda}(t) = \frac{1}{2} \int_0^{\infty} x^{\lambda-1} \exp\left[-\frac{1}{2}t(x + x^{-1})\right] dx \text{ (where, } K_{\lambda}(t) \text{ is the modified Bessel function of the third kind).}$$

Poisson-Lognormal [PLN] (Lognormal)

When the mixture distribution is selected as the lognormal distribution, the Poisson-Lognormal model is obtained. Here, $v \sim \log N(d, \sigma_{LN}^2)$. The mean of the lognormal distribution is expressed as follows (Connors et al., 2013):

$$E[v] = \exp\left(d + \frac{\sigma_{LN}^2}{2}\right) \quad (5)$$

The intercept identification condition $E[v] = 1$ is then obtained by requiring $d = -\sigma_{LN}^2/2$.

The variance of v is then obtained by substituting the aforementioned expression for d into the following equation for $Var[v]$:

$$\begin{aligned} Var[v] &= (e^{\sigma_{LN}^2} - 1)e^{2d + \sigma_{LN}^2} \\ &= e^{\sigma_{LN}^2} - 1 \end{aligned} \quad (6)$$

Poisson-Weibull [PW] (Weibull)

When the error term, v , follows the Weibull distribution, the Poisson-Weibull model is obtained. Specifically, $v \sim Weibull(\mu_W, \sigma_W)$. The mean of the Weibull distribution is expressed as follows (Cheng et al., 2013; Rigby and Stasinopolous 2009):

$$E[v] = \frac{1}{\mu_W^{1/\sigma_W}} \Gamma\left(\frac{1}{\sigma_W} + 1\right) \quad (7)$$

Thus, in order to meet the intercept identification condition of $E[v] = 1$, the following must hold:

$$\mu_W = \left(\Gamma\left(\frac{1}{\sigma_W} + 1\right)\right)^{\sigma_W} \quad (8)$$

With $E[v] = 1$, the variance of v can thus be obtained by substituting the aforementioned expression for μ_W into the following equation for $Var[v]$.

$$\begin{aligned} Var[v] &= \frac{1}{\mu_W^{2/\sigma_W}} \left[\Gamma\left(\frac{2}{\sigma_W} + 1\right) - \left(\Gamma\left(\frac{1}{\sigma_W} + 1\right)\right)^2 \right] \\ &= \frac{\Gamma\left(\frac{2}{\sigma_W} + 1\right)}{\left(\Gamma\left(\frac{1}{\sigma_W} + 1\right)\right)^2} - 1 \end{aligned} \quad (9)$$

Derivation of Confidence Intervals for Poisson Mean (True Mean Crash Frequency) (μ)

For this study, we consider a generalized linear model (GLM) for crash prediction, where each site of interest is a road segment, of form shown in Equations (10 a) and (10 b).

$$\eta = \log\left(\frac{\mu}{L * t}\right) = \log(\beta_0) + \sum_{i=1}^n x_{ij} * \beta_j \quad (10 a)$$

$$\eta = \log\left(\frac{\mu}{L * t}\right) = \beta'_0 + \sum_{i=1}^n x_{ij} * \beta_j \quad (10 b)$$

Where,

η = linear predictor;

β_j = j^{th} regression coefficient;

x_{ij} = j^{th} predictor for segment (site);

L = segment length (in miles); and

t = time period over which crash data was collected.

The product of segment length and the duration over which crash data was collected for each site are considered as an offset, leading to a reformulation in Equation (10) of the regression as follows:

$$\eta = \log(\mu) = \log(\beta_0) + \sum_{i=1}^n x_{ij} * \beta_j + \log(L * t) \quad (11)$$

Under Equation (11), if we consider x_1 as a traffic volume (F), we have:

$$\mu = \beta_0 F^{\beta_1} (L * t) * \exp\left(\sum_{i=2}^n x_{ij} * \beta_j\right)$$

In the GLM, estimators for the regression coefficients, $\hat{\beta}_j$, follow a multivariate normal distribution, $[\hat{\beta}'_0, \dots, \hat{\beta}'_n]' \sim N([\beta'_0, \dots, \beta'_n]', \Sigma)$. From, Equation (11), it is clear that $\mu = \exp(\eta)$. As done previously, the subscript i for the values of η and μ at site i are omitted without loss of generality. We can use this fact to derive an approximate $(1-\alpha)\%$ confidence interval (CI) for the Poisson mean (alternately, the true mean crash count), μ , as follows (where $Z_{1-\alpha/2}$ is the critical value for the $1 - \alpha/2$ quantile of the standard normal distribution) (Wood, 2005):

$$\begin{aligned} & \exp(\hat{\eta} \pm Z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\eta})}) \\ &= \exp(\hat{\eta}) * \exp(\pm Z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\eta})}) \\ &= \hat{\mu} * \exp(\pm Z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\eta})}) \end{aligned}$$

$$= \left[\frac{\hat{\mu}}{\exp(Z_{1-\alpha/2}\sqrt{\text{Var}(\hat{\eta}))}}, \hat{\mu} * \exp(Z_{1-\alpha/2}\sqrt{\text{Var}(\hat{\eta}))} \right] \quad (12)$$

Thus, regardless of the choice of mixture distribution, the approximate $(1 - \alpha)\%$ CI for the Poisson mean (alternately, the true mean crash count), μ , is as formulated in Equation (12). The reader is referred to Wood (2005) for the steps to calculate the variance of the linear predictor.

Derivation of Prediction Intervals for Poisson Parameter (m)

Here, the derivation of an approximate $(1 - \alpha)\%$ confidence interval for the Poisson parameter, alternately referred to as the safety, m is presented based on the procedure outlined in Wood (2005). Before beginning the derivation, it is important to note a useful result that is used in several subsequent calculations, that being that the distribution of the estimator for the Poisson mean ($\hat{\mu}$) although technically lognormal, can be approximated as normal (Wood, 2005). Hence, $\hat{\mu} \sim N(\mu_0 = \mu, \sigma_0^2 = \mu^2 \text{Var}(\hat{\eta}))$.

Following Wood (2005), the basic formulation for an approximate $(1 - \alpha)\%$ PI for the Poisson parameter, alternately the safety, m is presented in Equation (13).

$$\hat{\mu} \pm Z_{1-\alpha/2} * \sqrt{\text{Var}(m)} \quad (13)$$

Mathematically, it is possible for the lower bound of the PI in Equation (13) to be negative, though physically speaking, negative values of m are not sensible. Hence, the PI for m in equation (13) is reformulated as:

$$\left[\max\{0, \hat{\mu} - Z_{1-\frac{\alpha}{2}} * \sqrt{\text{Var}(m)}\}, \hat{\mu} + Z_{1-\frac{\alpha}{2}} * \sqrt{\text{Var}(m)} \right] \quad (14)$$

The variance of m is formulated as follows:

$$\begin{aligned}
 \text{Var}(m) &= \text{Var}(\mu v) \\
 &= E(\mu^2 v^2) - E(\mu v)^2 \\
 &= E(\mu^2)E(v^2) - E(\mu)^2 E(v)^2 \quad [\text{by independence of } \mu \text{ and } v] \\
 &= [\text{Var}(\mu) + E(\mu)^2] * [\text{Var}(v) + E(v)^2] - E(\mu)^2 E(v)^2
 \end{aligned} \tag{15}$$

The derived expressions for the variance of m for each of the mixture distributions considered in this study are presented in Table 1.

With formulations for $\text{Var}(m)$ in hand, and recalling that the distribution of $\hat{\mu}$ can be approximated as normal (which leads to the key substitution $\sigma_0^2 = \mu^2 \text{Var}(\hat{\eta})$), the derived PIs for m for each type of mixed-Poisson model considered in this study are presented in Table 2. For simplicity, 95% PIs are shown.

Derivation of Prediction Intervals for Predicted Crash Count (y)

The final type of prediction interval of interest for a mixed-Poisson model is that for the predicted crash count (y) at a new site. The formulation for the PI is developed based on Chebyshev's inequality and further assumes the following: (1) The lower bound for y is zero (make the PI more conservative and follows the convention of Wood (2005)); (2) y must be integer-valued (Wood, 2005). In general, a $(1 - \alpha)\%$ PI for y is shown in Equation (16), where the floor of the upper bound is taken to ensure it is integer-valued. As an example, to obtain a 95% PI for y , the expression under the first radical would evaluate to 19.

$$[0, [\hat{\mu} + \sqrt{\alpha^{-1} - 1} \sqrt{\text{Var}(y)}]] \quad (16)$$

The variance of Y is evaluated as follows:

$$\begin{aligned} \text{Var}(Y) &= E\{\text{Var}(Y | M)\} + \text{Var}\{E(Y|M)\} \\ &= E(M) + \text{Var}(M) \\ &= E(\mu v) + \text{Var}(M) \\ &= E(\mu) * E(v) + \text{Var}(M) \\ &= \mu_0 + \text{Var}(M) \end{aligned} \quad (17)$$

Hence, the $(1 - \alpha)\%$ PI for Y can be re-expressed as shown in Equation (18).

$$[0, [\hat{\mu} + \sqrt{\alpha^{-1} - 1} \sqrt{\hat{\mu} + \text{Var}(m)}]] \quad (18)$$

Thus, using the formulation for $\text{Var}(m)$ as provided in Equation (15), the 95% PIs for Y in the case of each of the five mixed-Poisson models are developed and shown in Table 3.

