

Chapman University

## Chapman University Digital Commons

---

Psychology Faculty Articles and Research

Psychology

---

11-20-2018

### **Crash Severity Analysis of Rear-End Crashes in California Using Statistical and Machine Learning Classification Methods**

Alidad Ahmadi

Arash Jahangiri

Vincent Berardi

Sahar Ghanipoor Machiani

Follow this and additional works at: [https://digitalcommons.chapman.edu/psychology\\_articles](https://digitalcommons.chapman.edu/psychology_articles)



Part of the [Automotive Engineering Commons](#), and the [Transportation Engineering Commons](#)

---

---

# Crash Severity Analysis of Rear-End Crashes in California Using Statistical and Machine Learning Classification Methods

## Comments

This is an Accepted Manuscript of an article published in *Journal of Transportation Safety & Security*, volume 12, issue 4, in 2020, available online at <https://doi.org/10.1080/19439962.2018.1505793>. It may differ slightly from the final version of record.

The Creative Commons license below applies only to this version of the article.

## Creative Commons License



This work is licensed under a [Creative Commons Attribution-Noncommercial 4.0 License](https://creativecommons.org/licenses/by-nc/4.0/)

## Copyright

Taylor & Francis

---

1  
2 **CRASH SEVERITY ANALYSIS OF REAR-END CRASHES IN CALIFORNIA**  
3 **USING STATISTICAL AND MACHINE LEARNING CLASSIFICATION METHODS**  
4

5  
6  
7 **Alidad Ahmadi**

8 Graduate Assistant

9 Department of Civil, Construction, and Environmental Engineering

10 San Diego State University

11 San Diego, CA 92182, United States

12 Tel: 760 798-5453; Email: aahmadi@sdsu.edu  
13

14 **Arash Jahangiri**

15 Assistant Professor

16 Department of Civil, Construction, and Environmental Engineering

17 San Diego State University

18 San Diego, CA 92182, United States

19 Tel: 540 200-7561; Email: ajahangiri@mail.sdsu.edu  
20

21 **Vincent Berardi**

22 Assistant Professor

23 Department of Psychology,

24 Chapman University

25 Orange, CA 92866, United States

26 Tel: 908 591-2948; Email: berardi@chapman.edu  
27

28 **Sahar Ghanipoor Machiani (corresponding author)**

29 Assistant Professor

30 Department of Civil, Construction, and Environmental Engineering

31 San Diego State University

32 San Diego, CA 92182, United States

33 Tel: 619 594-1937; Email: sghanipoor@mail.sdsu.edu  
34  
35  
36  
37  
38

1 **ABSTRACT**

2

3 Investigating drivers' injury level and detecting contributing factors that aggravate the damage  
4 level imposed on drivers and vehicles is a critical subject in the field of crash analysis. In this  
5 study, a comprehensive vehicle-by-vehicle crash dataset is developed by integrating five years of  
6 data from California crash, vehicles involved, and road databases. The dataset is used to model the  
7 severity of rear-end crashes for comparing three analytic techniques: multinomial logit (MNL),  
8 mixed multinomial logit (MMNL), and support vector machine (SVM). The results of the crash  
9 severity models and the role of contributing factors to the severity outcome of rear-end crashes are  
10 extensively discussed. In terms of prediction performance, all three models yielded comparable  
11 results; although, the SVM performed slightly better than the other two methods. The results from  
12 this study will inform aspects of our driver safety education and design, either vehicle or roadway  
13 design, required to be improved to alleviate the probability of severe injuries.

14

15 **Keywords:** traffic safety, crash severity classification, machine learning, mixed multinomial logit,  
16 support vector machine

1 **1. BACKGROUND**

2 According to crash statistics report presented annually by the Fatality Analysis Reporting  
3 System (FARS), in 2015, nearly 32 thousand people were killed in vehicle crashes throughout the  
4 United States (“Fatality Analysis Reporting System (FARS) Encyclopedia, National Highway  
5 Traffic Safety Administration (NHTSA),” n.d.). According to this report, since 1995, California,  
6 Texas, and Florida are among the states with the largest number of fatalities. The number of  
7 fatalities has decreased about 30 percent from the year the report was first initiated, from about  
8 47500 fatalities in 1994 to about 32000 fatalities in 2015; however, the current number is still  
9 extremely high, which reflects the need to find remedies to decrease the rate more quickly. In this  
10 regard, the crash analysis in the context of traffic safety has become one of the main areas of focus  
11 among the traffic engineers.

12 Crash databases are usually built by using police reports and are comprised of information such as  
13 the status of the crash, driver’s information, road segment detail, environmental factors, and traffic  
14 condition. Understanding crash models and identifying contributing factors and their significance  
15 are crucial as the outcomes can be used in higher-organizational and management-level actions to  
16 define countermeasures that could prevent future crashes, improve the standards for the roadway  
17 and network design, improve public health policies, provide better emergency services, alleviate  
18 driver’s injury severity, and nurture safer driving experience.

19 Crashes are naturally randomly-occurring incidents. Statistical models aid in better  
20 understanding the variability of these random events by examining the factors associated with  
21 them. There is a large body of research in the context of crash analysis. A review paper written by  
22 Lord and Mannering discusses the most common methodologies that have been used in studying  
23 crash frequency as well as the issues associated with them (Lord & Mannering, 2010). Savolainen

1 et al. (Savolainen, Mannering, Lord, & Quddus, 2011) published a similar paper discussing the  
2 most common methodologies for studying crash-severity analysis. A comprehensive study was  
3 carried out by Mannering and Bhat (Mannering & Bhat, 2014) in which an updated list of the most  
4 recent methodologies since the two previous studies [(Lord & Mannering, 2010) and (Savolainen  
5 et al., 2011)] was presented. Their analysis also focused on demonstrating how the crash analysis  
6 approach has evolved over time and how issues identified in previous studies have been addressed  
7 in more recently-introduced, advanced models (Mannering & Bhat, 2014). Most recently,  
8 Mannering et al. conducted a detailed discussion of how different statistical techniques can address  
9 the unobserved heterogeneity focusing on both injury-severity analysis and analysis of accident  
10 likelihood (Mannering, Shankar, & Bhat, 2016). In the remainder of this section, the most recent  
11 and relative studies concerning the modeling of injury-severity in crash databases are discussed.

12         Several studies that have developed injury-severity analysis models have focused on how  
13 the model's classification performance would change in absence or presence of various factors.  
14 For example, Li et al. (Z. Li, Wang, Liu, Bigham, & Ragland, 2013) have concentrated on  
15 demographic attributes which are believed to vary from one city to other cities. They showed the  
16 impact of implementing spatial attributes in development of crash prediction models in the  
17 framework of a Geographically Weighted Poisson Model (GWPM). Investigating the crash data  
18 and socio-demographic attributes of 58 counties in California, they compared the GWPM model's  
19 performance to the popular Generalized Linear Models (GLM) in predicting fatal crashes (Z. Li et  
20 al., 2013). In another study, Wu et al. (Q. Wu et al., 2014) used mixed logit models to analyze the  
21 injury-severity in single-vehicle crashes and crashes that involve two or more than two vehicles  
22 (Q. Wu et al., 2014). In a similar approach, focusing on truck-involved crashes, Chen et al. (C.  
23 Chen, Zhang, Tian, Bogus, & Yang, 2015) employed a modified Hierarchical Bayesian Random

1 Intercept model to predict the truck driver's injury severities based on a two year database of truck-  
2 involved crashes on rural roadways from New Mexico (C. Chen et al., 2015). Ye and Lord (Ye &  
3 Lord, 2014) investigated how the sample size can influence the performance of the Multinomial  
4 Logit, Ordered Probit, and Mixed Logit. Their results indicated that the sample size has a  
5 significant influence on the model performance. For example, a Mixed Logit model requires a  
6 much larger sample size compared to an Ordered Probit model, which can provide reasonable  
7 results even with small sample sizes (Ye & Lord, 2014).

8         The most basic crash severity modeling approaches, such as binary Logit and Probit  
9 models, have evolved into more advanced parametric and non-parametric models. These advanced  
10 models can address more of the unobserved characteristics of the data, which were not examined  
11 in earlier models. One such class of analytic approaches are mixed models (Milton, Shankar, &  
12 Mannering, 2008; Q. Wu et al., 2014; Yasmin & Eluru, 2013; Ye & Lord, 2014). For instance,  
13 Milton et al. (Milton et al., 2008) used the mixed logit model that is able to differentially account  
14 for the effect of independent variables on the level of severity over different road segments (Milton  
15 et al., 2008). Previous models had considered this effect to be the same on all road segments which  
16 results in biased results. The mixed feature of the model allows the coefficient of each variable  
17 affecting the injury-severity (explanatory variables) to vary across all individuals in the crash  
18 database in order to consider the heterogeneous effect and correlations of unobserved factors  
19 (Savolainen et al., 2011; K. E. Train, 2009). In another study, when studying single-vehicle crashes  
20 in the state of California from 2003 to 2004, Kim et al. (Kim, Ulfarsson, Kim, & Shankar, 2013)  
21 employed mixed logit models in injury-severity analysis. Their results indicated the importance of  
22 considering population heterogeneity according to the notable differences in results for different  
23 age groups, especially when comparing the older age group to younger drivers (Kim et al., 2013).

