

Modeling nonlinear relationship between crash frequency by severity and contributing factors by neural networks

Qiang Zeng^{a, b}, Helai Huang^{b, *}, Xin Pei^c, S.C. Wong^d

^a School of Civil Engineering and Transportation, South China University of Technology, Guangzhou, Guangdong, 510641 P.R. China

^b Urban Transport Research Center, School of Traffic and Transportation Engineering, Central South University, Changsha, Hunan, 410075 P.R. China

^c Department of Automation, Tsinghua University, Beijing, P.R. China

^d Department of Civil Engineering, The University of Hong Kong, Pokfulam Road, Hong Kong

ABSTRACT

This study develops neural network models to explore the nonlinear relationship between crash frequency by severity and risk factors. To eliminate the possibility of over-fitting and to deal with black-box characteristic, a network structure optimization and a rule extraction method are proposed. A case study compares the performance of the modified neural network models with that of the traditional multivariate Poisson-lognormal model for predicting crash frequency by severity on road segments in Hong Kong. The results indicate that the trained and optimized neural networks have better fitting and predictive performance than the multivariate Poisson-lognormal model. Moreover, the smaller differences between training and testing errors in the optimized neural networks with pruned input and hidden nodes demonstrate the ability of the structure optimization algorithm to identify insignificant factors and to improve the model's generalizability. Furthermore, two rule-sets are extracted from the optimized neural networks to explicitly reveal the exact effect of each significant explanatory variable on the crash frequency by severity under different conditions. The rules imply that there is a nonlinear relationship between risk factors and crash frequencies with each injury-severity outcome. With the structure optimization algorithm and rule extraction method, the modified neural network models have great potential for modeling crash frequency by severity, and should be considered a good alternative for road safety analysis.

Keywords: crash frequency by severity; neural network; over-fitting; structure optimization; rule extraction.

* Corresponding author

E-mail address: 641459622@qq.com (Q. Zeng), huanghelai@csu.edu.cn (H. Huang), peixin@mail.tsinghua.edu.cn (X. Pei), hhecwsc@hku.hk (S.C. Wong)

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

1. Introduction

In the past decade, there has been a quantity of research on predicting crash frequency by certain categories, such as injury severity (e.g., property damage only, possible injury, non-incapacitating injury, incapacitating injury or fatality) (Park and Lord, 2007), the number of vehicles involved (e.g., single vehicle, two vehicles, or three or more vehicles) (Venkataraman et al., 2013) or collision type (e.g., angle, head-on, rear-end, sideswipe or pedestrian-involved) (Ye et al., 2009). The first kind of classification covers most concerns, because crash injury severity is an important aspect in assessing safety performance, in addition to the crash frequency (AASHTO, 2010). Compared with conventional crash prediction models (referred to as “safety performance functions”), modeling crash frequency by severity identifies the effects of observed risk factors (such as the traffic, geometrical and environmental characteristics of sites) on the frequency of accidents with a particular injury-severity outcome. The expected crash frequencies at each level of severity provide deeper insights on the safety situation of a certain road entity (road segment, intersection, etc.). Therefore, while crash totals may not reveal a site deficiency, over exposure of a specific crash severity may uncover otherwise undetected deficiencies. Moreover, the models have been employed to rank road sites with promise for safety improvement, a critical step of network screening in the roadway safety management process (AASHTO, 2010), as injury severity and its associated costs are primary concerns in many programs (Miaou and Song, 2005).

Methodologically, there are mainly two groups of approaches to crash frequency by severity prediction: joint and separate modeling. In the former group, correlation between crash frequencies at various severity levels is the most important issue. To deal with it, a series of techniques have been investigated, such as multivariate regression models (Aguero-Valverde and Jovanis, 2009; Anastasopoulos et al., 2012; Barua et al., 2014, 2016; Bijleveld, 2005; El-Basyouny and Sayed, 2009; El-Basyouny et al., 2014; Ma and Kockelman, 2006; Ma et al., 2008; Park and Lord, 2007), simultaneous equations (Ye et al., 2009, 2013), a joint-probability approach (Pei et al., 2011), two-stage bivariate/multivariate models (Wang et al., 2011; Xu et al., 2014) and multinomial-generalized Poisson models (Chiou and Fu, 2013, 2015; Chiou et al., 2014). The multivariate Poisson regression proposed by Ma and Kockelman (2006) adds a common error term into the Poisson distributions of univariate regressions to account for their correlation, but it does not allow for the commonly observed over-dispersion, and it assumes the identical and positive covariances across crash frequencies (Park and Lord, 2007). In order to improve it, a multivariate Poisson-lognormal regression has been developed (Ma et al., 2008), which is able to accommodate over-dispersion and provides a fully general covariance structure. To account for the spatial correlation among neighboring sites, error terms with Gaussian conditional auto-regressive distribution have been introduced into the multivariate

1 Poisson-lognormal model (Barua et al., 2014). Based on it, Barua et al. (2016) have
2 proposed a multivariate random parameters count model to further capture
3 unobserved heterogeneity across observations.

4 Compared with multivariate regression models, the formulation of simultaneous
5 equations, the joint probability model and the two-stage bivariate/multivariate models
6 are less complicated (Pei et al., 2011; Wang et al., 2011; Xu et al., 2014; Ye et al.,
7 2009, 2013). Besides, the computation burden of simultaneous equations is lighter,
8 because their coefficients are calibrated by a simulated likelihood estimation method
9 (Ye et al., 2009), while the others are calibrated by Markov chain Monte Carlo
10 simulation, a typical Bayesian inference method. On the contrary, the
11 multinomial-generalized Poisson models (Chiou and Fu, 2013; Chiou et al., 2014),
12 especially the extension with accommodating spatio-temporal dependence (Chiou and
13 Fu, 2015), are even more complicated than multivariate count models.

14 Although the correlation across severity levels is significant in many studies, the
15 advantage of joint modeling over separate modeling is not “theoretical” but rather
16 “empirical”, as noted by Ma et al. (2008). In the comparative analysis conducted by
17 Lan and Persaud (2012), univariate models are found to fit the crash data better than
18 the multivariate model. Consequently, some researchers continue to separately model
19 crash frequencies at each severity level. For example, Venkataraman et al. (2013)
20 advocate univariate random parameter models to individually predict crash frequency
21 by severity, or other aggregation types, by accounting for heterogeneities across
22 unobserved or unobservable factors. All of the above-mentioned models are based on
23 a generalized linear function framework and certain assumed distributions of crash
24 data. However, in some cases, these assumptions may be violated and thereby result
25 in biased inferences (Li et al., 2008).

26 Relative to the statistical models, without any prior knowledge or assumption on
27 model structure, some artificial intelligence models can be used to approximate the
28 underlying nonlinear relationship between crash frequency by severity and safety
29 predictors (Haykin, 2009). As a common class of artificial intelligence models, neural
30 network models have been successfully used in many fields of transportation research
31 (Karlaftis and Vlahogianni, 2011). For highway safety analysis, a number of studies
32 have investigated the performance of neural network models in predicting crash
33 frequency or injury severity (Abdelwahab and Abdel-Aty, 2001; Chang, 2005; Huang
34 et al., 2016; Zeng and Huang, 2014b). The results show that neural network models
35 outperform some traditional statistical models, such as the negative binomial model of
36 crash frequency prediction and the ordered logit/probit models of crash injury severity
37 prediction. To the best of our knowledge, neural networks have not yet been employed
38 to predict crash frequency by severity.

39 Moreover, with the development of neural network techniques, the commonly
40 criticized weaknesses of crash prediction, the over-fitting problem and the black-box
41 characteristic, have been mostly eliminated. Advanced methods for network training
42 and structure optimization can establish generalized neural network models that

effectively approximate the relationship between crash frequency by severity and explanatory variables (Haykin, 2009). In addition, piecewise linear rules extracted from the developed neural networks are able to clearly illustrate the effects of risk factors (Setiono and Thong, 2004).

In summary, this study attempts to develop advanced neural networks for modeling the nonlinear relationship between crash frequency by severity and risk factors, and to clarify the effects of factors on the outcomes by extracting rules from the developed neural networks. To demonstrate the proposed methods, the neural network models are compared with the multivariate Poisson-lognormal model with regard to fitting and predictive performance. Accordingly, the remainder of this paper is organized as follows. The next section specifies the proposed models and methods. The collected data for model demonstration are described in Section 3. Section 4 introduces the detailed implementation of the proposed models and discusses the results. Finally, conclusions and recommendations for future research are presented in Section 5.

