# Evaluation of Injury Severity for Pedestrian-Vehicle Crashes in Jordan Using Extracted Rules

**Randa Oqab Mujalli, Ph.D [#, a], Laura Garach, Ph.D [b], Griselda López, Ph.D [c], Taleb Al-Rousan, Ph.D [d]**

[a] Assistant Professor at Department of Civil Engineering, The Hashemite University, 13115 Zarqa, Jordan
[b] Assistant Professor at Department of Civil Engineering, University of Granada, ETSI Caminos, Canales y Puertos, c/ Severo Ochoa, s/n, 18071 Granada, Spain, e-mail: lgarach@ugr.es
[c] Assistant Professor at Highway Engineering Research Group, Universitat Politècnica de València, Camino de Vera, s/n; 46022-Valencia, Spain, e-mail: grilomal@tra.upv.es
[d] Associate Professor at Department of Civil Engineering, The Hashemite University, 13115 Zarqa, Jordan, e-mail: taleb@hu.edu.jo
[#] Corresponding author, The Hashemite University, 13115 Zarqa, Jordan, Phone: +96253903333-4777
e-mail: randao@hu.edu.jo

**Abstract**

Safety of pedestrians is a major concern all over the world since pedestrians are considered to be the most vulnerable roadway users. This paper sought to determine the main factors in pedestrian-vehicle crashes that increase the risk of a fatality or severe injury. Pedestrian-vehicle crashes which occurred in urban and suburban areas in Jordan between 2009 and 2011 were investigated. Extracted rules from Bayesian networks were used to identify factors related to severity of pedestrian-vehicle crashes. To obtain as much information as possible about these factors, three subsets were used. The first and second subsets contain all types of collisions (pedestrians and non-pedestrians), where the first subset used collision type as a class variable and the second subset used injury severity. The third subset contains pedestrian collisions only and used injury severity as the class variable. The results indicate that when using collision type as the class variable, better performance was obtained and that the following variables increase the risk of fatality or severe injury: roadway type, number of lanes, speed limit, lighting, and adverse weather condition.

**Author Keywords**: pedestrians; Bayesian networks; rules extraction; imbalanced dataset; collisions; urban areas.

**Introduction**

Pedestrians are considered vulnerable road users who are more likely to suffer from severe injuries or even get killed than other road users if involved in a crash. In fact, the largest group of road fatalities worldwide is pedestrians hit by motorized vehicles (Peden et al. 2004; Naci et al. 2009; Rosén et al. 2011). In 2015, pedestrian deaths accounted for 22% of the traffic collision deaths worldwide (WHO 2015).

Economic level of countries highly affects the places or areas where pedestrians' collisions occur (WHO 2013). In high income countries pedestrians have high probability of being involved in collisions at urban areas, while pedestrians in some low and middle-income countries are expected be involved in collisions that occur in rural areas (WHO 2013). In Jordan, an increased vehicle ownership accompanied the rapid growth in population. According to the 2014 Jordan Statistics there are 6.675 million inhabitants (excluding visting labor and refugees) with 1,329,888 licensed vehicles (DOS 2014). 82.6% of the Jordanian population lives in urban or suburban areas, with 44% of this population living in Amman, the capital city (In Jordan, areas are defined as urban if the population is more than 10,000 inhabitants and as suburban for populations between 5,000 to 10,000 inhabitants) (Ababsa, 2013). Thus, rapid growth in both population and vehicle ownership has contributed, amongst other factors, to increasing the number of collisions. Data for 2014 indicate that in Jordan 102,441 traffic crashes occurred, of which 3,839 were pedestrian collisions. Although the number of pedestrian collisions does not seem to be high relative to the total number of crashes which accounts for 3.7% of total traffic crashes, however, the number of fatalities which resulted from pedestrian-vehicle collisions was 255 fatalities, which makes about (37%) of the total fatalities number which was (688) from traffic collisions (PTD 2014) as shown in Figure 1.

[Insert Figure 1]

Jordanian Traffic law imposes penalties on pedestrians crossing roads from illegal points i.e. not using defined crosswalks, pedestrian underpasses or overpasses. However, this penalty is not applied in practice because crosswalks sometimes do not provide the best crossing (MOI 2008). Often pedestrians are forced to cross at illegal locations because of bad driver behavior, for example stopping on crosswalks. Moreover, most crosswalks throughout Jordan are not clearly marked and pedestrians are forced to walk on the roads because of missing or inadequate sidewalks or because the sidewalks are obstructed by trees or merchants (JTI 2012). When pedestrian overpasses are available, about 60% of pedestrians do not use them because of high stairs, health related issues, extra walking distance, inconvenience and safety fears (Abojaradeh 2013).

67  Studying risk factors for pedestrian-vehicle collisions will help develop preventive measures and hence reducing

68  pedestrian-vehicle collisions rate (Zhang et al. 2014). Researchers have explored collisions involving

69  pedestrians using several statistical methods, such as logit or probit models (Zajac and Ivan 2003; Lee and

70  Abdel-Aty 2005; Xie et al. 2009; Tay et al. 2011; Islam and Hernandez 2013; Mohamed et al. 2013; Obeng and

71  Rokonuzzaman 2013;Olszewski et al. 2015). These models rely on a pre-assumption of the data distribution and

72  the relationship between explanatory variables and response variable (Chang and Wang 2006).