1 Developing separate models for two different age groups of older and younger drivers and for two  
2 gender groups on three different road surfaces, Morgan and Mannering (Morgan & Mannering,  
3 2011) utilized the mixed logit analysis to assess the effect of those factors on the severity of crashes  
4 in single-vehicle crashes. In a more recent study, Behnood et al. utilized mixed logit modeling and  
5 compared it to latent-class models using the eight-year pedestrian-injury database in Chicago  
6 (Behnood & Mannering, 2016a). A similar modeling approach has shown promising application  
7 in other studies as well (for examples see (Anastasopoulos & Mannering, 2011; Aziz, Ukkusuri,  
8 & Hasan, 2013; Behnood & Mannering, 2015, 2016b; Cerwick, Gkritza, Shaheed, & Hans, 2014;  
9 Malyshkina, Mannering, & Tarko, 2009; Manner & Wunsch-Ziegler, 2013; Moore, Schneider IV,  
10 Savolainen, & Farzaneh, 2011; Ye & Lord, 2011)).

11 More recently, non-parametric modeling methods such as SVM, Artificial Neural  
12 Networks, and Decision Table/Naïve Bayes (C. Chen, Zhang, Yang, Milton, & Alcántara, 2016)  
13 have become popular in the crash analysis studies. Non-parametric models are not built on the  
14 assumptions made based on the distribution properties of the data, which is usually the base for  
15 parametric models. An advantage of non-parametric models over parametric models is that they  
16 do not require the pre-defined relationship between the dependent and explanatory variables (Z.  
17 Li, Liu, Wang, & Xu, 2012).

18 SVM, a supervised machine learning technique, is one of the methods of  
19 classification/regression used in many different transportation and traffic safety-related areas (for  
20 examples see (Balali & Golparvar-Fard, 2016; Jahangiri & Rakha, 2015; Jahangiri, Rakha, &  
21 Dingus, 2016; X. Li, Lord, Zhang, & Xie, 2008; C.-H. Wu, Ho, & Lee, 2004)). Chen et al. (S.  
22 Chen, Wang, & van Zuylen, 2009) demonstrated the application of SVM as an incident detection  
23 tool to identify traffic incidents that reduce the capacity of the road (S. Chen et al., 2009). In the



1 context of active traffic management, Yu and Abdel-Aty (Yu & Abdel-Aty, 2013) showed the  
2 application of SVM in real-time risk analysis to predict crash occurrences (Yu & Abdel-Aty,  
3 2013). Studying 326 freeway sections around the state of Florida, Li et al. (Z. Li et al., 2012)  
4 employed SVM and Ordered Probit (OP) to predict the injury severity of individual crashes. Using  
5 the Radial Basic Function (RBF), they indicated that the SVM model performed better than the  
6 OP in terms of the percent of correct predictions. The result was achieved by comparing the two  
7 models with multi-class response (five injury-severity levels). It has also been demonstrated that  
8 classification results for a two-level SVM also resulted in a significant improvement in the  
9 prediction accuracy of SVM model (Z. Li et al., 2012).

10 Yu et al. (Yu & Abdel-Aty, 2013) also developed fixed parameter logit, SVM (with RBF),  
11 and random parameter logit for four-year data collected from a mountainous freeway section in  
12 Colorado. Comparing three models, they indicated that SVM and random parameter models  
13 provided a better fit than the fixed parameter logit models. In a more recent study, based on a two-  
14 year crash data from New Mexico, Chen et al. (C. Chen, Zhang, Qian, Tarefder, & Tian, 2016)  
15 studied the application of SVM in mapping the injury severity in rollover crashes. Their result  
16 indicated that the SVM model provided reasonable performance in terms of predicting the injury  
17 severity (C. Chen, Zhang, Qian, et al., 2016). Unlike the study conducted by Li et al. (Z. Li et al.,  
18 2012), where every single crash were taken as a single research unit, Chen et al. has taken each  
19 individual driver/vehicle as the research unit and taken into analysis a variable with the number of  
20 vehicles involved in the crash as an independent variable to the model (C. Chen, Zhang, Qian, et  
21 al., 2016).

22 In this paper, crash databases consisting of information on the crash, environment, vehicles,  
23 and occupants for five consecutive years (2007-2011) in the State of California were integrated.

1 To the best of our knowledge, this vehicle by vehicle integrated database has not been used to  
2 develop crash models. Three approaches for modeling the severity of rear-end crashes, -support  
3 vector machine (SVM), multinomial logit (MNL), and mixed multinomial logit (MMNL)-were  
4 applied to this database and compared. The remainder of this paper is structured as follows. The  
5 next section describes the methodology including data description and methods used to analyze  
6 the data. Then, results of the classification methods are presented and discussed for all three  
7 methods. The last section provides conclusions and future directions.

8

## 9 **2. METHODOLOGY**

### 10 **2.1. DATA DESCRIPTION**

11 The data used in this study was obtained from the Highway Safety Information System  
12 (HSIS) in the state of California for five consecutive years from 2007 to 2011. The database  
13 consists of three tables including the crash database, vehicle database, and road database. The rows  
14 in the crash database are built based on each case of a crash and includes information such as  
15 weather condition and lighting that is common among all vehicles involved in the crash. The  
16 vehicle database includes information specific to each vehicle such as driver's age, sex, and vehicle  
17 type. Finally, the road database is built based on information of each road segment such as a  
18 number of lanes and terrain level. In order to perform a vehicle by vehicle analysis for the purpose  
19 of this study, it was required to attach the information of each crash to every single vehicle involved  
20 in that specific crash, and then add the road information to each vehicle. This task was performed  
21 using a Matlab script.

22 The observations for the rear-end crashes were extracted from the database. The severity-  
23 injury levels in the dataset include 5 levels of (1) Property damage only, (2) Complaint of pain, (3)

1 Other visible injury, (4) Severe injury, and (5) Fatal. To select a subset of the independent variables  
2 for developing the models, all variables were individually examined to determine how well they  
3 classified severity levels using the Area Under the Curve (AUC) measure to rank the variables.  
4 Top variables that provided the best results in predicting the severity level were used in the full  
5 modeling procedure. The selected variables were age, sex, terrain level, weather condition,  
6 lighting, vehicle type, number of lanes, and crash cause. The models were developed on 9468  
7 individual vehicles with 70 percent of the data randomly selected to build the training model and  
8 the remaining 30 percent left as the test data to compare the performance of three models in terms  
9 of their prediction performance.

10

## 11 **2.2. MULTINOMIAL LOGIT (MNL) AND MIXED MULTINOMIAL LOGIT (MMNL)** 12 **REGRESSION**

13 Mixed Multinomial Logit regression, also known as random parameters logit model, is a  
14 generalized form of the Multinomial Logistic regression in which the coefficients of any of the  
15 variables are allowed to vary across the individuals and not be limited to a fixed value.  
16 Consequently, it allows the model to take into account the heterogeneity of the population. For the  
17 standard logistic regression (multinomial logit), the probability of individual  $i$  experiencing the  
18 severity level of  $l$  from the set of severity outcomes  $J$  is (Croissant, 2012):

$$19 \quad P_{il} = \frac{e^{\beta' x_{il}}}{\sum_{j=1}^J e^{\beta' x_{ij}}}$$

20 Here,  $x$  is the factor and  $\beta'$  is the fixed coefficient for all individuals. In the mixed  
21 multinomial logit, each individual has their own coefficient  $\beta'_i$ , and probabilities are described as  
22 probability of the individual  $i$ , conditional on the vector of individual-specific coefficient  $\beta_i$ ,  
23 experiencing severity level  $l$  as:

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

$$P_{il} | \beta_i = \frac{e^{\beta_i' x_{il}}}{\sum_{j=1}^J e^{\beta_i' x_{ij}}}$$

Since there are many observations, and finding an individual coefficient for each observation may not be of interest, the coefficients are considered to be random variables and the probabilities of each individual  $i$  is found conditional on the vector of random coefficients  $\beta$ . Later, the average of the probabilities for all values of  $\beta$  is found to obtain the unconditional probability. Given that  $\beta_i$  has the density of  $f(\beta, \theta)$  ( $\theta$  as distribution parameters of  $\beta$ ), for one individual coefficient, the unconditional probability for individual  $i$  experiencing injury level of  $l$  is:

$$P_{il} = \int_{\beta} (P_{il} | \beta_i) f(\beta | \theta) d\beta$$

Here, the function  $f$  indicates the density function for  $\beta$  with  $\theta$  defining the parameters of the density function. Solving this integral becomes more complicated when there is more than one parameter, which requires defining a separate  $\beta$  for each of the random variables. This necessitates simulation techniques. In this study, ‘mlogit’ package was utilized through the R software to perform the mixed logit analysis. Using 200 draws was found to be sufficient for results to converge. For more detail on how the simulation process is performed by the package and how to define simulation parameters, we refer the readers to (Croissant, 2012; K. Train & Croissant, 2012; K. E. Train, 2009).

