

Characteristics Based Heuristics to Select a Logical Distribution between the Poisson-Gamma and the Poisson-Lognormal for Crash Data Modeling

By:

Mohammadali Shirazi*

PhD. Candidate

Zachry Department of Civil Engineering

Texas A&M University, College Station, TX 77843, United States

Email: alishirazi@tamu.edu

and,

Dominique Lord, Ph.D.

Professor

Zachry Department of Civil Engineering

Texas A&M University, College Station, TX 77843, United States

Tel. (979) 458-3949

Email: d-lord@tamu.edu

Submission Date: August 1, 2017

Word Count: 4,568 (words) + 250×(5 Tables+5 Figures) = 7,068

*Corresponding Author

1 **ABSTRACT**

2 The Poisson-gamma (PG) and Poisson-lognormal (PLN) distributions are among the most popular
3 sampling distributions used by safety practitioners and analysts to model crash data. Several
4 studies have shown that the PLN offers a better alternative compared to the PG when data are
5 skewed while the PG is a more reliable option, otherwise. However, it is not explicitly clear when
6 the analyst needs to shift from the PG to the PLN – or vice versa, or what characteristics of data
7 should be observed a priori when deciding between these two alternative distributions. In addition,
8 in most research studies, the comparison between these two distributions or models and the
9 subsequent Model Selection decisions has usually been accomplished using the Goodness of Fit
10 (GoF) statistics or statistical tests. Such metrics rarely give any intuitions into why a specific
11 distribution or model is preferred over another (addressing the classical issue of Goodness-of-
12 Logic). This paper ponders into these topics by (1) designing characteristics based heuristics to
13 select a logical distribution between the PG and PLN, (2) prioritizing the most important
14 characteristics of the data under analysis to make a Model Selection decision between the PG and
15 the PLN. The proposed heuristics allows the analyst to select a logical distribution between the PG
16 and PLN, without any post-modeling efforts.

Keywords: Model Selection, Characteristics Based Heuristics, Classification, Poisson-gamma,
Poisson-lognormal

1 INTRODUCTION

2 Crash data modeling plays a pivotal role in most safety analyses or evaluations. Over last few
3 decades, safety scientists have placed significant efforts in introducing novel distributions or
4 models to study crash data (1, 2). However, among all potential modeling alternatives, the Poisson-
5 gamma (PG) and Poisson-lognormal (PLN) distributions still remain the most popular and
6 commonly used sampling distributions in the eyes of safety analysts and practitioners (1), mostly
7 due to their simplicity. Both of these distributions are classified as a member of the Poisson-
8 mixture family distributions, in which the Poisson distribution is mixed with another distribution
9 (known as a mixing distribution) to overcome the Poisson limitations in accounting over dispersion
10 or heterogeneity in data. In such mixture settings, it is assumed that the Poisson parameter is
11 randomly distributed by a logical mixing distribution. In the case of the PG mixture, the Poisson
12 parameter is distributed using a gamma distribution, while in the case of the PLN distribution, the
13 Poisson parameter is distributed using a lognormal distribution.

14 Although both are appropriate when data express a sign of over dispersion, each of these
15 distributions or models has its own positive and negative traits. As such, according to Lord and
16 Mannering (1), the PLN is more flexible than the PG to handle over dispersion and a better option
17 for modeling skewed data. In a more detailed examination of these two alternatives, Khazraee (3)
18 states that the thick tail of the lognormal distribution, theoretically, can give the PLN a substantial
19 boost when data are characterized by excessive large and/or unusual crash observations. The
20 comparison of the PG and PLN models is not limited to the safety literature. In a research that was
21 conducted to characterize the microbial counts in foods, Gonzales-Barron and Butler (4) showed
22 that the PLN is a better alternative when data include observations with large numbers, while the
23 PG outperforms the PLN for data with small count observations, and/or those with larger amount
24 of zero responses.

25 Overall, the previous studies indicate that the PLN is a better alternative for data with larger
26 skewness, and/or data that involve large count observations but fewer zero responses, while the
27 PG is a more suitable option for the opposing circumstances. However, it is not explicitly clear
28 when the analyst may need to switch from the PG to the PLN - or vice versa- and/or what
29 characteristics should be observed a priori to select a logical distribution between these two
30 alternatives. This paper addresses this topic and ponders into this issue by providing guidelines
31 and tools (or heuristics, to be exact) to select a logical distribution between the PG and PLN
32 distributions and recognizing the most important summary statistics to make a Model Selection
33 decision between these two sampling distributions.

