# The negative binomial–Lindley distribution as a tool for analyzing crash data characterized by a large amount of zeros

Dominique Lord [a,*], Srinivas Reddy Geedipally [b,1]

[a] *Zachry Department of Civil Engineering, Texas A&M University, 3136 TAMU, College Station, TX 77843-3136, USA*
[b] *Engineering Research Associate, Texas Transportation Institute, Texas A&M University, 3135 TAMU, College Station, TX 77843-3135, USA*

## ARTICLE INFO

## ABSTRACT

The modeling of crash count data is a very important topic in highway safety. As documented in the literature, given the characteristics associated with crash data, transportation safety analysts have proposed a significant number of analysis tools, statistical methods and models for analyzing such data. Among the data issues, we find the one related to crash data which have a large amount of zeros and a long or heavy tail. It has been found that using this kind of dataset could lead to erroneous results or conclusions if the wrong statistical tools or methods are used. Thus, the purpose of this paper is to introduce a new distribution, known as the negative binomial–Lindley (NB-L), which has very recently been introduced for analyzing data characterized by a large number of zeros. The NB–L offers the advantage of being able to handle this kind of datasets, while still maintaining similar characteristics as the traditional negative binomial (NB). In other words, the NB–L is a two-parameter distribution and the long-term mean is never equal to zero. To examine this distribution, simulated and observed data were used. The results show that the NB–L can provide a better statistical fit than the traditional NB for datasets that contain a large amount of zeros.

## 1. Introduction

The modeling of crash count data is a very important topic in highway safety. As documented in Lord and Mannering (2010), given the characteristics associated with crash data, transportation safety analysts have proposed a significant number of analysis tools and models for analyzing such data. Among the data issues documented in the paper, we find the one related to crash data which have a large amount of zeros and a long or heavy tail. For such datasets, the number of sites where no crash is observed is so large that traditional statistical distributions or models, such as the Poisson and Poisson-gamma or negative binomial (NB) distributions, cannot be used efficiently. The Poisson distribution tends to underestimate the number of zeros given the mean of the data, while the NB may over-estimate zeros, but under-estimate observations with a count. This is obviously dependent upon the characteristics of the tail.

The large amount of zeros observed in crash data have initially been attributed to observations or sites that can be categorized under two states: a safe state, where no crash can occur, and a

non-safe state (Miaou, 1994; Shankar et al., 1997, 2003). A portion of the zero counts come from the safe state, while the rest of the zero counts come from a Poisson or NB distribution. The observations classified under the first state are considered as 'added' zeros. The zero-inflated model (both used for the Poisson and NB) has been consequently proposed to analyze this kind of dataset, usually because they provide better statistical fit (Shankar et al., 1997, 2003; Kumara and Chin, 2003). Some researchers (Warton, 2005; Lord et al., 2005, 2007) however have raised important methodological issues about the use of such models, including the fact that the safe state has a distribution with a long-term mean equal to zero, which is theoretically impossible. Lord et al. (2005) noted that the large amount of zeros can be attributed to the following factors: (1) sites with a combination of low exposure, high heterogeneity, and sites categorized as high risk; (2) analyses conducted with small time or spatial scales; (3) data with a relatively high percentage of missing or mis-reported crashes; and (4) crash models with omitted important variables. More recently, Mayshkina and Mannering (2009) have proposed a zero-state Markov switching model, which overcomes some of the criticisms discussed above, for analyzing longitudinal datasets characterized by a large number of zeros.

