

Document downloaded from:

<http://hdl.handle.net/10251/120538>

This paper must be cited as:

López-Maldonado, G.; De Oña, J.; Garach, L.; Baena-Ruiz, L. (2016). Influence of deficiencies in traffic control devices in crashes on two-lane rural roads. *Accident Analysis & Prevention*. 96:130-139. <https://doi.org/10.1016/j.aap.2016.08.008>



The final publication is available at

<http://doi.org/10.1016/j.aap.2016.08.008>

Copyright Elsevier

Additional Information

1 **Abstract**

2

3 One of the main objectives of all public administrations is reducing traffic crashes. To this end, Road
4 Safety Inspections (RSI) stand out as a key measure. Signaling roads is one of the foremost tasks of
5 RSI. A road that is improperly or poorly signaled can lead to incorrect placement or maneuvers of
6 vehicles and ambiguous situations that can increase the risk of crashes. This paper analyses the
7 relationship between road crashes in two-lane rural highways and certain deficiencies in signaling. The
8 results show that deficiencies such as “incomplete removal of road works markings” or “no guide sign
9 or in incorrect position” are the ones associated with a higher probability of crashes in two-lane rural
10 highways. In view of these results, governmental agencies should verify that the original conditions of
11 a highway are re-established after any construction work is completed. They should also continuously
12 follow up on the signaling of this type of highway in order to maintain them in optimal conditions.

13

14 **Keywords.** Traffic crashes; road safety inspections; sign and marking; Decision Trees; Decision rules

15

16 **1. Introduction**

17

18 Traffic accidents are complex events involving the interaction of different contributory factors,
19 including road, driver, and vehicle. While it is well known that the human factor is the main cause of
20 traffic crashes, present in nearly 90% of them ([Siskind et al., 2011](#)), previous studies have shown that
21 the infrastructure also plays a significant role. Nearly 28% of crashes are due to infrastructure and, in
22 most cases a combination of human and road factors forms a major contribution in the road crashes
23 ([Odgen, 1996](#)).

24

25 In the literature, the crash contribution from human factors is usually analyzed in the context of driver
26 errors. The human error most often identified in the crashes is related to the perception and processing
27 of information presented by the road or traffic environment. Situations that cause problems with road

28 user perception, interpretation or judgment stages may lead to driver error or loss of control (Croft and
29 Schnerring, 2009). An estimated 30% of driver-distracted crashes derive from diverse sources outside
30 the vehicle (Regan et al., 2009). Hence, it is crucial to maintain the road features in optimal conditions
31 so that they have the least possible impact on the driver's performance.

32

33 Reducing highway crashes is one of the main aim of the Administrations. One means of reducing them
34 is to detect and correct roadway deficiencies. Road Safety Inspections (RSI) were established for this
35 purpose, they are an effective tool for the management of safety on existing roads. The European
36 Directive on Road Infrastructure Safety Management (EC, 2008) defines RSI as "an ordinary
37 periodical verification of the characteristics and defects that require maintenance work for reasons of
38 safety". Following the principle "Prevention is better than cure", the RSI are used to evaluate existing
39 road traffic facilities and to improve road safety performance (Alfredas et al., 2012). While some RSI
40 treatments will have a greater impact than others, as underlined by Elvik (SETRA 2008), significant
41 reductions in crashes can be expected as a result of a RSI and associated interventions.

42

43 After some years of experience with RSI, it is broadly recognized as one of the most important and
44 effective engineering tools available to improve road safety (Antov, 2011). This is why the European
45 Union makes RSI mandatory for trans-European Road Networks and they are recommended for the
46 rest of the transport infrastructures (EC, 2008). These inspections should be undertaken after
47 establishing a series of criteria to be articulated by means of checklists. The checklists are ordered lists
48 used to cover the most important issues which should be inspected during the RSI. The detected
49 hazards will be identified like Road Safety Deficiency (RSD).

50

51 Some aspects of a RSD that can be analyzed include those related with signaling. They are easier to
52 correct and involve a lower cost than other measures, such as the design of the road itself. It is
53 important for the highway to be properly signaled, and that the information provided is clear and
54 concise.

