

Intelligent Decision Support to Determine the Best Sensory Guardrail Locations

Noelia Rico, Irene Díaz, José R. Villar*, Enrique de la Cal
Computer Science Department, University of Oviedo, Oviedo, Spain

Abstract

Introducing intelligent safety devices in the roads would lead to enhance both the time of reaction and the traffic safety. Nevertheless, these intelligent devices are expensive, so choosing their location should be done carefully. This research is focused a decision support system to decide the placement of a specific safety device designed in a research project. This approach includes a feature selection stage, a model learning stage and the deployment stage. Decision models learn from real datasets with information related with accidents, classifying the samples as Fatal, Severe or Slight injury. Also, a case based risk index is proposed, so samples within the same label can be sorted. Therefore, in the deployment stage, each possible location is ranked and the user gets a feedback of the suitability of each of them to be considered for placing the intelligent safety device. The experimentation shows the proposal is valid provided the dataset for training includes enough granularity. However, it is shown that specific risk index should be designed for each road type and fork.

Keywords: Road Safety, Machine Learning, Sensory Guardrails

1. Introduction

2 This study focuses on the design of a decision support tool for the assess-
3 ment where safety intelligent road barriers should be located. These intelli-

*Corresponding author: José R. Villar, villarjose@uniovi.es
Email addresses: U0230790@uniovi.es (Noelia Rico), sirene@uniovi.es (Irene Díaz), villarjose@uniovi.es (José R. Villar*), delacal@uniovi.es (Enrique de la Cal)



Figure 1: An impact attenuator. These type of guardrail is placed at diversions in relatively high speed roads and motorways. *Image source: <http://www.hiasa.com/es/cargarFichaProducto.do?identificador=79>*

4 gent road barriers include sensory systems, on-crash video recording among
5 them, so they are complex and expensive; its location must be decided with
6 some sensibility. The focused road barriers -Fig. 1- are specifically designed
7 for intersections and diversions, and are usually known as impact attenua-
8 tors. Nevertheless, the solution to the problem was focused on a more generic
9 way, aiming to develop a tool that could propose the most interesting spots
10 for each type of road geometry.

11 The design of roads and the location of safety devices are two important
12 factors in reducing the accident rate in the network. As stated in [1, 2], a
13 safety focused road design leads to better performances in terms of traffic
14 security without high impact in the total budget. Concerning the safety
15 devices, the problem of deciding the correct spot to locate them is two-fold:
16 there are decisions to be made at design time and at the road exploiting time
17 [3, 4], when accident blackspots are found. In this latter case, the decision
18 is made based on the available traffic data; therefore, the examples of traffic
19 and crash data should be ranked and, somehow, a decision must be draw.
20 For sure, all these interventions and the decisions need to be evaluated [5, 6].

21 Ranking indexes have been widely use either in the design of guardrails
22 and their evaluation [7], in the design of road segments, the identification of
23 accident blackspots [8] or in the selection of road safety decisions. Some stud-
24 ies focus on the maintenance point of view, that is, the data from the road
25 maintenance companies; some others based their studies on crash experimen-
26 tal data, while other studies focuses on the data available from the National
27 Traffic Agency (NTA). Therefore, the ranking indexes vary enormously ac-
28 cording to the specific sub-field of interest. For instance, data coming from
29 crash experiments include physical variables, which are the basis for the in-

30 dexes [7]: the occupant impact velocity, the acceleration severity index or
31 the maximum change in the vehicle velocity, among others. These indexes
32 could be reused in other areas if there is mapping from other measures to
33 those physical variables. As an example, in those cases where the driver was
34 able to break marks on the asphalt can be used to estimate those physical
35 variables. However, these marks are not always measurable; therefore, this
36 type of indexes cannot be used as general measurements of the severity of an
37 accident.

38 The research on suitable ranking indexes shows that the most suitable
39 rank measurement directly relies on the available data. And each collection
40 may or not include a severity index. For instance, the KABCO injury sever-
41 ity scale is used by the police in the United States to classify each accident
42 as one of the five defined labels since 1966 [9]. However, the interpretation
43 of each label varies from one state to another [10]. Different indexes based
44 on the accident records were proposed in [5]. Chen et al. [11] proposed the
45 Road Safety Risk Index as a merged index of qualitative and quantitative
46 variables, including data from either the travellers, the roads, the vehicles,
47 the environment, the traffic fines, and the traffic accidents. These variables
48 were aggregated to obtain a single value that somehow reflects the risk of each
49 road; the entropy was proposed for determining the weights of each aggrega-
50 tion. Therefore, there are almost as many risk indexes different approaches
51 as contributions to the literature; each of them is valid for the specific focused
52 problem.