In cases in which the characteristics of the dataset cannot or is very difficult to be changed (as it will be discussed further below), the large number of zeros could still create a lot of difficulties for properly analyzing such dataset. This could obviously lead to erroneous results or conclusions if the wrong statistical tools or

---

methods are used. Thus, the purpose of this paper is to introduce a new distribution that has very recently been introduced for analyzing data characterized by a large number of zeros. This mixed distribution is known as the NB–Lindley (NB–L) distribution (Zamani and Ismail, 2010), which as the name implies, is a mixture of the NB and the Lindley distributions (Lindley, 1958; Ghitany et al., 2008). This two-parameter distribution has interesting and sound theoretical properties in which the distribution is characterized by a single long-term mean that is never equal to zero and a single variance function, similar to the traditional NB distribution. The properties of the NB–Lindley distribution are examined using simulated and observed data and a discussion is presented about the potential use of the NB–L distribution for traffic safety analyses. It important to point out that all documented distributions, such as the Poisson-gamma, Poisson-lognormal, Poisson-Pascal or the NB–L in highway safety research are in fact used as an approximation to describe the crash process. This process is known as the Poisson trials with unequal probability of events (See Lord et al., 2005, for additional details).

The paper is divided into five sections. Section 2 describes the characteristics of the NB–Lindley distribution. Section 3 presents the comparison analysis between the Poisson, NB and NB–L using simulated and observed data. Section 4 provides additional information for future work. Section 5 summarizes the study results.

## 2. Characteristics of the negative binomial–Lindley distribution

As discussed above, the NB–L distribution is a mixture of negative binomial and Lindley distributions. This mixed distribution has a thick tail and can be used when the data contains large number of zeros.

The negative binomial distribution is a mixture of Poisson and gamma distribution. The probability mass function (pmf) of the NB distribution can be given as:

$$P(Y = y; \phi, p) = \frac{\Gamma(\phi + y)}{\Gamma(\phi) \times y!}(1 - p)^{\phi}(p)^{y}; \quad \phi > 0, \quad 0 < p < 1 \quad (1)$$

The parameter '$p$' is defined as the probability of success in each trial and is given as:

$$p = \frac{\mu}{\mu + \phi} \quad (2)$$

where, $\mu = E(Y)$ = mean; and, $\phi$ = inverse dispersion parameter.

Then, it can be shown that the variance is (Casella and Berger, 1990):

$$\text{Var}(Y) = \phi \frac{p}{(1 - p)^2} = \frac{1}{\phi}\mu^2 + \mu \quad (3)$$

Using Eqs. (2) and (3), the pmf of the NB distribution can be re-parameterized this way:

$$P(Y = y; \mu, \phi) = \frac{\Gamma(\phi + y)}{\Gamma(\phi)\Gamma(y + 1)}\left(\frac{\phi}{\mu + \phi}\right)^{\phi}\left(\frac{\mu}{\mu + \phi}\right)^{y} \quad (4)$$

The pmf in Eq. (4) is the one normally used for crash count data.

The Lindley distribution is a mixture of exponential and gamma distribution (Lindley, 1958; Ghitany et al., 2008; Zamani and Ismail, 2010). The pmf of the Lindley distribution can be defined as follows:

$$f(X = x; \theta) = \frac{\theta^2}{\theta + 1}(1 + x)e^{-\theta x}; \quad \theta > 0, \quad x > 0 \quad (5)$$

A random variable $Z$ is assumed to follow a NB–L $(r, \theta)$ distribution when the following conditions satisfy:

$$Z \sim \text{NB}\left(r, P = 1 - e^{-\lambda}\right) \text{ and } \lambda \sim \text{Lindley}\left(\theta\right)$$

The pmf of the NB–L distribution is given as (Zamani and Ismail, 2010):

$$P(Z = z; r, \theta) = \frac{\Gamma(r + z)}{\Gamma(r) \times z!} \frac{\theta^2}{\theta + 1} \sum_{j=0}^{z} \frac{\Gamma(z + 1)}{\Gamma(j + 1) \times \Gamma(z + j + 1)}(-1)^j$$

$$\times \frac{\theta + r + j - 1}{(\theta + r + j)^2} \quad (6)$$

The parameter '$r$' is the shape parameter of NB–L distribution, similar to the inverse dispersion parameter '$\phi$' of the NB distribution. The parameter '$\theta$', in combination with shape parameter '$r$' dictates the mean and variance of the NB–L distribution.