55 Several authors stress the importance of correct signaling: [Miller \(1992\)](#) reported that existing
56 longitudinal pavement markings reduce crashes by 21%, and edge lines on rural two-lane highways
57 reduce crashes by 8%; and [Cho et al. \(2012\)](#) suggested that pavement markings provide guidance to
58 road travelers. An alteration of pavement color and/or texture or incomplete removal of pavement
59 markings during construction projects could confuse individuals driving through the construction work
60 zones. To make matters worse, under certain lighting and weather conditions the supposedly removed
61 markings may become more visible than the new ones. [Antov \(2011\)](#) highlighted as common
62 problems the missing, contradictory or incomprehensible signs/markings. [Croft and Schnerring \(2009\)](#)
63 pointed out that incorrect or poorly maintained pavement marking can lead to undue placement or
64 maneuvers of vehicles, thus increasing the risk of crashes. They also showed the influence of
65 delineation devices in road safety—poorly placed or missing delineation devices can transmit a false
66 picture of the way ahead, contributing to driver error ([Croft and Schnerring, 2009](#)). The
67 methodological approach for safety evaluation of two-lane rural highways segments put forth by
68 [Cafiso et al. \(2007\)](#) served to establish that daytime delineation of a road can be effectively
69 accomplished with pavement markings, whereas nighttime and rainy conditions may require a
70 different approach to provide long-range delineation of the roadway alignment. Supplementary
71 delineation is an important safety factor in any condition; but it may prove critical on horizontal
72 curves, especially on isolated curves with a short radius. [Croft and Schnerring \(2009\)](#) also indicated
73 that signs poorly located/incorrectly situated can cause confusion, increasing crash risk, just as
74 excessive signing can increase potential risk for road users. [Montella \(2005\)](#) described a systematic
75 process to determine which road features should be investigated and how each should be evaluated
76 during RSI. Accordingly, a safety improvement index was calculated and compared with the expected
77 collision frequency, and this procedure was carried out in 406 km of rural two-lane rolling highways in
78 Italy. The study revealed that for missing or ineffective curve warning signs on severe curves, the
79 relative risk factor could be assumed equal to 10% ([Montella, 2005](#)).

80

81 This study analyzes the relationship between crashes and certain deficiencies in signaling identified by
82 a previous RSI. The RSI was performed on two-lane rural highways in Andalusia (Spain). From the
83 Road Safety standpoint, it is vital that two-lane rural highways be studied, as they are the scenario of
84 most crashes. In Spain, 70% of crashes occur on this type of roads ([Ministerio del Interior, 2013](#)).

85

86 The analysis uses a data mining technique. This technique has been widely used in the road safety
87 field in recent years, giving satisfactory results ([Kuhnert et al., 2000](#); [Sohn and Shin, 2001](#); [Abdel
88 Wahab and Abdel-Aty, 2001](#); [Chang and Wang, 2006](#); [De Oña et al., 2011](#); [Kashani et al., 2011](#);
89 [Pakgohar et al., 2010](#); [Chang and Chien, 2013](#); [De Oña et al., 2013a](#); [De Oña et al., 2013b](#); [López et
90 al., 2014](#)). The main aim of this technique is the extraction of knowledge from large amounts of
91 previously unknown and indistinguishable data. In this case, Decision Trees (DTs) are employed. DTs
92 are appropriate for studying crashes because they are non-parametric techniques that do not require
93 prior probabilistic knowledge of the study phenomena. Further advantages of DTs with respect to
94 other methods having similar aims reside in the extraction of Decision Rules (DRs) ([De Oña et al.,
95 2013a](#)). Although each crash is the result of a unique chain of events, some specific factors are
96 common to several crash circumstances, and DRs can be used to identify these factors and their
97 interdependences ([Montella et al., 2011](#)). Safety analysts could use these rules to understand the events
98 leading up to an accident, and prioritize certain elements for actions intended to improve road safety.

99

100 In this paper, therefore, DRs extracted from DTs are used to analyze the relationship between the
101 actual occurrence of traffic crashes on two-lane rural highways and the deficiencies in roadway
102 signaling previously detected by means of RSI.

103

104 The paper is organized as follows: Section 2 presents a description of the data used, and also describes
105 the procedures for building DTs, extract DRs, and deriving the final rule set. Section 3 presents the
106 Results and a Discussion thereof. Finally, the last section succinctly presents some Conclusions.

107

108 2. MATERIALS AND METHODS

109

110 2.1 Description of the data

111

112 The data come from two different sources. The Andalusian Regional Government provided the Road
113 Inventory database and the Road Safety Inspections database, while the Spanish General Directorate
114 of Traffic provided the Spanish Road Crashes database.

115

116 The Road Inventory database contains a list of road sections with their geometrics and equipment
117 characteristics. Two-lane rural highways from the Complementary Road Network of Andalusia were
118 selected for this study. Urban segments, junctions and segments with road work places were removed
119 from the study; because the factors related to crashes taking place on these sections are different, they
120 should be analyzed separately (Moore et al., 2010). The total length of the investigated road network is
121 1,635 km.

122

123 The Road Safety Inspections database contains information about a RSI developed on two-lane rural
124 highways in the Complementary Road Network in Andalusia. In this RSI some risks associated with
125 RSD were identified. These risks were defined as Road Safety Deficiency Elements (RSD-E). The
126 risks related with the vertical signs and pavement markings identified during the RSI are denominated
127 Signaling Elements from Road Safety Inspection (RSD-SE). The main aim of this study is to
128 investigate the influence of RSD-SE on road crashes.

129

130 The Spanish Road Crashes database contains a description of the location and type of crashes that
131 occurred on Spanish roads. Information about the crashes in two-lane rural highways in the
132 Complementary Road Network of Andalusia was extracted from this database. The period of study is
133 three years (2006-2008), and during this period the total number of crashes with victims in these
134 segments was 1,454.

135 A global database with information about crashes, road characteristics, and RSD-SE was built using
136 the three databases. The following analysis is based on seven variables related to geometric and
137 environmental road characteristics (see Table 1) and eight RSD-SE (see Table 2).