CASE STUDY

This section provides a case study in which mixed-Poisson models are developed from a crash dataset collected in Washington State. Following model development, confidence and prediction intervals for the aforementioned crash rates are estimated and displayed graphically.

Data Description

The dataset used in this study was based upon that used in Lao, Wu, Corey, and Wang (2011). Specifically, it was collected to model animal-vehicle collisions along ten highways in Washington State over a total of 752 road segments. The dependent variable represents the number of animal carcasses removed from each road segment over a five-year period from 2002 to 2006. A summary of the data, including the explanatory variables and relevant summary statistics can be seen in Table 4. More detailed information on the dataset can be found in Lao et al. (2011).

Model Development

A total of five mixed-Poisson models were estimated from the animal-vehicle collision dataset. Each model included the set of variables that were found to be significant at the $\alpha = 0.05$ level in the Negative Binomial model, and all models included an offset term. The Negative Binomial, Poisson-Inverse-Gaussian, and Sichel models were estimated through a maximum likelihood (ML) approach in the GAMLSS package for the R statistical software (Rigby and Stasinopoulos, 2005). The Poisson-Lognormal and Poisson-Weibull likelihood functions do not have a closed form, hence these models had to be estimated through a Bayesian approach in the WinBUGS software package (Lunn, Thomas, Best, & Spiegelhalter, 2000) (using 20,000 iterations following 5000 iterations for burn in). Model parameters, associated standard errors (SE) (posterior values for the Bayesian models), p-values, and goodness-of-fit statistics including the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) are presented in Table 5 (a) and 5 (b).

Confidence and Prediction Intervals

In order to compare and contrast the prediction intervals for y and m , as well as the confidence intervals for μ associated with each of the five types of mixed-Poisson models considered, plots were made constructed to show the intervals for each model (Figure 1 (a)-(e)). The results in each plot are based upon a model using the coefficients shown in Table 5, in which AADT is allowed

1
2
3 to vary (in increments of 50) between 0 and 120,000 vehicles per day (as this was approximately
4 the range in the animal-vehicle collision dataset), and segment length and time period were fixed
5 the range in the animal-vehicle collision dataset), and segment length and time period were fixed
6 at one mile and 5 years, respectively, for the offset term. The rest of the variables were set to
7 constant values as follows (these values were the most common value of each variable,
8 respectively, in the dataset). That is to say, the intervals show numbers of animal-vehicle collisions
9 over a five-year period on a one-mile segment with varying AADT and marginalizing across all
10 other variables.
11
12
13
14
15
16
17
18

19 We first consider the 95% CI for the Poisson mean (μ). From Table 5 (a) and (b) it can be
20 seen that regardless of the type of model considered, the estimates for the model regression
21 coefficients are quite similar. Hence, the estimates for the Poisson mean were quite similar between
22 models, with maximum values ranging between 17.97 for the Poisson-Weibull model to as high
23 as 21.53 for the Sichel model, when AADT=120,000. From Figure 1 (a) through (e), it can be seen
24 that both the lower and upper bounds of the 95% CI for the Poisson mean are relatively similar in
25 value, respectively, across different models. Ultimately, at AADT=120,000, the tightest interval
26 around the estimate of μ resulted from the Poisson-Lognormal model ([11.94, 31.25],
27 width=19.31) and the widest interval resulted from the Sichel model ([11.62, 39.87],
28 width=28.25). As the true value of the Poisson mean is not known for each site, conclusions on
29 whether or not the narrowest interval is “best” cannot be made.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 For the plots of the 95% PIs for the safety, m , in Figure 1 (a) through (e), the lower bound
46 values are not shown as all values for each model, regardless of AADT were found to be zero. In
47 fact, such models produced negative values for the lower bounds when calculated, however, as
48 noted in Equation (14), negative values of the safety (m) are not sensible and hence the lowest
49 reasonable value is zero. For the Poisson-Inverse-Gaussian model, the width of the interval at an
50
51
52
53
54
55
56
57
58
59
60

1
2
3 AADT value of 120,000 was 149.10, the maximum value of m estimated from any of the 95%
4
5 CIs. The lowest value of an upper bound for the 95% CIs for m at 120,000 AADT, across all
6
7 models, was 72.89 as obtained from the Poisson-Weibull model. As was the case with the Poisson
8
9 mean, the true value of m is unknown, thus no comments on which interval is narrowest, while
10
11 still capturing the true parameter value can be made.
12
13

14
15 Regardless of model considered, the lower bound for the 95% prediction intervals for the
16
17 predicted response at a new site (y) is always zero, hence this interval is not shown in any of the
18
19 plots. Additionally, one will likely notice that the upper bounds for the PIs for y are much greater
20
21 (specifically, 1.92 to 2.06 times greater at 120,000 ADT) than the respective upper bounds for the
22
23 95% PIs for m . Besides yielding the largest values, the curves for the PIs for y are notably less
24
25 smooth than those representing the CIs for μ and PIs for m . This step-function appearance is a
26
27 result of the use of the floor function in calculation for the upper bound of the PI as is shown in
28
29 Equation (18), as the number of crashes predicted to occur at a new site should be integer-valued.
30
31 The upper bounds for the PIs for y ranged from as high as 307 for the Poisson-Inverse-Gaussian
32
33 model to as low as 141 for the Poisson-Weibull model. The upper bounds for the PIs for y as
34
35 predicted from the Negative Binomial, Sichel, and Poisson-Weibull models were found to be much
36
37 closer to each other than for the other models.
38
39
40
41

42
43 Table 7 shows a summary of key values for the different CIs and PIs (shown in Figure 1)
44
45 for each of the aforementioned mixed-Poisson models when the value of AADT is 120,000 (i.e.,
46
47 at the right end of the plots).
48
49
50

51
52 **As a final part of the case study, comparisons were made between the estimates of the**
53
54 **Poisson mean, μ , and the lower and upper bounds of the 95% CIs and PIs for each model**
55
56
57
58
59
60

1
2
3 across the range of values of covariates found in the dataset described in Table 5. While the
4 previous portion of the case study examined CIs and PIs based on varying the AADT and
5 keeping all other variable values fixed (to the values in Table 6), this portion computed all
6 estimates based on the full dataset described in Table 5 (i.e., all covariates covered a range
7 of values, none were fixed). For each of the mixed-Poisson models, the mean squared error
8 (MSE) was then computed between the estimated values of μ , and the lower and upper
9 bounds of the CI for μ and the PIs for the safety (m) and the predicted response (y)
10 considering all data points, respectively. The model coefficients in Table 5 were used to
11 estimate values of the Poisson mean (μ) for each data point in the dataset under each model.
12 Next, the 95% CI for the Poisson mean (μ) and the PIs for the safety (m) and the predicted
13 response (y) were estimated on a per-model basis for each data point. Finally, MSE values
14 were estimated between the estimates of μ and the lower and upper bounds of each of the
15 CIs and PIs, across all data points.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 The results of the MSE calculations between the estimated values of the Poisson mean
34 (μ) and the lower and upper bounds of the CIs and PIs, by model, are shown in Table 8.
35 From the table, it can be seen that, for the animal-vehicle collision dataset considered, the
36 Negative Binomial model yielded estimates for the Poisson mean (μ) that produced the
37 smallest MSE values for all confidence and prediction intervals. Thus, it appears that the
38 Negative Binomial model seems to provide less variation for the CIs and PIs than the other
39 models considered, based on the dataset investigated. On the other hand, the Poisson-
40 Lognormal model yielded the largest MSE values for the 95% CI for Poisson mean (μ), with
41 respect to both the lower and upper bounds of the interval. When considering the variation
42 between the estimate of μ and the lower bounds for the PI for the safety (m) and the PI for
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 the predicted response (y), across all models, the largest MSE values were observed under
4
5 the Poisson-Lognormal model. However, for the variation between the estimate of μ and the
6
7 upper bounds for the PI for the safety (m) and the PI for the predicted response (y), across
8
9 all models, the largest MSE values were observed under the Poisson-Inverse-Gaussian
10
11 model. For all models considered, the values of MSE for the lower bounds of the PIs for the
12
13 safety (m) and the predicted response (y) were equal on a per-model basis. This was due to
14
15 the fact that in all cases, the lower bound of the PIs for m and y was zero.
16
17
18
19
20
21