The first moment (i.e., the mean) of the NB–L $(r, \theta)$ is given as:

$$E(Z) = r\left[\frac{\theta^3}{(\theta + 1)(\theta - 1)^2} - 1\right] \quad (7)$$

It should be noted that $E[Z] = E[Y] = \mu$

The second moment of the NB–L $(r, \theta)$ is given as:

$$E(Z^2) = (r + r^2)\left[\frac{\theta^2(\theta - 1)}{(\theta + 1)(\theta - 2)^2}\right]$$

$$- (r + 2r^2)\left[\frac{\theta^3}{(\theta + 1)(\theta - 1)^2}\right] + r^2 \quad (8)$$

The variance of the NB–L $(r, \theta)$ is calculated as:

$$\text{Var}(Z) = E(Z^2) - (E(Z))^2 \quad (9)$$

As described above, the NB–L distribution is a two-parameter distribution, which implies that the mean is never equal to zero. The NB–L is in fact an extension of the NB distribution.

To estimate the parameters $r$ and $\theta$, Eqs. (7) and (8) need to be solved iteratively and both parameters should be greater than 0 ($r > 0$ and $\theta > 0$). The parameters can also be estimated by solving the likelihood function, but the function is difficult to manipulate, since the partial derivatives contain multiple solutions. Additional work is therefore needed for finding the optimal solution among all the possible ones. More detailed information can be found in Zamani and Ismail (2010).

## 3. Application of the negative binomial–Lindley distribution

This section presents the comparison analysis results between the Poisson, NB and NB–L distributions using simulated and observed data.

### 3.1. Simulated data

The simulation protocol is the same one used by Lord et al. (2005), which consisted in simulating a Poisson distribution and add observations with the value zero to 'simulate' a two-state process, one of which is characterized by a long-term equal to zero. For this example, count data with 100 observations were simulated using a Poisson distribution with a mean equal to 0.50. Then, 100, 150 and 200 observations with the value zero were added to the data. The simulated data are summarized in Fig. 1. The original simulated data produced 57 observations with the value zero, 34 with the value 1, 8 with the value 2, and 1 with a value above 3 (5 to be exact).

The Poisson, NB and NB–L distributions were fitted based on the simulated data. Using the mean and variance of the data, the parameters were estimated with Eqs. (2) and (3) for the NB distribution and Eqs. (7)–(9) for the NB–L distribution. After the parameters
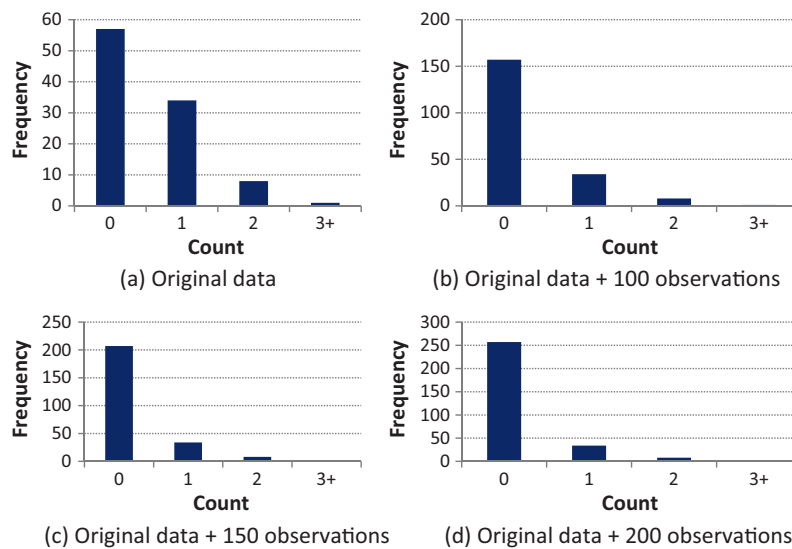
**Fig. 1.** Simulated data: Poisson distribution with a mean equal to 0.5 (and added 0s).