### 2.3. BINARY LOGIT

Binary Logit is the simple form of the multinomial logit (detailed above) with two instead of multiple outcomes. In this study, we build a binary model with the set of same variables for each of the severity levels. In each model, the first level is whether a specific severity occurs and the second level is whether that specific severity does not. For example, regarding the property damage

1 only level, the first outcome is when the severity was property damage only, and the second  
2 outcome is when the severity is not property damage only. In other words, the second outcome is  
3 the combination of the 4 remaining levels in the dataset (e.g. all possible outcomes except property  
4 damage only). The procedure is repeated for each level which makes a total of 5 binary models.  
5 It is expected that the binary logit does not yield results as efficient as multinomial logit; however,  
6 comparing results for all 5 models, we can investigate how disparate each of severity levels is from  
7 the rest of population.

8

#### 9 **2.4. SUPPORT VECTOR MACHINES (SVM)**

10 In addition to multinomial logit and mixed multinomial logit, SVM, a supervised machine  
11 learning algorithm, was employed to classify severity levels. SVM is known as a powerful method  
12 of classification (and also regression) problems as it tries to find the best possible decision  
13 boundaries between different classes. In model development, SVM applies the function  $\phi(\cdot)$  to  
14 transform the data from  $X$  space into some  $Z$  space. This transformation rearranges the data in such  
15 a way that the classification becomes an easier task. SVM was first introduced in (Boser, Guyon,  
16 & Vapnik, 1992) for separable data and was further expanded in (Cortes & Vapnik, 1995) for non-  
17 separable data. The SVM objective function maximizes the margin between different classes and  
18 at the same time accepts some errors that can be regulated through a penalty parameter. The  
19 formulation of SVM is as follows.

20

$$\min_{w,b,\xi} \left( \frac{1}{2} w^T w + C \sum_{n=1}^N \xi_n \right) \quad (1)$$

Subject to:

$$y_n(w^T \phi(x_n) + b) \geq 1 - \xi_n, n = 1, \dots, N \quad (2)$$

$$\xi_n \geq 0, n = 1, \dots, N \quad (3)$$

Where,

$w$	Parameters to define decision boundary between classes
$C$	Regularization (or penalty) parameter
$\xi_n$	Error parameter to denote margin violation
$b$	Intercept associated with decision boundaries
$\phi(x_n)$	Function to transform data from X space into some Z space
$y_n$	Target value for the $n^{th}$ observation

1 To develop a multi-class classification model using SVM, one versus one approach based on a  
 2 voting strategy has been employed; models were developed using only two classes (e.g. fatal vs  
 3 severe injury, fatal vs other visible injury, and so forth). To predict a new observation, all these  
 4 models are used to produce votes for different classes. The class with the highest vote is identified  
 5 as the predicted class. Since there are five classes for crash severity, combinations of two classes  
 6 from five classes results in ten different combinations, which led to the development of ten models.

7

### 8 **3. RESULTS**

#### 9 **3.1. MULTINOMIAL LOGIT AND MIXED MULTINOMIAL LOGIT REGRESSION**

10 Table 1 and Table 2, respectively, represent the estimation results for the coefficient of  
 11 each variable in the MNL and MMNL models as well as overall performance. Property damage  
 12 only severity level is selected as the base of comparison in the analysis and all the other severity

1 levels are compared to it. Levels of each independent variable are shown below the name of each  
2 variable with the first level representing the base level.

### 3 *Selecting Random Parameters*

4 To find the variables in which the coefficients vary across the individuals (random  
5 parameters in MMNL), all independent variables were checked to determine how strong the  
6 hypothesis of coefficients varying across individuals is. In the first run, all variables were put in  
7 the model assuming that the distribution of their coefficients was normal. The existence of  
8 variability among the coefficient (having the random effect) was signaled by a statistically  
9 significant variable that also had a standard deviation that is statistically different from zero (p-  
10 value lower than 0.05 in this study). For instance, according to Table 2, the coefficient associated  
11 to “4 lanes” level of the variable “Number of lanes” was found to be a significant factor for severe  
12 injury crashes and have a high standard deviation that is statistically different from zero (p-value  
13 = 0.00). Examining the normal distribution of coefficient associated with this variable, it was found  
14 that almost 74 percent of the sample has a negative coefficient and the rest are positive. This, in  
15 fact, indicates that this factor decreases the probability of crash resulting in a severe injury for 74  
16 percent of the population and increases the probability for the remaining 26 percent of the  
17 population. In a similar manner, significant variables with random effect were selected to remain  
18 in the model as variables with random effects and the rest of the variables that did not show the  
19 mixed effect w considered to have fixed coefficients.

### 20 *Variables Significance*

21 Only a few factors including, driving along a dark road with street light, driving on 4 lanes and 6  
22 lanes roads, and speeding show the heterogeneity effect. Results from MNL and MMNL are very  
23 similar with many variables showing significance in both MNL and MMNL. However, in order to

1 discuss the random effects as well, we used results found from MMNL (Table 2) as the basis for  
2 variable significance analysis.

3 Several variables were found to have a significant impact on different levels of severity. The results  
4 for different age groups show that aging has an increasing effect on the probability of getting killed  
5 in a crash. This, in fact, is consistent with the results from (Yasmin, Eluru, Pinjari, & Tay, 2014)  
6 in which they maintain the increase in the risk of fatal crashes for drivers of age 65 or above  
7 (compared to other age groups) for different types of crashes. Similar results were also achieved  
8 by (Q. Wu et al., 2014) in a study on the contribution of a variety of factors on different levels of  
9 injury in multi-vehicle crashes, where the probability of fatality was increased for drivers of 65  
10 years or older. Some other studies by (Kim et al., 2013; Xie, Zhang, & Liang, 2009; Yasmin &  
11 Eluru, 2013) have also confirmed the same finding.

12         Regarding the drivers' gender, it was found that compared to female drivers, male drivers  
13 are less likely to complain about pain after a crash. An explanation might be that men are physically  
14 stronger. There is also evidence that men are more likely to be killed compared to female drivers  
15 which raise the idea that male drivers might be more aggressive drivers than female drivers. This  
16 is in line with the finding by (Shankar, Mannering, & Barfield, 1996) where they maintain greater  
17 probability of fatal and disabling injury for crashes involved male drivers. Different results,  
18 however, have been concluded from the gender analysis. For instance, (Q. Wu et al., 2014) found  
19 that in multi-vehicle crashes, the coefficient varies across genders as for some female drivers it  
20 increases the probability of a fatal crash and for some, it reduces the probabilities. On the other  
21 hand, in some studies such as (Abdel-Aty, 2003; Xie et al., 2009), it was found that female drivers  
22 are more likely to get killed in a crash with the same circumstances.



1            Investigating driving on different terrain levels, the outcome of a crash is more likely to  
2 result in a severe injury or fatality when driving on mountainous and rolling terrains compared to  
3 flat terrains. This was expected since driving on mountainous or rolling terrains is usually more  
4 complex than driving on the flat terrains. This also confirms the results from (Shankar et al., 1996),  
5 where they indicated that high proportion of horizontal curves was found to increase the likelihood  
6 of a possible injury crash. As for driving in different weather condition, it was found that driving  
7 in foggy condition increases the probability of crashes leading to fatality or severe injuries  
8 compared to clear weather condition.

9            Results from the lighting condition of the street indicated that driving in a low light  
10 condition such as dusk or dark without street lights increases the probability of a fatal crash  
11 compared to driving in daylight. This is in line with the results from (Kim et al., 2013; Q. Wu et  
12 al., 2014; Xie et al., 2009), where they also found that undesirable lighting increases the propensity  
13 of a crash to be more severe or fatal. In this study; however, driving at night when there is street  
14 light found to have mixed effect. Evaluation of the distribution of the coefficient shows the  
15 existence of a statistically high standard deviation of 1.83 (p-value = 0.00) and a mean of 0.55 for  
16 the coefficients of the “Dark - street light” level of variable “Light”. This might be associated to  
17 the fact that some drivers tend to drive more carefully when their range of view is limited due to  
18 the lack of proper lighting while some drivers may not feel the need to compensate in this situation.

19            Comparing different vehicle types, it was found that motorcycle riders and truck drivers  
20 are more prone to higher levels of severity and fatality and less likely to experience property  
21 damage and surface injuries compared to passenger cars.

22            Increase in the number of lanes for most of the levels found to have a negative effect on  
23 the log odds of the crash to become fatal or severe. This, in fact, raises the idea that crashes are

1 more likely to have higher levels of severity on two-lane roads than roads and highways with a  
2 higher number of lanes.

3           Regarding crash causes, behavior such as following too closely and speeding are less likely  
4 to lead to fatality or severe injury compared to the situation where the driver is under influence of  
5 the alcohol. A similar result was achieved in other studies such as (Xie et al., 2009; Yasmin &  
6 Eluru, 2013).