22 SUMMARY AND CONCLUSIONS

23
24 Based upon the initial work of Wood (2005), confidence intervals for the Poisson mean (μ), safety
25
26 or Poisson parameter/safety (m), and predicted response (i.e., number of crashes at a new site, y)
27
28 for four types of mixed-Poisson models beyond that for the Negative Binomial (Poisson-Gamma)
29
30 model as given by Wood (2005) were developed. Formulae for these intervals are now available
31
32 for researchers and practitioners to use in order to obtain a window of uncertainty associated with
33
34 predictions, as compared to a sole point estimate. Specifically, the types of mixed-Poisson models
35
36 considered in this study were the Poisson-Inverse Gaussian, Sichel, Poisson-Lognormal, and
37
38 Poisson-Weibull models, all of which arise by allowing for a multiplicative error term following a
39
40 corresponding mixture distribution to enter the functional form of the Poisson parameter m . After
41
42 motivating mixed-Poisson models, derivations for the aforementioned confidence and prediction
43
44 intervals were provided.
45
46
47
48

49
50 Once the formulae for the CIs and PIs had been established, the theory was put into practice
51
52 by investigating the intervals associated with each of the five aforementioned types of mixed-
53
54 Poisson models. Since real-life, observed crash data was used, the true values of the Poisson mean
55
56
57
58
59
60

(μ) and safety or Poisson parameter (m) are unknown, hence comments cannot be made on which intervals perform “best.” Nonetheless, several important conclusions can be drawn from the case study considering the Texas data:

- (1) Regardless of the model considered, the estimates of the regression coefficients were relatively similar and hence values of the Poisson mean (μ) were quite similar regardless of AADT value;
- (2) Of the models developed in this study, the Sichel model yielded the widest interval for the Poisson mean (μ). The Poisson-Inverse-Gaussian model yielded the widest intervals for the safety (m) and predicted response at a new site (y). It further yielded the second greatest upper bound when considering all 95% PIs for the predicted response, y . That being said, there is no way to confirm narrower intervals on the μ , n , and y are necessarily better as the true values of these parameters is unknown;
- (3) For the Poisson mean (μ), the Poisson-Lognormal model yielded the narrowest 95% CI. The Poisson-Weibull model yielded the narrowest 95% PIs for m and for y of all models considered;
- (4) All models estimated yielded negative values for the lower bound on the safety (m) (before coercing them to be zero);
- (5) At the largest AADT values considered, the upper bounds on the PIs for y ranged from 1.92 to 2.06 times the values of the upper bounds of m at the same AADT;
- (6) MSE values were computed to study variation between the estimated values of the Poisson mean (μ) and the lower and upper bounds of the 95% CIs and PIs for each model, considering the datapoints in the animal-vehicle collision dataset. The Negative Binomial model yielded estimates for the Poisson mean (μ) that produced**

1
2
3 the smallest MSE values for all confidence and prediction intervals. The Poisson-
4 Lognormal model yielded the largest MSE values for the 95% CI for Poisson mean (μ
5
6
7
8), with respect to both the lower and upper bounds of the interval;
9

10 (7) When considering the variation between the estimate of μ and the lower bounds for
11 the PI for the safety (m) and the PI for the predicted response (y), across all models,
12 the largest MSE values were observed under the Poisson-Lognormal model;
13
14

15 (8) Finally, when examining the variation between the estimate of μ and the upper
16 bounds for the PI for the safety (m) and the PI for the predicted response (y), across
17 all models, the largest MSE values were observed under the Poisson-Inverse-
18 Gaussian model.
19
20
21
22
23
24
25
26
27

28 In terms of future work, this study introduces a few possibilities currently under
29 investigation by the authors. For example, a simulation study involving simulation of values for
30 the Poisson mean (μ), safety (m), and response (i.e., crash count, y) at a new site could help
31 determine which CIs and PIs best represent the true intervals. Finally, with the latest modeling
32 tools that have proposed, the estimated CIs and PIs should be developed for multiparameter models
33 (Geedipally, Lord, & Dhavala, 2012; Lord & Geedipally, 2018), random parameters models
34 (Anastasopoulos & Mannering, 2009; Rista et al., 2018; Shaon, Qin, Shirazi, Lord, & Geedipally,
35 2018), and semi-parametric models (Heydari, Fu, Lord, & Mallick, 2016; Shirazi, Lord, Dhavala,
36 & Geedipally, 2016; Ye, Wang, Zou, & Lord, 2018; Zou, Ash, Park, Lord, & Wu, 2018).
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 ACKNOWLEDGEMENT

52
53
54
55
56
57
58
59
60

This research is sponsored jointly by the National Natural Science Foundation of China (grant no. 51608386) and Shanghai Science and Technology Committee (grant no. 18510745400).

REFERENCES

- Aguero-Valverde, J., & Jovanis, P. (2008). Analysis of Road Crash Frequency with Spatial Models. *Transportation Research Record: Journal of the Transportation Research Board*, 2061, 55–63. <https://doi.org/10.3141/2061-07>
- Anastasopoulos, P. Ch., & Mannering, F. L. (2009). A note on modeling vehicle accident frequencies with random-parameters count models. *Accident Analysis & Prevention*, 41(1), 153–159. <https://doi.org/10.1016/j.aap.2008.10.005>
- Cameron, A. C., & Trivedi, P. K. (2013). *Regression Analysis of Count Data (Econometric Society Monographs)*. <https://doi.org/10.1017/CBO9781139013567>
- Casella, G., & Berger, R. L. (2001). *Statistical Inference* (2nd edition). Australia ; Pacific Grove, CA: Cengage Learning.
- Cheng, L., Geedipally, S. R., & Lord, D. (2013). The Poisson–Weibull generalized linear model for analyzing motor vehicle crash data. *Safety Science*, 54, 38–42. <https://doi.org/10.1016/j.ssci.2012.11.002>
- Connors, R. D., Maher, M., Wood, A., Mountain, L., & Ropkins, K. (2013). Methodology for fitting and updating predictive accident models with trend. *Accident Analysis & Prevention*, 56, 82–94. <https://doi.org/10.1016/j.aap.2013.03.009>
- Dean, C., Lawless, J. F., & Willmot, G. E. (1989). A Mixed Poisson-Inverse-Gaussian Regression Model. *The Canadian Journal of Statistics / La Revue Canadienne de Statistique*, 17(2), 171–181. <https://doi.org/10.2307/3314846>
- El-Basyouny, K., & Sayed, T. (2006). Comparison of Two Negative Binomial Regression Techniques in Developing Accident Prediction Models. *Transportation Research Record: Journal of the Transportation Research Board*, 1950, 9–16. <https://doi.org/10.3141/1950-02>
- Geedipally, S. R., & Lord, D. (2008). Effects of Varying Dispersion Parameter of Poisson–Gamma Models on Estimation of Confidence Intervals of Crash Prediction Models. *Transportation Research Record: Journal of the Transportation Research Board*, 2061(1), 46–54. <https://doi.org/10.3141/2061-06>
- Geedipally, S. R., & Lord, D. (2010). Investigating the effect of modeling single-vehicle and multi-vehicle crashes separately on confidence intervals of Poisson–gamma models. *Accident Analysis & Prevention*, 42(4), 1273–1282. <https://doi.org/10.1016/j.aap.2010.02.004>
- Geedipally, S. R., Lord, D., & Dhavala, S. S. (2012). The negative binomial-Lindley generalized linear model: Characteristics and application using crash data. *Accident Analysis & Prevention*, 45, 258–265. <https://doi.org/10.1016/j.aap.2011.07.012>
- Giuffrè, O., Granà, A., Roberta, M., & Corriere, F. (2011). Handling Underdispersion in Calibrating Safety Performance Function at Urban, Four-Leg, Signalized Intersections. *Journal of Transportation Safety & Security*, 3(3), 174–188. <https://doi.org/10.1080/19439962.2011.599014>
- Gustavsson, J. (1969). On the use of regression models in the study of road accidents. *Accident Analysis & Prevention*, 1(4), 315–321. [https://doi.org/10.1016/0001-4575\(69\)90077-3](https://doi.org/10.1016/0001-4575(69)90077-3)