2. Methodology

The multivariate Poisson-lognormal model, one of the most widely used statistical models for jointly predicting crash frequency and severity, is used as a benchmark in this study to compare its fitting and predictive performance with those of the proposed neural network models. In this section, the model architectures of the multivariate Poisson-lognormal and neural network models are specified. Then, the training, structure optimization, and rule extraction algorithms for the neural network models are described.

2.1. Model specification

2.1.1. Multivariate Poisson-lognormal model

In the multivariate Poisson-lognormal model, the crash count Y_{its} at site i during period t at injury severity degree s is assumed to follow a Poisson distribution (Ma et al., 2008), given λ_{its} , that is,

$$P(Y_{its} = y_{its} | \lambda_{its}) = \lambda_{its}^{y_{its}} \frac{e^{-\lambda_{its}}}{y_{its}!},$$

$$i = 1, 2, \dots, N, \quad t = 1, 2, \dots, T, \quad s = 1, 2, \dots, S, \quad y_{its} = 0, 1, 2, 3, \dots, \quad (1)$$

where N , T , and S are the number of observed sites, the periods and the categorized injury severity levels, respectively. The mean of Y_{its} , λ_{its} , is assumed to

1 have a generalized linear relationship with the explanatory variables, \mathbf{X}_{it} , such that

$$2 \quad \ln \lambda_{its} = \mathbf{X}'_{it} \boldsymbol{\beta}_s + \varepsilon_{its}, \quad (2)$$

3 in which $\boldsymbol{\beta}_s$ are the coefficients to be estimated. The error term ε_{its} accommodates
 4 the crash severity correlation and the common over-dispersion, which is
 5 multi-normally distributed as

$$6 \quad \boldsymbol{\varepsilon}_{it} \sim N_S(\mathbf{0}, \boldsymbol{\Sigma}), \quad \boldsymbol{\varepsilon}_{it} = \begin{pmatrix} \varepsilon_{it1} \\ \varepsilon_{it2} \\ \dots \\ \varepsilon_{its} \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1s} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2s} \\ \dots & \dots & \dots & \dots \\ \sigma_{s1} & \sigma_{s2} & \dots & \sigma_{ss} \end{pmatrix}. \quad (3)$$

7

8 2.1.2. Neural network model

9 Neural network models are information-processing mechanisms that are inspired
 10 by biological nervous systems (Haykin, 2009). Various categories of neural network
 11 models have been developed with different network architectures, such as the radial
 12 basis function and the self-organizing feature map. The multilayer perceptron, which
 13 is known as a universal approximator and is the most popular neural network for data
 14 mining, is used to model the underlying nonlinear relationship between crash
 15 frequency by severity and risk factors in this study. Although they can be predicted
 16 simultaneously in the same neural network, the crash frequencies at each injury
 17 severity are modeled separately herein to identify their respective pertinent predictors.
 18 Fig. 1 shows the structure of the developed multilayer perceptrons with fully
 19 connected neurons.

20 Consider a dataset containing N_1 continuous attributes and N_2 categorical
 21 attributes that may affect the crash frequency at the severity level $s(s=1,2,\dots,S)$. As
 22 in many statistical modeling methods, each categorical attribute $A_n(n=1,2,\dots,N_2)$
 23 is transformed into m_n-1 binary attribute(s) $a_1^n, \dots, a_j^n, \dots, a_{m_n-1}^n$, where m_n is the
 24 number of possible values for A_n . $a_j^n=1$ if A_n is equal to category j ; and
 25 $a_j^n=0$ otherwise. Each of the transformed attributes, together with the continuous
 26 attributes, is represented by a node $x_i(i=2,\dots,I)$ in the input layer. In addition, an
 27 input node, with $x_1=1$, is added. The weights of its connections with hidden neurons

1 are the biases. Therefore, the number of units I in the input layer is

$$2 \quad I = N_1 + \sum_{n=1}^{N_2} (m_n - 1) + 1. \quad (4)$$

3 To fit the training data well, the number of neurons in the hidden layer must be
 4 sufficiently large. If it is assumed to be J , then the connection weight between
 5 hidden node $j(j=1, \dots, J)$ and input node $i(i=1, \dots, I)$ is $w_{j,i}^{(1)}$. The hyperbolic
 6 function, $\tanh(\cdot)$, which is an odd sigmoid transfer function, is used for all of the
 7 hidden nodes. In the output layer, the only unit, ψ , represents the expected crash
 8 frequency at severity level s . $w_j^{(2)}$ denotes the weight of the connection between the
 9 output node and the hidden node $j(j=1, \dots, J)$. A linear function is employed as the
 10 transfer function for the output node. Then, the expected crash frequency at severity
 11 level s is given by,

$$12 \quad \psi = \sum_{j=1}^J w_j^{(2)} \tanh\left(\sum_{i=1}^I w_{j,i}^{(1)} x_i\right). \quad (5)$$

13

14 2.2. Network training

15

16 The conjugate gradient algorithm, which has a better learning performance than
 17 the popular back-propagation algorithm (Haykin, 2009), is adopted in this study to
 18 train the neural networks. For the collected samples $\{\mathbf{x}(m), o(m) | m=1, 2, \dots, M\}$ of
 19 severity $s(s=1, 2, \dots, S)$, where $\mathbf{x}(m)$ and $o(m)$ are a vector of risk factors and the
 20 corresponding observed crash frequency at severity s , respectively, and M is the
 21 number of samples, the conjugate gradient updates the connection weight vector \mathbf{w}
 22 as follows:

$$23 \quad \begin{aligned} \mathbf{w} &= (w_1, \dots, w_{(j-1)I+i}, \dots, w_{JI}, w_{JI+1}, \dots, w_{JI+j}, \dots, w_{J(1+I)}) \\ &= (w_{1,1}^{(1)}, \dots, w_{j,i}^{(1)}, \dots, w_{J,I}^{(1)}, w_1^{(2)}, \dots, w_j^{(2)}, \dots, w_J^{(2)}) \end{aligned}$$

24 1. Randomly select $w_{j,i}^{(1)}(j=2, \dots, J; i=1, \dots, I)$ and $w_j^{(2)}(j=1, \dots, J)$ from two
 25 uniform distributions. The means of both distributions are equal to 0, and their
 26 variances are $1/J$ and 1, respectively. Set the initial iteration $t=0$.

27 2. According to weight vector $\mathbf{w}(0)$, calculate the expected network outputs,

28 $\psi_m(m=1, 2, \dots, M)$, the derivative of outputs on all weights,

1 $\frac{\partial \psi(m)}{\partial \mathbf{w}(0)}$ ($m=1,2,\dots,M$), and the gradient vector, $\mathbf{g}(0)$:

$$2 \quad \frac{\partial \psi(m)}{\partial w_h} = \begin{cases} \frac{\partial \psi(m)}{\partial w_j^{(2)}} = \tanh\left(\sum_{i=1}^I w_{j,i}^{(1)} x_i(m)\right), & \text{if } w_h = w_j^{(2)} \\ \frac{\partial \psi_k(m)}{\partial w_{j,i}^{(1)}} = w_j^2 \tanh'\left(\sum_{i=1}^I w_{j,i}^{(1)} x_i(m)\right) x_i(m), & \text{if } w_h = w_{j,i}^{(1)} \end{cases}, \quad (6)$$

$$3 \quad \mathbf{g}(t) = \frac{1}{N} \sum_{m=1}^M [o(m) - \psi(m)] \frac{\partial \psi(m)}{\partial \mathbf{w}(t)}. \quad (7)$$

4 3. Set $\mathbf{s}(0) = \mathbf{r}(0) = -\mathbf{g}(0)$.

5 4. In iteration t , for the fixed $\mathbf{w}(t)$ and $\mathbf{s}(t)$, use the advance-and-retreat method
6 to linearly search the optimal $\eta(t)$ by minimizing the cost function,

7 $\xi_{av}(\mathbf{w}(t) + \eta \mathbf{s}(t))$:

$$8 \quad \xi_{av}(\mathbf{w}) = \frac{1}{2M} \sum_{m=1}^M [o(m) - \psi(m)]^2. \quad (8)$$

9 5. Check the convergence criteria. If the Euclidean norm of $\mathbf{r}(t)$ decreases to a
10 certain small portion, ε , of its initial value, $\|\mathbf{r}(0)\|$, or the iteration number meets
11 its maximum value, T , the algorithm is done:

$$12 \quad \|\mathbf{r}(t)\| \leq \varepsilon \|\mathbf{r}(0)\|, \text{ or } t = T.$$

13 6. Update the connection weight vector:

$$14 \quad \mathbf{w}(t+1) = \mathbf{w}(t) + \eta(t) \mathbf{s}(t). \quad (9)$$

15 7. Calculate the gradient vector $\mathbf{g}(t+1)$ by formulas (6)-(7) according to $\mathbf{w}(t+1)$.