were estimated, the predicted probability and frequency were calculated for each count. The results are summarized in Table 1. The goodness-of-fit was assessed using the Chi-squared test and the log-likelihood value. As discussed in Mitra and Washington (2007) and Lord and Park (2008), it is recommended to use more than one goodness-of-fit (GOF) measure for evaluating models and distributions. The log-likelihood value was calculated in two steps. First, the likelihood was calculated for each observation using the estimated parameters and pmf. Then, the natural logarithm of the likelihood was computed. For each dataset, the sum of the log-likelihood was calculated and compared for the Poisson, NB and NB–L distributions.

Table 1 shows that both the NB and the NB–L distributions provide adequate fit, but the NB–L offers an overall better statistical fit for all sample sizes with added zeros. Furthermore, the fit improves as the number or proportion of zeros increases. It should be pointed out that the fit using the log-likelihood shows that that all three distributions are very close. However, as discussed above about using different GOF measures, when a non-likelihood-based GOF criterion is used the NB–L shows a better fit. When 200 observations with zeros are added, the ratio between the two Chi-square values becomes larger than the ratio for the 150 added zeros.

As discussed in Miaou and Lord (2003) and Lord et al. (2005, 2007), the primary goal for analyzing distributions and regression models should not be solely based on finding the absolute best statistical fit. It is also very important to look at the data generating process, the relationship between the variables and whether the distribution or model is logically or theoretically sound. Miaou and Lord (2003) referred to the latter characteristic as "goodness-of-logic." Thus, although the NB–L offers a better statistical fit, it is still a rational distribution that can be used for characterizing crash data, given the fact that the long-term mean can never equal zero.

**Table 1**
Goodness-of-fit analysis – simulated results.

| Count | Observed frequency | Poisson distribution | NB distribution | NB–L distribution |
|---|---|---|---|---|
| Original data + 100 observations (% of zeros = 78.5) | | | | |
| 0 | 157 | 151.9 | 158.4 | 158.2 |
| 1 | 34 | 41.8 | 31.6 | 32.1 |
| 2 | 8 | 5.7 | 7.5 | 7.2 |
| 3+ | 1 | 0.6 | 2.5 | 2.5 |
| Parameters | | $\mu = 0.275$ | $\mu = 0.275\ \phi = 0.725$ | $\theta = 2.7113\ r = 6.925$ |
| Chi-square | | 2.77 | 0.30 | **0.17** |
| Log-likelihood | | −87.1970 | −87.1767 | **−87.1757** |
| Original data + 150 observations (% of zeros = 82.8) | | | | |
| 0 | 207 | 200.6 | 208.3 | 207.7 |
| 1 | 34 | 44.1 | 31.9 | 33.1 |
| 2 | 8 | 4.9 | 7.3 | 6.8 |
| 3+ | 1 | 0.4 | 2.5 | 2.4 |
| Parameters | | $\mu = 0.22$ | $\mu = 0.220\ \phi = 0.506$ | $\theta = 14.065\ r = 2.682$ |
| Chi-square | | 5.25 | 0.21 | **0.03** |
| Log-likelihood | | −108.9260 | −108.8877 | **−108.8858** |
| Original Data + 200 observations (% of zeros = 85.7) | | | | |
| 0 | 257 | 249.8 | 258.2 | 257.4 |
| 1 | 34 | 45.8 | 32.2 | 33.7 |
| 2 | 8 | 4.2 | 7.2 | 6.5 |
| 3+ | 1 | 0.3 | 2.5 | 2.4 |
| Parameters | | $\mu = 0.183$ | $\mu = 0.183\ \phi = 0.389$ | $\theta = 1.0889\ r = 1.659$ |
| Chi-square | | 7.85 | 0.16 | **0.01** |
| Log-likelihood | | −130.5944 | −130.5342 | **−130.5320** |

Bold value indicates best goodness-of-fit.