73  Some of the very interesting recent methods used in crashes analysis are data mining (DM), the methods used in

74  DM  do not rely on pre-assumptions of the data distribution nor the relationship between variables, and they

75  help in discovering some new relationships (Prato et al. 2010). In addition, patterns of collisions can be

76  discovered using DM techniques. One of the aims of Data Mining is identifying valid, novel, potentially useful

77  and ultimately understandable patterns in data (Fayyad et al. 1996). Different data mining methods have been

78  used to analyze road crashes, such as association rules (Pande and Abdel-Aty 2009; Montella et al. 2012),

79  decision trees (Prato et al. 2010; Abellán et al. 2013; De Oña et al. 2013; Jung et al. 2016), clustering analysis

80  (Mohamed et al. 2013; Thompson et al. 2013), neural networks (Abdel-Aty and Abdelwahab, 2004), Bayesian

81  Networks (BNs) (De Oña et al. 2011; Mujalli and De Oña 2011; Kwon et al. 2015; Mujalli et al. 2016). One of

82  the main advantages of BNs is that inference can be used to find out key variables which significantly affect the

83  dependent variable (class variable) in the resulting models (De Oña el al. 2012, Mujalli et al. 2016).

84  Traffic collision datasets are usually imbalanced, which means that the dataset has one or some of the classes

85  (categories) with much more samples in comparison to the others. The most prevalent class is called the

86  majority class, while the rarest class is called minority class (Yijing et al. 2016). The minority class usually

87  represents a concept with greater interest than the majority class. However, it is often outnumbered by the

88  majority class, and sometimes this scenario may be very severe (Sun et al. 2015). One of the most popular

89  methods for solving the class imbalance problems is sampling. The two most used sampling techniques are

90  undersampling and oversampling. In undersampling, instances of the minority and majority classes are selected

91  randomly in order to achieve a balanced stratified sample with equal class distributions, often using all instances

92  of the minority class and only a subset of the majority class, or undersampling both classes for even smaller

93  subsets with equal class sizes. In oversampling, the class distribution is approximately equalized between

94  minority and majority classes, where minority cases are replicated (Crone and Finlay 2012). Accordingly,

95  Mujalli et al. (2016) found that when using balanced datasets by oversampling, better classification of collisions

96  is achieved and less misclassification of the minority class occurs.

97 The main aim of this research is the exploratory analysis of pedestrian-vehicle collisions in Jordan using DM

98 techniques. DM allows knowledge extraction of previously hidden patterns (Fayyad et al. 1996). According to

99 Zhang et al (2004), using a single database limits the possibility of discovering specific relationships due to the

100 fact that knowledge might be hidden in multiple databases. In order to explore all possible patterns associated

101 with pedestrian-vehicle collisions and their severity, 3 subsets are developed and explored.Two subsets contain

102 all types of collisions, with one of them using collision type as a class variable and the other using injury

103 severity as the class variable, whereas the third subset is pedestrian-specific which only contains pedestrian

104 collisions. Both subsets were used in order to find out if different patterns of variables affecting severity could

105 be extracted when using a general subset or a pedestrian specific subset. To take into account the problem of

106 imbalanced data, sampling techniques will be applied prior to developing BN models and analysis.

107 The remainder of the paper is organized as follows: First, the method and data used are presented, followed by

108 results and their discussion. Finally, the main conclusions of this study and recommendations are stated.

109 **Material and Methods**

110 The procedure used in this study has been the following:

111    1. Dataset for collisions which occurred in urban and suburban areas in Jordan was obtained from

112       Jordanian Police Traffic Department (PTD), this dataset included 14 variables.

113    2.  All types of collisions were extracted from the dataset excluding property damage only collisions

114       (PDO); this subset is called hereafter general subset (see Table 1).

115                                         [Insert Table 1]

116    3. Subset Col-ACT and subset Col-SEV contain the same number of records as general subset; the only

117       difference between these subsets and the general subset is that the collision type variable (ACT) was

118       categorized into two categories (values): Pedestrian (PED) and others (OT) (see Table 1). Also, in

119       subset Col-ACT, the class variable was ACT and in subset Col-SEV the class variable was SEV. More

120       specifically, these subsets contain all records (pedestrian-vehicle collisions and non- pedestrian

121       collisions.

122    4. Subset Ped-SEV is obtained from subset Col-ACT. In subset PED-SEV, only records belonging to

123       (ACT=PED) were used and thus all records belonging to other collisions types (OT) were eliminated.

124       As a result, variable ACT was eliminated, and only 13 variables were used (see Table 1). More

125       specifically, this subset contains only pedestrian-vehicle collisions records.

126    5. All subsets are now characterized as follows:

4

127    a. Subsets Col-ACT and Col-SEV, in which all types of collisions were included: collisions with

128     vehicles and collisions with pedestrians.

129    b. Subset PED-SEV, in which only pedestrian-vehicle collisions were included.

130  6. In both subsets Col-ACT and PED-SEV, collisions were classified according to severity (SEV),

131   whereas in subset Col-ACT, collisions were classified according to collision type (ACT).

132  7. All subsets were balanced using oversampling balancing technique.

133  8. All balanced subsets were used to develop BNs using different search algorithms (K2 and simulated

134   Annealing) and scores (BDeu, MDL, AIC) (De Oña et al. 2011 and Mujalli et al. 2016).

135  9. Six Developed BNs were then compared (using 10-folds cross validation method) based on their

136   performance as obtained using performance measures.

137  10. The BN with the best performance using each subset was further analyzed using the conditional

138   probability tables to extract significant rules.

139  11. The extracted rules that obtained the highest probability of fatality or severe injury (FS) in each BN

140   from each subset were presented and analyzed.