16 Set $\mathbf{r}(t+1) = -\mathbf{g}(t+1)$.

17 8. Calculate $\beta(t+1)$ by the Polak-Ribiere method:

$$18 \quad \beta(t+1) = \max\left\{\frac{\mathbf{r}'(t+1)(\mathbf{r}(t+1) - \mathbf{r}(t))}{\mathbf{r}'(t)\mathbf{r}(t)}, 0\right\}. \quad (10)$$

19 9. Update the direction vector:

$$1 \quad \mathbf{s}(t+1) = \mathbf{r}(t+1) + \beta(t+1)\mathbf{s}(t). \quad (11)$$

2 10. Set $t = t+1$, and return to step 4.

3

4 2.3. Structure optimization

5

6 Following [Setiono and Leow \(2000\)](#), the structure optimization algorithm, which
7 has been successfully used to develop an optimized neural network model for crash
8 injury severity prediction ([Zeng and Huang, 2014b](#)), is proposed to improve the
9 generalization capacity of the neural network models and to identify the insignificant
10 explanatory variables. This method prunes the nodes that do not cause any significant
11 deterioration of the networks' accuracy. The mean absolute deviations of the training
12 set \mathbf{T} and testing set \mathbf{X} , that is, p and q , are used to evaluate the fitting and
13 predictive performance during network optimization:

$$14 \quad p = \frac{1}{M_1} \sum_{o(m) \in \mathbf{T}} |o(m) - \psi(m)|, \quad (12)$$

$$15 \quad q = \frac{1}{M_2} \sum_{o(m) \in \mathbf{X}} |o(m) - \psi(m)|, \quad (13)$$

16 where M_1 and M_2 are the number of samples in the training and testing sets,
17 respectively.

18 The following steps describe the detailed pruning process.

19 1. Train the network with a relatively large number of hidden nodes using the
20 conjugate gradient algorithm.

21 2. Calculate the p and q of the trained neural network, and set $p_b = p$,

22 $q_b = q$, and $ermax = \max\{p_b, q_b\}$.

23 3. For each $i(i=1, \dots, I)$, set $w_{i,j}^{(1)} = 0(j=1, \dots, J)$ and calculate the fitting errors p_i .

24 4. Retrain the network with $w_{i,j}^{(1)} = 0(j=1, \dots, J)$, where $p_i = \min_i p_i$, and compute
25 p and q for the retrained network.

26 5. If $p \leq (1+\sigma)ermax$ and $q \leq (1+\sigma)ermax$, then remove the input node l , set

27 $p_b = \min\{p, p_b\}$, $q_b = \min\{q, q_b\}$, $ermax = \max\{p_b, q_b\}$, $I = I - 1$,

28 and go back to step 3; otherwise, keep the previous weights of the network
29 connections.

30 6. For each $j(j=1, \dots, J)$, set $w_j^{(2)} = 0$ and calculate the fitting errors p_j .

1 7. Retrain the network with $w_h^{(2)} = 0$, where $p_h = \min_j p_j$, and compute p and q
 2 of the retrained network.

3 8. If $p \leq (1 + \sigma)ermax$ and $q \leq (1 + \sigma)ermax$, then remove the hidden node h . Set
 4 $p_b = \min\{p, p_b\}$, $q_b = \min\{q, q_b\}$, $ermax = \max\{p_b, q_b\}$, and
 5 $J = J - 1$, and go back to step 6; otherwise, keep the previous weights of the
 6 network connections.

7 In the above process, p_b and q_b represent, respectively, the minimal mean
 8 absolute deviations of the training and testing sets achieved so far. During the pruning
 9 process, generally, p_b increases while q_b decreases. $ermax$ is used to
 10 determine whether or not a node can be removed to remove as many insignificant
 11 nodes as possible without sacrificing the generalization accuracy. In addition, σ is
 12 the margin by which the error is allowed to increase when pruning a certain node.

13

14 2.4. Rule extraction

15

16 The rule extraction method developed by [Setiono and Thong \(2004\)](#) is modified to
 17 generate exact and comprehensible rules from the pruned neural network to illustrate
 18 the effects of significant explanatory variables. In the next subsections, a particle
 19 swarm optimization algorithm-based approach to approximating the transfer functions
 20 of hidden units is introduced as a critical step in the method, and then the rule
 21 extraction process is described.

22

23 2.4.1. Approximating transfer functions

24 The transfer functions of the hidden nodes can be approximated by piecewise
 25 functions. Theoretically, the more pieces fit the function, the more accurate the rule
 26 set, and the more rules may be extracted. To balance the two aspects, a three-piece
 27 linear function suggested by [Setiono and Thong \(2004\)](#) is used to approximate the
 28 transfer function of each hidden node $j(j = 1, \dots, J)$, $\tanh(\cdot)$, as shown in [Fig. 2](#). The
 29 slopes, β_{j_0} and β_{j_1} , and the cut-off point, ξ_{j_0} , are three undetermined parameters
 30 that minimize the sum of the squared deviations,

$$31 \min \sum_{m=1}^M (\tanh(v_j(m)) - L_j(v_j(m)))^2, \quad (14)$$

32 where

$$33 L_j(x) = \begin{cases} -\alpha_{j_1} + \beta_{j_1}x & \text{if } x < -\xi_{j_0} \\ \beta_{j_0}x & \text{if } -\xi_{j_0} \leq x \leq \xi_{j_0} \\ \alpha_{j_1} + \beta_{j_1}x & \text{if } x > \xi_{j_0} \end{cases}, \quad (15)$$

$$1 \quad v_j(m) = \sum_{i=1}^I w_{j,i}^{(1)} x_i(m), \quad (16)$$

$$2 \quad \alpha_{j1} = (\beta_{j0} - \beta_{j1}) \xi_{j0}. \quad (17)$$

3

4 *2.4.2. Searching the optimal parameters*

5 To approximate the transfer function accurately, the particle swarm optimization
6 algorithm, an efficient global search method, is used to solve the preceding nonlinear
7 optimization problem. The particle swarm optimization algorithm is well-known for
8 its exploration capacity, its exploitation capacity and its easy implementation (Poli et
9 al., 2012). In the algorithm, each feasible solution $(\beta_{j0}, \beta_{j1}, \xi_{j0})$ is referred to as a

10 “particle”, \mathbf{U} , and each particle flies around the three-dimensional search space with
11 a velocity \mathbf{V} , which is updated iteratively according to the best solution of the
12 particle achieved so far (particle best, **pbest**) and the best solution obtained by all of
13 the particles in the swarm so far (global best, **gbest**):

$$14 \quad \mathbf{V}_s^{r+1} = \mathbf{V}_s^r + c_1 \lambda_1 (\mathbf{pbest}_s^r - \mathbf{U}_s^r) + c_2 \lambda_2 (\mathbf{gbest}_s^r - \mathbf{U}_s^r), \quad (18)$$

$$15 \quad \mathbf{U}_s^{r+1} = \mathbf{U}_s^r + \mathbf{V}_s^{r+1}, \quad (19)$$

16

$$\mathbf{U} = (\beta_{j0}, \beta_{j1}, \xi_{j0}); r = 1, 2, \dots, R; s = 1, 2, \dots, S.$$

17 where \mathbf{U}_s^r is the s th particle at the r th iteration, and \mathbf{V}_s^{r+1} is its flying velocity to
18 the $r+1$ th iteration. c_1 and c_2 are two acceleration constants, while λ_1 and λ_2
19 are two uniform random numbers in $[0,1]$. R is the maximum iteration number, and
20 S is the number of particles used for searching the optimal solution.

21

22 *2.4.3. Generating regression rules*

23 Once the transfer functions of the hidden units have been approximated, the
24 relationship between the network inputs and outputs can be formulated with piecewise
25 linear functions. The detailed steps for extracting rules from the optimized neural
26 network are as follows:

- 27 1. For each hidden unit $j(j = 1, \dots, J)$, generate a three-piece linear function $L_j(x)$
28 with the approach previously described.
- 29 2. According to the pair cut-off points in $L_j(x)$, $-\xi_{j0}$ and ξ_{j0} , a certain input can
30 be located in one of three sections of hidden node j . Then, J hidden nodes will

1 result in $\underbrace{3 \times 3 \times \dots \times 3}_J$ locations for inputs. Consequently, the whole input space can
 2 be separated into 3^J subspaces.