138

139

[Insert here Table 1]

140

141

[Insert here Table 2]

142

143 The following criteria were corroborated in order to identify the RSD-SE:

144

- 145 • RSD-SE1: The length of the passing zone was at least the minimum indicated by Spanish
146 National Standards. For a speed of 100 km/h, the minimal distance is 250 meters. For 60 km/h
147 the minimum is 75 meters. For speeds in-between, intermediate distances are established.
- 148 • RSD-SE2: The regulatory signs are present and correctly positioned (e.g., speed limit or no
149 passing zone).
- 150 • RSD-SE3: Signs indicating danger/precaution are present and correctly positioned (e.g., road
151 narrows, dangerous curve, animal crossing, etc.)
- 152 • RSD-SE4: Guide signs are properly situated.
- 153 • RSD-SE5: Road markings are clear and visible.
- 154 • RSD-SE6: There are no contradictions between vertical signs and road markings at a given
155 point.
- 156 • RSD-SE7: This deficiency is considered to exist in segments where road markings have not
157 been adequately eliminated.
- 158 • RSD-SE8: This deficiency is considered to exist in segments where the road width is greater
159 than 7 meters and there are no post-mounted delineators or they present damage amounting to
160 over 50%.

161 **[Insert here Figure 1]**

162

163 **2.2. Classification and Decisions Trees (CART)**

164

165 DTs are one of the most widely used data mining techniques for classifying and predicting class
166 variables. When the target variable is discrete, a classification tree is developed, whereas a regression
167 tree is developed for continuous variables. CART can be used for both kinds of target variables. In this
168 study, the target variable is the occurrence of the accident (ACC: YES or NO) and, therefore, a
169 classification tree is developed.

170

171 A DT is an oriented graph formed by a finite number of nodes departing from the root node. DTs are
172 built recursively, following a descending strategy, starting with the full data set (made by the root
173 node). Using specific split criteria, the full set of data is then split into even smaller subsets. Each
174 subset is split recursively until all of them are pure (when all cases in each subset present the same
175 class) or their “purity” cannot be increased. Thus the tree’s terminal nodes are formed, obtained
176 according to the answer values of the target variable (De Oña et al., 2013a).

177

178 The CART method is a particular methodology for building binary Decision Trees in which the Gini
179 Index is used as the splitting criterion. The development of a CART model generally consists of three
180 steps: (1) growth of the tree; (2) the pruning process; and (3) selection of an optimal tree from the
181 pruned trees. Tree growing entails recursive partitions of the target variable to maximize “purity” in
182 the two subsequent child nodes. By definition, the terminal nodes present a lower degree of impurity
183 compared to the root node. In tree growing, predictors generate candidate partitions (or splits) at each
184 internal node of the tree, so that a suitable criterion needs to be defined in order to choose the best
185 partition (or the best split) of the objects. The Gini reduction criteria is applied to measure the “worth”
186 of each split in terms of its contribution toward maximizing the homogeneity through the resulting

187 split. If a split results in the splitting of one parent node into B branches, the “worth” of that split may
188 be measured as follows:

189

$$190 \quad \text{Worth} = \text{Impurity (Parent node)} - \sum_{n=1}^N P(n) * \text{Impurity}(n), \quad (1)$$

191

192 where Impurity (Parent node) denotes the Gini measure for the impurity (i.e., non-homogeneity) of the
193 parent node, and P(b) denotes the proportion of observations in the node assigned to branch b. The
194 impurity measure, Impurity (node), may be defined as follows:

195

$$196 \quad \text{Impurity (node)} = 1 - \sum_{i=1}^I \left(\frac{\text{number of class } i \text{ cases}}{\text{all cases in the node}} \right)^2. \quad (2)$$

197

198 When a node is ‘pure’, Eq. (2) gives the minimum value, and its value will be higher for less
199 homogeneous nodes. If one considers the definition of “worth” according to Eq. (1), a split resulting in
200 more homogeneous branches (Child nodes) will have more “worth”.

201

202 While developing a CART this criterion is applied recursively to the descendants to achieve Child
203 nodes having maximum worth which, in turn, become the parents to successive splits, and so on. The
204 splitting process goes on until there is no (or less than a pre-specified minimum) reduction in impurity
205 and/or the limit for a minimum number of observations in a leaf node is reached. Following this
206 process, a saturated tree is obtained. The saturated tree provides the best fit for the used database, but
207 overfits the information contained within the database, and this overfitting does not help in classifying
208 other databases. Therefore, when developing a CART model data is usually divided into two subsets:
209 one for learning (or training) and the other for testing (or validation).

210

211 The learning sample is used to split nodes, while the testing sample is used to compare the
212 misclassification. The saturated tree is obtained from the learning data. Overly large trees could result

213 in higher misclassification when applied to classify new databases. To decrease its complexity, the tree
214 is pruned in a second step according to the cost-complexity algorithm, which is based on removing the
215 branches that add little to the predictive value of the tree. The cost-complexity measure combines the
216 precision criteria as opposed to complexity in the number of nodes and processing speed, searching for
217 the tree that obtains the lowest value for this parameter. The final step gives rise to the optimal tree. A
218 more detailed description of the CART method can be found in [Breiman et al. \(1984\)](#).