141 The procedure is shown in Figure 2.

142          [Insert Figure 2]

143 **Data**

144 Collision data were obtained from PTD with records for traffic collisions that occurred on urban and suburban

145 roads in Jordan from year 2009 to 2011.Unrealistic data was removed from the dataset, as well as, PDO

146 collisions. The final dataset contained 21,852 records (general subset). Table 1 includes detailed explanation of

147 the subsets used.

148 Fourteen (14) variables were used to find out which factors significantly affect the occurrence of a severe or

149 fatal outcome in pedestrian-vehicle collisions. These variables were related to road characteristics, prevailing

150 conditions, prevailing weather conditions, collision type and the resulting injury severity of the worst injured

151 (Chang and Wang 2006; Abellán et al. 2013, De Oña et al. 2013; López et al. 2014).

152 Table 2 includes the classification of records between the classes of classification, where in subset ACT and

153 subset 1: there are 17,769 records with slight injuries and 4,083records with fatal and severe injuries. In subset

154 Col-ACT the class variable ACT is imbalanced with (94.4%) cases classified as OT and 5.6% cases classified as

155 PED. In subsets Col-SEV the class variable SEV is imbalanced, with 81.32% cases classified as SI, and 18.68%

156 cases classified as FS. Also, in subset 2 the class variable SEV was imbalanced with 82.9% cases classified as

157 slight injuries and 17.1% cases classified as fatality or severe injuries.

158 [Insert Table 2]

159 *Balancing subsets*

160 Datasets in which records belonging to categories of the class variable do not have almost the same distribution

161 between the categories of the class variable are called imbalanced datasets (Crone and Finlay 2012).The

162 problem of imbalanced datasets is more apparent in cases where the main interest of the classification task is to

163 distinguish cases belonging to the class having rare occurrence in nature, such as collisions resulting in severe or

164 fatal outcomes. The possible problem is misclassifying a collision as minor when in reality it should have been

165 classified as severe injury or fatal. To deal with the problem of imbalanced datasets, sampling prior to

166 classification is most commonly used, in which class distribution of records is altered so that the minority class

167 (rare cases) records in the dataset are increased in size (Thammasiri et al. 2014).

168 In this research work, and based on the results obtained by Mujalli et al. (2016), oversampling balancing

169 technique in Weka software was used (Witten and Frank 2005).

170 *Oversampling*

171 The oversampling technique used herein is Synthetic Minority Oversampling Technique (SMOTE), in which a

172 subset of the original dataset is created by creating synthetic minority examples; each minority class sample is

173 taken and synthetic examples along the line of the segments joining any/all of the (k) minority class nearest

174 neighbors are introduced. Neighbors from the (k) nearest neighbors are randomly chosen and one sample is

175 generated in the direction of each, based on the amount of oversampling prescribed (Chawla et al. 2002). For

176 further details on sampling techniques used in the analysis of traffic crashes, interested readers are referred to

177 Mujalli et al. (2016)

178 *Bayesian networks*

179 A BN is a Directed Acyclic Graph (DAG) over a set of variables: $U=\{x_1, x_2, ..., x_n\}$, $n \geq 1$ and a set of probability

180 tables $B_p = p(x_i/pa\ (x_i)), x_i \in U$ where $pa\ (x_i)$ is the set of parents of $x_i$ in BN and $i = (1,2,3,....,n)$. A BN represents

181 joint probability distributions $P(U) = \prod_{x_i \in U} p\left(x_i \mid pa(x_i)\right)$. The arcs are interpreted as direct dependence

182 relationships between the linked variables, however indirect dependence relationships between variables could

183 exist (Acid et al., 2004).

6

184  The classification task consists in classifying a variable $y=x_0$ called the class variable, given a set of variables

185  $U=x_1 \ldots x_n$, called attribute variables. A classifier $h : U \rightarrow y$ is a function that classify an instance of $U$ to a value

186  of $y$. A dataset $D$ consisting of samples over $(U, y)$ is used to learn a classifier where the learning task consists of

187  finding an appropriate BN given a data set $D$ over $U$.

188  In this study, BNs are used in order to develop different models and to compare their results in terms of their

189  ability to correctly classify collisions according to either of their injury severity or collision type and also to

190  extract significant rules. Two search methods were used: hill climber algorithm restricted by an order on the

191  variables (K2) and simulated annealing search algorithm. Also, three different score functions were used: BDe

192  score metric (BDeu); Minimum Description Length (MDL); and the Akaike Information Criterion (AIC). The

193  search algorithms and the scores were applied in this study because they are widely used, relatively quick, and

194  they produce good results in terms of network overall performance (Madden 2009; Mujalli et al. 2016).

195  *Performance evaluation measures*

196  To compare the different BNs, the following performance measures were used: accuracy, sensitivity, specificity,

197  and area under receiver operating characteristic curve (AUC). Their equations and explanation of AUC are

198  given below:

199
$$Accuracy = \frac{TSI + TFS}{TSI + FSI + TFS + FFS} \tag{1}$$

200
$$Sensitivity = \frac{TSI}{TSI + FFS} \tag{2}$$

201
$$Specificity = \frac{TFS}{TFS + FSI} \tag{3}$$

202  Where,

203  TSI: number of (SI) instances correctly classified,

204  TFS: number of (FS) instances correctly classified,

205  FSI: number of (SI) instances incorrectly classified, and

206  FFS: number of (FS) instances incorrectly classified.

207  AUC is the area below a curve where sensitivity (plotted on y-axis) and 1-specificity (Bradley 1997).

208  **Results**

209  The aim herein is to investigate the main factors that affect the occurrence of a specific outcome in pedestrian-

210  vehicle collisions. To achieve this aim three subsets have been analyzed:

211      -    Subset Col-ACT: contains all records (pedestrian-vehicle collisions and non- pedestrian collisions.