3 3. For each non-empty subspace, the rule consequence is $\tilde{y} = \sum_{j=1}^J w_j^{(2)} \cdot L_j(v_j)$, where

4 $v_j = \sum_{i=1}^I w_{i,j}^{(1)} \cdot x_i$, and the rule condition is $C_1 \& C_2 \& \dots \& C_J$, where C_j is either

5 $v_j < -\xi_{j0}$, $-\xi_{j0} \leq v_j \leq \xi_{j0}$ or $v_j > \xi_{j0}$.

6

7 **3. Data preparation and preliminary analysis**

8

9 A crash dataset obtained from the Traffic Information System maintained by the
 10 Transport Department of Hong Kong is used to demonstrate the proposed neural
 11 network models and to compare them with the multivariate Poisson-lognormal model.
 12 This dataset contains 211 road segments that are evenly and widely distributed across
 13 Hong Kong. Geographical information system techniques are used to map crashes to
 14 these segments. The injury severity outcomes of the crashes are divided into two
 15 levels, fatality or serious injury and slight injury, and the annual crash numbers at
 16 each severity level at each site during 2002 to 2006 are obtained. The road geometric
 17 and traffic information is also included in the dataset. [Table 1](#) illustrates the
 18 definitions and descriptive statistics of the variables used in the model development.

19 The lane changing opportunity (LCO) variable refers to the different types of
 20 central lane marking, with values 0, 1 and 2 representing, respectively, double
 21 continuous lines, double lines with one continuous line and one dashed line, and a
 22 single dashed line. For those sub-segments with more than one type of central lane
 23 marking, the length-weighted average values are used. [Pei et al. \(2012\)](#) provides a
 24 more detailed description of the lane changing opportunity.

25 According to [Table 1](#), the mean and variance of crash frequency at the slight
 26 injury level are 6.04 and 25.81, respectively, indicating a possible over-dispersion. A
 27 similar characteristic is found in the fatality or serious injury crash frequency. In the
 28 multivariate Poisson-lognormal model, to account for the potential nonlinear
 29 relationship between crash frequencies and traffic volumes, the natural logarithm of
 30 AADT and Length, $\ln(\text{AADT})$ and $\ln(\text{Length})$, are modeled as other factors ([Zeng and
 31 Huang, 2014a](#)).

32 Correlation tests and multi-collinearity diagnoses for the risk factors are then
 33 conducted. According to the results of the Pearson correlation tests, we find that
 34 $\ln(\text{AADT})$ and Lane, $\ln(\text{AADT})$ and Park, Lane and LCO, SL and Shoulder, SL and
 35 Park are significantly correlated with correlation coefficients greater than 0.6. To
 36 reduce the model complexity, Lane, Park, and Shoulder are therefore excluded from

1 the models. The results of the diagnoses indicate that there is no significant
2 collinearity in the remaining factors.

4. Model implementation and result analysis

4.1. Model implementation

8 The multivariate Poisson-lognormal model is estimated with the freeware
9 WinBUGS, which is a popular platform to make Bayesian inference and is
10 well-known for its flexible programming environment (Zeng and Huang, 2014a). In
11 the absence of sufficient prior knowledge, non-informative priors are specified for the
12 parameters and the hyper-parameters. Specifically, a diffused normal distribution
13 $N(0, 10^4)$ is used as the priors of all elements of β_s ($s=1,2$), while a Wishart prior

14 $W(\mathbf{P}, r)$ is used for Σ^{-1} , where $\mathbf{P} = \begin{bmatrix} 1, & 0 \\ 0, & 1 \end{bmatrix}$ represents the scale matrix and $r = 2$

15 is the degrees of freedom (El-Basyouny and Sayed, 2009; Park and Lord, 2007). Five
16 hundred thousand iterations of the Markov chain Monte Carlo simulation are made,
17 with the first 4000 iterations acting as burn-ins. After ensuring the Markov chain
18 Monte Carlo convergence by the Gelman-Rubin statistics available in WinBUGS,
19 another 50,000 iterations are set to make summaries for the (hyper-) parameters.

20 The training, the structure optimization, and the rule extraction algorithms of the
21 neural network models are programmed in MATLAB. All of the variables are
22 normalized for the convenience of network training. To compare the performance of
23 the models fully, a 5-fold cross validation is conducted, where the dataset is randomly
24 divided into five parts with equal number of observations/patterns. Each time, the
25 sub-dataset of any four parts is input for training the models while the rest is used for
26 testing the predictive performance. Based on the collected data, $I = 14$, we first set
27 $J = 10$ for all networks. In the network training, $\varepsilon = 0.001$ and $T = 50$. We assume
28 that $\sigma = 0.05$ in the structure optimization algorithm, while $R = 300$ and $S = 700$
29 in the particle swarm optimization algorithm.

4.2. Model comparison

33 The results of the model comparison are summarized in Table 2. With regard to
34 the five folds of model comparison, in terms of the mean absolute deviation criteria,
35 all of the trained and optimized neural network models have lower fitting and
36 predictive errors for the training and testing datasets than the multivariate
37 Poisson-lognormal models, at both the fatality or serious injury and the slight injury
38 levels. This demonstrates that neural network models of crash frequency prediction
39 may give a better approximation performance than certain traditional statistical

1 models, which is probably due to the neural network's capacity for approximating
2 arbitrary nonlinear functions.

3 After pruning the network structure with the structure optimization algorithm, the
4 model fitting is generally expected to be degraded to some extent but the model
5 prediction should be improved as discussed in the [section 2.3](#). But in the results as
6 shown in [Table 2](#), it is surprisingly found that both the fitting and predictive errors of
7 the neural network models are reduced by the proposed model structure optimization
8 algorithm. As generally known, like other training algorithms, the proposed conjugate
9 gradient algorithm may sometimes be locally converged ([Haykin, 2009](#)). Therefore, a
10 presumable cause for the reduced model-fitting errors may be that pruning nodes and
11 retraining network could help to escape from local minima and to search for better
12 solutions. As a result, we may argue that the model generalization performance
13 associated with the proposed algorithm is improved as reflected by the reduced model
14 fitting and predictive errors.

15 Moreover, certain numbers of input and hidden nodes are removed from the
16 trained neural networks in all of the five folds, which indicates that the original
17 models have redundant nodes, and that the factors corresponding to those removed
18 input nodes may have no significant effects on the crash frequency by severity.

19 It is also noticeable that the five pairs of optimized neural networks end up with
20 slight distinctions in their mean absolute deviation values and the final sets of input
21 and hidden nodes. This instability is presumably attributable to the small sample size
22 ([Xie et al., 2007](#)), given the important impact of sample size on a model's
23 generalizability ([Haykin, 2009](#)).

24 25 *4.3. Interpretation of the explanatory variables*

26
27 The specific conditions and consequences of the rules extracted from the
28 optimized neural networks are shown in [Tables 3-6](#). In these tables, we can clearly see
29 the effects of the significant factors on crash frequencies at the two levels of injury
30 severity, under diverse conditions. For the purpose of comparison, the estimation
31 results of the parameters and the hyper-parameters in the multivariate
32 Poisson-lognormal model are shown in [Table 7](#) and [Table 8](#), respectively. According
33 to the results in [Table 8](#), we see that both the fatality or serious injury and the slight
34 injury crash data are over-dispersed, as their extra-Poisson variations (σ_{11} and σ_{22})
35 are significantly positive at the 95% credible level. Moreover, the correlation
36 coefficient $\rho (= \sigma_{12} / \sqrt{\sigma_{11}\sigma_{22}})$ reaches 0.763, showing that the crash frequencies at
37 the two injury levels are highly correlated.

38 In this section, we analyze mainly the rule consequences in [Tables 5 and 6](#), as the
39 rule conditions in [Tables 3 and 4](#) may be difficult to understand. Instead, we employ
40 the characteristics of the road segments involved at certain particular rules to illustrate

1 the effects of the risk factors. Even so, it is noticeable that, based on the conditions,
2 the rule to which each observation in the analysis should be assigned can be
3 determined accurately. Comparing the results in [Tables 5 and 6](#) with those in [Table 7](#),
4 we find that the coefficients of all of the identified factors in the optimized neural
5 networks are significant at the 95% credible level in the multivariate
6 Poisson-lognormal model, except Rainfall in the fatality or serious injury neural
7 network.