34 Recently, Shirazi et al. (5) introduced a systematic methodology to design heuristics to
35 select the 'most-likely-true' and logical distribution among potential alternative distributions for
36 modeling count data. The authors demonstrated the application of the methodology by designing
37 heuristics to select a model between the Negative Binomial (NB) and Negative Binomial Lindley
38 (NB-L) (6, 7, 8) distributions. The proposed heuristics for the NB vs. NB-L recently has been
39 successfully examined in the context of Model Selection with covariates as well (9). The
40 methodology described in Shirazi et al.'s study (5) is used as a benchmark to address our topic. As

1 noted by Shirazi et al. (5), when designed, such heuristics have notable advantages to typical Model
2 Selection metrics, such as:

- 3 • Unlike the Goodness of Fit (GoF) metrics or typical statistical tests, these heuristics
4 examine the characteristics of data – addressing the classical issue of Goodness of Logic
5 (GoL) - for model recommendation.
- 6 • They can be used before fitting the distributions since only the characteristics of data, in
7 terms of the summary statistics, are considered to come up with the model
8 recommendation.
- 9 • They can be used as quick characteristics based guidelines for the safety analysts or
10 practitioners to select a model between the potential alternatives.
- 11 • The complexity of the potential alternatives is considered implicitly in such Model
12 Selection perspective.
- 13 • They can be used as quick heuristics when the analyst deals with high velocity of big data
14 and prompt Model Selection decisions are needed periodically.

15 The objectives of this paper consequently are: (1) provide simple guidelines or heuristics
16 to select a logically-sound distribution between the PG and PLN sampling distributions, given a
17 set of summary statistics of data, and (2) determine and prioritize the most important characteristics
18 of data, reflected into the summary statistics, to make a decision between these two distributions.
19 The objectives are accomplished by applying the two-steps, i.e.: (1) Monte Carlo simulations and
20 (2) Classifications, systematic methodology described by Shirazi et al. (5).

21 MIXED-POISSON FAMILY MODELS

22 Both of the PG and PLN distributions are classified as a member of the mixed-Poisson family
23 distributions, where the Poisson parameter is mixed with a distribution to accommodate the over-
24 dispersed data. The PG and PLN are two common models used to analyze crash data in safety
25 literature (1, 10, 11, 12). The characteristics of the PG and PLN distributions are described in this
26 section.

27 The probability mass function (pmf) of the Poisson distribution is defined as follows:

$$\text{Poisson}(\lambda) \equiv P(Y = y | \lambda) = \frac{\lambda^y \times e^{-\lambda}}{y!} \quad (1)$$

28 where the mean (m), variance (VAR) and variance-to-mean ratio (VMR) of the observations are
29 equal to:

$$E(y) = m = \lambda \quad (2a)$$

$$V(y) = \text{VAR} = \lambda \quad (2b)$$

$$\text{VMR}(y) = \text{VMR} = 1 \quad (2c)$$

1 The PG distribution is a mixture of the Poisson and gamma distributions, which can be
 2 structured as the following hierarchical representation:

$$y | \lambda \sim \text{Poisson}(\lambda) \quad (3a)$$

$$\lambda | \mu, \phi \sim \text{gamma}\left(\phi, \frac{\phi}{\mu}\right) \quad (3b)$$

3 The above mixture would result in a closed form NB distribution. The pmf of the NB distribution
 4 is defined as follows:

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

5 where μ = mean response of observations, and ϕ = inverse dispersion parameter. The mean (m),
 6 variance (VAR) and variance-to-mean ratio (VMR) of the PG distribution are defined as:

$$E(y) = \text{mean} = \mu \quad (5a)$$

$$V(y) = \text{VAR} = \mu + \frac{\mu^2}{\phi} \quad (5b)$$

$$\text{VMR}(y) = \text{VMR} = 1 + \frac{\mu}{\phi} \quad (5c)$$

7 The PLN distribution is a mixture of the Poisson and lognormal distributions, which can
 8 be structured as the following hierarchical representation:

$$y | \lambda \sim \text{Poisson}(\lambda) \quad (6a)$$

$$\log(\lambda) | \nu, \sigma^2 \sim \text{normal}(\nu, \sigma^2) \quad (6b)$$

10 Note that the mean (μ_λ) and variance (V_λ) of the lognormal distribution with parameters ν, σ^2 are
 11 equal to:

$$E(\lambda) = \mu_\lambda = e^{\nu + \sigma^2} \quad (7a)$$

$$\text{Var}(\lambda) = V_\lambda = \frac{e^{\sigma^2 - 1}}{e^{2\nu + \sigma^2}} \quad (7b)$$

12 Therefore, the mean (m), variance (VAR), and variance-to-mean ratio (VMR) of the PLN
 13 distribution are defined as:

$$E(y) = m = \mu_\lambda \quad (8a)$$

$$V(y) = \text{VAR} = \mu_\lambda + V_\lambda \quad (8b)$$

$$\text{VMR}(y) = \text{VMR} = 1 + \frac{V_\lambda}{\mu_\lambda} \quad (8c)$$

1 MODEL SELECTION HEURISTICS

2 Shirazi et al. (5) documented and discussed a systematic methodology to design simple
 3 *characteristics* based heuristics to predict the label of the most-likely-true distribution to model
 4 the data under analysis. In such perspective, the Model Selection problem is treated as a
 5 classification problem. The key to this approach are (1) simulating datasets that closely represent
 6 the population under consideration and recording the summary statistics of each dataset, and (2)
 7 training a classifier over the summary statistics to learn the patterns in the data to discriminate one
 8 distribution from another. For more information on rationales behind such Model Selection
 9 perspective and detailed steps of the methodology, the readers are referred to the work of Shirazi
 10 et al. (5).