**Table 2**
Single-vehicle fatal crashes on divided multilane rural highways in Texas between 1997 and 2001.

| Crashes | Observed frequency | Poisson | NB | NB–L |
|---|---|---|---|---|
| 0 | 1532 | 1509.2 | 1534.4 | 1532.9 |
| 1 | 162 | 198.2 | 154.7 | 158.3 |
| 2 | 19 | 13.0 | 25.8 | 23.7 |
| 3 | 6 | 0.6 | 4.9 | 4.6 |
| 4+ | 2 | 0.0 | 1.2 | 1.4 |
| Parameters | | $\mu = 0.131$ | $\mu = 0.131\ \phi = 0.434$ | $\theta = 15.984\ r = 1.851$ |
| Chi-square | | 102.99 | 2.73 | **1.68** |
| Log-likelihood | | −715.1 | −696.1 | **−695.6** |

Bold value indicates best goodness-of-fit.

**Table 3**
Single-vehicle roadway departure crashes on rural two-lane horizontal curves in Texas between 2003 and 2008.

| Crashes | Observed frequency | Poisson | NB | NB–L |
|---|---|---|---|---|
| 0 | 29087 | 28471.6 | 29204.8 | 29133.6 |
| 1 | 2952 | 3918.0 | 2706.0 | 2855.5 |
| 2 | 464 | 269.6 | 567.4 | 503.1 |
| 3 | 108 | 12.4 | 141.1 | 120.9 |
| 4 | 40 | 0.4 | 37.8 | 35.9 |
| 5 | 9 | 0.0 | 10.6 | 13.1 |
| 6 | 5 | 0.0 | 3.0 | 3.3 |
| 7 | 2 | 0.0 | 0.9 | 3.3 |
| 8 | 3 | 0.0 | 0.3 | 0.0 |
| 9 | 1 | 0.0 | 0.1 | 0.0 |
| 10+ | 1 | 0.0 | 0.0 | 3.3 |
| Parameters | | $\mu = 0.138$ | $\mu = 0.138\ \phi = 0.284$ | $\theta = 9.212\ r = 1.018$ |
| Chi-square | | 2297.31 | 57.47 | **11.68** |
| Log-likelihood | | −14,208.1 | −13,557.7 | **−13,529.8** |

Bold value indicates best goodness-of-fit.

The next section describes the results for the observed data.

### 3.2. Observed data

Two datasets were used for this part of the analysis. The first one included single-vehicle fatal crashes that occurred on divided multilane rural highways between 1997 and 2001. The data were collected as a part of NCHRP 17-29 research project titled "*Methodology for estimating the safety performance of multilane rural highways*" (Lord et al., 2008). The data contained 1721 segments that varied from 0.10 mile to 11.21 miles, with an average equal to 1.01 miles. The sample mean was equal to 0.13. About 89% of the segments had no fatal crash. As discussed above, one way to reduce the number of zeros would be to change the spatial scale by aggregating shorter segments and create a new sample that contains longer segments. However, it may not always possible to aggregate the data due to loss of homogeneity of various geometric elements among different segments. For the purpose of this example, we will leave this dataset as is.

The second dataset included single-vehicle roadway departure fatal crashes that occurred on 32,672 rural two-lane horizontal curves between 2003 and 2008. The sample mean is equal to 0.14. For this dataset, about 90% of the data experienced no crash during the 5-year period. Given the proportion of zeros, one may believe that horizontal curves are very safe, although previous research indicate that they are not when compared to tangent segments (AAHSTO, 2010). The large number of zeros for this dataset can be explained by the sample that contains very short segments (about 90% are less than 0.3 miles). It would consequently be very difficult to change the spatial scale to reduce the number of zeros. One could perhaps increase the number of years to increase the sample size, if the data are available.