8 Regarding the main effects of the risk factors identified, most of the risk factors
9 have consistent signs, as shown in [Tables 5 and 6](#), which also conform to the signs in
10 the multivariate Poisson-lognormal model results shown in [Table 7](#). The signs of the
11 coefficients of the factors AADT, Length, SL, BS and Rainfall in the slight injury
12 neural network, and AADT and Diverge in the fatality or serious injury neural
13 network are identical at all rules. As for the other factors, it is interesting to find a few
14 different signs in several specific rules. Moreover, it is observed that the estimated
15 coefficient values are also distinct for several specific rules. This implies that those
16 risk factors probably have variable safety effects under different road conditions. This
17 could be important evidence of nonlinear relationship between crash frequency by
18 severity and the risk factors, which cannot be identified and modeled with the
19 traditional generalized linear regression models, such as the multivariate
20 Poisson-lognormal model.

21 According to the results in [Table 5](#), more slight injury crashes tend to occur on
22 longer road segments with more daily traffic, as observed by the positive coefficient
23 estimations associated with all of the eleven rules for AADT and Length. This is a
24 reasonable conclusion, given that AADT and segment length are always used as two
25 of the crash exposure variables in highway safety analysis ([AASHTO, 2010](#); [Zeng
26 and Huang, 2014a](#)). Nonetheless, the proposed neural network model presents specific
27 values for varied safety effects under different conditions. For example, increasing
28 one unit of AADT is expected to increase only 0.08 crashes (based on the normalized
29 data) under Condition 10, but 2.11 crashes (almost 26 times the former) under
30 Conditions 1 and 2.

31 Slight injury crash frequencies are found lower on road segments with higher
32 speed limits. It may be attributed to two reasons: (1) roadway segments designed for
33 higher speeds are usually well planned, constructed, and managed, features that
34 promote road safety, as argued by some previous researchers ([Milton and Mannering,
35 1998](#)), and (2) given a collision occurs, higher speed usually increases the likelihood
36 of severe injury and fatality while decreases that of slight injury ([Zeng and Huang,
37 2014b](#)).

38 The presence of median barriers is found to reduce slight injury crash occurrence
39 at most rules. A number of existing studies have also found that median barriers can
40 effectively prevent cross-median crashes ([Donnell and Mason, 2006](#)). However, the
41 estimated coefficients are positive at Rules 6 and 9. For those observations at these
42 rules, about 90 % of the road segments have median barriers, of which most are

1 inner-city highways with heavy daily traffic (mean = 36,538 vehicles) and many
2 merging ramps (mean = 1.74). These factors may hinder safe driving and bring about
3 more slight injury collisions related to median barriers. Under all conditions, the
4 presence of bus stops decreases slight injury crash frequencies, which may be
5 attributed to the increased interaction between buses and other vehicles when entering
6 or leaving bus bays (Pei et al., 2012).

7 The negative coefficients of Gradient under most conditions indicate that more
8 slight crashes are expected to occur on road segments with steeper downgrade slopes,
9 which is generally consistent with engineering experience. Besides, Gradient is found
10 to decrease the crash frequencies at Rules 6 and 10. Most of the involved road
11 segments are very long, such as Tsing Long Highway (9.07 km), Shek O Road (7.75
12 km) and Tolo Highway (5.60 km). Driving on the downgrade directions of these long
13 highways, drivers may be more careful, thus reducing the crash risk.

14 Slight injury crash frequencies usually increase with more lane changing
15 opportunities. Lane-cutting maneuvers often increase vehicle interaction, such as
16 overtaking, thereby raising the incidence of traffic conflict (Pei et al., 2012). It is
17 interesting to find that more lane changing opportunities could bring about more slight
18 injury collisions under Condition 10. The referred roadways mainly consists of
19 freeways, such as the longest segment in the dataset— Tsing Long Highway (9.07
20 km). Lane changing maneuver is less frequent on these freeways than on those busy
21 inner-city roadways, which may reduce the vehicle speed variance. This may possibly
22 explain why LCO negatively affects the slight injury crash frequency on them.

23 Generally, rainfall impairs visibility and makes road surfaces slippery, thereby
24 reducing skidding resistance, which raises the probability of crash occurrence. This is
25 why Rainfall has positive model coefficients in Table 5, which indicates that rainfall
26 may lead to more slight injury crashes (Pei et al., 2012).

27 Based on the results in Table 6, we find that more fatality or serious injury crashes
28 are associated with longer roadway segments, more daily traffic, no median barrier,
29 presence of bus stop, steeper downgrades and more precipitation under most or all
30 conditions, which is similar to the results of slight injury crashes. The negative
31 coefficients for the variable Diverge indicate that more diverging ramps give rise to a
32 higher fatality or serious injury crash risk, which may be attributable to more conflicts
33 at the sites approaching diverging ramps.

34 At Rule 21, the length is found negatively related to fatality or serious injury crash
35 frequency, in which all fatality or serious injury crashes occurred on Tsing Long
36 Highway. For the longest road segment, its average annual fatality or serious injury
37 crash number is only 1.3, smaller than the mean of the whole population (1.8). A
38 possible reason for the negative coefficient of Length may be that some unobserved
39 factors (such as well design and maintenance) associated with this highway greatly
40 promote the safety situation. Under the same Condition, Gradient is found positively
41 related to the fatality or serious injury crash frequency. It may be a result of a similar
42 reason to the corresponding findings in slight injury crashes, that is, drivers are

1 usually more cautious when driving on the downgrade of the so long (9.07 km)
2 segment.

3 Regarding the observations at Rules 16 and 17, most of the road segments are also
4 inner-city highways with heavy daily traffic (mean = 35,462 vehicles) and many
5 merging ramps (mean = 1.82). Like the situation at Rules 6 and 9, these factors may
6 impede safe driving and result in more fatality or serious injury median-related
7 crashes. Meanwhile, there are bus stops on all of the roadway segments at Rules 16
8 and 17. The decreased travel speed of buses entering or leaving bus bays could reduce
9 the probability of severe crashes. As a consequence, the presence of bus stops
10 decreases the fatality or serious injury crash frequency on the segments.

11 Probably due to the same reason as for the slight injury crashes, Rainfall has
12 positive coefficients under most conditions. However, drivers tend to be more careful
13 and reduce their speed when driving on rainy areas. That may be why Rainfall is
14 negatively related to the fatality or serious injury crash frequency under Conditions 14
15 and 15, since the annual precipitation of most involved observations are over 3000
16 mm.

18 **5. Conclusions and future research**

19
20 This study develops advanced neural networks for modeling the nonlinear
21 relationship between crash frequency by severity and the related factors. To improve
22 the generalization capacity and to handle the black-box characteristic of neural
23 networks, a structure optimization algorithm and a modified rule extraction algorithm
24 are proposed. A crash dataset obtained from the Traffic Information System
25 maintained by the Transport Department of Hong Kong, where crashes are classified
26 into slight injury and fatality or serious injury severity degrees, is used to demonstrate
27 the proposed methods and to compare them with the results of a multivariate
28 Poisson-lognormal model.

29 Despite the over-dispersed crash data and the high correlation between the crash
30 frequencies of the different injury degrees, the results show that both the trained and
31 the optimized neural networks outperform the multivariate Poisson-lognormal model
32 in fitting and predictive performance. It indicates the neural network's superiority
33 over the multivariate Poisson-lognormal model in modeling crash frequency by
34 severity. When several input and hidden nodes are deleted from the original neural
35 networks, better approximation performance is achieved, demonstrating the structure
36 optimization algorithm's ability to identify insignificant factors and to improve the
37 model's generalization capacity. The optimized neural networks generate two rule-sets
38 in which the coefficients of the explanatory variables are different, which confirms
39 that they are nonlinearly related to the crash frequencies. The signs of these
40 coefficients have identical directions under most conditions, and are consistent with
41 those in the multivariate Poisson-lognormal model. Moreover, most of the results for
42 the explanatory variables are reasonable and conform to traffic engineering

1 experience or the findings of previous studies, which further validates the proposed
2 methods.

3 It is worth noting that the other aforementioned statistical model may have better
4 performance than the multivariate Poisson-lognormal model for the collected data in
5 this study, although the latter is the most popular method for jointly modeling crash
6 frequency and severity. For example, the identified nonlinear relationship between
7 crash frequency by severity and risk factors could be viewed as unobserved
8 heterogeneities across observations. The heterogeneities could be accommodated in a
9 multivariate random parameters Poisson-lognormal model, and the empirical analysis
10 based on our collected dataset indicates that it is potentially a better fitting approach.
11 Further research efforts could be made to compare the proposed neural network
12 models with the emerging advanced statistical models on more field datasets. Further,
13 as mentioned above, the developed neural network models can be employed as an
14 alternative approach for identifying sites with promise for improving safety. In the
15 absence of the averaged crash cost of each injury degree level, this part of the
16 application has not been conducted. More comprehensive crash data are needed to
17 compare the proposed neural network techniques with the state-of-the-art methods,
18 such as the Bayesian hierarchical models, for site ranking.