53 Data coming from NTA have been analyzed in several studies. For in-
54 stance, [12] proposed to cluster the data from the Indian NTA and to use
55 Association Rule Mining. This study shows the typical rules that can be
56 extracted due to the inherent categorical type of the provided variables and
57 its relatively reduced granularity and cardinality. What the studies in the
58 literature clearly remark is the need of more detailed data, so better data
59 mining can be performed on the given collections [8, 12]. Additionally, the
60 merge of datasets coming from the NTA with the datasets coming from other
61 sources -like the road maintenance agencies- can improve the benefits of ap-
62 plying these methods. Similar studies can be analysed from the literature
63 [13, 14].

64 Al-Badairi et al analysed the relationships between the input features of
65 the dataset from the Oregon Department of Traffic and the casualties of the
66 big truck crashes in run-off the road accidents [15]. They proposed ordered
67 probit models to try to find out the relevant knowledge. Basically, the main

68 factors were mainly due to the familiarity of the drivers with the network, but
69 no real hint concerning safety devices were found. Similarly, Mussone et al
70 studied the relationships among the variables in the urban traffic in the city
71 of Turin (Italy) [16] using traffic and weather data. After a pre-processing
72 and feature selection process based on correlations and SOM clustering, PCA
73 was performed to extract up to 8 features to a cumulative representation of
74 93%; the most linked original features were then chosen as the input variables
75 for the second stage. This second step made use of Artificial Neural Networks
76 to predict the level of severity of an accident.

77 From both analysis -design of ranking indexes and learning and deploying
78 the decision models-, it is clear that the most important thing is to obtain
79 high quality data. As stated in [17], 'The most serious data quality issues
80 appear to be: inaccuracies in crash location and time, difficulties in data
81 linkage (e.g. with traffic data) due to inconsistencies in databases, severity
82 misclassification, inaccuracies and incompleteness of involved users demo-
83 graphics and inaccurate identification of crash contributory factors.' With
84 all these lacks in the available data, obtaining suitable decision models to
85 locate safety devices becomes a real challenge.

86 This work was initially inspired in [13], where the Andalusia Region
87 database of Susceptible Elements of Improvement were used together with
88 the database of accidents from the Spanish Traffic National Agency (DGT
89 [18], Dirección General de Tráfico) to extract rules in order to find relation-
90 ships between the crashes and those Susceptible Elements of Improvement.
91 The authors proposed several machine learning methods for that task. How-
92 ever, no results were reported, mainly because there are no correspondences
93 between the data from both databases.

94 This research tackles with a decision support system to assess the loca-
95 tion of intelligent barriers and safety devices. A method is proposed to learn
96 models able to classify and sort the traffic locations according to the pre-
97 dicted risk of accidents and their severity if data is provided with the enough
98 granularity. The study proposes a generalized solution for the problem, al-
99 though it has to be adapted to the specific data available for each case. To
100 our knowledge, this is the first hybridized solution for this type of problems.

101 Therefore, the safety devices -the impact attenuators among them- loca-
102 tions can be sorted, so an assessment in the maintenance investment can be
103 delivered. To do so, a four stage process is detailed, including the database
104 generation, the modeling stage, the ranking of the locations stage and the
105 final deployment. The current implementation has been developed for the

106 available datasets from the DGT. Results show that, despite having datasets
107 with a very poor granularity, the models can perform properly and that the
108 suggested rank function fits to those constrains.

109 This paper is organized as follows. Next section details the four stage
110 method, including the theoretical backgrounds and the specific adaptations
111 to the available data. In Section 3 details the experimentation and the pa-
112 rameter setting, discussing the obtained results as well. Finally, conclusions
113 are drawn.

114 2. Materials and Methods

115 This solution hybridizes two different artificial intelligent techniques. On
116 the one hand, a C5.0 decision tree is performed to label new location can-
117 didates with the more suitable crash severity label. On the other hand, a
118 risk index, computed using a retrieval and a reuse stages from Case Based
119 Reasoning (CBR), is proposed to sort the location candidates as a function
120 of the similar spots in the historical database.