- 1
2
3 Gustavsson, J., & Svensson, Å. (1976). A Poisson Regression Model Applied to Classes of Road
4 Accidents with Small Frequencies. *Scandinavian Journal of Statistics*, 3(2), 49–60.
- 5 Hauer, E. (1992). Empirical bayes approach to the estimation of “unsafety”: The multivariate
6 regression method. *Accident Analysis & Prevention*, 24(5), 457–477.
7 [https://doi.org/10.1016/0001-4575\(92\)90056-O](https://doi.org/10.1016/0001-4575(92)90056-O)
- 8
9 Hauer, E. (1997). *Observational Before/After Studies in Road Safety. Estimating the Effect of*
10 *Highway and Traffic Engineering Measures on Road Safety*. Retrieved from
11 <https://trid.trb.org/view/1136089>
- 12 Hauer, E., Ng, J. C. N., & Lovell, J. (1988). ESTIMATION OF SAFETY AT SIGNALIZED
13 INTERSECTIONS (WITH DISCUSSION AND CLOSURE). *Transportation Research*
14 *Record*, (1185). Retrieved from <https://trid.trb.org/view/301420>
- 15 Heydari, S., Fu, L., Lord, D., & Mallick, B. K. (2016). Multilevel Dirichlet process mixture
16 analysis of railway grade crossing crash data. *Analytic Methods in Accident Research*, 9,
17 27–43. <https://doi.org/10.1016/j.amar.2016.02.001>
- 18
19 Jovanis, P. P., & Chang, H.-L. (1986). MODELING THE RELATIONSHIP OF ACCIDENTS TO
20 MILES TRAVELED. *Transportation Research Record*, (1068). Retrieved from
21 <https://trid.trb.org/view/288394>
- 22
23 Lan, B., & Persaud, B. (2012). Evaluation of Multivariate Poisson Log Normal Bayesian Methods
24 for Before-After Road Safety Evaluations. *Journal of Transportation Safety & Security*,
25 4(3), 193–210. <https://doi.org/10.1080/19439962.2011.649194>
- 26 Lao, Y., Wu, Y.-J., Corey, J., & Wang, Y. (2011). Modeling animal-vehicle collisions using
27 diagonal inflated bivariate Poisson regression. *Accident Analysis & Prevention*, 43(1),
28 220–227. <https://doi.org/10.1016/j.aap.2010.08.013>
- 29
30 Lawless, J. F. (1987). Negative binomial and mixed poisson regression. *Canadian Journal of*
31 *Statistics*, 15(3), 209–225. <https://doi.org/10.2307/3314912>
- 32 Lord, D. (2008). Methodology for estimating the variance and confidence intervals for the estimate
33 of the product of baseline models and AMFs. *Accident Analysis & Prevention*, 40(3),
34 1013–1017. <https://doi.org/10.1016/j.aap.2007.11.008>
- 35 Lord, D., & Geedipally, S. R. (2018). Chapter 14. Safety Prediction with Datasets Characterised
36 with Excess Zero Responses and Long Tails. In *Transport and Sustainability: Vol. 11. Safe*
37 *Mobility: Challenges, Methodology and Solutions* (Vol. 11, pp. 297–323).
38 <https://doi.org/10.1108/S2044-994120180000011016>
- 39
40 Lord, D., Geedipally, S. R., & Guikema, S. D. (2010). Extension of the Application of Conway-
41 Maxwell-Poisson Models: Analyzing Traffic Crash Data Exhibiting Underdispersion. *Risk*
42 *Analysis*, 30(8), 1268–1276. <https://doi.org/10.1111/j.1539-6924.2010.01417.x>
- 43
44 Lord, D., Kuo, P.-F., & Geedipally, S. R. (2010). Comparison of Application of Product of
45 Baseline Models and Accident-Modification Factors and Models with Covariates:
46 Predicted Mean Values and Variance. *Transportation Research Record: Journal of the*
47 *Transportation Research Board*, (2147). Retrieved from <https://trid.trb.org/view/909798>
- 48 Lord, D., & Mannering, F. (2010). The statistical analysis of crash-frequency data: A review and
49 assessment of methodological alternatives. *Transportation Research Part A: Policy and*
50 *Practice*, 44(5), 291–305. <https://doi.org/10.1016/j.tra.2010.02.001>
- 51
52 Lord, D., & Miranda-Moreno, L. F. (2008). Effects of low sample mean values and small sample
53 size on the estimation of the fixed dispersion parameter of Poisson-gamma models for
54 modeling motor vehicle crashes: A Bayesian perspective. *Safety Science*, 46(5), 751–770.
55 <https://doi.org/10.1016/j.ssci.2007.03.005>
- 56
57
58
59
60

- 1
2
3 Lord, D., Washington, S. P., & Ivan, J. N. (2005). Poisson, Poisson-gamma and zero-inflated
4 regression models of motor vehicle crashes: Balancing statistical fit and theory. *Accident*
5 *Analysis & Prevention*, 37(1), 35–46. <https://doi.org/10.1016/j.aap.2004.02.004>
6
7 Lunn, D. J., Thomas, A., Best, N., & Spiegelhalter, D. (2000). WinBUGS - A Bayesian modelling
8 framework: Concepts, structure, and extensibility. *Statistics and Computing*, 10(4), 325–
9 337. <https://doi.org/10.1023/A:1008929526011>
10
11 Mannering, F. L., & Bhat, C. R. (2014). Analytic methods in accident research: Methodological
12 frontier and future directions. *Analytic Methods in Accident Research*, 1, 1–22.
13 <https://doi.org/10.1016/j.amar.2013.09.001>
14
15 Maycock, G., & Hall, R. D. (1984). *ACCIDENTS AT 4-ARM ROUNDABOUTS*. Presented at the
16 Planning & Transport Res & Comp, Sum Ann Mtg, Proc. Retrieved from
17 <https://trid.trb.org/view/269602>
18
19 Park, E. S., Carlson, P. J., Porter, R. J., & Andersen, C. K. (2012). Safety effects of wider edge
20 lines on rural, two-lane highways. *Accident Analysis & Prevention*, 48, 317–325.
21 <https://doi.org/10.1016/j.aap.2012.01.028>
22
23 Rigby, R. A., & Stasinopoulos, D. M. (2005). Generalized additive models for location, scale and
24 shape. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 54(3), 507–
25 554. <https://doi.org/10.1111/j.1467-9876.2005.00510.x>
26
27 Rigby, R. A., Stasinopoulos, D. M., & Akantziliotou, C. (2008). A framework for modelling
28 overdispersed count data, including the Poisson-shifted generalized inverse Gaussian
29 distribution. *Computational Statistics & Data Analysis*, 53(2), 381–393.
30 <https://doi.org/10.1016/j.csda.2008.07.043>
31
32 Rigby, R., & Stasinopolous, D. (2009). *A flexible regression approach using GAMLSS in R*.
33 London: London Metropolitan University.
34
35 Rista, E., Goswamy, A., Wang, B., Barrette, T., Hamzeie, R., Russo, B., ... Savolainen, P. T.
36 (2018). Examining the safety impacts of narrow lane widths on urban/suburban arterials:
37 Estimation of a panel data random parameters negative binomial model. *Journal of*
38 *Transportation Safety & Security*, 10(3), 213–228.
39 <https://doi.org/10.1080/19439962.2016.1273291>
40
41 Shaon, M. R. R., Qin, X., Shirazi, M., Lord, D., & Geedipally, S. R. (2018). Developing a Random
42 Parameters Negative Binomial-Lindley Model to analyze highly over-dispersed crash
43 count data. *Analytic Methods in Accident Research*, 18, 33–44.
44 <https://doi.org/10.1016/j.amar.2018.04.002>
45
46 Shirazi, M., Lord, D., Dhavala, S. S., & Geedipally, S. R. (2016). A semiparametric negative
47 binomial generalized linear model for modeling over-dispersed count data with a heavy
48 tail: Characteristics and applications to crash data. *Accident Analysis & Prevention*, 91,
49 10–18. <https://doi.org/10.1016/j.aap.2016.02.020>
50
51 Srinivasan, R., Baek, J., & Council, F. (2010). Safety Evaluation of Transverse Rumble Strips on
52 Approaches to Stop-Controlled Intersections in Rural Areas. *Journal of Transportation*
53 *Safety & Security*, 2(3), 261–278. <https://doi.org/10.1080/19439962.2010.508571>
54
55 Wood, G. R. (2005). Confidence and prediction intervals for generalised linear accident models.
56 *Accident Analysis & Prevention*, 37(2), 267–273.
57 <https://doi.org/10.1016/j.aap.2004.10.005>
58
59 Ye, X., Pendyala, R. M., Shankar, V., & Konduri, K. C. (2013). A simultaneous equations model
60 of crash frequency by severity level for freeway sections. *Accident Analysis & Prevention*,
57, 140–149. <https://doi.org/10.1016/j.aap.2013.03.025>