11 This section is divided into three parts. First, the detailed steps of the simulation design are
 12 described. In the second part, a Decision Tree (DT) (13) classifier is used to design simple and
 13 straightforward heuristics to select a distribution for modeling between the PG and PLN
 14 distributions. The results of this section can be used as straightforward guidelines to select a logical
 15 distribution between these two alternatives. In the third part, a Random Forest (RF) (14) classifier
 16 is trained to design a more accurate Model Selection tool to predict the ‘most-likely-true’
 17 distribution between the PG and PLN distributions, as well as prioritizing the key summary
 18 statistics to discriminate these two distributions.

19 Simulation Design

20 Simulation is a key step in designing such Model Selection heuristics. It is essential to first make
 21 sure that the simulated datasets represent the characteristics of the target population, and then
 22 ensure that the alternative distributions have fair representations among simulated data (5). The
 23 first concern can be addressed by simulating data given the most common range observed in
 24 context population, in our case, the crash data population. The second concern can be addressed
 25 by ensuring that some summary statistics (referred to as control factors) are distributed similarly
 26 among the simulated datasets from alternative distributions (5). In other words, the analyst seeks
 27 to discriminate the distributions based on factors such as the ‘Kurtosis’ and/or ‘Skewness’, while
 28 the control factors such as the ‘mean’ or the ‘VMR’ are distributed similarly among simulated
 29 datasets.

30 In our problem design, we ensure that the ‘mean’ and the ‘VMR’ of data are uniformly
 31 distributed among the generated datasets from both of these distributions, simply, by simulating
 32 the mean (m) and the VMR from a uniform distribution with a range that is the most common
 33 observed range in crash data, as shown in Eqs. (9a) and (9b).

$$m \sim \text{uniform}(0.1, 20) \quad (9a)$$

$$\text{VAR} \sim \text{uniform}(1, 25) \quad (9b)$$

34 Next, given Eqs. (5a) and (5c), the parameters of the PG distribution can be estimated as:

$$\mu = m \quad (10a)$$

$$\phi = \frac{\mu}{\text{VMR} - 1} \quad (10b)$$

1 Similarly, given the Eqs. (8a) and (8c), first, we have:

$$\mu_\lambda = \mu \quad (11a)$$

$$V_\lambda = (\text{VMR} - 1) \times \mu_\lambda \quad (11b)$$

2 Then, given the Eqs. (7a) and (7b), the parameters of the PLN distribution can be derived as:

$$v = \log \left(\frac{\mu_\lambda^2}{\sqrt{V_\lambda + \mu_\lambda^2}} \right) \quad (12a)$$

$$\sigma = \sqrt{\log \left(\frac{V_\lambda}{\mu_\lambda^2} + 1 \right)} \quad (12b)$$

3 Now, it is possible to simulate a dataset with a size of $n=5,000$ from the PG distribution given
 4 parameters derived in Eq. 10, and from the PLN distribution given the parameters derived in Eq.
 5 12. The above procedure can be repeated for $N=100,000$ iterations, for each one of these
 6 distributions. Each time, m -types of summary statistics are recorded. We used 22 summary
 7 statistics in our analysis. These summary statistics include the mean (μ), variance (σ^2), standard
 8 deviation (σ), variance-to-mean ratio (VMR), coefficient-of-variation (CV), Skewness (skew),
 9 Kurtosis (K), percentage-of-zeros (Zeros), percentiles in 10% increments, the 10-th, 20-th, 30-th
 10 and 40-th inter-percentiles, and the range (R).

11 The detailed steps of the simulation protocol are described as follows:

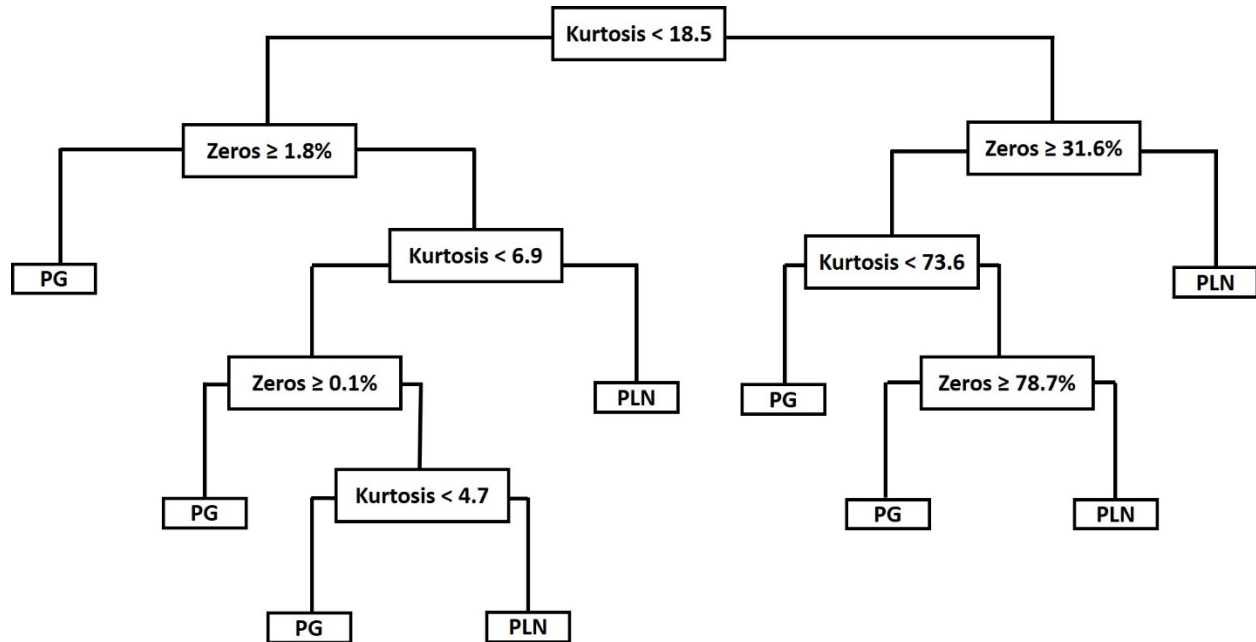
12 Repeat the following steps for 'N' iterations:

- 13 1. Simulate the mean (m) and the VMR from the Eqs. (8a) and (8b)
- 14 2. Find the parameters of the PG distribution from the Eqs. (9a) and (9b) and the PLN
 15 distribution from Eqs. (11a) and (11b).
- 16 3. Simulate a dataset with a size of 'n' given the parameters derived in Step 2, from both of
 17 the PG and the PLN distributions.
- 18 4. Record all the 22 types of summary statistics described above for the simulated datasets.

19 Decision Tree Heuristic

20 A Decision Tree classifier was used as a tool to partition the 22-dimensional predictor space that
 21 is created by the simulated summary statistics, and assign a label (either the PG or the PLN) to
 22 each partition. Fig. 1 shows the outcome of the Decision Tree classifier. As shown in Fig. 1, the
 23 population Kurtosis and the percentage-of-zeros play a substantial role in making a decision
 24 between the PG and PLN distributions. As seen in this figure, overall, the PLN is recommended
 25 for situations when data are more skewed but has fewer zero responses, while the PG distribution

1 is a better option otherwise; these results confirm the trends observed and/or reported in previous
 2 studies in the literature (1, 3, 4). Unlike previous studies, however, Fig. 1 provides a more
 3 perspicuous characteristics based guidance on selecting a sampling distribution between these two
 4 alternatives.



5
 6 **FIGURE 1: Characteristics Based Heuristic to select a Model between the PG and PLN Distributions.** (Tree
 7 can be used for data with the characteristics of $0.1 < \text{mean} < 20$ and $1 < \text{VMR} < 25$)

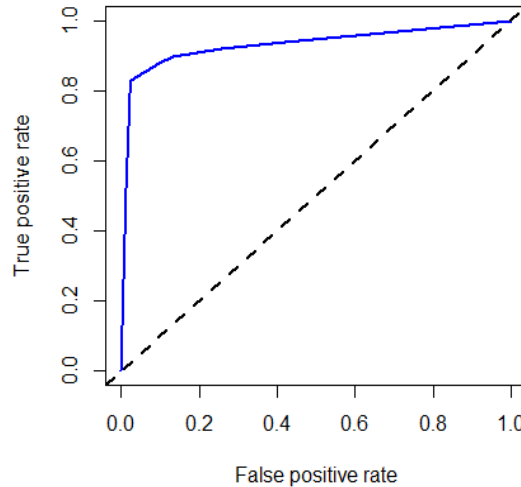
8 The output of a binary classifier can be either True (T) when it correctly classifies the label
 9 of the distribution, or False (F) when it misclassifies the label of the correct distribution. Let the
 10 PLN and PG distributions, respectively, be labeled as the positive (P) and negative (N) outputs of
 11 the binary classification. Such definitions represent a test when the analyst assumes the PG
 12 distribution as a base model, while he or she seeks to know when a shift to the PLN distribution is
 13 recommended. Table 1 shows the confusion matrix of the binary classification given such
 14 assumptions. The overall misclassification error is equal to 9.68% and the sensitivity [Note:
 15 $\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN})$] and specificity [Note: $\text{Specificity} = \text{TN} / (\text{TN} + \text{FP})$]
 16 are equal to 97.24% and 85.12%, respectively. The sensitivity of the classification is very high
 17 indicating that when the outcome of the binary classifier is the PLN distribution, there is a very
 18 high chance that the classifier has correctly detected the label of the distribution. However, the
 19 specificity of the classification is not as high as its sensitivity, meaning that when the outcome of
 20 the classifier is the PG distribution, there are still some chances that the output label was detected
 21 incorrectly. When the output of the classifier is the PG distribution, the analyst may consider other
 22 tests as well to decide between these two distributions and/or can decide to choose an alter
 23 tolerance threshold to decide between the PG and PLN. In the next section, we use a Random
 24 Forest classifier for a more accurate classification. Note that when the sample Kurtosis and the
 25 percentage-of-zeros deviate further away from the discriminating threshold, the ‘most-likely-true’

1 label can potentially be selected with greater confidence. Although not reported here, the DT
 2 heuristic was tested for simulated test data and the misclassification error was less than 10% for
 3 the test data.