7

### 8           **3.2. BINARY LOGIT**

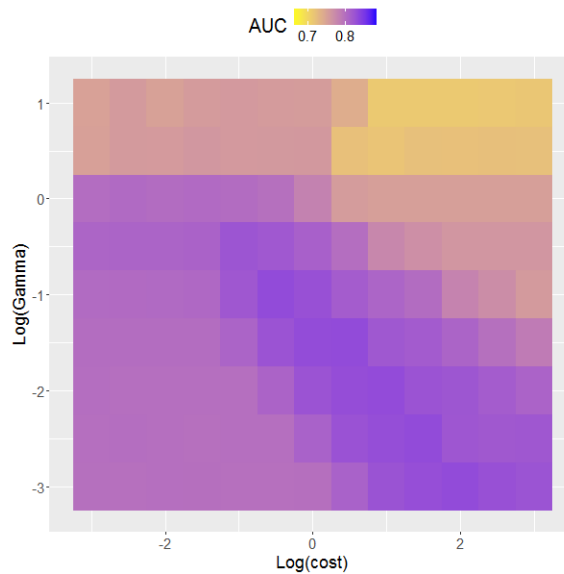
9 Table 3 demonstrates the summary of the results for the binary logit model. Results from both  
10 Table 1 and Table 3 indicate the same effects from selected factors in different ways. For example,  
11 looking at the motorcycle data, from Table 1 (multinomial approach), we see that motorcycle  
12 drivers are more prone to being killed in accidents compared to passenger car drivers. The same  
13 result could be ascertained by comparing the estimated coefficients/significances retrieved from  
14 each of the five binary logits. Starting from the first binary model (property damage only vs. all  
15 levels) there is a low negative coefficient that increases and eventually becomes positive as you  
16 move toward the fifth binary model (fatal vs. all levels). The same concept applies to other  
17 variables as well.

18

### 19           **3.3. SUPPORT VECTOR MACHINE (SVM)**

20           The same train and test datasets that were used for the MNL and MMNL models were  
21 applied to develop and validate the SVM model using Gaussian kernels. All attributes (variables)  
22 were scaled as SVM usually performs poorly without feature scaling. This is a simple  
23 preprocessing step that could significantly impact the SVM model performance. In order to find

1 the optimal SVM model, two parameters, the regularization parameter, and the Gamma parameter,  
2 must be tuned. This parameter tuning was conducted for all ten models, one of which is shown in  
3 Fig. 1. Focusing on the fatal vs property damage only model, Fig. 1 illustrates how different  
4 parameter settings impact the area under curve (AUC) value. This figure shows a grid search of  
5 the Gamma and Cost parameters to achieve the best possible performance of one of the SVM  
6 models. The best AUC was achieved with cost and Gamma parameters being equal to 0.32 and  
7 0.1, respectively. The darker regions in the figure represent higher values of AUC, which  
8 correspond to better performance. The optimal performance of the other SVM models was  
9 achieved in the same fashion.



10  
11 Figure 1 SVM parameter tuning for the fatal vs property damage only.

#### 12 4. DISCUSSION OF PREDICTION RESULTS

13 In this section, the SVM, MNL (both binary and multi-class), and MMNL are compared in  
14 terms of their predicting accuracy on the test data. Table 4 demonstrates the number of true and  
15 false predictions in each model as well as for each class.

1           While both binary and multinomial logit models indicate similar information regarding the  
2 effect of factors on the outcome of the accident, they do not maintain the same accuracy when used  
3 for prediction. Based on Table 4, the binary models built for the first three levels (property damage  
4 only, complaint of pain, and other visible injury) did not maintain good prediction results based  
5 on the test data, where every single prediction for the desired level (e.g. property damage only in  
6 the first model) is incorrect. Although binary models built for last two levels (severity and fatal)  
7 maintained better prediction results, the accuracy was still not very good. While the prediction  
8 results from binary models display estimation inefficiency, it also demonstrates how discrete the  
9 last two levels (severity and fatal) are from the other levels.

10           Regression and support vector approaches with all five levels in the model led to more  
11 accurate predictions. The prediction results are also included in Table 4 for comparison. Both MNL  
12 and MMNL models maintained a comparable accuracy, which is expected since most parameters  
13 were not found to have random effects. Yet it is important that random parameters are included to  
14 better report the effect of these variables. It should be noted that most machine learning algorithms  
15 such as SVM are considered as black-box methods that are difficult to interpret. In general, they  
16 could perform better than classical statistical methods such as logit models but have limited  
17 explanatory power. In the present study, the SVM performed slightly better in specific cases as  
18 explained below.

19           According to Table 4, the SVM has maintained greater accuracy for “other visible injury”  
20 and “fatal” class of the database. On the other hand, MNL and MMNL had greater accuracy for  
21 “property damage only”, “complaint of pain”, and “severe injury”. Although, in total, the number  
22 of true predictions by the SVM is slightly higher than the total for MNL and MMNL.

1 More detailed results regarding the prediction of each model in the form of confusion  
2 matrix are shown in Table 5. According to Table 5, the largest number of true predictions belongs  
3 to “Complaint of pain” level. On the other hand, the lowest number of true predictions belong to  
4 “Other visible injury” level. This is a broad, and somewhat vague, a category which is overlapping  
5 with the “Complaint of pain” category, so it is expected that it would be hard to predict. The  
6 remaining “Property damage only”, “Severe injury”, and “Fatal” levels have similar accuracy  
7 levels. The last row under each confusion matrix indicates the percent of true prediction for each  
8 individual level of severity.

9  
10

## 11 **5. CONCLUSION**

12 Considering the significant losses resulting from crashes, predicting the likelihood of crash  
13 severity has been a vital line of research. Studying crashes and their contributing factors allow  
14 traffic engineers and practitioners to determine and prioritize dangerous crash prone situations,  
15 which in turn aids efforts to implement suitable countermeasures and enhance the level of safety.  
16 In this study, crash severity models were developed using three statistical approaches of MNL  
17 (binary and multi-class), MMNL, and SVM. The database used for the model training is an  
18 integrated database of the California crash database from 2007 to 2011 including information on  
19 driver, vehicle, and environment.

20 Studying effects of different factors on severity of accidents, results from both MNL and  
21 MMNL models showed that older and male drivers are more likely to have a fatal crash while  
22 complain of pain is observed more among female drivers compared to males. As expected, driving  
23 in foggy and mountainous or rolling terrains is associated with more fatal and severe injury crashes.

1 Similarly, low light conditions are associated to more fatal crashes. Regarding the vehicle type,  
2 motorcycle riders and truck drivers are more likely to experience a higher level of crash severity.  
3 Also, a fewer number of lanes was found to be associated with higher level of severity. Regarding  
4 the cause of crashes, improper turning and driving under influence were found to cause severe  
5 injuries and fatal crashes. In terms of random parameters, factors including driving on dark roads  
6 with a street light, driving on 4 and 6 lane roads, and speeding could have various impact on the  
7 severity outcome of accidents as for some people they have been a contributing factor and for the  
8 other, they have acted as the opposite.

9 Studying the performance of the models based on their prediction accuracy, the MNL and  
10 MMNL maintained comparable results. This result was expected since only a few of the variables  
11 were found to have a mixed effect, and those that found to have mixed effect were not dominant  
12 enough to make a significant change in the overall outcome of the model. The SVM approach  
13 maintained a slightly better prediction accuracy using the test data.

14 In the future, the crash analysis will be conducted for other crash types than rear-end  
15 crashes. Also, other techniques such as neural network and discriminant analysis method will be  
16 considered on the same dataset and results will be compared in terms of model performance.

17

## 18 **ACKNOWLEDGEMENT**

19 The authors would like to thank the Highway Safety Information System (HSIS) for their  
20 support and provision of the crash dataset. The authors wish to also thank the University Grants  
21 Program (UGP) at San Diego State University for partially sponsoring this study.

22

23

1 **REFERENCES**

- 2 Abdel-Aty, M. (2003). Analysis of driver injury severity levels at multiple locations using  
3 ordered probit models. *Journal of Safety Research*, 34(5), 597–603.  
4 <https://doi.org/10.1016/j.jsr.2003.05.009>
- 5 Anastasopoulos, P. C., & Mannering, F. L. (2011). An empirical assessment of fixed and random  
6 parameter logit models using crash- and non-crash-specific injury data. *Accident Analysis  
7 & Prevention*, 43(3), 1140–1147. <https://doi.org/10.1016/j.aap.2010.12.024>
- 8 Aziz, H. M. A., Ukkusuri, S. V., & Hasan, S. (2013). Exploring the determinants of pedestrian–  
9 vehicle crash severity in New York City. *Accident Analysis & Prevention*, 50, 1298–  
10 1309. <https://doi.org/10.1016/j.aap.2012.09.034>
- 11 Balali, V., & Golparvar-Fard, M. (2016). Evaluation of Multiclass Traffic Sign Detection and  
12 Classification Methods for U.S. Roadway Asset Inventory Management. *Journal of  
13 Computing in Civil Engineering*, 30(2), 04015022.  
14 [https://doi.org/10.1061/\(ASCE\)CP.1943-5487.0000491](https://doi.org/10.1061/(ASCE)CP.1943-5487.0000491)
- 15 Behnood, A., & Mannering, F. L. (2015). The temporal stability of factors affecting driver-injury  
16 severities in single-vehicle crashes: Some empirical evidence. *Analytic Methods in  
17 Accident Research*, 8, 7–32. <https://doi.org/10.1016/j.amar.2015.08.001>
- 18 Behnood, A., & Mannering, F. L. (2016a). An empirical assessment of the effects of economic  
19 recessions on pedestrian-injury crashes using mixed and latent-class models. *Analytic  
20 Methods in Accident Research*, 12, 1–17. <https://doi.org/10.1016/j.amar.2016.07.002>
- 21 Behnood, A., & Mannering, F. L. (2016b). The effects of drug and alcohol consumption on  
22 driver injury severities in single-vehicle crashes. *Traffic Injury Prevention*, 18(5), 456–  
23 462. <https://doi.org/10.1080/15389588.2016.1262540>