- 1
2
3 Ye, X., Wang, K., Zou, Y., & Lord, D. (2018). A semi-nonparametric Poisson regression model
4 for analyzing motor vehicle crash data. *PLOS ONE*, *13*(5), e0197338.
5 <https://doi.org/10.1371/journal.pone.0197338>
6
7 Zha, L., Lord, D., & Zou, Y. (2016). The Poisson inverse Gaussian (PIG) generalized linear
8 regression model for analyzing motor vehicle crash data. *Journal of Transportation Safety*
9 *& Security*, *8*(1), 18–35. <https://doi.org/10.1080/19439962.2014.977502>
10
11 Zhao, M., Liu, C., Li, W., & Sharma, A. (2018). Multivariate Poisson-lognormal model for
12 analysis of crashes on urban signalized intersections approach. *Journal of Transportation*
13 *Safety & Security*, *10*(3), 251–265. <https://doi.org/10.1080/19439962.2017.1323059>
14
15 Zou, Y., Ash, J. E., Park, B.-J., Lord, D., & Wu, L. (2018). Empirical Bayes estimates of finite
16 mixture of negative binomial regression models and its application to highway safety.
17 *Journal of Applied Statistics*, *45*(9), 1652–1669.
18 <https://doi.org/10.1080/02664763.2017.1389863>
19
20 Zou, Y., Wu, L., & Lord, D. (2015). Modeling over-dispersed crash data with a long tail:
21 Examining the accuracy of the dispersion parameter in Negative Binomial models. *Analytic*
22 *Methods in Accident Research*, *5–6*, 1–16. <https://doi.org/10.1016/j.amar.2014.12.002>
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. Variance of m for Mixture Distributions

Mixture Distribution	Var(m)
Gamma	$\alpha * (\sigma_0^2 + \mu_0^2) + \sigma_0^2$ [Note: $\alpha = \frac{1}{\phi}$]
Inverse Gaussian	$\frac{1}{\lambda} * (\sigma_0^2 + \mu_0^2) + \sigma_0^2$
Generalized Inverse Gaussian	$[\sigma_0^2 + \mu_0^2] * \left(\frac{2\sigma_{GIG}(\nu_{GIG} + 1)}{c} + \frac{1}{c^2} \right) - \mu_0^2$
Lognormal	$e^{\sigma_{LN}^2} * [\sigma_0^2 + \mu_0^2] - \mu_0^2$
Weibull	$[\sigma_0^2 + \mu_0^2] * \left[\frac{\Gamma\left(\frac{2}{\sigma_W} + 1\right)}{\left(\Gamma\left(\frac{1}{\sigma_W} + 1\right)\right)^2} \right] - \mu_0^2$

Table 2. 95% Prediction Intervals for m

Model	95% PI for m
Negative Binomial (NB)	$\left[\max \left(0, \hat{\mu} - 1.96 * \sqrt{\hat{\mu}^2 [\hat{\alpha}(\text{Var}(\hat{\eta}) + 1) + \text{Var}(\hat{\eta})]} \right), \hat{\mu} + 1.96 * \sqrt{\hat{\mu}^2 [\hat{\alpha}(\text{Var}(\hat{\eta}) + 1) + \text{Var}(\hat{\eta})]} \right]$
Poisson-Inverse-Gaussian (PIG)	$\left[\max \left(0, \hat{\mu} - 1.96 * \sqrt{\hat{\mu}^2 \left[\frac{1}{\hat{\lambda}} (\text{Var}(\hat{\eta}) + 1) + \text{Var}(\hat{\eta}) \right]} \right), \hat{\mu} + 1.96 * \sqrt{\hat{\mu}^2 \left[\frac{1}{\hat{\lambda}} (\text{Var}(\hat{\eta}) + 1) + \text{Var}(\hat{\eta}) \right]} \right]$
Sichel (SI)	$\left[\max \left(0, \hat{\mu} - 1.96 * \sqrt{\hat{\mu}^2 \left\{ [\text{Var}(\hat{\eta}) + 1] * \left(\frac{2\hat{\sigma}_{GIG}(\hat{\nu}_{GIG} + 1)}{c} + \frac{1}{c^2} \right) - 1 \right\}} \right), \hat{\mu} + 1.96 * \sqrt{\hat{\mu}^2 \left\{ [\text{Var}(\hat{\eta}) + 1] * \left(\frac{2\hat{\sigma}_{GIG}(\hat{\nu}_{GIG} + 1)}{c} + \frac{1}{c^2} \right) - 1 \right\}} \right]$
Poisson-Lognormal (PLN)	$\left[\max \left(0, \hat{\mu} - 1.96 * \sqrt{\hat{\mu}^2 [e^{\hat{\sigma}_{LN}^2} (\text{Var}(\hat{\eta}) + 1) - 1]} \right), \hat{\mu} + 1.96 * \sqrt{\hat{\mu}^2 [e^{\hat{\sigma}_{LN}^2} (\text{Var}(\hat{\eta}) + 1) - 1]} \right]$
Poisson-Weibull (PW)	$\left[\max \left(0, \hat{\mu} - 1.96 * \sqrt{\hat{\mu}^2 \left([\text{Var}(\hat{\eta}) + 1] * \frac{\Gamma\left(\frac{2}{\hat{\sigma}_W} + 1\right)}{\left(\Gamma\left(\frac{1}{\hat{\sigma}_W} + 1\right)\right)^2} - 1 \right)} \right), \hat{\mu} + 1.96 * \sqrt{\hat{\mu}^2 \left([\text{Var}(\hat{\eta}) + 1] * \frac{\Gamma\left(\frac{2}{\hat{\sigma}_W} + 1\right)}{\left(\Gamma\left(\frac{1}{\hat{\sigma}_W} + 1\right)\right)^2} - 1 \right)} \right]$

Table 3. 95% Prediction Intervals for y

Model	95% PI for y
Negative Binomial (NB)	$\left[0, \left[\hat{\mu} + \sqrt{19} \sqrt{\hat{\mu} + \hat{\mu}^2 [\hat{\alpha}(\text{Var}(\hat{\eta}) + 1) + \text{Var}(\hat{\eta})]}\right]\right]$
Poisson-Inverse-Gaussian (PIG)	$\left[0, \left[\hat{\mu} + \sqrt{19} \sqrt{\hat{\mu} + \hat{\mu}^2 \left[\frac{1}{\hat{\lambda}}(\text{Var}(\hat{\eta}) + 1) + \text{Var}(\hat{\eta})\right]}\right]\right]$
Sichel (SI)	$\left[0, \left[\hat{\mu} + \sqrt{19} \sqrt{\hat{\mu} + \hat{\mu}^2 \left\{[\text{Var}(\hat{\eta}) + 1] * \left(\frac{2\hat{\sigma}_{\text{GIG}}(\hat{\nu} + 1)}{c} + \frac{1}{c^2}\right) - 1\right\}}\right]\right]$
Poisson-Lognormal (PLN)	$\left[0, \left[\hat{\mu} + \sqrt{19} \sqrt{\hat{\mu} + \hat{\mu}^2 \left[e^{\hat{\sigma}_{\text{LN}}^2}(\text{Var}(\hat{\eta}) + 1) - 1\right]}\right]\right]$
Poisson-Weibull (PW)	$\left[0, \left[\hat{\mu} + \sqrt{19} \sqrt{\hat{\mu} + \hat{\mu}^2 \left([\text{Var}(\hat{\eta}) + 1] * \left[\frac{\Gamma\left(\frac{2}{\hat{\sigma}_W} + 1\right)}{\left(\Gamma\left(\frac{1}{\hat{\sigma}_W} + 1\right)\right)^2} - 1\right]}\right)\right]\right]$

Table 4. Summary Statistics for Animal-Vehicle Collision Dataset

Variable	Description	Min	Max	Mean	SD
Carcass	Number of carcasses per segment	0	53	3.47	6.86
AADT	Annual average daily traffic	612	120173	7721.85	10820.29
Access	Restrictive access control (0=No, 1=Yes)			0.15	
Spd_limt	Speed limit (miles per hour)	25	70	58.60	6.81
Trkpets	Truck percentage (%)	0	54.16	15.54	8.88
Nolanes	Number of lanes	2	7	2.48	0.95
Seg_lng	Segment length (miles)	0.5	1	0.69	0.14
TerRol	Terrain type rolling (0=No, 1=Yes)			0.76	
TerMou	Terrain type mountainous (0=No, 1=Yes)			0.13	
Lanewid	Lane width (feet)	10	17	11.78	0.55
Lshlw	Left shoulder width (feet)	0	20	6.02	2.77
Rshlw	Right shoulder width (feet)	0	26	8.64	6.24
White	White-tailed deer habitat (0=No, 1=Yes)			0.36	
Elk	Elk deer habitat (0=No, 1=Yes)			0.36	
Mule	Mule deer habitat (0=No, 1=Yes)			0.60	