219

220 **2.3. Decision Rules (DRs)**

221

222 The DT's structure can be transformed into rules in order to extract its potentially useful information.
223 A DR is a logical, conditional structure of the type if $A \rightarrow B$, in which A is the antecedent of the rule
224 and B is the consequent, with all the splits of the parent nodes being the antecedent and the class of the
225 terminal node being the consequent.

226

227 Each rule starts at the root node and each variable that is included in tree division makes an IF of the
228 rule, which ends in terminal node with a value of THEN (which is associated with the state resulting
229 from the terminal node). The class of a node is the status that shows the highest number of cases.
230 Thus, a priori, the number of rules can be identified with the number of terminal nodes in the tree.

231

232 Due to the fact that the occurrence of crashes is infrequent in comparison with the non-occurrence of
233 crashes, the class of the terminal node —and therefore the class resulting from the rule— will usually
234 be the non-occurrence of an accident ($ACC=NO$). Notwithstanding, from the road safety perspective,
235 the rules of interest are the rules involving crashes. To identify this type of rule, and following
236 previous studies ([Montella et al., 2012](#); [López et al., 2014](#)), we use the posterior classification ratio
237 (PCR) in order to re-assign a response class (the consequent) to each rule extracted. PCR compares the
238 classification of the terminal nodes of the tree with the classification of the root node (Eq. 3):

239

240
$$PCR(j|t) = \frac{p(j|t)}{p(j|t_{raiz})} \quad (3)$$

241

242 where:

243 $p(j|t)$ = Proportion of observations in node “t” that belong to the class “j”, where class “j” is
 244 “YES”; t_{root} = Root node of the tree.

245

246 The assignment of the class to each rule was performed selecting the class j^* with the greatest value of
 247 PCR. In addition, we will analyze only rules in which the consequent of the class variable (ACC) is
 248 the accident occurrence (YES). For each rule, then, two parameters are calculated: Support and
 249 Probability of accident (in three years’ time).

250

251 Support: The support of the rule (S) is the percentage of the data set for which both A (antecedent) and
 252 B (consequent) appear, that is, the number of cases in which the following rule is fulfilled:

253

254
$$S(A \rightarrow B) = \frac{|(A \rightarrow B)_t|}{N} \quad (4)$$

255 Where $(A \rightarrow B)_t$ is the number of crashes for which both conditions A and B are verified; N is the total
 256 number of crashes.

257

258 Probability of accident (in 3 years): Indicates the probability that an accident will occur in three years’
 259 time as a consequence of the circumstances given in the rule.

260

261
$$\text{Prob. accident 3 years} = PCR * \text{Prob. global acc} \quad (5)$$

262

263 Where $\text{Prob. global acc} = \frac{\text{crashes}}{\text{km network}}$

264

265 Because of the large number of patterns considered, DTs may suffer from an extreme risk of Type I
266 error, that is, of finding patterns that appear only by chance to satisfy constraints on the sample data
267 ([Webb, 2007](#)). To reduce the risk of Type I error, and following other authors ([Montella et al., 2012](#);
268 [Kashani and Mohaymany, 2011](#)), the dataset was split randomly in two parts: a training set (70%) and
269 a testing set (30%). The rules extracted on the training set were validated using the testing set. The
270 application of the tree structure obtained in the training set to the testing sample produced the testing
271 tree that was used for validation. To reduce the risk that results were overfitted to the sample, at each
272 node of the testing tree the assignment of the class was compared with the assignment performed in
273 the training tree. As a result, only nodes with the same class in both the training and the testing trees
274 were validated.

275

276 **2.4. Decision Rules obtained from a Decision Tree: The global DRs set**

277

278 The extraction of knowledge with DRs extracted from a DT has some limitations. The rules depend on
279 the DT's structure because they are extracted from each tree branch from the root node to the terminal
280 node. Therefore, knowledge is extracted only in the direction dictated from the root node to the
281 terminal node even if other possible important rules could exist. To extract all the possible patterns
282 from a particular data set, [Abellán et al. \(2013\)](#) proposed a method called information root node
283 variation (IRNV). The main characteristic of the IRNV method is that a set of DTs is built by varying
284 the root node. Thus, every possible set of DRs is obtained from each tree, providing a set of rules with
285 potentially useful information.

286

287 The first step in order to obtain DRs from the different DTs built varying root node was to randomly
288 split the dataset into the training set (70% of the data) and the testing set (30%). Then, based in the
289 IRNV method, a total of 15 models of DTs varying the root node are developed (i.e., a different model
290 for each one of the seven variables and the eight RSD-SE considered for studying). All the rules in
291 which the consequent is the occurrence of the accident are extracted from these models. The main

292 problem with this method is that most rules are extracted from DTs to which the root node has been
 293 imposed, and this node could not be essential for the pattern that describes the rule. To overcome this
 294 issue a procedure of verification of the root node is performed (López et al., 2014) to determine
 295 whether the rule should be simplified.