- 1 Boser, B. E., Guyon, I. M., & Vapnik, V. N. (1992). A Training Algorithm for Optimal Margin  
2 Classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning*  
3 *Theory* (pp. 144–152). New York, NY, USA: ACM.  
4 <https://doi.org/10.1145/130385.130401>
- 5 Cerwick, D. M., Gkritza, K., Shaheed, M. S., & Hans, Z. (2014). A comparison of the mixed  
6 logit and latent class methods for crash severity analysis. *Analytic Methods in Accident*  
7 *Research*, 3–4, 11–27. <https://doi.org/10.1016/j.amar.2014.09.002>
- 8 Chen, C., Zhang, G., Qian, Z., Tarefder, R. A., & Tian, Z. (2016). Investigating driver injury  
9 severity patterns in rollover crashes using support vector machine models. *Accident*  
10 *Analysis & Prevention*, 90, 128–139. <https://doi.org/10.1016/j.aap.2016.02.011>
- 11 Chen, C., Zhang, G., Tian, Z., Bogus, S. M., & Yang, Y. (2015). Hierarchical Bayesian random  
12 intercept model-based cross-level interaction decomposition for truck driver injury  
13 severity investigations. *Accident Analysis & Prevention*, 85, 186–198.  
14 <https://doi.org/10.1016/j.aap.2015.09.005>
- 15 Chen, C., Zhang, G., Yang, J., Milton, J. C., & Alcántara, A. “Dely.” (2016). An explanatory  
16 analysis of driver injury severity in rear-end crashes using a decision table/Naïve Bayes  
17 (DTNB) hybrid classifier. *Accident Analysis & Prevention*, 90, 95–107.  
18 <https://doi.org/10.1016/j.aap.2016.02.002>
- 19 Chen, S., Wang, W., & van Zuylen, H. (2009). Construct support vector machine ensemble to  
20 detect traffic incident. *Expert Systems with Applications*, 36(8), 10976–10986.  
21 <https://doi.org/10.1016/j.eswa.2009.02.039>
- 22 Cortes, C., & Vapnik, V. (1995). Support-Vector Networks. *Machine Learning*, 20(3), 273–297.  
23 <https://doi.org/10.1023/A:1022627411411>



- 1 Croissant, Y. (2012). Estimation of multinomial logit models in R: The mlogit Packages. *R*  
2 *Package Version 0.2-2*. URL: [Http://Cran. R-Project.](http://Cran.R-Project.Org/Web/Packages/Mlogit/Vignettes/Mlogit.Pdf)  
3 *Org/Web/Packages/Mlogit/Vignettes/Mlogit. Pdf*. Retrieved from  
4 <ftp://193.1.193.75/disk1/cran.r-project.org/web/packages/mlogit/vignettes/mlogit.pdf>
- 5 Fatality Analysis Reporting System (FARS) Encyclopedia, National Highway Traffic Safety  
6 Administration (NHTSA). (n.d.). Retrieved May 1, 2017, from [http://www-](http://www-fars.nhtsa.dot.gov/Main/index.aspx)  
7 [fars.nhtsa.dot.gov/Main/index.aspx](http://www-fars.nhtsa.dot.gov/Main/index.aspx)
- 8 Jahangiri, A., & Rakha, H. A. (2015). Applying machine learning techniques to transportation  
9 mode recognition using mobile phone sensor data. *IEEE Transactions on Intelligent*  
10 *Transportation Systems, 16*(5), 2406–2417.
- 11 Jahangiri, A., Rakha, H., & Dingus, T. A. (2016). Red-light running violation prediction using  
12 observational and simulator data. *Accident Analysis & Prevention*. Retrieved from  
13 <http://www.sciencedirect.com/science/article/pii/S0001457516302056>
- 14 Kim, J.-K., Ulfarsson, G. F., Kim, S., & Shankar, V. N. (2013). Driver-injury severity in single-  
15 vehicle crashes in California: A mixed logit analysis of heterogeneity due to age and  
16 gender. *Accident Analysis & Prevention, 50*, 1073–1081.  
17 <https://doi.org/10.1016/j.aap.2012.08.011>
- 18 Li, X., Lord, D., Zhang, Y., & Xie, Y. (2008). Predicting motor vehicle crashes using Support  
19 Vector Machine models. *Accident Analysis & Prevention, 40*(4), 1611–1618.  
20 <https://doi.org/10.1016/j.aap.2008.04.010>
- 21 Li, Z., Liu, P., Wang, W., & Xu, C. (2012). Using support vector machine models for crash  
22 injury severity analysis. *Accident Analysis & Prevention, 45*, 478–486.  
23 <https://doi.org/10.1016/j.aap.2011.08.016>

- 1 Li, Z., Wang, W., Liu, P., Bigham, J. M., & Ragland, D. R. (2013). Using Geographically  
2 Weighted Poisson Regression for county-level crash modeling in California. *Safety*  
3 *Science*, 58, 89–97. <https://doi.org/10.1016/j.ssci.2013.04.005>
- 4 Lord, D., & Mannering, F. (2010). The statistical analysis of crash-frequency data: A review and  
5 assessment of methodological alternatives. *Transportation Research Part A: Policy and*  
6 *Practice*, 44(5), 291–305. <https://doi.org/10.1016/j.tra.2010.02.001>
- 7 Malyshkina, N. V., Mannering, F. L., & Tarko, A. P. (2009). Markov switching negative  
8 binomial models: An application to vehicle accident frequencies. *Accident Analysis &*  
9 *Prevention*, 41(2), 217–226. <https://doi.org/10.1016/j.aap.2008.11.001>
- 10 Manner, H., & Wunsch-Ziegler, L. (2013). Analyzing the severity of accidents on the German  
11 Autobahn. *Accident Analysis & Prevention*, 57, 40–48.  
12 <https://doi.org/10.1016/j.aap.2013.03.022>
- 13 Mannering, F. L., & Bhat, C. R. (2014). Analytic methods in accident research: Methodological  
14 frontier and future directions. *Analytic Methods in Accident Research*, 1, 1–22.  
15 <https://doi.org/10.1016/j.amar.2013.09.001>
- 16 Mannering, F. L., Shankar, V., & Bhat, C. R. (2016). Unobserved heterogeneity and the  
17 statistical analysis of highway accident data. *Analytic Methods in Accident Research*, 11,  
18 1–16. <https://doi.org/10.1016/j.amar.2016.04.001>
- 19 Milton, J. C., Shankar, V. N., & Mannering, F. L. (2008). Highway accident severities and the  
20 mixed logit model: An exploratory empirical analysis. *Accident Analysis & Prevention*,  
21 40(1), 260–266. <https://doi.org/10.1016/j.aap.2007.06.006>
- 22 Moore, D. N., Schneider IV, W. H., Savolainen, P. T., & Farzaneh, M. (2011). Mixed logit  
23 analysis of bicyclist injury severity resulting from motor vehicle crashes at intersection

1 and non-intersection locations. *Accident Analysis & Prevention*, 43(3), 621–630.  
2 <https://doi.org/10.1016/j.aap.2010.09.015>

3 Morgan, A., & Mannering, F. L. (2011). The effects of road-surface conditions, age, and gender  
4 on driver-injury severities. *Accident Analysis & Prevention*, 43(5), 1852–1863.  
5 <https://doi.org/10.1016/j.aap.2011.04.024>

6 Savolainen, P. T., Mannering, F. L., Lord, D., & Quddus, M. A. (2011). The statistical analysis  
7 of highway crash-injury severities: A review and assessment of methodological  
8 alternatives. *Accident Analysis & Prevention*, 43(5), 1666–1676.  
9 <https://doi.org/10.1016/j.aap.2011.03.025>

10 Shankar, V., Mannering, F., & Barfield, W. (1996). Statistical analysis of accident severity on  
11 rural freeways. *Accident Analysis & Prevention*, 28(3), 391–401.  
12 [https://doi.org/10.1016/0001-4575\(96\)00009-7](https://doi.org/10.1016/0001-4575(96)00009-7)

13 Train, K., & Croissant, Y. (2012). Kenneth Train’s exercises using the mlogit package for R. *R*,  
14 25, 0–2.

15 Train, K. E. (2009). *Discrete Choice Methods with Simulation*. Cambridge University Press.

16 Wu, C.-H., Ho, J.-M., & Lee, D. T. (2004). Travel-time prediction with support vector  
17 regression. *IEEE Transactions on Intelligent Transportation Systems*, 5(4), 276–281.  
18 <https://doi.org/10.1109/TITS.2004.837813>

19 Wu, Q., Chen, F., Zhang, G., Liu, X. C., Wang, H., & Bogus, S. M. (2014). Mixed logit model-  
20 based driver injury severity investigations in single- and multi-vehicle crashes on rural  
21 two-lane highways. *Accident Analysis & Prevention*, 72, 105–115.  
22 <https://doi.org/10.1016/j.aap.2014.06.014>