121 Besides, the performance of the final solution completely relies on the
122 data quality. Thus, the first subsection describes the available dataset and
123 its main features. The design of the classifier is detailed in subsection 2.2,
124 followed by the CBR based index computation. Finally, the integration of
125 both solutions is explained.

126 2.1. A description of the DGT dataset

127 From the study published in [13], we contacted the DGT, which is the
128 Spanish Agency for guaranteeing the traffic security. The DGT publishes
129 every year a report of the accidents in the Spanish roads; these data sets
130 are publicly available at [18]. In this work we have considered the data of
131 accidents in Spain from 2008 to 2013 (both included), Table 1 shows the
132 number of accidents per year.

Year	2008	2009	2010	2011	2012	2013
Number of Accidents	93161	88251	85503	83027	83115	89519

Table 1: Accidents per year, from the DGT data set.

133 There are 36 features included in the main DGT dataset, plus 10 features
134 related with the vehicles involved in the crash and 26 features related with

135 the injuries. The main part of these features are categorical, that is, there
136 are just a finite set of values allowed for each of them. The road location
137 is barely identified; although there keep a field for the council, it is rarely
138 filled: the location of the crash is identified by the road and the province
139 in the vast majority of the examples. Of course, there are also numerical
140 values: the number of vehicles involved or the number of fatal injuries are
141 clear examples. Table 2 describes the most important variables.

Table 2: Description of the variables in the datasets.

Variable	Values
Hour	1 h periods from a day
Week day	Monday(1) to Sunday(7)
Province	52 Spanish provinces
Region	18 Spanish regions
Area	Road(1), Urban Area (2), Side Street (3), Detour (4)
Grouped area	Intercity road(1), City road (2)
Road	Road identifier
Road owner	National(1), Regional(2), Provincial(3), Municipal(4), Other(5)
Road kind	Motorway(1), Highway(2), Motored-vehicle road(3) Road with slow lane(4), Road without slow lane(5) Byway(5), Side road(7), Road fork(8), other(9)
Road elements	No available data(0), Nothing remarkable(1), Zebra-cross or island(3), Middle road island(4), Central stop lane(5), Left-turn traffic circle(6), Other (7)
Priority	No available data(0), Traffic officer(1), Traffic Lights(2), STOP sign(3), GIVE WAY sign(4), Road markings only(5), Zebra-crossing sign(6), Other signalization(7), Other(8)
Roadbed Conditions	Dry and Clean(1), Shaded(2), Damped(3), Frozen(4), Snowed(5), Muddy(6), Loose gravel(7), Oily (8), Other(9)
Road Stretch	Straight line(1), Soft curves(2), Strong curves without traffic signs(3)
Luminosity	Day light(1), Twilight(2), Night: good lighting(3), Night: bad lighting (4), Night: No lighting (5)
Traffic volume	Low(1), High(2), Congested(3), None(4)

Continued on next page

Table 2 – continued from previous page

Variable	Values
Special action	Reversible lane(1), Roadside set up(2), Other measure(3), None (4)
Sidewalk	No, Yes
Additional context	14 different contexts
Atmospheric conditions	Good weather(1), Dripping fog(2), Fog(3), Rain(4), Strong rain(5), Snowing(6), Strong wind(7)
Restricted visibility	Buildings(1), Layout of nature(2), Vegetation(3), Blinding(5), Dust or smoke(6), A different constraint(7), No restriction(8)

142

143 After analyzing the data, it was found several irrelevant variables with
 144 respect to the accident rate or the security device needed. Nevertheless, the
 145 worst result was that the dataset itself was rather incomplete and suffering
 146 from granularity. Consequently, the dataset was filtered, keeping only those
 147 features for which there was no evidence of being irrelevant.

148 As this research is focused on learning the severity of a location, the
 149 locations provided by variable *Road fork kind* are considered. Thus, the goal
 150 of this learning process is to study the severity for each road fork identified
 151 by the variable *Road fork kind* (See Table 3).

Road fork type
T or Y shape
X or + shape
Acceleration lane
Diverting lane
Roundabout

Table 3: Different road forks considered in the dataset.