296

297 DR is the rule extracted from a DT in which the root node is imposed (called in this study the extended
 298 rule); and DR^- is the rule without the root node (called in this study the simple rule); A is the
 299 antecedent of the DR and is formed by n variables ($X'1, X'2, \dots, X'n$); A^- is the antecedent of the DR^-
 300 and is formed by $n-1$ variables ($X'2, \dots, X'n$). In this way, we have to compare $DR: A (X'1, X'2, \dots,$
 301 $X'n) \rightarrow B$ vs. $DR^-: A^- (X'2, \dots, X'n) \rightarrow B$, where B is the consequent. The extended rule (rule with n
 302 items) is selected over a simple rule (rule with $n-1$ items) if it verifies two conditions (López et al.,
 303 2014):

304 Condition 1:
$$\frac{PCR(A \rightarrow B)}{PCR(A^- \rightarrow B)} \geq 1.03 \quad (6)$$

305

306 Condition 2:
$$\frac{S(A \rightarrow B)}{S(A^- \rightarrow B)} \geq 0.2 \quad (7)$$

307

308 Condition 1 establishes that the increase of PCR in the DR should be over 3%; and Condition 2
 309 indicates that the support of the DR with respect to the DR^- should be, at least, 20%. Thus, the global
 310 DRs set is formed by extended rules (DRs) when conditions 1 and 2 are verified simultaneously, or
 311 simple rules (DRs^-) if one of the conditions is not verified.

312

313 Once the simple rule or the extended rule has been selected, the chosen rules are validated in the
 314 testing set. The PCR is calculated again, and the rules fulfilling $PCR \geq 1$ (the rules whose consequents
 315 are the occurrence of the accident) are the validated rules. The rules that are validated become part of
 316 the final set of RDs and should be analyzed from the road safety standpoint.

317

318 **3. RESULTS**

319

320 In the first step, the dataset was randomly split into training and testing sets: 1,174 km formed the
321 training set, having 738 road sections with crashes and 10,989 road sections without crashes.

322

323 The different models of DTs are built varying the root node using the training set. DT₁ is the model
324 obtained directly, without imposing the root node, whereas DT₂ to DT₁₅ are the models obtained
325 varying the root node. Table 3 shows the main results for the 15 models. 106 rules were extracted from
326 the different models. From these rules, only rules with $PCR \geq 1$ in the training set (i.e., rules whose
327 consequent is the occurrence of the accident) are selected (62 rules).

328

329 **[Insert here Table 3]**

330

331 **[Insert here Table 4]**

332

333 In following, the root node is verified. This verification is only necessary for patterns obtained from
334 DT₂ to DT₁₅ (DTs in which the root node was imposed). Rules for DT₁ do not call for such verification
335 because they do not come from a DT whose root node is imposed. The procedure compares the
336 extended rule (*DR*) and the simple rule (*DR*⁻). Altogether, 62 rules were analyzed: 20 rules are *DR*
337 (verify conditions 1 and 2 simultaneously), and all the others (42) are *DR*⁻ (do not verify one
338 condition).

339

340 Finally, the rules were validated in the testing set and a total of 61 rules were obtained. Given that
341 most of the rules are simplified (they are *DR*⁻) some appear more than once. After this process, only 17
342 rules remain forming the global *DR* set.

343 Table 4 shows the rules grouped in four sets to show their common patterns. In the first group (three
344 rules), the rules only have one RSD-SE as a deficiency; in the second group (five rules), the rules are

345 formed by RSD-SE with some deficiencies and geometric or environmental variables; in the third
346 group (six rules), the rules are formed by RSD-SE without deficiencies (RSD-SE=NO) and geometric
347 or environmental variables; and in the fourth group (three rules), the rules are formed by RSD-SE
348 (with or without deficiencies) and geometric or environmental variables.

349

350 Table 4 shows the values of PCR and the probability of accident (in three years) in the training and
351 testing set (for each rule). The average probability of accident in the network is 6.28%.

352

353 Rules in the first group show a direct relationship between some signaling deficiencies and the
354 occurrence of crashes. Rule 1 shows that the incomplete removal of road works markings (RSD-SE
355 7=Y) presents a probability of accident between 20% (value in the training set) to 22% (values in the
356 testing set). This translates as an increased probability of the order 220% to 254% (see value of PCR)
357 with regard to the mean values of accident probability in the network analyzed (6.28%).