19 **Acknowledgements**

20 This research was jointly supported by the Natural Science Foundation of China
21 (No. 71371192, 71301083), the Hong Kong Research Grants Council of the Hong
22 Kong Special Administrative Region, China (No. 717512) and a grant from the Joint
23 Research Scheme of National Natural Science Foundation of China/Research Grants
24 Council of Hong Kong (No. 71561167001 & N_HKU707/15).

25 **References**

- 26 Highway Safety Manual, 1st edition, 2010. AASHTO, Washington D.C.
- 27 Abdelwahab, H.T., Abdel-Aty, M.A., 2001. Development of artificial neural network
28 models to predict driver injury severity in traffic accidents at signalized
29 intersections. *Transportation Research Record* 1746, 6-13.
- 30 Aguero-Valverde, J., Jovanis, P.P., 2009. Bayesian multivariate Poisson lognormal
31 models for crash severity modeling and site ranking. *Transportation Research*
32 *Record* 2136, 82-91.
- 33 Anastasopoulos, P. C., Shankar, V. N., Haddock, J. E., Mannering, F. L., 2012. A
34 multivariate tobit analysis of highway accident-injury-severity rates. *Accident*
35 *Analysis and Prevention* 45, 110-119.
- 36 Barua, S., El-Basyouny, K., Islam, M. T., 2014. A full Bayesian multivariate count
37 data model of collision severity with spatial correlation. *Analytic Methods in*
38 *Accident Research* 3, 28-43.

- 1 Barua, S., El-Basyouny, K., Islam, M. T., 2016. Multivariate random parameters
2 collision count data models with spatial heterogeneity. *Analytic Methods in*
3 *Accident Research* 9, 1-15.
- 4 Bijleveld, F. D. 2005. The covariance between the number of accidents and the
5 number of victims in multivariate analysis of accident related outcomes. *Accident*
6 *Analysis and Prevention* 37 (4), 591-600.
- 7 Chang, L., 2005. Analysis of freeway accident frequencies: negative binomial
8 regression versus artificial neural network. *Safety Science* 43 (8), 541-557.
- 9 Chiou, Y. C., Fu, C., 2013. Modeling crash frequency and severity using
10 multinomial-generalized Poisson model with error components. *Accident Analysis*
11 *and Prevention* 50, 73-82.
- 12 Chiou, Y. C., Fu, C., Chih-Wei, H., 2014. Incorporating spatial dependence in
13 simultaneously modeling crash frequency and severity. *Analytic methods in*
14 *accident research* 2, 1-11.
- 15 Chiou, Y. C., Fu, C., 2015. Modeling crash frequency and severity with
16 spatiotemporal dependence. *Analytic Methods in Accident Research* 5, 43-58.
- 17 Donnell, E.T., Mason Jr, J.M., 2006. Predicting the frequency of median barrier
18 crashes on Pennsylvania interstate highways. *Accident Analysis and Prevention*
19 38(3), 590-599.
- 20 El-Basyouny, K., Sayed, T., 2009. Collision prediction models using multivariate
21 Poisson-lognormal regression. *Accident Analysis and Prevention* 41 (4), 820-828.
- 22 El-Basyouny, K., Barua, S., Islam, M. T., 2014. Investigation of time and weather
23 effects on crash types using full Bayesian multivariate Poisson lognormal models.
24 *Accident Analysis and Prevention* 73, 91-99.
- 25 Haykin, S.S., 2009. *Neural networks and learning machines*, 3rd edition. Prentice Hall,
26 New York.
- 27 Huang, H., Zeng, Q., Pei, X., Wong, S.C., Xu, P., 2016. Predicting crash frequency
28 using an optimized radial basis function neural network model. *Transportmetrica*
29 *A* 12 (4): 330-345.
- 30 Karlaftis, M.G., Vlahogianni, E.I., 2011. Statistical methods versus neural networks in
31 transportation research: Differences, similarities and some insights. *Transportation*
32 *Research Part C* 19 (3), 387-399.
- 33 Labi, S., 2011. Efficacies of roadway safety improvements across functional
34 subclasses of rural two-lane highways. *Journal of Safety Research* 42 (4),
35 231-239.
- 36 Lan, B., Persaud, B., 2012. Evaluation of multivariate Poisson log normal Bayesian
37 methods for before-after road safety evaluations. *Journal of Transportation Safety*
38 *and Security* 4 (3), 193-210.
- 39 Li, X., Lord, D., Zhang, Y., Xie, Y., 2008. Predicting motor vehicle crashes using
40 support vector machine models. *Accident Analysis and Prevention* 40 (4),
41 1611-1618.
- 42 Ma, J., Kockelman, K.M., 2006. Bayesian multivariate Poisson regression for models

1 of injury count, by severity. *Transportation Research Record* 1950, 24-34.

2 Ma, J., Kockelman, K.M., Damien, P., 2008. A multivariate Poisson-lognormal
3 regression model for prediction of crash counts by severity, using Bayesian
4 methods. *Accident Analysis and Prevention* 40 (3), 964-975.

5 Miaou, S. P., Song, J. J., 2005. Bayesian ranking of sites for engineering safety
6 improvements: decision parameter, treatability concept, statistical criterion, and
7 spatial dependence. *Accident Analysis and Prevention* 37 (4), 699-720.

8 Milton, J., Mannering, F., 1998. The relationship among highway geometrics,
9 traffic-related elements and motor-vehicle accident frequencies. *Transportation* 25
10 (4), 395-413.

11 Park, E.S., Lord, D., 2007. Multivariate Poisson-lognormal models for jointly
12 modeling crash frequency by severity. *Transportation Research Record* 2019, 1-6.

13 Pei, X., Wong, S.C., Sze, N.N., 2011. A joint-probability approach to crash prediction
14 models. *Accident Analysis and Prevention* 43 (3), 1160-1166.

15 Pei, X., Wong, S.C., Sze, N.N., 2012. The roles of exposure and speed in road safety
16 analysis. *Accident Analysis and Prevention* 48, 464-471.

17 Poli, R., Kennedy, J., Blackwell, T., 2007. Particle swarm optimization. *Swarm*
18 *intelligence* 1, 33-57.

19 Setiono, R., Leow, W.K., 2000. Pruned neural networks for regression. *PRICAI 2000*
20 *Topics in Artificial Intelligence*. Springer Berlin, Heidelberg, pp. 500-509.

21 Setiono, R., Thong, J.Y., 2004. An approach to generate rules from neural networks
22 for regression problems. *European Journal of Operational Research* 155, 239-250.

23 Venkataraman, N., Ulfarsson, G.F., Shankar, V.N., 2013. Random parameter models
24 of interstate crash frequencies by severity, number of vehicles involved, collision
25 and location type. *Accident Analysis and Prevention* 59, 309-318.

26 Wang, C., Quddus, M. A., Ison, S. G., 2011. Predicting accident frequency at their
27 severity levels and its application in site ranking using a two-stage mixed
28 multivariate model. *Accident Analysis and Prevention* 43 (6), 1979-1990.

29 Xie, Y., Lord, D., Zhang, Y., 2007. Predicting motor vehicle collisions using Bayesian
30 neural networks: an empirical analysis. *Accident Analysis and Prevention* 39 (5),
31 922-933.

32 Xu, X., Wong, S.C., Choi, K., 2014. A two-stage bivariate logistic-Tobit model for the
33 safety analysis of signalized intersections. *Analytic Methods in Accident Research*
34 3-4, 1-10.

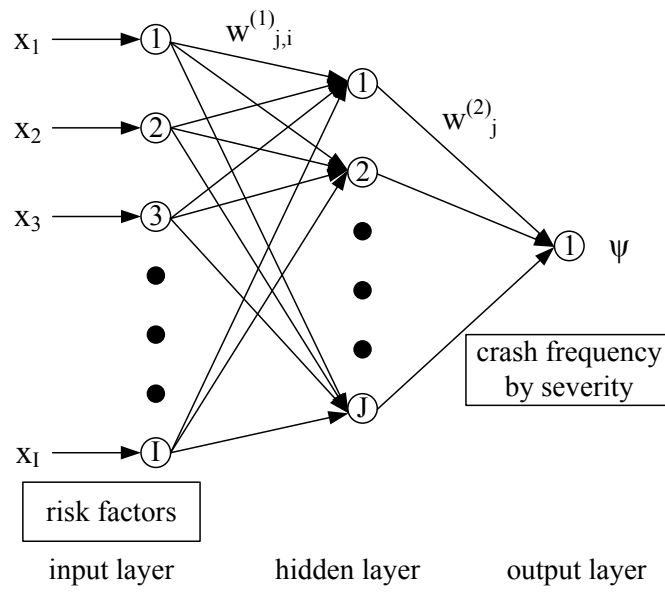
35 Ye, X., Pendyala, R.M., Washington, S.P., Konduri, K., Oh, J., 2009. A simultaneous
36 equations model of crash frequency by collision type for rural intersections. *Safety*
37 *Science* 47 (3), 443-452.