152 To analyze the severity associated to each road fork we must identify the
 153 variable associated to severity. According to [19], it is possible to classify
 154 accidents as *slight injury*, *serious injury* and *fatal injury*. Fatal injury in-
 155 cludes the cases where death occurs in less than 30 days as a result of the
 156 accident. Serious injuries are those where either immediate or later detention

157 in hospital as an in-patient, was required. Data provided by DGT includes
158 three variables related to accident severity: *Number of deaths (D)*, *Number*
159 *of serious injured (SE)*, *Number of slightly injured (SI)*.

160 Using these three variables it is possible to construct a new variable called
161 *Accident severity* in the following way.

- 162 • Fatality: $D > 0$.
- 163 • Serious injury: $SE > 0$ and $D = 0$.
- 164 • Slight injury: $SI > 0$, $SE = 0$ and $D = 0$

165 Accident severity restricted to each road fork, which takes the values
166 {Fatality, Serious injury, Slight injury}, is the goal of this research.

167 2.2. Learning the severity of a location

168 At a first sight, the problem seemed to be easily addressable using frequent
169 patterns and association rule mining (ASM). Therefore, the initial stage of
170 this phase was performing ASM on the data [20]. This method produced a
171 fairly vast amount of rules, typical for the ASM, that need further filtering
172 and processing. However, either no suitable method for filtering and finding
173 meaningful rules was found, or the set of rules included only elemental ones.
174 Consequently, an alternative based on machine learning is proposed.

175 The strategy followed to extract knowledge from the accident databases
176 is based on first studying the main factors affecting accident severity at an
177 intersection year by year. In a second step, the rules associated to each
178 different intersection across the years are mined in order to obtain the most
179 frequent rule sets using a voting strategy across years. The basic algorithm
180 is described below.

```
181 For each road fork
182     For each year from 2008 to 2013
183         Return the best classifier
184     End for
185     Return the most frequent rule set per road fork
186 End for
```


187 *2.2.1. Classifier selection*

188 The prediction of accident severity on crosses and intersections is modeled
189 in this problem as a Machine Learning (ML) problem. As the final goal of the
190 work is to provide an understandable output to be used by traffic analysts,
191 we focus our attention on approaches based on tree models, due to their
192 performance and great interpretability. At this concern, tree based models
193 has been chosen to build classification trees.

194 In general, tree-based ML models and algorithms work in the following
195 procedure [21]. They grow a tree using forward selection using a top-down
196 approach from root to leaves until some stopping condition is reached. At
197 each step they find the best split according to some impurity measure. The
198 node associated to the maximal impurity reduction is then selected. Finally,
199 most method prune the tree back and obtain a rule per path from the root
200 to each leaf.

201 Different combinations of metrics, splits, stopping conditions and pruning
202 methods lead to different approaches. In this work, considering both inter-
203 pretability and performance, C5.0 ([22]), recursive partitioning (PART) and
204 Random Forests ([23] are selected as classifiers. The performance of these
205 classifiers is measured in terms of the well-known measures Precision, Recall,
206 F_1 and Accuracy [24]. Precision is the fraction of relevant examples among
207 the retrieved instances. Recall is the fraction of relevant instances that have
208 been retrieved over the total amount of relevant instances. F-score is the
209 harmonic mean of Precision and Recall. It is a quite common measure used
210 to weight the existing trade-off between Precision and Recall.

211 *2.3. Measuring road hazardousness*

212 Several different indexes have been proposed in the literature to rank
213 the hazardousness of either roads or current driving conditions, etc. As an
214 example, [25] proposed a ranking index of the segments of road in order to
215 alert drivers about arising contingencies in the traffic or in the road.

216 The problem of designing a ranking index was analyzed in [11], finding
217 that they are problem specific and also data specific. This means that de-
218 signing a generalizable ranking index is a real challenge, which should be
219 kept open to fit the specificity of the problem and the data available. There
220 are several reasons for this challenge. There is no clear relationship among
221 the different factors involved in road hazardous indexes; for instance, it has
222 been found that there is no clear relationship between the crash frequency

223 and the traffic flow [26]. Furthermore, there are differences in the relation-
224 ships when data from different countries are analyzed [27]. Therefore, each
225 problem would lead to the most suited set of factors and their corresponding
226 aggregation [28].