- 1 Xie, Y., Zhang, Y., & Liang, F. (2009). Crash Injury Severity Analysis Using Bayesian Ordered  
2 Probit Models. *Journal of Transportation Engineering*, 135(1), 18–25.  
3 [https://doi.org/10.1061/\(ASCE\)0733-947X\(2009\)135:1\(18\)](https://doi.org/10.1061/(ASCE)0733-947X(2009)135:1(18))
- 4 Yasmin, S., & Eluru, N. (2013). Evaluating alternate discrete outcome frameworks for modeling  
5 crash injury severity. *Accident Analysis & Prevention*, 59, 506–521.  
6 <https://doi.org/10.1016/j.aap.2013.06.040>
- 7 Yasmin, S., Eluru, N., Pinjari, A. R., & Tay, R. (2014). Examining driver injury severity in two  
8 vehicle crashes – A copula based approach. *Accident Analysis & Prevention*, 66, 120–  
9 135. <https://doi.org/10.1016/j.aap.2014.01.018>
- 10 Ye, F., & Lord, D. (2011). Investigation of Effects of Underreporting Crash Data on Three  
11 Commonly Used Traffic Crash Severity Models. *Transportation Research Record:  
12 Journal of the Transportation Research Board*, 2241, 51–58.  
13 <https://doi.org/10.3141/2241-06>
- 14 Ye, F., & Lord, D. (2014). Comparing three commonly used crash severity models on sample  
15 size requirements: multinomial logit, ordered probit and mixed logit models. *Analytic  
16 Methods in Accident Research*, 1, 72–85.
- 17 Yu, R., & Abdel-Aty, M. (2013). Utilizing support vector machine in real-time crash risk  
18 evaluation. *Accident Analysis & Prevention*, 51, 252–259.  
19 <https://doi.org/10.1016/j.aap.2012.11.027>  
20

1 Table 1 - Summary of results for multinomial logit model

Variable	Complaint of pain				Other visible injuries				Severe injury				Fatal			
	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )
:(intercept)	0.26	0.31	0.84	0.40	1.39	0.27	5.11	0.00***	1.94	0.26	7.40	0.00***	1.22	0.28	4.44	0.00***
<b>Age</b>																
:Young adult																
:Adult	0.06	0.10	0.58	0.56	0.04	0.10	0.45	0.65	0.02	0.10	0.16	0.88	0.13	0.12	1.09	0.27
:Middle aged	0.03	0.11	0.31	0.76	0.11	0.11	1.03	0.30	0.09	0.12	0.78	0.44	0.27	0.13	2.10	0.04*
:Old	0.03	0.18	0.18	0.86	0.16	0.18	0.86	0.39	0.20	0.19	1.03	0.30	0.63	0.20	3.09	0.00**
<b>Sex</b>																
:Female																
:Male	-0.30	0.08	-3.70	0.00***	-0.06	0.08	-0.66	0.51	0.10	0.09	1.14	0.26	0.25	0.10	2.47	0.01*
<b>Terrain</b>																
:Flat																
:Mountainous	0.06	0.21	0.26	0.79	0.35	0.20	1.76	0.08 <sup>ˆ</sup>	0.42	0.20	2.08	0.04*	0.29	0.22	1.31	0.19
:Rolling	0.11	0.09	1.14	0.26	0.07	0.10	0.76	0.45	0.16	0.10	1.64	0.10 <sup>ˆ</sup>	0.34	0.11	3.24	0.00**
<b>Weather</b>																
:Clear																
:Cloudy	0.11	0.11	1.08	0.28	-0.19	0.11	-1.67	0.09 <sup>ˆ</sup>	-0.21	0.12	-1.76	0.08 <sup>ˆ</sup>	-0.15	0.13	-1.16	0.25
:Raining	-0.02	0.24	-0.10	0.92	0.12	0.23	0.53	0.60	-0.27	0.26	-1.04	0.30	-0.36	0.28	-1.28	0.20
:Fog	0.89	0.84	1.06	0.29	0.62	0.87	0.71	0.48	2.33	0.75	3.10	0.00**	2.22	0.77	2.87	0.00**
<b>Light</b>																
:Daylight																
:Dusk - Dawn	-0.43	0.22	-1.93	0.05 <sup>ˆ</sup>	-0.31	0.22	-1.38	0.17	0.14	0.22	0.66	0.51	0.87	0.22	3.93	0.00***
:Dark - Street Lights	-0.01	0.12	-0.08	0.93	-0.08	0.12	-0.65	0.52	0.57	0.12	4.81	0.00***	1.06	0.13	8.29	0.00***
:Dark - No Street Lights	-0.13	0.14	-0.92	0.36	0.29	0.13	2.17	0.03*	1.00	0.13	7.78	0.00***	1.78	0.13	13.56	0.00***
<b>Vehicle type</b>																
:Passenger Car																
:Truck	-0.17	0.24	-0.69	0.49	0.36	0.21	1.72	0.09 <sup>ˆ</sup>	1.13	0.19	5.88	0.00***	1.62	0.19	8.40	0.00***
:Motorcycle	1.72	0.78	2.21	0.03*	3.47	0.72	4.80	0.00***	4.73	0.72	6.62	0.00***	4.74	0.72	6.60	0.00***
:Pickup Truck	-0.01	0.11	-0.05	0.96	0.02	0.11	0.17	0.87	0.08	0.12	0.71	0.48	-0.04	0.13	-0.30	0.76
:Bus	1.04	1.16	0.90	0.37	1.11	1.16	0.96	0.34	2.39	1.07	2.23	0.03*	2.06	1.14	1.80	0.07 <sup>ˆ</sup>
:Emergency Vehicle	0.22	0.55	0.40	0.69	-0.96	0.83	-1.17	0.24	-0.44	0.73	-0.60	0.55	-0.51	0.81	-0.63	0.53
<b>Number of lanes</b>																
:2																
:3	0.46	0.71	0.64	0.52	0.38	0.72	0.53	0.59	-1.00	0.85	-1.18	0.24	0.60	0.70	0.86	0.39

Variable	Complaint of pain				Other visible injuries				Severe injury				Fatal			
	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )
:4	-0.09	0.20	-0.46	0.65	-0.05	0.20	-0.25	0.80	-0.51	0.20	-2.57	0.01*	-0.79	0.21	-3.82	0.00***
:5	0.33	0.35	0.92	0.36	0.28	0.37	0.76	0.45	-0.18	0.36	-0.49	0.63	-0.52	0.38	-1.39	0.16
:6	-0.55	0.20	-2.76	0.01**	-0.36	0.20	-1.78	0.07*	-1.15	0.20	-5.81	0.00***	-1.34	0.21	-6.44	0.00***
:7	-0.76	0.30	-2.53	0.01*	-0.32	0.29	-1.11	0.27	-0.84	0.29	-2.95	0.00**	-1.31	0.32	-4.13	0.00***
:8	-0.26	0.19	-1.39	0.16	-0.34	0.20	-1.73	0.08*	-0.86	0.19	-4.56	0.00***	-1.22	0.20	-6.14	0.00***
:> 8	-0.46	0.19	-2.48	0.01*	-0.27	0.19	-1.44	0.15	-1.01	0.18	-5.46	0.00***	-1.33	0.19	-6.85	0.00***
<b>Cause</b>																
:Under influence																
:Following too closely	0.35	0.31	1.15	0.25	-1.40	0.28	-5.05	0.00***	-2.51	0.32	-7.94	0.00***	-3.07	0.43	-7.14	0.00***
:Improper Turn	0.34	0.40	0.86	0.39	-0.81	0.35	-2.31	0.02*	-0.45	0.32	-1.42	0.16	0.16	0.32	0.51	0.61
:Speeding	0.26	0.24	1.04	0.30	-1.31	0.19	-6.91	0.00***	-1.90	0.18	-10.45	0.00***	-1.80	0.19	-9.67	0.00***
:Other Violations	-0.01	0.29	-0.05	0.96	-1.30	0.24	-5.50	0.00***	-1.61	0.23	-7.00	0.00***	-2.00	0.26	-7.84	0.00***
:Other Than Driving	0.90	0.90	1.00	0.32	-0.75	0.93	-0.80	0.43	-0.32	0.84	-0.38	0.70	1.69	0.77	2.20	0.03*
:Unknown	-0.77	0.64	-1.19	0.23	-2.46	0.63	-3.87	0.00***	-3.35	0.72	-4.64	0.00***	-1.84	0.55	-3.33	0.00***
<b>Model overall</b>																
Log-Likelihood: -9693.1 McFadden R <sup>2</sup> : 0.09 Likelihood ratio test : chisq = 1894.7 (p.value = < 2.22e-16)																