358

359 Table 4 shows that RSD-SE7 is present in the rules entailing a greater probability of accident (rules 1,
360 4, 5, 8 and 16). Rules 4 and 5, in which there are incomplete removal of road works markings, on
361 roads with AADT less than or equal to 5000 veh/day, have similar values for probability of accident
362 (between 18% and 23%). With the same values of AADT, if the terrain is flat or rolling, the probability
363 of accident is between 21% and 25% (see Rule 8). Rule 16 reflects another pattern for roads with the
364 RSD-SE7. Although it does not involve deficiencies with warning signs (RSD-SE3=N), the values of
365 probability of traffic crashes are also high, varying between 20% and 25%. Some researches ([Cafiso et
366 al., 2007](#); [Miller, 1992](#)) have described the involvement of deficient road marking in crashes, showing
367 that their improvement is likely to be cost-effective. For example, on roads with edge lines missing, a
368 relative increase in injury accident risk of 8% could be assumed; and when the center line is missing
369 the risk increases to 13% ([Safety Audits of Existing Roads, 2003](#)). [Ellis and Pyeon \(2006\)](#) indicated
370 that pavement work markings not properly removed may confuse or distract drivers. Alteration of
371 pavement color and/or texture, as well as incomplete removal of pavement markings, has been

372 identified as a particular problem for motorists; they can be mistaken for navigable lanes through
373 construction work zones. Because motorists or drivers heavily rely on pavement markings for roadway
374 guidance, it is imperative to remove old markings to reduce crashes owing to lane confusion ([Cho et](#)
375 [al., 2012](#)).

376

377 Rule 2 shows that lack of correspondence between vertical signs and road markings (RSD-SE6)
378 presents a direct relationship with crashes. When this deficiency appears on the analyzed roads, the
379 probability of accident reaches 14% to 21%. This means an increased probability that is 127% to
380 235% greater than the mean values for accident probability. According to the results, several road
381 safety problems identified with the lack of correspondence between vertical signs and road markings
382 can be tied to accident risk. The study by [Antov \(2011\)](#) highlights problems stemming from missing,
383 contradictory or difficult to read signs/markings, but further typical deficiencies are incomplete or
384 misleading signs/road markings, or an “overload” of information.

385

386 Rule 3 shows a direct relationship between crashes and instances when the guide sign does not exist or
387 it is in an incorrect position (RSD-SE4). In this case, the probability of accident is similar to Rule 2
388 (between 15% and 21%). This stands as an increased probability ranging from 135% to 238% beyond
389 the mean values. Some investigations reveal the importance of vertical signs on traffic crashes. [Cafiso](#)
390 [et al. \(2007\)](#) evaluated two-lane rural highways and established that regulatory signs, such as speed
391 limits, could affect road safety by conveying essential information on safe behavior. For missing or
392 ineffective signs, the relative risk factor was assumed as equal to 20%. [Croft and Schnerring \(2009\)](#)
393 likewise established that poorly located or incorrect signs could lead to a confusing and ambiguous
394 situation, increasing crash risk.

395

396 Rules 6 and 7 show the relationship between deficiencies in RSD-SE 2 and crashes on road with
397 AADT higher than 5,000 veh/day, and with road width between 5 and 6.5 meters. Rule 7 shows that
398 the probability of accident varies from 20% to 21%, meaning respective increased probabilities of the

399 order 215% to 237% over the mean values. The influence of regulatory signs on accident occurrence
400 has been investigated in two-lane rural highways by some researchers. [Cafiso et al. \(2007\)](#) reported
401 explanations of the relative increase in accident risk for some safety issues as vertical sign,
402 determining that regulatory signs such as for speed limits could affect road safety by conveying
403 essential information on safe behavior.

404

405 Rule 6 the same pattern, adding deficiencies in road markings (not exist or were deleted - RSD-SE5).
406 In this case, the probability of an accident in 3 years is around 20%. This would be an increased
407 probability between 214% 217% greater than the average values for probability of accident in the
408 network. Previous studies have shown that incorrect or poorly maintained pavement markings can lead
409 to incorrect placement or maneuvers of vehicles, and increase the risk of crashes ([Croft and
410 Schnerring, 2009](#)).

411

412 In the third group, the rules are formed by RSD-SE without deficiencies and geometric or
413 environmental variables. In Rule 14 the probability of accident is similar to the probability of accident
414 of the network analyzed (6.28%). The values of probability in Rules 9, 11 and 12 range between
415 6.49% and 14.89%. Only in two rules (Rules 10 and 13) the probabilities are greater, increasing to
416 values of 17%-21%; these rules are identified on roads with high values of AADT (>5,000 veh/day)
417 and roadway width between 5 and 6.5 meters. Such findings underline that two-lane rural highways
418 with major traffic flow (AADT) entail an increased risk of accident.

419

420 In the fourth group, the rules are formed by RSD-SE (with or without deficiencies) and geometric or
421 environmental variables. As seen for Rules 15 and 17, when there are no deficiencies with the RSD-
422 SE7 (incomplete removal of road work markings) on roads with values of AADT between 1,000 to
423 5,000 veh/day, even if other elements fail, the probability of accident is low, very similar to the
424 probability of the network.

425

426 4. CONCLUSIONS

427 This paper presents an analysis of deficiencies in signaling with regard to crashes on rural highways.
428 In addition, some variables related to geometric and environmental road characteristic were used, and
429 the Data Mining technique of Decision Trees was applied. In order to derive all the information
430 possible from the database analyzed, different DT models were built, varying the root node, and from
431 each of the models the DRs of interest were extracted (rules whose consequence is the occurrence of
432 an accident). As a result, 62 rules were obtained, and 61 of them were validated. After elimination of
433 the rules that were the same, a total of 26 rules made up the final set.