7

212      -      Subset Col-SEV: contains all records (pedestrian-vehicle collisions and non- pedestrian collisions).

213      -      Subset Ped-SEV: contains only pedestrian-vehicle collisions records.

214 In subset Col-ACT the class variable was the collision type (ACT) and the collisions are classified according to

215 ACT: PED, OT; whereas in subsets Col-SEV and Ped-SEV, the class variable was the severity and thus

216 collisions were classified according to SEV: FS or SI.

217 The following sections include a presentation of the subsets used in analysis and the BNs developed using these

218 subsets. Next, a detailed analysis of the developed BNs along with the findings of this research work is shown.

219 *Subsets used*

220 As previously indicated subsets were imbalanced as shown in Table 2. To solve this issue, oversampling was

221 used to develop balanced subsets. Details of the balanced subsets are shown in Table 2.

222 It is illustrated in Table 2 that the number of the instances in the resulting subsets was increased to the size of

223 the majority class (20,636 instances for PED in subset Col-ACT; 17,761 instances for FS in subset Col-SEV and

224 1,008 for FS in subset Ped-SEV).

225 *Bayesian networks and subsets*

226 For each balanced subset, six models using six different algorithms (as shown in Figure 1) were developed.

227 Subsets were first divided into ten subsamples, training subsamples and testing subsamples. In each run, nine

228 subsamples were used to train the model and the last was used to test the model. The process was repeated for

229 ten runs for each subsample (ten-fold cross-validation).

230 To compare the resulting models, performance evaluation measures were used. The results shown in Table 3 are

231 the averages of the ten runs of the test subsamples. It should be mentioned that the statistical significance of

232 each model obtained was tested using a corrected paired t-test.

233                                           [Insert Table 3]

234 Results extracted from the developed models using the balanced subsets indicated that accuracy results for both

235 subset Col-SEV and subset Ped-SEV were in range of 59% to 61% and 51% to 61%, respectively, where

236 accuracy for subset Col-ACT was 85.7% to 86.8% indicating a better overall performance obtained by subset

237 ACT.

238 Regarding sensitivity, the values ranges in subset Col-SEV were between 58% and 61.6%; whereas in subset

239 Ped-SEV, the ranges increased from 1.9% to 58.7% indicating that subset COl-SEV overpowered subset Ped-

240 SEV in classifying instances belonging to SI class. In subset Col-ACT sensitivity results ranged from 75.9% to

241 76.4%.

8

242 In terms of specificity, in subset Col-SEV the values ranged from 56.3% to 64.2%, where in subset Ped-SEV the

243 values ranged from 64.2% to 99.7%. It is evident here that all classifiers in both subsets were remarkably

244 capable of classifying the minority class (FS) correctly. In subset Col-ACT, the values ranged from 95.5% to

245 97.3% indicating that all algorithms used were capable of classifying PED correctly.

246 Finally, with respect to AUC, results for subset Col-SEV ranged from 63% to 65.6%, for subset Ped-SEV from

247 50.8% to 64.7% and for subset Col-ACT the results ranged from 90.4% to 91.6%.

248 Based on the aforementioned results, subset Col-ACT had more wins (21 wins) than subset Col-SEV with no

249 wins and subset Ped-SEV with 2 wins only. Regarding the best algorithms used, results in Table 4 shows the

250 averages for the test subsamples using the ten-fold cross-validation and testing the statistical significance of

251 each algorithm using a corrected paired t-test. It is shown than the best algorithm used K2+AIC with 11 wins.

252                                                      [Insert Table 4]

253 As shown in Table 4, algorithm K2+AIC was the best in the balanced subsets. In subsets Col-ACT and Col-SEV

254 it obtained the highest values in accuracy, specificity and AUC; whereas in subset Ped-SEV it achieved the

255 highest values in accuracy, sensitivity and AUC. To this end, BNs obtained using K2+AIC will be used to

256 interpret the contributing factors into pedestrian-vehicle collisions.

257 *Interpretation of Bayesian networks*

258 BNs developed using K2+AIC algorithm were used for further analysis and to determine which variables are

259 responsible for FS injuries in pedestrian-vehicle collisions as shown in Table 5.

260                                                      [Insert Table 5]

261 As illustrated in Table 5, the developed BN using Col-ACT subset had the following variables directly

262 connected to the variable severity (SEV) which is the variable of interest: collision type (ACT), road type (DIR),

263 number of lanes (LANE), traffic control (CONT), and speed (SPE). The severity presents a direct relationship

264 with both pedestrian and vehicle crashes, with road type, number of lanes and traffic control and speed limits in

265 all types of crashes included in subset ACT. When using subset Col-SEV it was found that all the independent

266 variables were connected to severity (SEV) except for traffic control (CONT), whereas in subset Ped-SEV (only

267 pedestrian collisions) seven variables were connected to severity (SEV): vehicles involved (VEH), collision

268 pattern (PAT), road type (DIR), number of lanes (LANE), road grade (GRADE), lighting condition (LIG) and

269 speed limit (SPE). On the other hand, the following variables were found to be common in all BNs using all

270      subsets and with arcs connecting them to severity (SEV): number of lanes (LANE), road type (DIR) and speed

271      limit (SPE); indicating the significance of these variables in pedestrian-vehicle collisions.