Table 5. (a) Model Results Estimated with ML Approach

	NB			PIG			SI		
	Estimate	SE	p-val	Estimate	SE	p-val	Estimate	SE	p-val
Intercept $\log(\beta_0)$	-10.36	1.11	< 2.00E-16	-10.38	1.13	< 2.00E-16	-10.42	1.12	< 2.00E-16
$\log(\text{AADT}) \beta_1$	0.55	0.08	1.45E-10	0.59	0.09	2.12E-11	0.56	0.09	2.00E-10
Access β_2	-1.12	0.30	1.98E-04	-1.02	0.30	5.92E-04	-1.11	0.30	2.45E-04
Spd_limt β_3	0.09	0.02	1.15E-07	0.08	0.02	5.49E-07	0.09	0.02	1.22E-07
Nolanes β_4	-0.40	0.11	1.65E-04	-0.39	0.12	8.22E-04	-0.41	0.11	1.71E-04
Lshlw β_5	0.12	0.03	4.07E-06	0.11	0.03	1.07E-04	0.12	0.03	7.06E-06
White β_6	1.39	0.13	< 2.00E-16	1.59	0.13	< 2.00E-16	1.41	0.13	< 2.00E-16
Elk β_7	0.37	0.13	3.54E-03	0.50	0.14	2.43E-04	0.38	0.13	3.47E-03
Distribution Parameter(s)	$\alpha = 1/\varphi = 1.85$			$\lambda = 0.106$			$v_{\text{GIG}} = 0.4716, \sigma_{\text{GIG}} = 271$		
Global Deviance	2986.46			3010.00			2986.27		
AIC	3004.46			3028.00			3006.27		
BIC	3046.06			3069.60			3052.49		

Table 5. (b) Model Results Estimated with Bayesian Approach

	PLN					PW				
	Mean	SE	2.5%	50%	97.5%	Mean	SE	2.5%	50%	97.5%
Intercept $\log(\beta_0)$	-10.28	0.60	-11.20	-10.38	-8.87	-9.08	1.37	-10.75	-9.42	-6.87
$\log(\text{AADT}) \beta_1$	0.58	0.06	0.44	0.59	0.67	0.53	0.08	0.38	0.53	0.66
Access β_2	-1.05	0.32	-1.69	-1.04	-0.45	-0.79	0.31	-1.43	-0.79	-0.21
Spd_limt β_3	0.08	0.01	0.06	0.08	0.10	0.07	0.02	0.04	0.07	0.10
Nolanes β_4	-0.35	0.14	-0.57	-0.35	-0.10	-0.41	0.10	-0.61	-0.41	-0.22
Lshlw β_5	0.12	0.03	0.07	0.13	0.17	0.13	0.03	0.08	0.13	0.19
White β_6	1.69	0.15	1.40	1.69	2.00	1.51	0.13	1.25	1.51	1.74
Elk β_7	0.60	0.14	0.32	0.60	0.88	0.39	0.14	0.13	0.39	0.66
Distribution Parameter(s)	$\sigma_{\text{LN}} = 1.35$					$\sigma_{\text{W}} = 0.70$				

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 6. Default Values of Variables used in Models for Interval Construction

Variable	Value
Access β_2	0
Spd_limt β_3	60
Nolanes β_4	2
Lshlw β_5	8
White β_6	0
Elk β_7	0

For Peer Review Only

Table 7. Summary of Values for Mixed-Poisson CIs and PIs at AADT=120,000

	NB	PIG	SI	PLN	PW
μ Lower Bound	11.55	10.82	11.62	11.94	10.24
μ Max	21.23	20.27	21.53	19.32	17.97
μ Upper Bound	39.02	38.00	39.87	31.25	31.53
m Lower Bound	0.00	0.00	0.00	0.00	0.00
m Upper Bound	81.92	149.10	87.06	109.07	72.89
y Lower Bound	0.00	0.00	0.00	0.00	0.00
y Upper Bound	157	307	168	219	141

Table 8. MSE Values Calculated between Estimates of μ and 95% CI and PI Lower and Upper Bounds for the Animal-Vehicle Collision Dataset

	NB	PIG	SI	PLN	PW
μ Lower Bound	2.24	3.28	2.41	4.09	2.37
μ Upper Bound	4.49	6.96	4.95	8.83	4.76
m Lower Bound	30.56	38.35	31.37	48.26	32.68
m Upper Bound	226.46	1439.31	264.18	1008.96	282.80
y Lower Bound	30.56	38.35	31.37	48.26	32.68
y Upper Bound	1165.45	7141.38	1351.75	5025.57	1443.71