434

435 In order to perform a safety analysis of the rules, they are grouped in four sets: rules directly relating
436 crashes with signaling deficiencies; rules relating crashes with deficiencies in signaling and roadway
437 characteristics; rules that do not involve deficient signaling or highway characteristics, but under
438 certain geometric and/or environmental conditions bear a relation with crashes; and rules that present
439 deficiencies in some elements, in others no, and have geometric and/or environmental variables
440 present.

441

442 In general, the element RSD-SE7 (incomplete removal of road works markings) appears in the rules
443 with the greatest probability of accident in 3 years (Rules 4, 5, 8 and 16), producing in turn a greater
444 probability of accident. RSD-SE7 appears with AADT less than or equal to 5,000 veh/day, which may
445 indicate that the pavement markings are not properly re-established in this type of roadway. This result
446 shows that the government agencies or local administration should verify that after construction is
447 finished, the original conditions of a roadway must be re-established quickly and efficiently.

448 Deficiencies in RSD-SE4 (sign does not exist or it is not correctly situated) are also associated with a
449 greatly increased probability of accident in the network analyzed. This finding serves to accentuate the
450 importance of maintaining signaling. Indeed, it is recommended that administrations make vigilance
451 and follow up of roadway signs and signals a priority, ensuring that they are in optimal conditions.

452 **Acknowledgments**

453 The authors would like to acknowledge FEDER funding by the European Union for financial support
454 via project “Análisis de la relación entre Elementos Susceptibles de Mejora, Accidentes y TCA” of the
455 “Programa Operativo FEDER de Andalucía 2007-2013”. We also thank the Public Works Agency and
456 Regional Ministry of Public Works and Housing of the Regional Government of Andalusia. The
457 authors are grateful to the Spanish General Directorate of Traffic (DGT) for providing the data
458 necessary for this research. Griselda López wishes to express her acknowledgement to the regional
459 ministry of Economy, Innovation and Science of the regional government of Andalusia (Spain) for a
460 scholarship to train teachers and researchers in Deficit Areas.

461

462 **REFERENCES**

463 Abdelwahab, H., Abdel-Aty, M. (2001). Development of artificial neural network models to predict
464 driver injury severity in traffic accidents at signalized intersections. Transportation Research Record:
465 Journal of the Transportation Research Board, (1746), pp. 6-13.

466

467 Abellán, J., G. López, J. de Oña (2013). Analysis of Traffic Accident Severity Using Decision Rules
468 via Decision Trees. Expert Systems with Applications, 40, pp.6047-6054.

469

470 Alfredas Laurinavičius , Lina Juknevičiūtė-Žilinskienė , Kornelija Ratkevičiūtė , Ineta Lingytė , Laura
471 Čygaitė , Vytautas Grigonis , Rasa Ušpalytė-Vitkūnienė , Dago Antov , Tiit Metsvahi , Zsuzsanna
472 Toth-Szabo, András Várhelyi (2012) Policy instruments for managing road safety on EU-roads,
473 Transport, 27:4, 397-404, DOI: 10.3846/16484142.2012.751934

474

475 Antov, D. (2011). Road Safety Inspection Guidelines and Checklists. BALTRIS, 43 pp.

476

477 Breiman, L., Friedman, J., Olshen, R., Stone, C. (1984). Classification and Regression Trees. Belmont,
478 CA: Chapman and Hall.
479

480 Cafiso, S., Cava, G., Montella, A. (2007). Safety index for evaluation of two-lane rural highways.
481 Transportation Research Record: Journal of the Transportation Research Board, 2019, pp. 136-145.
482

483 Chang, L. Y., Wang, H. W. (2006). Analysis of traffic injury severity: An application of non-parametric
484 classification tree techniques. Accident Analysis and Prevention, 38, pp.1019–1027.
485

486 Chang, L. Y., Chien, J. T. (2013). Analysis of driver injury severity in truck-involved accidents using a
487 non-parametric classification tree model. Safety science, 51(1), pp. 17-22.
488

489 Croft, P., Schnerring, F. (2009). Guide to Traffic Management. Part 13: Road Environment Safety.
490 Sydney, Australia: Austroads.
491

492 DaCoTA (2012) Roads, Deliverable 4.8q of the EC FP7 project DaCoTA
493

494 De Oña, J., López, G., Abellán, J. (2013a). Extracting Decision Rules from police accident reports
495 through Decision Trees. Accident Analysis and Prevention, 50, pp. 1151–1160.
496

497 De Oña, J., López, G., Mujalli, R. O., Calvo, F. J. (2013b). Analysis of traffic accidents on rural
498 highways using Latent Class Clustering and Bayesian Networks. Accident Analysis and Prevention,
499 51, pp. 1–10.
500

501 De Oña, J., Mujalli, R. O., Calvo, F. J. (2011). Analysis of traffic accident injury severity on Spanish
502 rural highways using Bayesian networks. Accident Analysis & Prevention, 43(1), pp. 402-411.
503

504 EC. 2008. Directive 2008/96/EC of the European Parliament and of the Council of 19 November 2008
505 on Road Infrastructure Safety Management. European Commission. Available from Internet:
506 <http://eur-lex.europa.eu/LexUriServ/LexUriServ>.
507