1

1 Table 2 - Summary of results for mixed multinomial logit model

Variable	Complaint of pain				Other visible injuries				Severe injury				Fatal			
	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )
:(intercept)	0.25	0.31	0.82	0.41	1.37	0.32	4.33	0.00***	1.97	0.27	7.34	0.00***	1.29	0.29	4.48	0.00***
<b>Age</b>																
:Young adult																
:Adult	0.06	0.10	0.61	0.54	0.07	0.13	0.55	0.58	0.02	0.11	0.18	0.86	0.17	0.13	1.30	0.19
:Middle aged	0.04	0.11	0.33	0.74	0.16	0.14	1.11	0.27	0.09	0.12	0.73	0.46	0.33	0.14	2.31	0.02*
:Old	0.03	0.19	0.17	0.86	0.10	0.24	0.40	0.69	0.18	0.21	0.89	0.37	0.75	0.22	3.34	0.00**
<b>Sex</b>																
:Female																
:Male	-0.30	0.08	-3.68	0.00***	-0.08	0.11	-0.75	0.46	0.13	0.10	1.35	0.18	0.28	0.11	2.48	0.01*
<b>Terrain</b>																
:Flat																
:Mountainous	0.06	0.22	0.26	0.80	0.46	0.25	1.81	0.07'	0.47	0.22	2.15	0.03*	0.29	0.23	1.27	0.20
:Rolling	0.11	0.10	1.15	0.25	0.02	0.12	0.19	0.85	0.13	0.11	1.24	0.22	0.40	0.12	3.43	0.00***
<b>Weather</b>																
:Clear																
:Cloudy	0.11	0.11	1.04	0.30	-0.26	0.15	-1.74	0.08'	-0.19	0.13	-1.49	0.14	-0.17	0.14	-1.17	0.24
:Raining	-0.04	0.24	-0.15	0.88	0.26	0.30	0.87	0.38	-0.26	0.28	-0.95	0.34	-0.32	0.30	-1.05	0.30
:Fog	0.91	0.85	1.07	0.28	0.30	1.00	0.30	0.77	2.59	0.80	3.23	0.00**	2.17	0.85	2.55	0.01*
<b>Light</b>																
:Daylight																
:Dusk - Dawn	-0.43	0.23	-1.89	0.06'	-0.39	0.29	-1.33	0.18	0.20	0.23	0.89	0.37	1.00	0.23	4.29	0.00***
:Dark - Street Lights	-0.01	0.12	-0.11	0.91	-0.22	0.17	-1.28	0.20	0.61	0.13	4.75	0.00***	0.55	0.36	1.51	0.13
:Dark - No Street Lights	-0.13	0.15	-0.87	0.38	0.20	0.17	1.16	0.25	1.07	0.14	7.68	0.00***	1.91	0.15	13.04	0.00***
<b>Vehicle type</b>																
:Passenger Car																
:Truck	-0.16	0.24	-0.64	0.52	0.32	0.26	1.24	0.21	1.23	0.22	5.71	0.00***	1.81	0.22	8.34	0.00***
:Motorcycle	1.73	0.79	2.19	0.03*	3.43	0.74	4.61	0.00***	4.94	0.73	6.78	0.00***	4.88	0.73	6.67	0.00***
:Pickup Truck	0.00	0.11	-0.02	0.98	-0.01	0.15	-0.05	0.96	0.09	0.13	0.68	0.50	-0.08	0.15	-0.53	0.60
:Bus	1.00	1.19	0.84	0.40	1.00	1.35	0.73	0.46	2.58	1.20	2.15	0.03*	2.18	1.27	1.72	0.09'
:Emergency Vehicle	0.23	0.55	0.42	0.67	-1.40	1.24	-1.13	0.26	-0.54	0.86	-0.63	0.53	-0.81	1.20	-0.67	0.50
<b>Number of lanes</b>																
:2																
:3	0.45	0.72	0.63	0.53	0.63	0.81	0.77	0.44	-1.01	0.88	-1.15	0.25	0.66	0.68	0.97	0.33

Variable	Complaint of pain				Other visible injuries				Severe injury				Fatal			
	Parameter Estimate	Standard Error	t-value	Pr(< t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(< t )	Parameter Estimate	Standard Error	t-value	Pr(> t )
:4	-0.09	0.20	-0.45	0.66	0.03	0.27	0.11	0.91	-1.05	0.36	-2.92	0.00**	-1.28	0.30	-4.21	0.00***
:5	0.33	0.36	0.92	0.36	0.35	0.45	0.77	0.44	-0.19	0.37	-0.52	0.61	-0.72	0.42	-1.72	0.09'
:6	-0.54	0.20	-2.71	0.01**	-0.25	0.27	-0.92	0.36	-1.82	0.39	-4.62	0.00***	-1.82	0.32	-5.65	0.00***
:7	-0.75	0.30	-2.51	0.01*	-0.14	0.38	-0.37	0.71	-0.87	0.30	-2.88	0.00**	-1.37	0.34	-4.07	0.00***
:8	-0.26	0.19	-1.36	0.17	-0.27	0.25	-1.06	0.29	-0.89	0.19	-4.70	0.00***	-1.30	0.20	-6.47	0.00***
:> 8	-0.46	0.19	-2.44	0.01*	-0.12	0.26	-0.47	0.64	-1.04	0.18	-5.64	0.00***	-1.47	0.20	-7.47	0.00***
<b>Cause</b>																
:Under influence																
:Following too closely	0.35	0.31	1.14	0.26	-1.44	0.29	-5.05	0.00***	-2.67	0.35	-7.64	0.00***	-3.28	0.44	-7.49	0.00***
:Improper Turn	0.35	0.40	0.86	0.39	-0.85	0.36	-2.39	0.02*	-0.41	0.34	-1.22	0.22	0.19	0.34	0.55	0.58
:Speeding	0.26	0.25	1.04	0.30	-2.21	0.63	-3.53	0.00***	-1.97	0.19	-10.14	0.00***	-1.95	0.20	-9.60	0.00***
:Other Violations	-0.01	0.29	-0.04	0.97	-1.35	0.24	-5.55	0.00***	-1.63	0.24	-6.73	0.00***	-2.20	0.28	-7.85	0.00***
:Other Than Driving	0.88	0.95	0.92	0.36	-0.81	0.94	-0.86	0.39	-0.57	0.91	-0.63	0.53	1.83	0.79	2.33	0.02*
:Unknown	-0.77	0.64	-1.21	0.23	-2.52	0.66	-3.84	0.00***	-3.53	0.78	-4.56	0.00***	-2.04	0.67	-3.04	0.00**
<b>Random parameters</b>																
Light :Dark - Street Lights													1.83	0.61	2.98	0.00**
Number of lanes :4									1.66	0.59	2.82	0.00**	1.36	0.40	3.41	0.00***
Number of lanes :6									1.66	0.53	3.10	0.00**	1.23	0.44	2.78	0.01**
Cause :Speeding					-2.14	0.92	-2.34	0.02**								
<b>Model overall</b>																
Log-Likelihood: -9684.8																
McFadden R^2: 0.09																
Likelihood ratio test :																
chisq = 1911.3 (p.value = < 2.22e-16)																

1

2



1 Table 3 - Summary of results for binary logit model (each column belongs to one of the five binary logits)

Variable	Property damage only Versus All levels				Complaint of pain Versus All levels				Other visible injuries Versus All levels				Severe injury Versus All levels				Fatal Versus All levels			
	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )
:(intercept)	-2.90	0.23	-12.38	0.00***	-2.51	0.23	-10.71	0.00***	-1.35	0.18	-7.55	0.00***	-0.62	0.16	-3.88	0.00***	-1.55	0.18	-8.62	0.00***
<b>Age</b>																				
:Young adult																				
:Adult	-0.06	0.08	-0.73	0.46	0.02	0.08	0.22	0.82	0.00	0.08	-0.05	0.96	-0.05	0.08	-0.65	0.52	0.09	0.10	0.95	0.34
:Middle aged	-0.11	0.09	-1.20	0.23	-0.06	0.09	-0.71	0.48	0.03	0.09	0.34	0.74	-0.04	0.09	-0.41	0.68	0.19	0.10	1.84	0.07
:Old	-0.20	0.15	-1.35	0.18	-0.16	0.15	-1.09	0.27	-0.01	0.15	-0.10	0.92	-0.02	0.15	-0.16	0.88	0.51	0.16	3.10	0.00**
<b>Sex</b>																				
:Female																				
:Male	0.07	0.07	1.03	0.30	-0.34	0.07	-5.16	0.00***	-0.02	0.07	-0.27	0.79	0.15	0.07	2.09	0.04*	0.31	0.09	3.61	0.00***
<b>Terrain</b>																				
:Flat																				
:Mountainous	-0.26	0.17	-1.53	0.13	-0.20	0.16	-1.21	0.23	0.16	0.15	1.13	0.26	0.25	0.14	1.76	0.08	0.02	0.16	0.12	0.90
:Rolling	-0.15	0.08	-1.88	0.06	0.00	0.08	-0.06	0.95	-0.06	0.07	-0.83	0.41	0.02	0.07	0.29	0.77	0.24	0.08	2.96	0.00**
<b>Weather</b>																				
:Clear																				
:Cloudy	0.08	0.09	0.92	0.36	0.24	0.09	2.77	0.01**	-0.15	0.09	-1.59	0.11	-0.15	0.09	-1.59	0.11	-0.06	0.10	-0.59	0.56
:Raining	0.09	0.20	0.48	0.63	0.06	0.20	0.29	0.77	0.28	0.18	1.53	0.13	-0.17	0.21	-0.82	0.41	-0.28	0.23	-1.23	0.22
:Fog	-1.63	0.73	-2.24	0.03*	-0.64	0.49	-1.32	0.19	-1.01	0.53	-1.92	0.05	1.07	0.31	3.51	0.00***	0.74	0.35	2.12	0.03*
<b>Light</b>																				
:Daylight																				
:Dusk - Dawn	0.05	0.17	0.27	0.78	-0.53	0.19	-2.76	0.01**	-0.39	0.19	-2.12	0.03*	0.08	0.18	0.46	0.64	0.97	0.18	5.51	0.00***
:Dark - Street Lights	-0.27	0.10	-2.86	0.00**	-0.29	0.10	-3.00	0.00**	-0.39	0.10	-4.13	0.00***	0.35	0.09	4.00	0.00***	0.91	0.10	9.16	0.00***
:Dark - No Street Lights	-0.66	0.11	-6.05	0.00***	-0.83	0.11	-7.31	0.00***	-0.40	0.10	-4.12	0.00***	0.37	0.09	4.17	0.00***	1.37	0.09	14.94	0.00***
<b>Vehicle type</b>																				
:Passenger Car																				
:Truck	-0.81	0.17	-4.75	0.00***	-0.99	0.19	-5.29	0.00***	-0.45	0.15	-3.04	0.00**	0.40	0.12	3.27	0.00**	1.09	0.12	8.99	0.00***