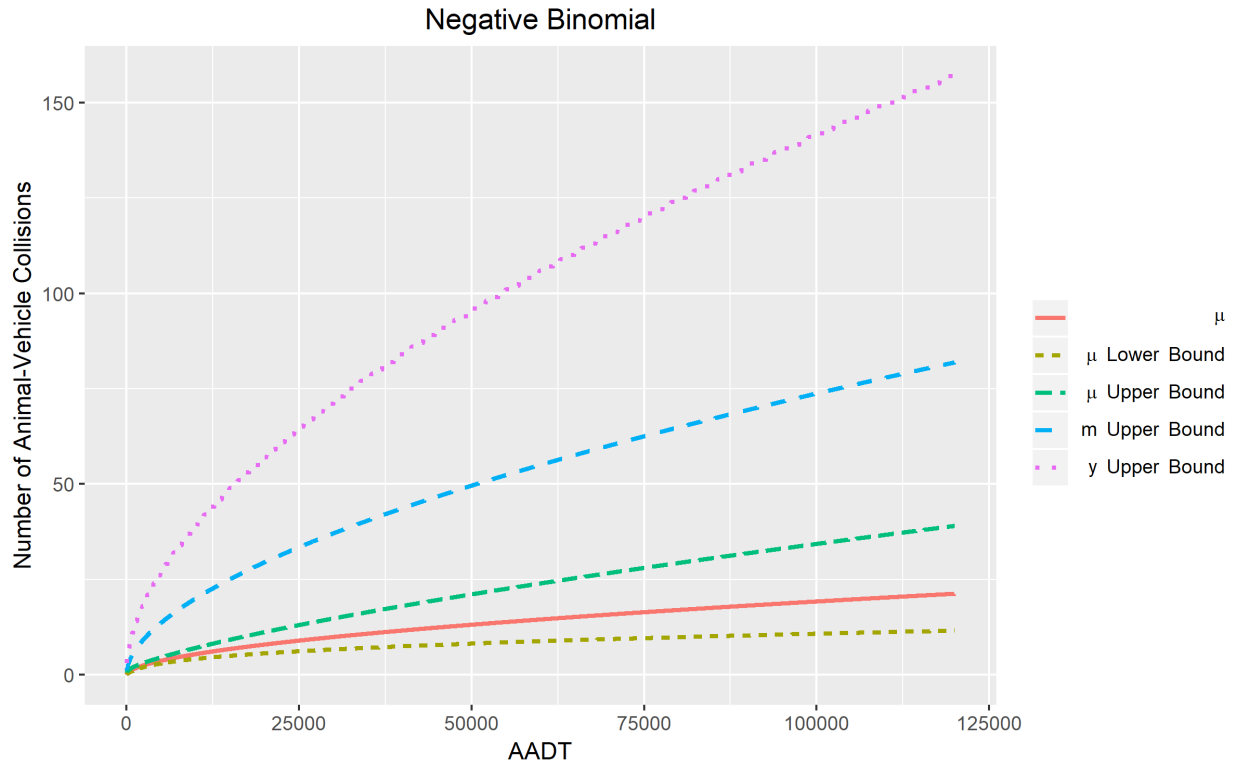


Figure 1 (a) 95% CIs and PI for Negative Binomial Model

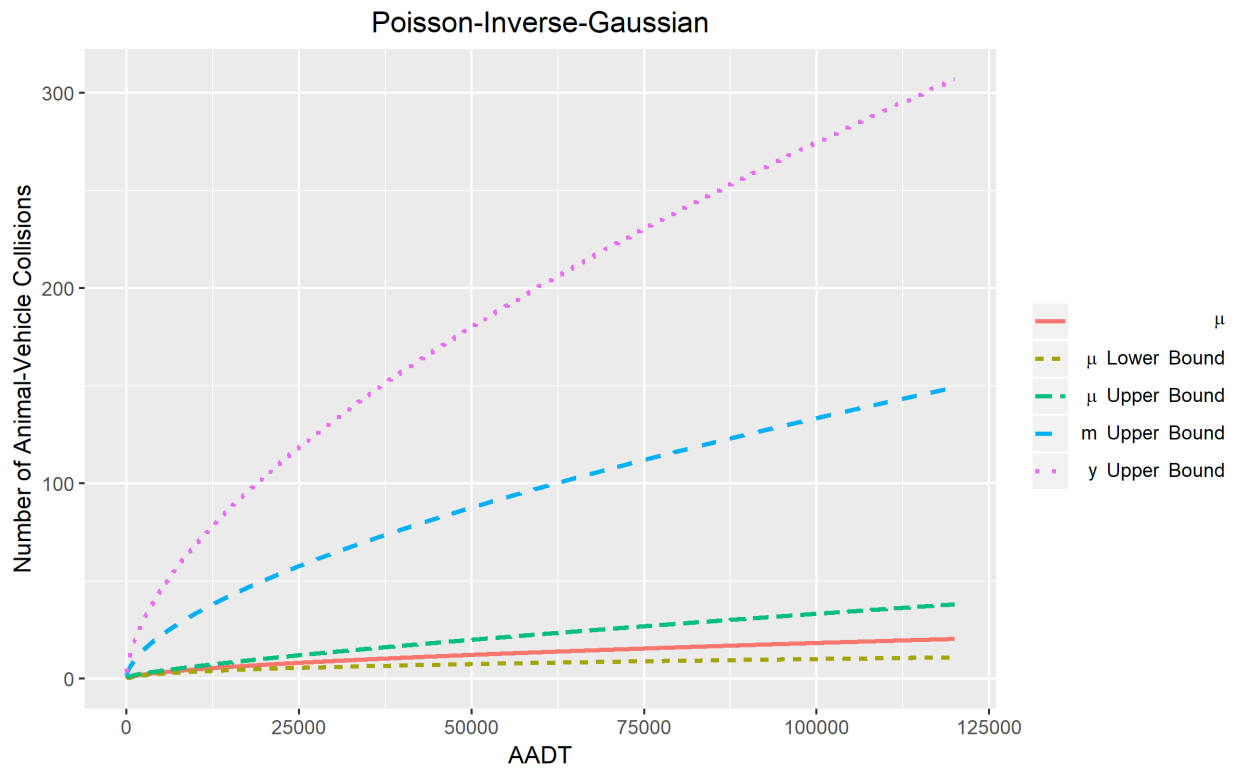


Figure 1 (b) 95% CIs and PI for Poisson-Inverse-Gaussian Model

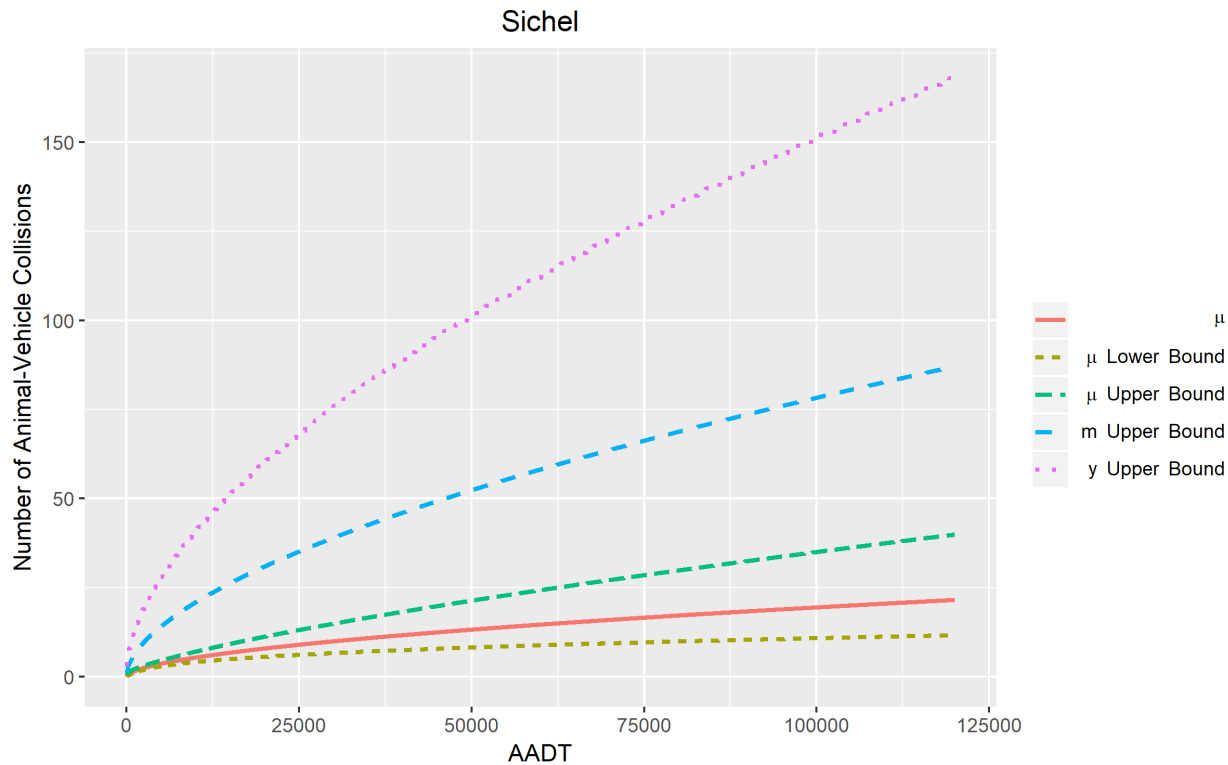


Figure 1 (c) 95% CIs and PI for Sichel Model

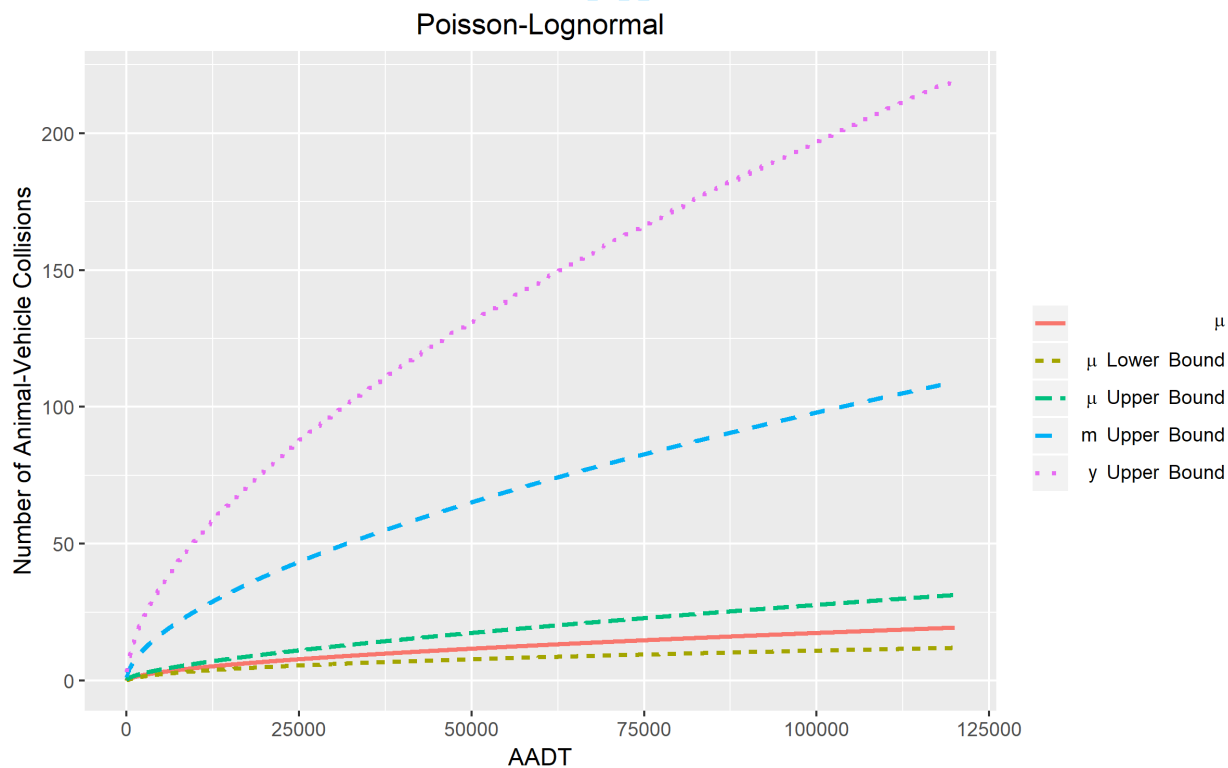


Figure 1 (d) 95% CIs and PI for Poisson-Lognormal Model

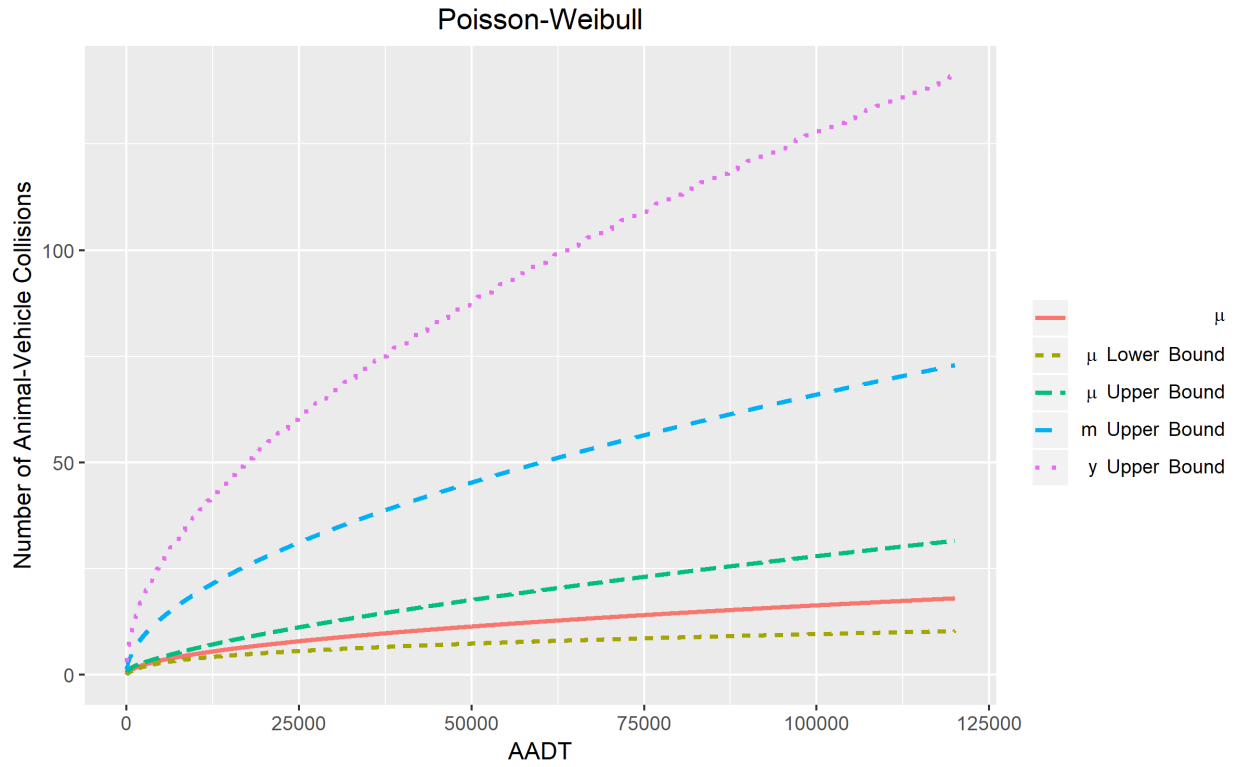


Figure 1 (e) 95% CIs and PI for Poisson-Weibull Model

Review Only

Response to comments

The authors would like to thank the reviewers for their comments on this manuscript. We have tried to address them all in the revised paper. Please see our responses to each comment in the following. **All new text in the document is indicated in bold and colored purple.**

Comment #1:

“The paper's focus is very clear and the authors have undertaken the work reasonably. My only comments are regarding how they have used these values. For example, in table 7, the authors provide a summary of each of the limits. However, i find this too simplistic. How is this any different than simply providing the upper and lower limits of the counts in the data. Basically, what i am saying is, the presentation does not provide adequate idea of the actual ranges provided on an individual record level. My suggestion would be compute the lower and upper limits at a record level and then provide summary statistics of the deviance across each record in the data as a table. This is where the value of the contribution be beneficial”

Response:

Thank you for this comment. We have since added a new table, Table 8 showing mean squared error values calculated between the linear predictor, μ , and the lower and upper limits of each of the CIs and PIs, across all datapoints in the dataset. We have further added an explanation of the results in Table 8 in two newly added paragraphs, directly after Table 7. The conclusions section has also been updated to reflect these changes. Please see Table 8 and the describing text added to the revised manuscript below...

“As a final part of the case study, comparisons were made between the estimates of the Poisson mean, μ , and the lower and upper bounds of the 95% CIs and PIs for each model across the range of values of covariates found in the dataset described in Table 5. While the previous portion of the case study examined CIs and PIs based on varying the AADT and keeping all other variable values fixed (to the values in Table 6), this portion computed all estimates based on the full dataset described in Table 5 (i.e., all covariates covered a range of values, none were fixed). For each of the mixed-Poisson models, the mean squared error (MSE) was then computed between the estimated values of μ , and the lower and upper bounds of the CI for μ and the PIs for the safety (m) and the predicted response (y) considering all data points, respectively. The model coefficients in Table 5 were used to estimate values of the Poisson mean (μ) for each data point in the dataset under each model. Next, the 95% CI for the Poisson mean (μ) and the PIs for the safety (m) and the predicted response (y) were estimated on a per-model basis for each data point. Finally, MSE values were estimated between the estimates of μ and the lower and upper bounds of each of the CIs and PIs, across all data points.

The results of the MSE calculations between the estimated values of the Poisson mean (μ) and the lower and upper bounds of the CIs and PIs, by model, are shown in Table 8. From the table, it can be seen that, for the animal-vehicle collision dataset considered, the Negative Binomial model yielded estimates for the Poisson mean (μ) that produced the smallest MSE values for all confidence and prediction intervals. Thus, it appears that the Negative Binomial model seems to provide less variation for the CIs and PIs than the other models considered, based on the dataset