508 Ellis, R., Pyeon, J. (2006). Development of improved procedures for removing temporary pavement
509 markings during highway construction. Proc., Transportation Research Board 85th Annual Meeting,
510 Transportation Research Board, Washington, DC.
511

512 Kashani, A. T., Mohaymany, A. S. (2011). Analysis of the traffic injury severity on two-lane, two-way
513 rural roads based on classification tree models. *Safety science*, 49, pp. 1314-1320.
514

515 Kuhnert, P.M., Do, K.A., McClure, R. (2000). Combining non-parametric models with logistic
516 regression: an application to motor vehicle injury data. *Computational Statistics & Data Analysis* 34,
517 pp. 371–386.
518

519 López, G., Abellán, J., Montella, A., de Oña, J. (2014). Patterns of Single-Vehicle Crashes on Two-
520 Lane Rural Highways in Granada Province, Spain: In-Depth Analysis Through Decision Rules.
521 *Transportation Research Record: Journal of the Transportation Research Board*, 2432, pp. 133-141.
522

523 Miller, T.R. Benefit–cost analysis of lane marking. In *Transportation Research Record*, No. 1334,
524 TRB, National Research Council, Washington, D.C., 1992, pp. 38- 45.
525

526 Ministerio del Interior (2013). Las principales cifras de la Siniestralidad Vial en España 2013.
527 Dirección General de Tráfico, Madrid. Disponible en: [http://www.dgt.es/Galerias/seguridad-](http://www.dgt.es/Galerias/seguridad-vial/estadisticas-e-indicadores/publicaciones/principales-cifras-siniestralidad/Siniestralidad_Vial_2013.pdf)
528 [vial/estadisticas-e-indicadores/publicaciones/principales-cifras-](http://www.dgt.es/Galerias/seguridad-vial/estadisticas-e-indicadores/publicaciones/principales-cifras-siniestralidad/Siniestralidad_Vial_2013.pdf)
529 [siniestralidad/Siniestralidad_Vial_2013.pdf](http://www.dgt.es/Galerias/seguridad-vial/estadisticas-e-indicadores/publicaciones/principales-cifras-siniestralidad/Siniestralidad_Vial_2013.pdf)
530

531 Montella, A., Aria, M., D'Ambrosio, A., Mauriello, F. (2012). Analysis of powered two-wheeler
532 crashes in Italy by classification trees and rules discovery. *Accident Analysis and Prevention*, 49, pp.
533 58–72.

534

535 Moore, D.N., Schneider IV, W.H., Savolainenb, P.T., Farzaneh, M. (2010). Mixed logit analysis of
536 bicyclist injury severity resulting from motor vehicle crashes at intersection and nonintersection
537 locations. *Accident Analysis and Prevention*, *Accident Analysis and Prevention* 43(3), pp. 621-630.

538

539 Ogden, KW 1996, *Safer roads: a guide to road safety engineering*, Avebury Technical, Aldershot, UK

540

541 Pakgohar, A., Tabrizi, R.S., Khalilli, M., Esmaeili, A. (2010). The role of human factor in incidence
542 and severity of road crashes based on the CART and LR regression: a data mining approach. *Procedia*
543 *Computer Science* 3, pp. 764–769.

544

545 Regan, M.A., Trent, W.V., Lee, J.D., Young, K.L. (2009) Driver distraction injury prevention
546 countermeasures – Part 3: Vehicle, technology and road design. In: Regan, M.A., Lee, J.D., Young,
547 K.L. (Ed.), *Driver Distraction: Theory, Effects and Mitigation*, pp. 579-601. CRC Press, Boca Raton,
548 Florida.

549

550 SETRA. 2008. *Methodological Guide: Road Safety Inspections*. Service d'étudessur les transports, les
551 routes et leursa menagements (SETRA). 70 p. Available from Internet: [http://](http://www.setra.equipement.gouv.fr/IMG/pdf/US_0928A_RSI_-_Road_Safety_Inspections.pdf)
552 www.setra.equipement.gouv.fr/IMG/pdf/US_0928A_RSI_-_Road_Safety_Inspections.pdf

553

554 Siskind, V., Steinhardt, D., Sheehan, M., O'Connor, T., Hanks, H. (2011). Risk factors for fatal crashes
555 in rural Australia. *Accident Analysis & Prevention*, 43(3), pp. 1082-1088.

556

557 Sohn, S.Y., Shin, H.W. (2001). Data mining for road traffic accident type classification. *Ergonomics*
558 44, pp. 107–117.
559
560 Webb, G. I. (2007). Discovering significant patterns. *Machine g significant patterns. Machine*
561 *Learning*, 68(1), 1-33 pp.
562

563 **List of Tables:**

564

565 Table 1. Description, values and codes of the geometric and environmental road characteristics

566

567 Table 2. Description, values and codes of the RSD-SE

568

569 Table 3. Description of the rules extracted

570

571 Table 4. Final rule set

572

573 **List of Figures:**

574

575 Figure 1. Examples of road safety deficiencies (RSD)