272      In view of the above results, it is of interest for safety researchers to find out what differences exist between

273      collisions if pedestrian only subset was used and those if all types of collisions subsets were used, as well as,

274      what differences exist if the root node (class variable) was altered. In light of this, the developed BNs need to be

275      further analyzed using more deep knowledge extraction analysis. The following section describes the main

276      differences and similarities between the developed BNs using the different subsets.

277      ***Rules extraction using BNs***

278      In order to identify the most significant values (categories) of variables that affect the occurrence of fatality or

279      severe injury in a pedestrian-vehicle collision, the most significant rules were extracted from the conditional

280      probability tables of the developed BNs using each subset. The extracted rules allow finding the particular

281      variables and their respective values which are associated with the highest probability (confidence) of

282      occurrence for fatality or severe pedestrian-vehicle collision. The results for the extracted rules are shown in

283      Table 6.

284      [Insert Table 6]

285      Based on results in Table 6, it would be possible to identify the factors that affect the occurrence of fatal or

286      severe pedestrian-vehicle collisions. In which the combinations of factors that define a pedestrian collision to be

287      classified as FS were determined. To this end the rules were divided into three groups:

288      1. High confidence rules group: the rules which achieved a confidence of at least 90% are listed.

289      2. Moderate confidence rules group: the rules which achieved a confidence between 70-89% are listed.

290      3. Low confidence rules group: the rules which achieved a confidence between 50-69% are listed.

291      It should be noted that the categorization into these three categories was based on the number of rules and the

292      confidence level obtained by these rules where it was found in literature that researchers used similar thresholds

293      to extract knowledge from used datasets. For example both Abellán et al. (2013) and López et al. (2014) used a

294      confidence threshold of 60% in order to extract significant rules when using decision trees.

295      In the high confidence group, the results indicated that there were only two rules belonging to this group: the

296      first most important rule having a confidence of 92.1% was obtained when using subset Col-ACT. This rule has

297      the parent node severity (SEV) with its state (FS), the rule strongly suggests that the occurrence of FS

298      pedestrian-vehicle collision is almost certain if the following conditions exist: the collision occurring on roads

299      with speed limit of 30 km/hr on a 4 lanes divided roadway. This points out that as the width of the crossed roads

300 increases the risk of fatal or severe pedestrian collisions increase even if the posted speed was as low as 30

301 km/hr.

302 The second rule which had a confidence of 92.1% of fatal or severe collision occurrence was obtained when

303 using subset Ped-SEV. The extracted rule has the parent node number of lanes (LANE) with its state (1) where

304 it involved the existence of the following factors at time of collision: the collision occurs on 2 lanes one way

305 roads, in which the rule basically states that the risk of crossing wider roads is associated with almost a certain

306 fatal or severe pedestrian-vehicle collision occurrence.

307 In the moderate confidence rules group, there were six extracted rules belonging to all subsets, the first rule with

308 the highest confidence in this group was extracted from subset Col-SEV. This rule has the parent node lighting

309 condition (LIG) with its state (DARK) and it indicates that the occurrence of fatal or severe pedestrian collisions

310 is associated with the existence of the following conditions at time of collision: if the collision occurs at night

311 without lighting on 2 ways undivided roadway. This rule does not clearly identify the prevailing weather

312 condition, even though the atmospheric weather exists as a contributory factor but it says (other). If we refer to

313 the original data, (other) refers to one of the following states: snow, storm, wind, fog, or dust which means that

314 the prevailing conditions were not clear. This rule ascertains the importance of the type of roadway and adds

315 two other factors: dark and unclear weather.

316 The second highest confidence rule was extracted from subset Ped-SEV with a 76.1% confidence. This rule has

317 the parent node number of lanes (LANE) with its state (2) and it states that if the collision occurs on 2 lanes

318 undivided roads, the probability of a pedestrian being involved in severe or fatal collision is moderate;

319 highlighting the significance of the type of roadway in pedestrian collisions.

320 The third highest confidence rule was extracted from subset Col-SEV with a 74% confidence, where a new

321 factor was added to the contributory factors in pedestrian-vehicle collisions which is the number of involved

322 vehicles. The rule has the parent node speed (SPE) with its state (40) and it states that if collision occurs on one

323 way road with a speed limit of 40 km/hr and there was only one vehicle involved in this pedestrian-vehicle

324 collision, then there is a moderate probability that the collision will be classified as fatal or severe collision. The

325 last rule in this group was extracted from subset Ped-SEV with a confidence of 73.9%, in which the rule has the

326 parent node number of lanes (LANE) with its state (2) and it states that the following combination of factors

327 have a moderate confidence that the collision would be classified as fatal or severe injury: if the collision occurs

328 on a four lanes-tow-2 way divided road.

In the low confidence group, there were 5 extracted rules belonging to this group. The first highest confidence rule was extracted from subset Ped-SEV, in which the rule has the parent node speed (SPE) with its state (40) and it states that if collision occurs on two-ways undivided road with a speed limit of 40 km/hr, then the probability is relatively low that the collision will be classified as fatal or severe injury. It should be mentioned that the absence of the number of vehicles factor (as shown from rule no. 2 from subset Col-SEV) and the change of the type of roadway from one way (as shown in rule no. 2 from subset Col-SEV) to two-ways has decreased the confidence from 74% to 66% indicating the significance of the factors which were present in rule Ped-SEV from subset Col-SEV in fatal or severe injury collisions.