Tables 2 and 3 summarize the goodness-of-fit results for the two datasets. Similar to the results shown in the previous section, the two tables show that the NB–L distribution provides a better fit.

In summary, the NB–L provided a better fit than the traditional NB distribution for all simulated and observed datasets that included a large number of zeros. However, as pointed out by Zamani and Ismail (2010), the NB–L distribution works very well only when the dataset contains many observations with zero count. Based on the additional simulated data (not shown here), when the proportion of zeros is below 80%, the traditional NB distribution offers a performance that is equal to that of the NB–L. This is probably explained by the fact that the variance may not be large enough for the NB–L distribution to properly capture the variation. Since the distribution is also influenced by the length of the tail, it is suggested to evaluate both the NB and NB–L distributions when the percentage of zeros is above 70% (to be conservative) and select the one that provides the best goodness-of-fit statistic, as both are logically sound.

## 4. Further work

Since this is a newly introduced distribution, there are a lot of avenues for further work, both in terms of theoretical and practical applications. This paper focused on the distribution, but most of the research work in highway safety research is related to the development and application regression or statistical models that link crashes to covariates, such as traffic flow and geometric design characteristics. This means that the data should be analyzed in the context of a generalized linear model (GLM) similar to what Guikema and Coffelt (2008) and Sellers and Shmueli (2010) did for the Conway-Maxwell-Poisson (COM-Poisson) model.

Like the COM-Poisson distribution, the NB–L is a complex distribution and the pmf and likelihood need to be further manipulated in order to use it for modeling crash data. Unlike the NB model, the mean cannot be easily linked to the covariates when the MLE is used for estimating the coefficients (note: the centric parameter $\theta$ is not equal or the same as the mean, $\mu$, of the NB distribution). Furthermore, this type of distribution falls under the category of

hierarchical models, where the intermediate parameters need to be integrated separately (Booth et al., 2001). Although Zamani and Ismail (2010) have provided an approach to maximize the log-likelihood, Bayesian methods may be needed for developing the GLM. This is governed by the function that uses intermediate parameters; this is in addition to the limitations described at the end of Section 2.

Once the GLM is fully developed, the NB–L model should be evaluated and compared to other existing models that could potentially be used for handling data characterized by a large number of zeros. Those include random parameter (Anastasopoulos and Mannering, 2009) and zero-state Markov switching models (Mayshkina and Mannering, 2009) among others. This comparison will be the true test for the applicability of the NB–L for handling data with a lot of zeros. Finally, the NB–L model should also be evaluated for the small sample size and low sample mean problem (Lord, 2006).

## 5. Summary and conclusions

This paper has described the application of the NB–L distribution to datasets characterized by a large number of zeros and a heavy tail. Traditional statistical methods that have been proposed for analyzing such datasets have been found to suffer from important numerical and methodological problems. The newly introduced NB–L offers the advantage of being able to handle datasets with a large number of zeros, while still maintaining similar characteristics as the traditional NB. In other words, the NB–L is a two-parameter distribution and the long-term mean is never equal to zero. To examine this distribution, simulated and observed data were used. The results have shown that the NB–L always provided better statistical fit than the NB for datasets having a large number of zeros. However, since the distribution is also influenced by the length of the tail, it is suggested to evaluate both the NB and NB–L distributions and select the one that provides the best goodness-of-fit statistic. In conclusion, it is believed that this new distribution and subsequent GLM may offer a very useful tool for analyzing data characterized with a large amount of zeros.

## Acknowledgements

## References

AAHSTO, 2010. Highway Safety Manual, 1st ed. American Association of State Highway and Transportation Officials, Washington, DC.

Anastasopoulos, P.C., Mannering, F.L., 2009. A note on modeling vehicle accident frequencies with random-parameters count models. Accident Analysis and Prevention 41 (1), 153–159.