38 Ye, X., Pendyala, R. M., Shankar, V., Konduri, K. C., 2013. A simultaneous equations
39 model of crash frequency by severity level for freeway sections. *Accident*
40 *Analysis and Prevention* 57, 140-149.

41 Zeng, Q., Huang, H., 2014a. Bayesian spatial joint modeling of traffic crashes on an
42 urban road network. *Accident Analysis and Prevention* 67, 105-112.

- 1 Zeng, Q., Huang, H., 2014b. A stable and optimized neural network model for crash
- 2 injury severity prediction. *Accident Analysis and Prevention* 73, 351-358.

1



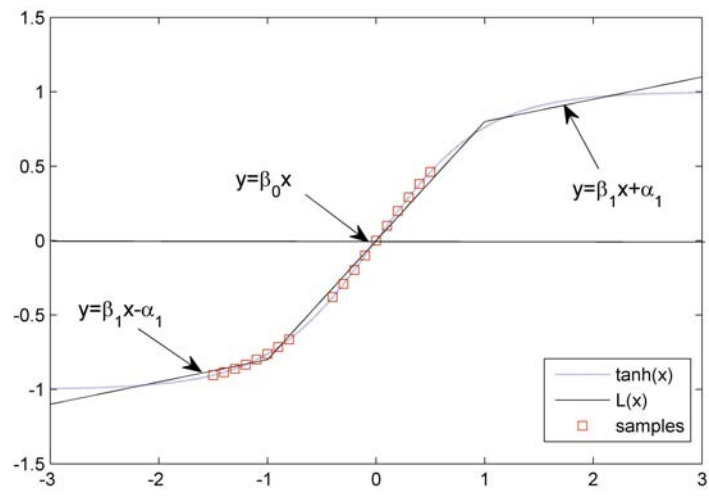
2

3

4

Fig. 1. Developed multilayer perceptron structure

1



2

3

4

Fig. 2. Three-piece linear approximation of $\tanh(\cdot)$

1 **Table 1** Descriptive statistics of the variables

Variable	Description	Mean	SD	Min.	Max.
<i>Response variable</i>					
Slight	Slightly injured crash count per segment per year	6.04	5.08	0	40
KSI	Killed and seriously injured crash count per segment per year	1.60	1.90	0	12
<i>Explanatory variables</i>					
AADT	Average annual daily traffic (vehicles)	22077	19945	1164	101632
Length	Segment length (km)	1.47	1.55	0.15	9.07
Lane	Number of lanes	2.41	1.18	1	7
Width	Average width of each lane (m)	3.63	0.64	2.40	7.30
SL	Posted speed limit (km/h)	60.3	14.7	50	110
Merge	Number of merging ramps	0.84	1.00	0	4
Diverge	Number of diverging ramps	1.75	2.27	0	17
Inter	Number of intersections	1.90	2.37	0	16
Gradient	Average segment gradient (10^{-2})	0.04	2.74	-11	11
Curvature	Average segment curvature	21.9	17.5	0	85
LCO	Lane changing opportunity	2.43	1.61	0	7.85
Median	Presence of median barrier: yes = 1, no = 0	0.70	0.46	0	1
BS	Presence of bus stop: yes = 1, no = 0	0.64	0.48	0	1
Shoulder	Presence of hard shoulder: yes = 1, no = 0	0.13	0.34	0	1
Park	Presence of on-street parking: yes = 1, no = 0	0.51	0.49	0	1
Rainfall	Annual precipitation (mm)	2279	565	761	3215

2

1 **Table 2 Model comparison**

Injury severity	Model	Training mean absolute deviation	Testing mean absolute deviation	Number of input nodes	Number of hidden nodes	
1	Multivariate	2.87	3.03	—	—	
	Poisson-lognormal					
	Slight injury	Trained neural network	2.85	2.99	14	10
		Optimized neural network	2.71	2.93	9	4
1	Multivariate	1.17	1.00	—	—	
	Poisson-lognormal					
	Fatality or serious injury	Trained neural network	1.09	1.07	14	10
		Optimized neural network	1.08	1.03	8	4
2	Multivariate	2.85	3.00	—	—	
	Poisson-lognormal					
	Slight injury	Trained neural network	2.77	2.95	14	10
		Optimized neural network	2.76	2.88	7	4
2	Multivariate	1.09	1.30	—	—	
	Poisson-lognormal					
	Fatality or serious injury	Trained neural network	1.06	1.26	14	10
		Optimized neural network	1.03	1.25	8	4
3	Multivariate	2.89	2.86	—	—	
	Poisson-lognormal					
	Slight injury	Trained neural network	2.84	2.80	14	10
		Optimized neural network	2.63	2.71	11	4
3	Multivariate	1.13	1.15	—	—	
	Poisson-lognormal					
	Fatality or serious injury	Trained neural network	1.08	1.14	14	10
		Optimized neural network	1.06	1.13	10	4

2

1 **Table 2 (continued)** Model comparison

	Multivariate	2.87	2.95	—	—	
	Poisson-lognormal					
	Slight injury	Trained neural network	2.78	2.85	14	10
4		Optimized neural network	2.71	2.76	9	3
	Fatality or serious injury	Multivariate	1.14	1.13	—	—
		Poisson-lognormal				
Trained neural network		1.09	1.10	14	10	
	Optimized neural network	1.06	1.07	7	4	
	Multivariate	2.90	2.82	—	—	
	Poisson-lognormal					
	Slight injury	Trained neural network	2.87	2.76	14	10
5		Optimized neural network	2.75	2.70	11	4
	Fatality or serious injury	Multivariate	1.12	1.21	—	—
		Poisson-lognormal				
Trained neural network		1.07	1.12	14	10	
	Optimized neural network	1.02	1.08	11	6	

2

1 **Table 3** Rule conditions for slight injury crashes

Rule	Condition ^a
1	$v_1 < -1.498 \ \& \ -0.87 \leq v_2 \leq 0.87 \ \& \ -1.042 \leq v_3 \leq 1.042 \ \& \ v_4 < -0.52$
2	$v_1 > 1.498 \ \& \ -0.87 \leq v_2 \leq 0.87 \ \& \ -1.042 \leq v_3 \leq 1.042 \ \& \ v_4 < -0.52$
3	$v_1 > 1.498 \ \& \ v_2 > 0.87 \ \& \ -1.042 \leq v_3 \leq 1.042 \ \& \ v_4 < -0.52$
4	$v_1 < -1.498 \ \& \ -0.87 \leq v_2 \leq 0.87 \ \& \ v_3 > 1.042 \ \& \ v_4 < -0.52$
5	$v_1 > 1.498 \ \& \ -0.87 \leq v_2 \leq 0.87 \ \& \ v_3 > 1.042 \ \& \ v_4 < -0.52$
6	$v_1 > 1.498 \ \& \ v_2 > 0.87 \ \& \ v_3 > 1.042 \ \& \ v_4 < -0.52$
7	$v_1 < -1.498 \ \& \ v_2 > 0.87 \ \& \ -1.042 \leq v_3 \leq 1.042 \ \& \ -0.52 \leq v_4 \leq 0.52$
8	$v_1 > 1.498 \ \& \ v_2 > 0.87 \ \& \ -1.042 \leq v_3 \leq 1.042 \ \& \ -0.52 \leq v_4 \leq 0.52$
9	$v_1 > 1.498 \ \& \ -0.87 \leq v_2 \leq 0.87 \ \& \ v_3 > 1.042 \ \& \ -0.52 \leq v_4 \leq 0.52$
10	$v_1 > 1.498 \ \& \ v_2 > 0.87 \ \& \ v_3 > 1.042 \ \& \ -0.52 \leq v_4 \leq 0.52$
11	$v_1 > 1.498 \ \& \ v_2 > 0.87 \ \& \ -1.042 \leq v_3 \leq 1.042 \ \& \ v_4 > 0.52$
2	$v_1 = -0.819 + 0.344AADT - 1.061Length + 0.159SL - 0.357Median$ $+ 0.032BS - 0.443Gradient + 0.076LCO + 0.760Rainfall$
3	$v_2 = 0.668 + 1.555AADT - 0.532Length - 0.019SL - 0.179Median$ $+ 0.163BS - 0.134Gradient + 0.238LCO - 0.029Rainfall$
4	$v_3 = 1.809 - 1.502AADT - 0.763Length - 0.030SL + 0.736Median$ $- 0.305BS + 0.635Gradient - 0.878LCO - 0.323Rainfall$
5	$v_4 = -0.504 + 1.271AADT - 0.647Length + 0.285SL - 0.364Median$ $+ 0.174BS - 0.177Gradient + 0.504LCO - 0.440Rainfall$
6	