The second two highest confidence rules were extracted from subsets Col-ACT and Col-SEV. The extracted rule from subset Col-ACT has the parent node severity (SEV) and it state (FS) and it states that if the collision occurs on four lanes -two- ways undivided road with a speed limit of 90 km/hr, then there is a low probability that the collision will be classified as fatal or severe injury. This result might indicate the degree of alert that the pedestrian pay when crossing roads with high speed limits, in which the probability of the pedestrians being involved in such collisions becomes relatively low. The extracted rule from subset Col-SEV has the parent node road type (DIR) and its state (2_DIV) and it states that: if the collision occurs on two ways divided road and two vehicles were involved in the collision, then the probability of having fatal or severe injury is relatively low.

The third highest rule in this group was extracted from subset Ped-SEV, in which the rule has the parent node speed (SPE) and its state (40) and it states that: if the collision occurs on a one way road having a speed limit of 40 km/hr then the probability is relatively low that the collision will be classified as fatal or severe injury. It should be mentioned that the change of road type from two ways undivided (as shown in rule no. 4 from subset Ped-SEV) into one way in the current rule has decreased the probability of fatal or severe injury collisions having the speed limit the same in both rules.

The last extracted rule in this group was obtained from subset Ped-SEV, this rule has the parent node lighting conditions (LIG) and its state (DARK) and it states that: if the collision occurs at night under dark lighting conditions when the prevailing weather condition is other (snow, storm, wind, fog, or dust) then the probability is relatively low that the collision will be classified as fatal or severe injury. It should be mentioned that the absence of the road type factor (as shown in rule 1 in subset Col-ACT) has decreased the obtained probability in the current rule, indicating the significance of road type in such collisions once more.

357   Absence of control type from all significant extracted rules raised some concerns since it is believed

358   to have some effect on severity outcome, consequently inference was performed in order to caputre

359   the effect of control type (if any) on fatality or severe injury outcome. For more details on BN

360   inference, readers are referred to (De Oña et al., 2011).

361   The evidence was set to certainty for the following variables/subset combination as follows: (Subsets

362   Col-ACT and Col-SEV: Probability (ACT=PED, SEV=FS); Subset PED-SEV: Probability

363   (SEV=FS)). The probabilities of control type (CONT) were then calculated and the results indicated

364   that in subset Col-ACT the probability of CONT=NO_CONT was found to be 0.8950, in subset Col-

365   SEV the probability of CONT=NO-CONT was found to be 0.8758, and in subset PED_SEV the

366   probability of CONT=NO-CONT was found to be 0.8765. The inference results indicat that fatal or

367   severe pedestrian-vehicle crashes are significantly associated with locations where no control exists;

368   this result should be interpreted with caution because this category of the variable CONT includes

369   both uncontrolled intersections and non-intersections. In a study prepared by Chong (2018) risk of

370   fatality was found to increase at non-intersections as compared to intersections. On the other hand,

371   intersections without signals were found to have increased risk of fatal or injury pedestrian crashes

372   (Moudon et al., 2011).

373   **Discussion**

374   The research work presented in this paper analyzed the severity of pedestrian-vehicle collisions in urban and

375   sub-urban areas in Jordan using BNs. In this study, in order to explore all the main contributory factors affecting

376   these collisions, three subsets were used. Two subsets included all types of collisions, pedestrian and no

377   pedestrian collisions (subsets Col-ACT and Col-SEV), whereas the third subset included only pedestrian

378   collisions (subset Ped-SEV). It was found that using collision type as the class variable in subset Col-ACT

379   enhanced the performance of the developed BNs in terms of better classification of cases as compared to the

380   other two subsets. On the other hand, it was found that using subset Ped-SEV which included only pedestrian

381   collisions enhanced the performance with respect to subset Col-SEV which included all collision types. The

382   enhanced performance of subset Col-ACT increases its reliability in terms of defining the factors that are highly

383   associated with pedestrian involved collisions.

384     Table 7 shows a summary of the number of times each factor was present in the BNs resulting from using each

385     subset. Road type was present 10 times in the resulting rules indicating the importance of this factor. Speed and

386     number of lanes were present 5 times each in the resulting extracted rules. On the other hand, number of

387     vehicles, lighting and weather were the least present.

388                                                 [Insert Table 7]

389     Regarding to which states of factors caused FS pedestrian-vehicle collisions, higher risks of FS pedestrian

390     collision were associated with collisions involving single vehicle as shown in rule number 2 from subset Col-

391     SEV. This result contradicts with the results obtained by Verzosa and Miles (2016) where they found that the

392     odds of fatal collisions involving pedestrians are three times higher when there are heavy vehicles or multiple

393     vehicles involved than when there are light vehicles. The results herein indicated that single vehicle crashes

394     rather than multiple vehicles, as found by Verzosa and Miles (2016), have high risk of resulting in fatality when

395     the speed limit is 40km/hr and the road is one way. However, both studies are not completely comparable since

396     in this paper the type of vehicle involved was missing and thus the results obtained for single vehicle crashes

397     might be for single heavy vehicles crashes. In absence of any reserach work on injury severity of vehicle-

398     pedestrian crashes that included number of vehicles as a predictor without reference to type of

399     vehicle, authors of this work recommend that the result obtained herein should be interpreted with

400     caution.