Booth, J.G., Casella, G., Fried, H., Hobert, J.P., 2001. Negative Binomial Loglinear Mixed Models. Working Paper. Department of Statistics, University of Florida, Gainesville, FL.

Casella, G., Berger, R.L., 1990. Statistical Inference. Wadsworth Brooks/Cole, Pacific Grove, CA.

Ghitany, M.E., Atieh, B., Nadarajah, S., 2008. Lindley distribution and its application. Mathematics and Computers in Simulation 78, 39–49.

Guikema, S.D., Coffelt, J.P., 2008. A flexible count data regression model for risk analysis. Risk Analysis 28 (1), 213–223.

Kumara, S.S.P., Chin, H.C., 2003. Modeling accident occurrence at signalized tee intersections with special emphasis on excess zeros. Traffic Injury Prevention 3 (4), 53–57.

Lindley, D.V., 1958. Fiducial distributions and Bayes' theorem. J.R. Stat. Soc. (20), 102–107. http://www.jstor.org/stable/2983909.

Lord, D., 2006. Modeling motor vehicle crashes using Poisson-gamma models: examining the effects of low sample mean values and small sample size on the estimation of the fixed dispersion parameter. Accident Analysis & Prevention 38 (4), 751–766.

Lord, D., Geedipally, S.R., Persaud, B.N., Washington, S.P., van Schalkwyk, I., Ivan, J.N., Lyon, C., Jonsson, T., 2008. Methodology for Estimating the Safety Performance of Multilane Rural Highways NCHRP Web-Only Document 126. National Cooperation Highway Research Program, Washington, DC.

Lord, D., Mannering, F.L., 2010. The statistical analysis of crash-frequency data: a review and assessment of methodological alternatives. Transportation Research - Part A 44 (5), 291–305.

Lord, D., Park, P.Y-J., 2008. Investigating the effects of the fixed and varying dispersion parameters of Poisson-gamma models on empirical Bayes estimates. Accident Analysis & Prevention 40 (4), 1441–1457.

Lord, D., Washington, S.P., Ivan, J.N., 2005. Poisson Poisson-gamma and zero inflated regression models of motor vehicle crashes: balancing statistical fit and theory. Accident Analysis & Prevention 37 (1), 35–46.

Lord, D., Washington, S.P., Ivan, J.N., 2007. Further notes on the application of zero inflated models in highway safety. Accident Analysis & Prevention 39 (1), 53–57.

Mayshkina, N.V., Mannering, F.L., 2009. Zero-state Markov switching count-data models: an empirical assessment. Accident Analysis and Prevention Vol. 42 (1), 122–130.

Miaou, S.-P., 1994. The relationship between truck accidents and geometric design of road sections: Poisson versus negative binomial regressions. Accident Analysis & Prevention 26 (4), 471–482.

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

Mitra, S., Washington, S., 2007. On the nature of over-dispersion in motor vehicle crash prediction models. Accident Analysis & Prevention 39 (3), 459–468.

Shankar, V., Milton, J., Mannering, F.L., 1997. Modeling accident frequency as zero-altered probability processes: an empirical inquiry. Accident Analysis & Prevention 29 (6), 829–837.

Shankar, V.N., Ulfarsson, G.F., Pendyala, R.M., Nebergal, M.B., 2003. Modeling crashes involving pedestrians and motorized traffic. Safety Science 41 (7), 627–640.

Sellers, K.F., Shmueli, G., 2010. A flexible regression model for count data. Annals of Applied Statistics 4, 943–961.

Warton, D.I., 2005. Many zeros does not mean zero inflation: comparing the goodness-of-fit of parametric models to multivariate abundance data. Environmetrics 16 (2), 275–289.

Zamani, H., Ismail, N., 2010. Negative binomial–Lindley distribution and its application. Journal of Mathematics and Statistics 6 (1), 4–9.