227 Examples of such indexes are those propose in [29, 13]. In [29], three dif-
228 ferent indexes are proposed to evaluate the design and maintenance strategies
229 of roads. One index is focused on density of crashes on each road segments,
230 a second index focuses on the number of injuries in the road sections, while a
231 third one deals with the number of fatalities occurred in the segments. These
232 indexes were estimated for the design of the roads; however, they can also
233 be obtained from historical data. With the density of crashes, injuries and
234 fatalities for each section, and the length of the segment, the three indexes
235 are computed for each road segment and -by means of aggregation- for the
236 complete road.

237 Similarly, Martin et al. [13] proposed two indexes -one referred to the
238 number of crashes and the second referred to the number of injuries and
239 fatalities- computed using historical data from i) 5 previous years, ii) last
240 two years, iii) two years back and iv) last year. With all these indexes, the
241 authors proposed a set of four rules to classify a road section as hazardous
242 or not.

243 In this study a different approach is proposed considering not the histor-
244 ical data of the segment but from similar segments. The idea underneath is
245 to apply similar retrieval concepts than those used in Case Based Reasoning
246 (CBR) [30] to find those road segments that better match the current posi-
247 tion, using these cases as the historical data for computing the hazardousness
248 index. CBR typically includes four stages: Retrieval, Reuse, Revise and Re-
249 tain. The retrieval is concerned with finding the most similar cases, assigning
250 a similarity degree to each one. Reuse deals with the selection of the most
251 interesting retrieved cases, even partially, to generate a new outcome. Revise
252 makes use of reasoning to generate completely new proposals based on the
253 retrieved and chosen information. Finally, Retain estimates when a new gen-
254 erated case is found worth to be included in the case base for future reuse.
255 For the extend of this study, only the first two stages are needed.

256 Each segment must be described not only using the geometry design in-
257 formation -like the slope, cant, radius, gps location, etc.- but also including
258 the environmental factors -tar type, surrounding nature, etc.-. All these fac-
259 tors must be assessed, that is, the implication of each factor in the similarity
260 of the cases is fundamental to obtain good results. This assessment can be

261 implemented by means of weights, whose values can be fixed a priori [29, 13]
262 or using any other method [11]. Actually, if variables related to the envi-
263 ronmental factors are introduced then it may be interesting to cluster these
264 factors and to define different set of weights to each cluster.

265 Unfortunately, only the DGT statistical crash data collection was avail-
266 able for this research. As mentioned before, this dataset has a very poor
267 granularity, and there were no information concerning with the factors men-
268 tioned above. The granularity was given in terms of road identification,
269 crashing council and the type of segment where the crash took place -straight
270 segment, Y or X intersection, etc.- Therefore, for experimentation purposes,
271 similarity of the cases were restricted to these variables, assigning them the
272 same weight. The reason for this weight selection is the poor granularity of
273 the DGT dataset; a candidate is compared to all the accident cases for the
274 same road, council and segment type. Nevertheless, these variables do not
275 tell much about the similarities in the context of these spots, so there is no
276 much reason to think a better set of weights can be set.

277 Once the similar segments have been retrieved, then they are reused.
278 Reusing cases means considering the factors that are found relevant in the
279 severity of a crash. These factors are the number of fatalities (fn), the
280 number of severe injured individuals (si), the number of lightly injured in-
281 dividuals (li) and the number of involved vehicles (iv); all these factors are
282 integer numbers. All of them have to be aggregated in order to obtain a
283 single scalar ranking index, so we propose to scale each factor to the interval
284 $[0.0, 1.0]$ and then to aggregated them using a weighted sum.

285 Each factor is lineally scaled from 0 to an upper limit, corresponding from
286 0.0 to 1.0. However, the factors have no real upper limit, so the scaled factors
287 may surpass the 1.0 scaled limit. To illustrate this problem, let's focus on
288 the number of fatalities; let's say the upper limit is 2, so this factor is scaled
289 0 to 2 to the interval $[0.0, 1.0]$. Sadly, in an accident there can be more than
290 two fatalities, meaning the scaled factor surpasses the value of 1.0.

291 The weights, for this research, have been manually chosen according to
292 the relevance we think each one has. For instance, the number of fatalities has
293 the higher weight ($w_{fn} = 0.4$), then the number of severe injuries ($w_{si} = 0.3$);
294 the remaining weights are $w_{li} = 0.1$ and $w_{iv} = 0.2$ for the light injuries and
295 the number of involved vehicles, respectively. The weight selection must be
296 carefully defined according to the problem to solve and the data quality.