4 **TABLE 1: PG vs. PLN: Confusion Matrix Based on the Results of the Decision Tree Classifier**

Predicted	Actual	
	PLN	PG
PLN	41.50% (TP)	1.18% (FN)
PG	8.50% (FP)	48.82% (TN)

5 The performance of a binary classifier can also be evaluated by the Receiver-Operating-
 6 Characteristics (ROC) plot (15, 16). The ROC plot is created by plotting the true positive rate
 7 (sensitivity) against the false positive rate (1-specificity) by varying the discriminating threshold.
 8 The performance of the classifier then is evaluated by the area under the curve (AUC) which can
 9 be any value from 0.5 to 1. A higher value represents a better classification. The ROC plot is shown
 10 in Fig. 2 and the value of the AUC is equal to 0.93.



11 **FIGURE 2: ROC Plot of the Classification between PG and PLN Based on the Decision Tree Results**

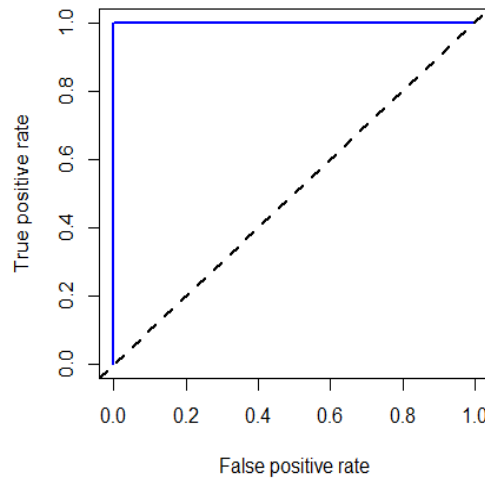
13 **Random Forest Heuristic**

14 Although they are easy to interpret and use, decision trees may not be as accurate as other
 15 classifiers (say Random Forest) and can be non-robust (15, 16). This means that a potential change
 16 in data could possibly result in altering in the final decision tree. The Random Forest classifier
 17 tries to overcome this issue by building many trees, instead of one, to substantially improve the
 18 performance of the classification (15, 16).

1 In our Random Forest classification, the number of trees was set to 100. Unlike the
 2 Decision Tree classification, the outcome of a Random Frost classification cannot be shown
 3 graphically. However, the trained forest can be recorded and still be used as an easy
 4 Characteristics-Based Model Selection tool to select a distribution between the PG and PLN
 5 distributions, without any post-molding efforts. Table 2 shows the confusion matrix of the binary
 6 classification between the PG and PLN, based on the results of the Random Forest classifier. The
 7 misclassification error is equal to 0.01%, and the sensitivity and specificity of the classifier are
 8 almost equal to 100%. The ROC plot is shown in Fig. 3 and the value of the AUC is almost equal
 9 to 1. Although not reported here, the Random Forest heuristic was tested for simulated test data
 10 and the misclassification error was less than 1.5% for the test data.

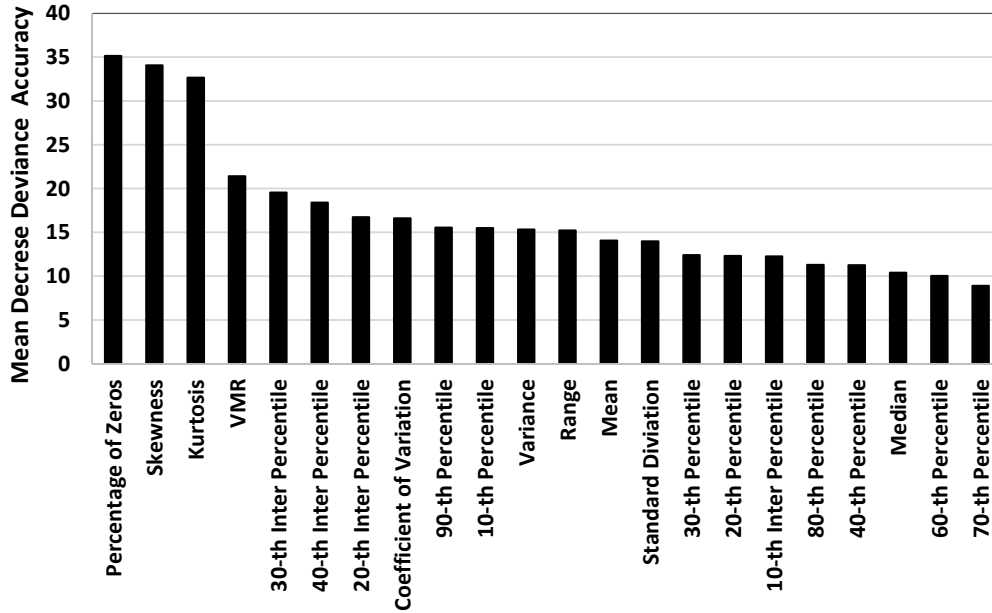
11 **TABLE 2: PG vs. PLN: Confusion Matrix Based on the Results of the Random Forest Classifier**