investigated. On the other hand, the Poisson-Lognormal model yielded the largest MSE values for the 95% CI for Poisson mean (μ), with respect to both the lower and upper bounds of the interval. When considering the variation between the estimate of μ and the lower bounds for the PI for the safety (m) and the PI for the predicted response (y), across all models, the largest MSE values were observed under the Poisson-Lognormal model. However, for the variation between the estimate of μ and the upper bounds for the PI for the safety (m) and the PI for the predicted response (y), across all models, the largest MSE values were observed under the Poisson-Inverse-Gaussian model. For all models considered, the values of MSE for the lower bounds of the PIs for the safety (m) and the predicted response (y) were equal on a per-model basis. This was due to the fact that in all cases, the lower bound of the PIs for m and y was zero.”

Table 8. MSE Values Calculated between Estimates of μ and 95% CI and PI Lower and Upper Bounds for the Animal-Vehicle Collision Dataset

	NB	PIG	SI	PLN	PW
μ Lower Bound	2.24	3.28	2.41	4.09	2.37
μ Upper Bound	4.49	6.96	4.95	8.83	4.76
m Lower Bound	30.56	38.35	31.37	48.26	32.68
m Upper Bound	226.46	1439.31	264.18	1008.96	282.80
y Lower Bound	30.56	38.35	31.37	48.26	32.68
y Upper Bound	1165.45	7141.38	1351.75	5025.57	1443.71

Comment #2:

During revising the paper, authors are suggested to cite the topic-related articles published in Journal of Transportation Safety and Security in recent years.

Response:

Thank you for the comment. References to the following five Journal of Transportation Safety and Security (JTSS) papers have been added. We had one reference already, so there are now six JTSS papers cited.

Orazio Giuffrè, Anna Granà, Marino Roberta & Ferdinando Corriere (2011) Handling Underdispersion in Calibrating Safety Performance Function at Urban, Four-Leg, Signalized Intersections, Journal of Transportation Safety & Security, 3:3, 174-188, DOI: [10.1080/19439962.2011.599014](https://doi.org/10.1080/19439962.2011.599014)

Bo Lan & Bhagwant Persaud (2012) Evaluation of Multivariate Poisson Log Normal Bayesian Methods for Before-After Road Safety Evaluations, Journal of Transportation Safety & Security, 4:3, 193-210, DOI: [10.1080/19439962.2011.649194](https://doi.org/10.1080/19439962.2011.649194)

1
2
3
4 Emira Rista, Amrita Goswamy, Bo Wang, Timothy Barrette, Raha Hamzeie, Brendan Russo,
5 Georges Bou-Saab & Peter T. Savolainen (2018) Examining the safety impacts of narrow lane
6 widths on urban/suburban arterials: Estimation of a panel data random parameters negative
7 binomial model, Journal of Transportation Safety & Security, 10:3, 213-228, DOI:
8 10.1080/19439962.2016.1273291
9

10
11 Raghavan Srinivasan, Jongdae Baek & Forrest Council (2010) Safety Evaluation of Transverse
12 Rumble Strips on Approaches to Stop-Controlled Intersections in Rural Areas, Journal of
13 Transportation Safety & Security, 2:3, 261-278, DOI: [10.1080/19439962.2010.508571](https://doi.org/10.1080/19439962.2010.508571)
14

15
16 Mo Zhao, Chenhui Liu, Wei Li & Anuj Sharma (2018) Multivariate Poisson-lognormal model
17 for analysis of crashes on urban signalized intersections approach, Journal of Transportation
18 Safety & Security, 10:3, 251-265, DOI: [10.1080/19439962.2017.1323059](https://doi.org/10.1080/19439962.2017.1323059)
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60