297 The index is computed in three steps. Firstly, compute the risk index for
298 each of the retrieved cases; let NC be the number of retrieved cases. Each

299 case includes four mentioned factors. These factors represent the information
300 the dataset from the DGT includes for each crash sample. Each of these
301 feature values is mapped to the $[0.0, 1.0]$ interval using linear functions, a 0.0
302 in each feature value was assigned to a 0.0 in the mapping. An upper limit
303 (UL) was fixed to each feature, so this UL was mapped to a 1.0. The UL for
304 fn , si , li and iv were set to $UL_{fn} = 2$, $UL_{si} = 2$, $UL_{li} = 5$ and $UL_{iv} = 3.0$.
305 In the case of iv an extra linearity segment was introduced, so a value of 2.0
306 for iv was mapped to 0.75. Finally, an aggregation of the scaled factors is
307 obtained as a weighted sum with a-priori fixed weights.

308 Secondly, the maximum (MAX) and the mean (MN) values of this
309 index among all the retrieved cases should be calculated. Thirdly, in order
310 to shift the final index to the worst scenario, the average of MAX and MN
311 is powered to $\frac{1}{NC}$. This last computed value is the risk index of the current
312 sample based on the most similar cases in the database.

313 2.4. Ranking the locations and deployment

314 Finally, the two approaches need to be merged so the analyzed locations
315 could be sorted. In this approach we perform the classification and the
316 ranking independently. Then, the calculated risk index is assigned to each
317 classified sample. Finally, the examples are sorted by the class severity and
318 then by the risk index.

319 3. Experiments and Results

320 3.1. Selection of the classifier and parameter tuning

321 Before starting with parameter tuning, it is necessary to focus on the
322 distribution of *Accident severity* variable. As Figure 2 shows, this variable is
323 extremely imbalanced. In fact the number of slight injuries is almost 9 times
324 of the total amount of serious injuries and this is 10 times the number of fatal
325 accidents (independently on the year and road fork kind considered). That
326 makes the problem difficult to solve so that different re-sampling as well as
327 learning strategies have been tested.

328 Thus, the default configuration of the algorithms selected in Section 2.2.1
329 was considered to select the most suitable learning strategy among multi-
330 category, one to all and one to one learning strategy. In addition, given
331 that the dataset is quite unbalanced, resampling must be considered. Under-
332 sampling, Over-sampling, SMOTE and ROSE strategies were also tested
333 and compared to no resampling. Thus, each classifier was trained using cross



Figure 2: Box-plot graph for the different classes of accident severity

334 validation with k folds for each learning strategy and resampling technique,
 335 obtaining that the best combination in terms of F_1 is one to one as learning
 336 strategy and ROSE ([31]) as resampling strategy. ROSE provides a unified
 337 framework to deal simultaneously with the problem of model estimation and
 338 accuracy evaluation in imbalanced learning. It builds on the generation of
 339 new artificial examples from the classes, according to a smoothed bootstrap
 340 approach.

341 Once the learning (one to one) and resampling (ROSE) strategies are
 342 set, each classifier is optimized via parameter tuning. In particular, the
 343 parameters analyzed are detailed below.

- 344 • Random forests: Number of predictors sampled for splitting at each
 345 node (ranging from 1 to 6).
- 346 • C5.0. In this case several parameters have been studied: Feature
 347 selection, number of boosting iterations (trials), ranging in the set
 348 $\{1, 2, 5, 10\}$ and output model (tree or a rule set).
- 349 • Recursive partitioning: Complexity parameter (ranging from 1 to 10).

350 Again, the learning process is tested using cross validation with 10 folds.
 351 R version 3.3.3 as well as the caret package [32] were the tools used to perform
 352 these experiments.

353 Table 4 shows the method that obtains the best F-score. Each value
 354 represents the percentage of experiments the method obtained the best F-
 355 score. As it can be seen, C5.0 is the method performing better most of times.
 356 Thus, it is the one selected for predicting accident severity.