401     Type of roadway; whether it was divided or one-way were found to have the highest risk in FS collisions as seen

402     in rule no. 1 in subset Col-ACT and rule no. 1 in subset Ped-SEV. Olszewski et al. (2015) found that divided

403     roadways were associated with increased risk of pedestrian death at un-signalized crosswalks which is

404     consistent with the results herein. It should be mentioned that divided roadways are riskier when the speed limit

405     is 30 km/h and the number of lanes per direction is two lanes as shown in rule no. 1 in subset Col-ACT, this

406     might give an indication of the impaired decision of pedestrians when crossing wide divided roadways with low

407     speed limit and their false perception of safety. On the other hand, rule no. 2 in subset Col-ACT has the same

408     combination of factors as that in rule no. 1, but speed limit was 90 km/hr which decreased the risk of FS

409     collisions, which might also indicate the limited accessibility the pedestrians have on higher speed roads, where

410     the accessibility is provided through pedestrian facilities; and even if there were some violations made by

411     pedestrians, they will be probably more careful when crossing roadways with higher speed limits. Another result

412     is that single lane one-way roads have the same risk as four lanes divided roadways as shown in rule no. 1 in

413 subset Ped-SEV which may indicate once again the false perception of safety that encourage pedestrians to cross

414 roads without taking necessary precautions.

415 Number of lanes, whether one or two lanes, was found to be associated with FS collisions in both rules in subset

416 Col-ACT and in rules 1 to 3 in subset Ped-SEV. Sze and Wong (2007) found that more lanes are significantly

417 associated with the higher likelihood of a pedestrian fatality. Also, Sherony and Zhang (2015) found that in light

418 vehicle collisions which occurred in the United Sates, more pedestrian crashes occurred on principle arterial

419 than any other type of roads. In addition, Amoh-Gyimah et al. (2016) found that wider roads are significantly

420 associated with fatalities in pedestrian-vehicle crashes. The interpretation of such results would be that as

421 number of lanes in roadways increases the ability of the driver to violate speed limits increases and hence

422 increases her/his odds of being involved in collisions.

423 As for speed, it was found that when the posted speed is 30 km/hr, the probability of FS collisions will be as

424 high as 92.1% with less probability of FS collisions with higher speeds (rule no. 2 from subset Col-ACT, rule

425 no. 2 from subset Col-SEV, rules no. 3 and 4 from subset Ped-SEV). According to Johansson et al. (2004), the

426 90[th]-percentile speed should not exceed 30 km/h for higher safety. Other researchers found that the risk of a

427 pedestrian dying of a collision with a vehicle is estimated to jump from 5% to 45% when the speed of the

428 striking vehicle increases from 20 (32 km/hr) to 30 mph (48 km/hr); it is 85% when vehicular speed reaches 40

429 mph (64 km/hr) (Limpert 1984; Leaf and Preusser 1999; Evaluating Safety and Health Impacts 2014). In

430 addition, results reported by Al-Omari and Obaidat (2013) supported the results found herein; they indicated

431 that 65% of pedestrian collisions which occurred at Irbid city in Jordan were at speed limit of 40 km/hr.

432 Moreover, Knowles et al. (2012) found that around 90% of fatal pedestrian collisions on London's roads

433 occurred on roads with speed limits of 48 km/hour or lower.

434 Higher levels of severity of pedestrian collisions were also found to be highly associated with night and dark

435 conditions as shown in rule 6 from subset 2. Sherony and Zhang (2015) found that in the United States for light

436 vehicles collisions involving pedestrians, most of the injury crashes occurred under daylight condition (62%)

437 while most of the fatal crashes occurred in dark condition (36% for dark but lighted, 34% for dark). In addition,

438 many other researchers reported that the majority of pedestrian collisions and majority of severe injuries have

439 occurred during nighttime (Campbell et al. 2004; Kim et al. 2008 and 2010; Verzosa and Miles 2016 and Al-

440 Ghamedi 2002, Amoh-Gyimah et al. 2016).

441 According to the results obtained herein, a snowy, stormy, windy, fogy, or dusty day was found to be mostly

442 associated with FS. DiMaggio and Durkin (2002) found that weather conditions have not been directly related to

443 the risk of injury. Kim et al. (2010) found that inclement weather (rain, snow, fog, and smog) decreased the

444 probability of fatal injury (−37%), but increased the probability of incapacitating injury (5%), non-

445 incapacitating injury (5%), and possible or no injury (5%). Moreover, Sherony and Zhang (2015) indicated that

446 the majority of collisions in the United States occurred under clear / cloudy condition (no adverse conditions), in

447 which 84.6% pedestrian injury crashes, 85.4% pedestrian fatal crashes occurred during no adverse weather

448 condition. However, Amoh-Gyimah et al. (2016) found that there is a significantly higher probability of

449 pedestrian-vehicle crashes to occur during unclear weather conditions which is consistent with the result found

450 herein.

451 **Limitations of the study**

452 First of all, this study was aimed at finding out variables that increase the probability of having a

453 fatality or severe injury resulting from a vehicle-pedestrian crash, and should not be confused with

454 other studies that aim at finding predictor variables that are associated with frequency of pedestrian

455 crashes. That is said, it should be mentioned that the crash data collection is based on the standard

456 Jordanian Police Traffic Department reporting system, where crashes were not geocoded by using

457 milepost data. Consequently, some variables that might have significant effect on vehicle-pedestrian

458 crashes were either unavailable or could not be linked to the data obtained and hence were not

459 included herein.

460 This study did not include studying the effect of Average Annual Daily Traffic (AADT) on injury

461 severity of vehicle-pedestrian crashes since it was not available at any transportation agency. Studies

462 found in literature used AADT as a predictor variable in frequency of vehicle pedestrian crashes, and

463 it was found to be a significant predictor in many of these studies (Brugge et al., 2002; LaScala et al.,

464 2000; Lee and Abdel-Aty, 2005; Loukaitou-Sideris et al., 2007). However, a study on injury severity

465 of vehicle-pedestrian crashes found that AADT was a relatively a weak predictor; they referred the

466 reason of this result to the fact that state routes with high AADT have lower speeds (Moudona et al.