Predicted	Actual	
	PLN	PG
PLN	50.00% (TP)	0.01% (FN)
PG	0.00% (FP)	49.99% (TN)



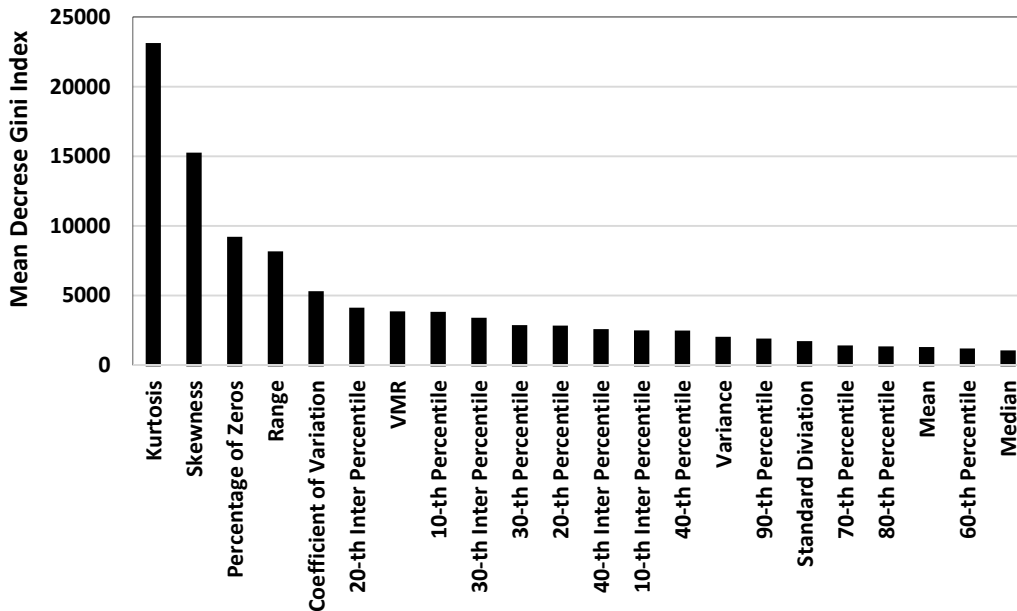
12
 13 **FIGURE 3: ROC Plot of the Classification between the PG and PLN based the Random Forest Results**

14 As a byproduct of the Random Forest classifier, the predictors (summary statistics) can be
 15 ranked by their importance. Fig. 4 and 5 show the importance of the summary statistics based on
 16 two criteria: (1) mean decrees Deviance Accuracy and (2) mean decrease Gini index (15, 16). As
 17 shown in these figures, Kurtosis, Skewness and the percentage-of-zeros are among the most
 18 important summary statistics to select a model between the PG and PLN distributions.



1

2 **FIGURE 4: Importance of the Summary Statistics to Select a Distribution between the PG and PLN Based on**
 3 **the Mean Decrease Deviance Accuracy Given the Results of the Random Forest Classifier**



4

5 **FIGURE 5: Importance of the Summary Statistics to Select a Distribution between the PG and PLN Based on**
 6 **the Mean Decrease Gini, Given the Results of the Random Forest Classifier**

7 **APPLICATION OF PROPOSED HEURISTICS TO OBSERVED DATA**

8 In this section, two datasets are used to evaluate the proposed heuristics. The first dataset includes
 9 information related to single-vehicle crashes that occurred on Michigan rural two-lane highway in
 10 2006. This dataset was utilized in several previous studies (17, 7, 8). The dataset includes 33,970

1 segments, and the mean, variance, VMR, Kurtosis, and the percentage-of-zeros of data are equal
 2 to: 0.68, 3.15, 4.62 123.6 and 69.7%, respectively. The second dataset contains crash data that
 3 occurred between 2012 and 2014 on Texas urban four-lane arterials. This dataset also has been
 4 used in several studies (18, 19, 20) in the past. The dataset includes 4,264 segments, and the mean,
 5 variance, VMR, Kurtosis, and the percentage-of-zeros of data are equal to: 2.26, 45.53, 19.27, 92.8
 6 and 56.5%, respectively. The detailed summary statistics of the two datasets are shown in Table
 7 3.