Road Fork			
	Serious vs Slight	Fatality vs Slight	Fatality vs Serious
T or Y	C5.0 (100%)	C5.0 (100%)	C5.0 (100%)
X or +	C5.0 (100%)	C5.0 (100%)	C5.0 (100%)
Acceleration lane	C5.0 (100%)	rf (50%)	rf (68%)
		C5.0 (33%)	C5.0 (16%)
		rpart (17%)	rpart (16%)
Diverting lane	C5.0 (100%)	C5.0 (50%)	rpart (50%)
		rf (33%)	rf (33%)
		rpart (17%)	C5.0 (17%)
Roundabout	C5.0 (100%)	C5.0 (50%)	C5.0 (100%)
		rf (33%)	
		rpart (17%)	

Table 4: Methods obtaining the best performance for each combination of road fork and crash severity. In brackets the percentage of times that each method was the best

357 3.2. Training and validation of the proposal

358 Thus, C5.0 is selected as base classifier to predict accident severity. Ac-
 359 cording to preliminary experiments detailed in Section 3.3, the best configu-
 360 ration for C5.0 is winnowing and a tree based structure instead of rules. The
 361 number of boosting iterations depends on the data set. The performance of
 362 C5.0 measured in terms of F-score is shown in Table 5. F-score ranges from
 363 0.62 to 0.8. As it can be seen, the algorithm performance is quite similar
 364 across years and road forks kinds. F-score is about 0.7 in average. Although
 365 it is not so bad, it is a challenge to improve classifier efficiency.

366 When studying accident severity at a certain road fork, the decision trees
 367 provided by C5.0 are different across years as the data sets are obviously
 368 different. However, it is interesting to check if there is any consistent rule
 369 to predict accident severity across the years. To do that, tree models are
 370 translated into rules, obtaining 6 different rule bases for each road fork. Thus,
 371 to obtain the final rule base, only rules occurring half of years are considered.
 372 Table 6 shows the number of rules obtained by road fork and accident severity.

Year	T or Y	X or +	Accelerating lane	Diverting lane	Roundabout
2008	0.67	0.71	0.67	0.70	0.69
2009	0.67	0.74	0.73	0.73	0.73
2010	0.68	0.74	0.70	0.62	0.72
2011	0.66	0.71	0.74	0.65	0.75
2012	0.64	0.75	0.72	0.70	0.73
2013	0.65	0.76	0.80	0.65	0.77

Table 5: F-score obtained when accident severity is predicted with C5.0

373 The rule base associated to accident severity in roundabouts is shown in
374 Table 7 as example. As it can be seen, accident severity depends on priority
375 signs, existence of sidewalks, road luminosity, traffic volume and other road
376 properties.

Road fork type	Fatal	Serious injury	Slight injury
T or Y shape	19	2	5
X or + shape	6	4	4
Acceleration lane	11	2	2
Diverting lane	1	0	6
Roundabout	6	6	7

Table 6: Number of rules occurring during at least 3 years.

377 3.3. Measuring risk

378 As the final goal of this research is to predict where to place sensory
379 Guardrail Locations, it is necessary to establish a priority among the different
380 places or road profiles. Thus, once a test example is classified as *fatal*, *serious*
381 *injury* or *slight injury* the risk index is computed and the examples are
382 ranked according to it. As the risk factor is computed over the data set,
383 we have checked the risk index for the different kinds of road forks studied
384 in the classification problem. The obtained factor risks are shown in Figure
385 3 where each risk value is the average of the factor risks obtained during the
386 different years. As it can be seen, the riskiest forks are roundabouts while
387 acceleration lanes are the least ones. Note that this is a measure independent
388 of the accident severity classification.