467 2011). The results of this study did not differ than the results obtained by Ewing (2006) and  NHTSA

468 (1999) in which they found that  injury severity is strongly associated with vehicle speed.

469 Population density  which is usually used to define area type (urban or rural) was also found by many

470 presvious research works to have a significant effect on pedestrians' crashes. However, its effect on

471   higher injury severity levels is reduced as population density increases (Goel et al. 2018). The reason

472   behind this result is referred to the fact that high numbers of pedestrians are expected to be in highly

473   populated areas, and in the absence of dedicated facilities for pedestrians, pedestrians tend to use lanes

474   next to curb-side which slows down vehicles' speeds and decreases the risk of high severity of

475   pedestrian crashes, this phenomenon is called safety-in-numbers (Jacobsen, 2003;Elvik and

476   Bjørnskau, 2017).

477   Vehicle type was also missing from the dataset, eventhough it is considered as an important predictor

478   that affects the severity outcome of a vehicle-pedestrian crash. Ballesteros et al. (2004) found that

479   size, weight, and design of vehicles involved in vehicle-pedestrian crashes are related to injuries of

480   some types.

481   **Conclusions**

482   Most research on traffic crashes injury classification use injury severity as the class variable without taking into

483   account the possibility of having better performance if another class variable was chosen for the same purpose.

484   In order to find out if this assumption is correct, a comparison was made between Bayesian networks developed

485   using injury severity as a class variable and that using collision type as a class variable. The results in this

486   research work indicated that the developed Bayesian networks using collision type as a class variable obtained

487   better performance (more reliable results).

488   On the other hand, one of the main challenges encountered in safety research that affect the classification of

489   crashes is imbalanced datasets, as a rule of thumb, traffic crashes datasets are usually imbalanced and thus the

490   results will be biased towards the class with higher number of records (majority class i.e. slight injuries class)

491   (Mujalli et al. 2016). In the present work, a dataset was obtained from Jordanian Police Traffic Department,

492   where the data that was obtained was imbalanced and hence needed to be balanced prior to developing models

493   in order to enhance the classification of collisions.

494   For analysis and modeling purposes, three subsets were extracted from the original dataset; the first subset

495   included all crash types but used collision type as a class variable, the second subset also included all crash

496   types but used injury severity as a class variable, and the third subset included only pedestrian crash type and

497   used injury severity as a class variable. All subsets were first balanced, and Bayesian networks were developed

498   combining different search algorithms (K2 and simulated Annealing) and scores (BDeu, MDL, AIC). This

499   process resulted in developing a total of six Bayesian networks per subset.

500 In order to measure the reliability of each developed Bayesian network using each subset, their classification

501 capabilities were tested and the results indicated that the most reliable Bayesian networks were those which used

502 collision type as the class variable.

503 The Bayesian network with the most reliable results from each subset was used to extract significant rules which

504 if existed in the same crash, the crash will end up in a fatal or a severe injury. The extracted rules were grouped

505 into three levels of confidence, as a result the following variables were found to significantly increase the

506 probability of a fatal or severe injury outcome in pedestrian-vehicle collisions: speed limit, road type, and

507 number of lanes where all these variables were present in the rules belonging to high confidence level group.

508 Other variables that were found to be significant in fatal or severe injury collisions included: number of involved

509 vehicles in a crash, lighting and adverse weather conditions. Consequently, using all subsets enhanced the

510 extraction of rules and possibly using more variables (if available) would help extracting more valuable

511 knowledge when more than one subset is employed. These results, which were consistent with previous studies,

512 could be of a great interest in case they were used to enhance the existing shortcomings in the applied norms

513 that govern the pedestrian safety on roads.

514 **Recommendations**

515 Authors recommend that some prevention strategies should be followed by both governmental and private

516 agencies to improve road safety through actions aimed at decreasing the occurrence of pedestrian-vehicle

517 collisions and resulting injuries and fatalities. Actions should be directed towards monitoring both attitude and

518 behavior of drivers and pedestrians especially in urban areas where roads are wide and divided, without proper

519 lighting, and with speed limits of 30-40 km/hr since these locations were found to be more hazardous and

520 awareness must be raised to help reduce both fatalities and severity of collisions at these locations. Also,

521 pedestrian sidewalks should be maintained and be only used by pedestrians, where crosswalks should be clearly

522 marked and penalties should be imposed on drivers stopping on them so that pedestrians' penalties crossing

523 from illegal points are accordingly activated. Speed limit compliance enforcement should be done using speed

524 cameras especially at pedestrian "black spots", and traffic calming measures as well should be used more

525 frequently at these locations to lessen the problem.

526 Authors believe that important knowledge could be extraced if the flow rates or average anual daily traffic

527 (AADT) on travelled roadways where crashes occured were available, and hence a recommendation for Traffic

528 Police Department is given to collaborate with the ministry of public Works and Housing and local

529 muncipalities in order to collect and maintain this type of data.

530    Future work in this field might be directed towards using other data mining techniques in order to compare the

531    results obtained using Bayesian networks. Moreover, to the knowledge of authors, no research is currently

532    available about run-off-road crashes and property damage only crashes that occur in Jordan, and hence this

533    might be a future line of investigation.

534    **Acknowledgements**