8

TABLE 3: Summary Statistics of the Datasets

Summary Statistics	Michigan Dataset	Texas Dataset
Mean	0.68	2.36
Variance	3.15	45.53
Standard Deviation (Sd.)	1.77	6.75
Variance-to-Mean-Ratio (VMR)	4.62	19.27
Coefficient-of-Variation (CV)	2.60	2.86
Skewness (skew)	7.76	7.92
Kurtosis (K)	123.59	92.67
Percentage-of-Zeros (Z)	69.6%	56.5%
10-th Percentile	0	0
20-th Percentile	0	0
30-th Percentile	0	0
40-th Percentile	0	0
50-th Percentile (Median)	0	0
60-th Percentile	0	1
70-th Percentile	1	1
80-th Percentile	1	3
90-th Percentile	2	6
10-th Inter-Percentile	1	1
20-th Inter-Percentile	1	1
30-th Inter-Percentile	1	3
40-th Inter-Percentile	2	6
Range	61	120

9 Table 4 and 5, respectively, show the recommended models for the Michigan and Texas
 10 data based on the proposed heuristics and the log-likelihood metric. While the classical metrics
 11 require the distributions to be fitted to the data before coming up with the model recommendation,
 12 the proposed heuristics can be used without any post-modeling inputs and/or efforts. The decision
 13 based on the proposed heuristics solely rely on characteristics of data. For both datasets, the PLN
 14 distribution is the favored distribution to model data, based on the classical log-likelihood metric
 15 and the proposed heuristics. Classical metrics, such as the log-likelihood, do not give any intuitions
 16 into why the PLN is preferred to the PG (addressing the Goodness-of-Logic issue). On the other
 17 hand, the proposed heuristics come up with the model recommendation by considering the
 18 characteristics of data; hence, in this case, the analyst can select a logical distribution to model
 19 data. For example, a large kurtosis value in both datasets plays a substantial role in choosing the
 20 PLN over the PG.

1 **TABLE 4: Model Selection for the Michigan Data.**

Method	PG	PLN	Criteria	Favored Distribution
Log-Likelihood (LL) ¹	-36332.85 ($\phi = 0.30, \mu = 0.68$)	-36117.54 ($\nu = -1.48, \sigma = 1.50$)	$LL_{PLN} > LL_{PG}$	PLN
Decision Tree Heuristic ²	Kurtosis= 123.6 Zeros=69.7%		Kurtosis > 73.6 Zeros < 78.7%	PLN
Random Forest Heuristic ²	Using All 22 Summary Statistics		Using the RF Heuristic	PLN

2 ¹Requires fitting the distributions.3 ²Do not require fitting the distributions4 **TABLE 5: Model Selection for the Texas Data.**

Method	PG	PLN	Criteria	Favored Distribution
Log-Likelihood (LL) ¹	-7462.91 ($\phi = 0.23, \mu = 2.36$)	-7432.35 ($\nu = -0.82, \sigma = 1.95$)	$LL_{PLN} > LL_{PG}$	PLN
Decision Tree Heuristic ²	Kurtosis= 92.8 Zeros= 56.5%		Kurtosis > 73.6 Zeros < 78.7%	PLN
Random Forest Heuristic ²	Using All 22 Summary Statistics		Using the RF Heuristic	PLN

5 ¹Requires fitting the distributions.6 ²Do not require fitting the distributions7 **SUMMARY AND CONCLUSIONS**

8 The Poisson-gamma and Poisson-lognormal are the most popular sampling distributions used in
9 safety analyses and evaluations as a means to analyze crash data. This study investigated, in details,
10 under what circumstances the PLN is preferred over the PG, and vice versa, based on
11 characteristics of data, reflected in the summary statistics. A decision tree was constructed and
12 proposed as quick guidelines to select a distribution between these two alternatives. The Kurtosis
13 and percentage-of-zeros were the only summary statistics used by the classifier in the decision
14 tree. Although Decision Tree classifiers are non-robust and potentially provide different tree splits,
15 the results shown in Fig.1 can be used by practitioners as useful guidelines for selecting a sampling
16 distribution between the PG and PLN. We used a Random Forest classifier to design a more
17 accurate tool to select a distribution among these two options. As a byproduct of a Random Forest
18 classifier, the summary statistics can be ranked by their importance. Among the 22 types of
19 summary statistics used in the analysis, Kurtosis, Skewness and the percentage-of-zeros were
20 found the most important and critical summary statistics to select a Model between the PG and
21 PLN. The next step should compare the PLN and NB-L, to decide when the percentage-of-zeros
22 favors a model over the other. Further analysis in context of heuristics is also needed to consider
23 the effect of the sample-size (21, 22, 23) on proposed heuristics.

24 **Acknowledgments**

25 The authors would like to thank the Safe-D UTC program for the support obtained for this
26 research. The opinions expressed by the authors in this research do not necessarily reflect those
27 from the Safe-D UTC program. We also would like to thank Dr. Soma Dhavala for his comments
28 and sharing his valuable insights with us.