1 **Table 4** Rule conditions for fatality or serious injury crashes

Rule	Condition ^a
12	$v_1 < -1.148 \ \& \ -0.407 \leq v_2 \leq 0.407 \ \& \ v_3 < -0.746 \ \& \ v_4 < -0.139$
13	$v_1 > 1.148 \ \& \ -0.407 \leq v_2 \leq 0.407 \ \& \ v_3 < -0.746 \ \& \ v_4 < -0.139$
14	$v_1 < -1.148 \ \& \ -0.407 \leq v_2 \leq 0.407 \ \& \ -0.746 \leq v_3 \leq 0.746 \ \& \ v_4 < -0.139$
15	$v_1 > 1.148 \ \& \ -0.407 \leq v_2 \leq 0.407 \ \& \ -0.746 \leq v_3 \leq 0.746 \ \& \ v_4 < -0.139$
16	$v_1 < -1.148 \ \& \ v_2 > 0.407 \ \& \ -0.746 \leq v_3 \leq 0.746 \ \& \ v_4 < -0.139$
17	$v_1 > 1.148 \ \& \ v_2 > 0.407 \ \& \ -0.746 \leq v_3 \leq 0.746 \ \& \ v_4 < -0.139$
18	$v_1 > 1.148 \ \& \ -0.407 \leq v_2 \leq 0.407 \ \& \ v_3 > 0.746 \ \& \ v_4 < -0.139$
19	$v_1 < -1.148 \ \& \ v_2 > 0.407 \ \& \ v_3 > 0.746 \ \& \ v_4 < -0.139$
20	$v_1 > 1.148 \ \& \ v_2 > 0.407 \ \& \ v_3 > 0.746 \ \& \ v_4 < -0.139$
21	$v_1 > 1.148 \ \& \ -0.407 \leq v_2 \leq 0.407 \ \& \ v_3 < -0.746 \ \& \ -0.139 \leq v_4 \leq 0.139$
22	$v_1 > 1.148 \ \& \ -0.407 \leq v_2 \leq 0.407 \ \& \ v_3 < -0.746 \ \& \ v_4 > 0.139$
2	^a $v_1 = -1.126 + 0.789AADT + 1.082Length + 0.923Diverge - 0.353Median - 0.066BS - 0.281Gradient - 0.317Rainfall$
3	$v_2 = 0.348 + 0.352AADT + 0.127Length - 0.6Diverge - 0.483Median + 0.473BS - 0.06Gradient - 0.111Rainfall$
4	$v_3 = 0.177 + 0.966AADT - 1.36Length - 0.383Diverge - 1.079Median + 0.751BS - 0.034Gradient + 0.188Rainfall$
5	$v_4 = -0.69 - 0.424AADT + 1.047Length - 0.496Diverge + 0.084Median + 0.06BS - 0.126Gradient - 0.587Rainfall$
6	

1 **Table 5** Rule consequences for slight injury crashes

Rule	Coefficient of the variable in the consequence (linear function)								
	Constant	AADT	Length	SL	Median	BS	Gradient	LCO	Rainfall
1	-0.618	2.110	0.746	-0.105	-0.626	0.347	-0.529	0.830	0.257
2	-0.830	2.110	0.746	-0.105	-0.626	0.347	-0.529	0.830	0.257
3	-0.751	1.322	1.015	-0.095	-0.535	0.265	-0.461	0.710	0.242
4	-0.038	0.975	0.170	-0.127	-0.069	0.117	-0.049	0.167	0.012
5	-0.250	0.975	0.170	-0.127	-0.069	0.117	-0.049	0.167	0.012
6	-0.172	0.188	0.439	-0.117	0.021	0.034	0.019	0.046	0.002
7	-0.543	1.055	1.152	-0.155	-0.458	0.228	-0.424	0.604	0.335
8	-0.755	1.055	1.152	-0.155	-0.458	0.228	-0.424	0.604	0.335
9	-0.254	0.707	0.306	-0.187	0.007	0.080	-0.012	0.061	0.105
10	-0.175	0.080	0.576	-0.177	-0.098	0.003	0.056	-0.060	0.091
11	-0.971	1.322	1.015	-0.095	-0.535	0.265	-0.461	0.710	0.242

2

1

Table 6 Rule consequences for fatality or serious injury crashes

Rule	Coefficient of the variable in the consequence (linear function)							
	Constant	AADT	Length	Diverge	Median	BS	Gradient	Rainfall
12	0.336	0.546	0.174	0.049	-0.366	0.225	-0.071	0.038
13	0.994	0.546	0.174	0.049	-0.366	0.225	-0.071	0.038
14	0.0667	0.264	0.570	0.161	-0.051	0.006	-0.061	-0.017
15	0.725	0.264	0.570	0.161	-0.051	0.006	-0.061	-0.017
16	0.080	0.186	0.542	0.294	0.056	-0.099	-0.048	0.008
17	0.738	0.186	0.542	0.294	0.056	-0.099	-0.048	0.008
18	0.559	0.546	0.174	0.049	-0.366	0.225	-0.071	0.038
19	-0.086	0.467	0.146	0.182	-0.258	0.120	-0.058	0.063
20	0.572	0.467	0.146	0.182	-0.258	0.120	-0.058	0.063
21	1.913	1.253	-1.574	0.876	-0.505	0.125	0.138	1.017
22	0.531	0.546	0.174	0.049	-0.366	0.225	-0.071	0.038

2

1 **Table 7** Parameter estimation in the multivariate Poisson-lognormal model

Variable	Slight injury			Fatality or serious injury		
	Mean	S.D.	95 % Credible interval	Mean	S.D.	95 % Credible interval
Constant	1.136	0.201	(0.774, 1.481)	-0.474	0.307	(-1.052, 0.126)
ln(AADT)	0.563	0.035	(0.496, 0.633)	0.334	0.056	(0.228, 0.446)
ln(Length)	0.557	0.030	(0.498, 0.616)	0.670	0.049	(0.574, 0.765)
SL	-0.027	0.002	(-0.034, -0.021)	-0.017	0.004	(-0.026, -0.009)
Merge	-0.046	0.024	(-0.093, 0.001)	-0.064	0.036	(-0.135, 0.006)
Diverge	0.023	0.012	(0.0003, 0.046)	0.057	0.017	(0.023, 0.090)
Inter	0.015	0.012	(-0.010, 0.040)	-0.010	0.019	(-0.047, 0.026)
Median	-0.185	0.068	(-0.313, -0.050)	-0.316	0.111	(-0.533, -0.099)
BS	0.384	0.054	(0.282, 0.490)	0.300	0.085	(0.135, 0.464)
Gradient	-1.559	0.771	(-3.069, -0.058)	-2.53	1.16	(-4.816, -0.253)
Curvature	-0.002	0.001	(-0.007, 0.004)	-0.001	0.002	(-0.007, 0.006)
LCO	0.104	0.017	(0.071, 0.136)	0.149	0.026	(0.099, 0.201)
Width	0.006	0.035	(-0.064, 0.075)	0.089	0.056	(-0.020, 0.197)
Rainfall	0.083	0.037	(0.015, 0.156)	0.062	0.056	(-0.045, 0.172)

2 The values in bold are those significantly positive or negative with the 95 % credible intervals
 3 bounded away from zero.

4

1 **Table 8** Hyper-parameter estimation in multivariate Poisson-lognormal

Hyper-parameter	Mean	S.D.	95 % Credible interval	
			2.5 %	97.5 %
σ_{11}	0.228	0.020	0.191	0.269
$\sigma_{21}(=\sigma_{12})$	0.199	0.022	0.198	0.243
σ_{22}	0.299	0.042	0.222	0.388
ρ^a	0.763	0.047	0.667	0.847

2 The values in bold are those significantly positive or negative with the 95 % credible intervals
 3 bounded away from zero.

4 ^{a:} $\rho = \sigma_{12} / \sqrt{\sigma_{11}\sigma_{22}}$.

5