Variable	Property damage only Versus All levels				Complaint of pain Versus All levels				Other visible injuries Versus All levels				Severe injury Versus All levels				Fatal Versus All levels			
	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )	Parameter Estimate	Standard Error	t-value	Pr(> t )
:Motorcycle	-3.95	0.71	-5.55	0.00***	-2.18	0.33	-6.70	0.00***	-0.32	0.15	-2.09	0.04*	1.32	0.12	10.86	0.00***	1.10	0.14	8.12	0.00***
:Pickup Truck	-0.02	0.09	-0.18	0.86	-0.02	0.09	-0.25	0.80	0.01	0.09	0.11	0.91	0.10	0.09	1.06	0.29	-0.06	0.11	-0.58	0.56
:Bus	-1.64	1.03	-1.59	0.11	-0.50	0.64	-0.78	0.43	-0.45	0.63	-0.72	0.47	1.23	0.46	2.65	0.01**	0.54	0.58	0.93	0.35
:Emergency Vehicle	0.26	0.50	0.53	0.60	0.62	0.46	1.34	0.18	-0.89	0.75	-1.19	0.24	-0.25	0.63	-0.40	0.69	-0.35	0.72	-0.49	0.63
<b>Number of lanes</b>																				
:2																				
:3	-0.24	0.63	-0.38	0.70	0.33	0.46	0.71	0.48	0.26	0.45	0.57	0.57	-1.44	0.62	-2.31	0.02*	0.73	0.40	1.84	0.07 <sup>ˆ</sup>
:4	0.31	0.17	1.81	0.07 <sup>ˆ</sup>	0.21	0.15	1.44	0.15	0.30	0.15	2.00	0.05*	-0.23	0.13	-1.70	0.09 <sup>ˆ</sup>	-0.57	0.15	-3.89	0.00***
:5	-0.05	0.31	-0.17	0.86	0.38	0.25	1.52	0.13	0.31	0.25	1.25	0.21	-0.21	0.23	-0.88	0.38	-0.62	0.26	-2.40	0.02*
:6	0.77	0.16	4.69	0.00***	0.07	0.15	0.50	0.62	0.38	0.15	2.57	0.01*	-0.51	0.14	-3.72	0.00***	-0.71	0.15	-4.71	0.00***
:7	0.75	0.23	3.19	0.00**	-0.21	0.24	-0.88	0.38	0.39	0.22	1.75	0.08 <sup>ˆ</sup>	-0.18	0.21	-0.85	0.40	-0.77	0.25	-3.06	0.00**
:8	0.59	0.16	3.68	0.00***	0.27	0.14	1.94	0.05 <sup>ˆ</sup>	0.23	0.14	1.62	0.10	-0.33	0.13	-2.58	0.01**	-0.76	0.14	-5.32	0.00***
:> 8	0.69	0.16	4.40	0.00***	0.10	0.14	0.74	0.46	0.41	0.14	2.93	0.00**	-0.42	0.13	-3.34	0.00***	-0.79	0.14	-5.70	0.00***
<b>Cause</b>																				
:Under influence																				
:Following too closely	1.58	0.24	6.68	0.00***	2.01	0.24	8.30	0.00***	0.00	0.20	0.00	1.00	-1.30	0.24	-5.36	0.00***	-1.76	0.38	-4.64	0.00***
:Improper Turn	0.35	0.30	1.18	0.24	0.72	0.30	2.44	0.01*	-0.58	0.22	-2.59	0.01**	-0.27	0.18	-1.52	0.13	0.65	0.18	3.67	0.00***
:Speeding	1.38	0.17	8.04	0.00***	1.65	0.19	8.78	0.00***	-0.11	0.11	-0.99	0.32	-0.85	0.09	-9.10	0.00***	-0.62	0.10	-6.20	0.00***
:Other Violations	1.41	0.21	6.84	0.00***	1.34	0.23	5.89	0.00***	-0.07	0.16	-0.43	0.67	-0.42	0.14	-3.01	0.00**	-0.86	0.18	-4.84	0.00
:Other Than Driving	-0.44	0.75	-0.59	0.56	0.54	0.56	0.95	0.34	-1.35	0.61	-2.21	0.03*	-1.19	0.46	-2.60	0.01**	2.05	0.35	5.79	0.00
:Unknown	2.27	0.45	5.03	0.00***	1.14	0.58	1.96	0.05 <sup>ˆ</sup>	-0.68	0.55	-1.24	0.21	-1.74	0.63	-2.78	0.01**	0.28	0.43	0.64	0.52
<b>Model Overall</b>	Log-Likelihood: -3127.9 McFadden R^2: 0.072 Likelihood ratio test : chisq = 482.72 (p.value = < 2.22e-16)				Log-Likelihood: -3171.6 McFadden R^2: 0.075 Likelihood ratio test : chisq = 514.85 (p.value = < 2.22e-16)				Log-Likelihood: -3171.6 McFadden R^2: 0.075 Likelihood ratio test : chisq = 514.85 (p.value = < 2.22e-16)				Log-Likelihood: -3207.9 McFadden R^2: 0.057 Likelihood ratio test : chisq = 386.92 (p.value = < 2.22e-16)				Log-Likelihood: -2570.5 McFadden R^2: 0.134 Likelihood ratio test : chisq = 797.17 (p.value = < 2.22e-16)			

1 Table 4 - summary of number of true and false predictions for each model

		Binary outcome Model										Multi-level outcome Model							
		Property damage only vs. all other levels		Complaint of pain vs. all other levels		Other visible injury vs. all other levels		Severe injury vs. all other levels		Fatal vs. all other levels		Multinomial model with all five levels		Mixed multinomial model with all five levels		Support vector machine with all five levels			
		True	False	True	False	True	False	True	False	True	False	True	False	True	False	True	False		
Prediction result for individual levels	Property damage only	0	620									216	404	215	405	208	412		
	Complaint of pain											0	579	275	304	280	299	270	310
	Other visible injuries											0	609	17	592	17	592	30	579
	Severe injury											29	563	225	367	222	370	213	379
	Fatal											48	392	133	307	132	308	158	282
Total prediction	2219			621	2258	582	2231	609	2244	596	2391	449	<b>866</b>	1974	<b>866</b>	1974	<b>879</b>	1962	
	AUC	0.65		0.66		0.54		0.67		0.74									

2 \* Shaded areas indicate the combination of all remaining levels for each model.

3

1 Table 5 - Confusion matrix for multinomial logit (1), mixed multinomial logit (2), and support vector machine (3)

		Reference				
		Property damage only	Complaint of pain	Other visible injury	Severe injury	Fatal
Prediction	Property damage only	216	197	168	115	65
	Complaint of pain	275	275	236	138	68
	Other visible injury	24	12	17	19	10
	Severe injury	67	53	122	225	164
	Fatal	38	42	66	95	133
	<b>True prediction in %</b>	<b>35 %</b>	<b>47 %</b>	<b>3 %</b>	<b>38 %</b>	<b>30 %</b>

2 (1) Multinomial Logit

		Reference				
		Property damage only	Complaint of pain	Other visible injury	Severe injury	fatal
Prediction	Property damage only	215	192	166	115	64
	Complaint of pain	278	280	238	136	68
	Other visible injury	24	13	17	25	14
	Severe injury	67	50	117	222	162
	Fatal	36	44	71	94	132
	<b>True predictions in %</b>	<b>35 %</b>	<b>48 %</b>	<b>3 %</b>	<b>38 %</b>	<b>30 %</b>

3 (2) Mixed Multinomial Logit

		Reference				
		Property damage only	Complaint of pain	Other visible injury	Severe injury	fatal
Prediction	Property damage only	208	178	147	108	62
	Complaint of pain	266	270	229	121	66
	Other visible injury	20	21	30	19	15
	Severe injury	64	53	119	213	139
	Fatal	62	58	84	131	158
	<b>True predictions in %</b>	<b>34 %</b>	<b>47 %</b>	<b>5 %</b>	<b>36 %</b>	<b>36 %</b>

4 (3) Support Vector Machine