REFERENCES

- 1 [1] Lord, D., and Mannering, F. The statistical analysis of crash-frequency data: a review and assessment of
2 methodological alternatives. *Transportation Research Part A: Policy and Practice*, 44(5), 2010, pp. 291-305.
- 3 [2] Mannering, F. L., and Bhat, C. R. Analytic methods in accident research: methodological frontier and future
4 directions. *Analytic Methods in Accident Research*, 1, 2014, pp. 1-22.
- 5 [3] Khazraee S.H. Full Bayesian Poisson-hierarchical models for crash data analysis: investigating the impact of model
6 choice on site-specific predictions. PhD Dissertation. Department of Civil Engineering, Texas A&M University,
7 College Station, Texas, 2016.
- 8 [4] Gonzales-Barron, U., & Butler, F. A comparison between the discrete Poisson-gamma and Poisson-lognormal
9 distributions to characterise microbial counts in foods. *Food Control*, 22(8), 2011, pp. 1279-1286.
- 10 [5] Shirazi, M., Dhavala, S.S., Lord, D., Geedipally, S.R. A methodology to design heuristics for model selection
11 based on characteristics of data: application to investigate when the negative binomial Lindley (NB-L) is preferred
12 over the negative binomial (NB). *Accident Analysis & Prevention*, 2017, in press. (<http://dx.doi.org/10.1016/j.aap.2017.07.002>).
- 13 [6] Lord, D., and Geedipally, S. R. The negative binomial–Lindley distribution as a tool for analyzing crash data
14 characterized by a large amount of zeros. *Accident Analysis & Prevention*, 43(5), 2011, pp. 1738-1742.
- 15 [7] Geedipally, S. R., Lord, D., and Dhavala, S. S. The negative binomial-Lindley generalized linear model:
16 Characteristics and application using crash data. *Accident Analysis & Prevention*, 45, 2012, pp. 258-265.
- 17 [8] Shirazi, M., Lord, D., Dhavala, S.S., Geedipally, S.R. A semiparametric negative binomial generalized linear
18 model for modeling over dispersed count data with a heavy tail: characteristics and applications to crash data. *Accident*
19 *Analysis & Prevention*, 91, 2016, pp. 10-18.
- 20 [9] Shaon, M.R.R, Qin, X., Shirazi, M., Lord, D., Geedipally, S.R. A random parameters negative binomial-Lindley
21 generalized linear model to analyze over-dispersed data. Working paper to be submitted for publication.
- 22 [10] Lord, D., and Miranda-Moreno, L.F. Effects of Low Sample Mean Values and Small Sample Size on the
23 Estimation of the Fixed Dispersion Parameter of Poisson-gamma Models for Modeling Motor Vehicle Crashes: A
24 Bayesian Perspective. *Safety Science*, 46 (5), 2008, pp. 751-770.
- 25 [11] Aguero-Valverde, J., Jovanis, P.P. Analysis of road crash frequency with spatial models. *Transportation Research*
26 *Record*, 2061, 2008, 55–63.
- 27 [12] Aguero-Valverde, J. Full Bayes Poisson gamma, Poisson lognormal, and zero inflated random effects models:
28 Comparing the precision of crash frequency estimates. *Accident Analysis & Prevention*, 50, 2013, 289-297.
- 29 [13] Breiman, L, Friedman, J. H., Olshen, R. A., Stone, C. J. Classification and regression trees. Monterey, CA:
30 Wadsworth & Brooks/Cole Advanced Books & Software. ISBN 978-0-412-04841-8, 1984.
- 31 [14] Breiman, L. Random forests. *Machine learning*, 45(1), 2001, pp. 5-32.
- 32 [15] Hastie, T., Tibshirani, R. and Friedman, J. The elements of statistical learning. *NY Springer*, 2001
- 33 [16] James, G., Witten, D., Hastie, T., and Tibshirani, R. An introduction to statistical learning (Vol. 6). New York:
34 springer. 2013.
- 35 [17] Qin, X., Ivan, J.N., Ravishanker, N. Selecting exposure measures in crash rate prediction for two-lane highway
36 segments. *Accident Analysis and Prevention* 36 (2), 2004, pp. 183-191.
- 37 [18] Lord, D., Geedipally, S.R., Shirazi, M. Improved Guidelines for Estimating the Highway Safety Manual
38 Calibration Factors. ATLAS-2015-10, 2016.
- 39 [19] Geedipally, S. R., Shirazi, M., & Lord, D. Exploring the Need for Having Region-Specific Calibration Factors
40 *Transportation Research Record*, 2017, in press.
- 41 [20] Shirazi, M., Geedipally, S. R., & Lord, D. A procedure to determine when safety performance functions should
42 be recalibrated. *Journal of Transportation Safety & Security*, 9(4), 2017, pp. 457-469.

- 1 [21] Lord, D. Modeling motor vehicle crashes using Poisson-gamma models: Examining the effects of low sample
2 mean values and small sample size on the estimation of the fixed dispersion parameter. *Accident Analysis &*
3 *Prevention*, 38(4), 2006, 751-766.
- 4 [22] Shirazi, M., Lord, D., & Geedipally, S. R. Sample-size guidelines for recalibrating crash prediction models:
5 recommendations for the Highway Safety Manual. *Accident Analysis & Prevention*, 93, 2016, pp. 160-168.
- 6 [23] Shirazi, M., Geedipally, S. R., & Lord, D. A Monte-Carlo simulation analysis for evaluating the severity
7 distribution functions (SDFs) calibration methodology and determining the minimum sample-size requirements.
8 *Accident Analysis & Prevention*, 98, 2017, pp. 303-311.