<i>Factors for predicting Slight injuries</i>
{Roadbed conditions != Dry and clean, Luminosity !=Twilight, Sidewalk=No, Traffic volume=Low}
{Road kind != Highway, Roadbed conditions = Dry and clean, Priority=Give Way sign, Sidewalk=No, Traffic volume=Low}
{Priority=STOP sign, Traffic volume=Low}
{Luminosity=Night: No lighting, Traffic volume=Low}
{Restricted visibility != No restriction, Traffic volume=Low}
{Road kind != Motored-vehicle road, Priority= Other signals, Traffic volume=Low}
{Priority != Give way sign OR Other signals, Grouped area= City road, Luminosity= Night: good lighting, Traffic volume=Low}
<i>Factors for predicting Serious injuries</i>
{Luminosity =Night: good lighting, Sidewalk=Yes, Special action=None}
{Priority=Give Way sign, Traffic volume=Low}
{Priority!=Zebra-crossing sign, Traffic volume!=Low, Sidewalk=Yes}
{Roadbed conditions != Dry and clean, Traffic volume=Low}
{Sidewalk = No, Traffic volume!=Low}
{Sidewalk = No}
<i>Factors for predicting Fatal injuries</i>
{Priority != GIVE WAY sign Luminosity !=Night: good lighting, Roadbed conditions = Dry and clean, Sidewalk=Yes}
{Luminosity!= Night: No lighting, Traffic volume!=Low}
{Priority != STOP or Give way signs or Zebra-crossing sign, Roadbed conditions = Dry and clean}
{Priority != STOP sign, Road kind=Highway, Roadbed conditions = Dry and clean, Traffic volume=Low, Sidewalk=No, Restricted visibility != No}
{Priority != STOP sign or or Zebra-crossing sign, Roadbed conditions = Dry and clean, Traffic volume=Low, Sidewalk=Yes, Restricted visibility != No}
{Traffic volume!=Low, Sidewalk=No}

Table 7: Factors for predicting accident severity for a X or + Road Fork

389 As can be seen, the proposed risk index lacks in generalization as the
390 values obtained clearly differs from one type of crossing to another. This
391 problem is present in any ad-hoc risk index because the same set of variables
392 are used in computing the risk index for all type of crash.

393 For instance, the number of fatalities in accidents in a roundabout is
394 higher because of the speed and the traffic density. However, the number

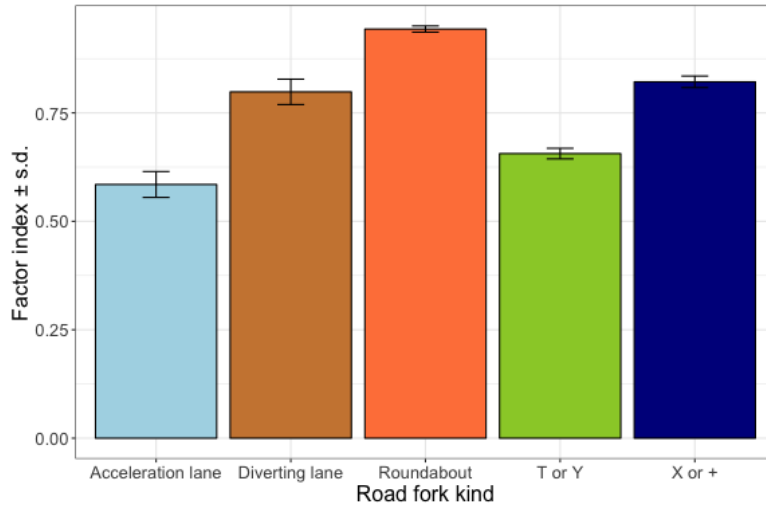


Figure 3: Factor risk by Road Fork. Each bar represents the factor risk in average associated to each road fork during the years under study (2008 to 2013)

395 of fatalities when diverting of a motorway is smaller because the front-lateral
 396 crashes might produce more injuries and less fatalities.

397 The obtained results suggest that a different risk index should be devel-
 398 oped for each type of fork and road. Therefore, more research is needed in
 399 order to discriminate the most relevant features involved in the risk index
 400 calculation for each type of fork and road.

401 4. Conclusions

402 In this research, a proposal of a decision support system to help in choos-
 403 ing the location where to place an intelligent sensory guardrail among several
 404 candidates. The decision support system suggests an order of the candidates
 405 according to a label of the accident severity and a risk factor. The intelli-
 406 gent model use to classify the candidates is learned based on historical data
 407 from a accident database. Moreover, the risk index is also obtained based on
 408 similarly retrieved cases within the historical data, in a case-based reason-
 409 ing fashion. The experimentation carried out in this research made use the
 410 public accident dataset published by the DGT. This dataset includes plenty
 411 of nominal and discrete features. A study of the relationships and the most
 412 interesting features has been performed before the process of learning models
 413 for the classification.

414 A comparison of the performance of different models, C5.0 has been found
415 the most interesting and robust model for each type of crossing type. To-
416 gether with the risk index, the proposal has been found interesting and with
417 a good performance on the data. However, one of the main problems found
418 is the high data granularity. For instance, the risk index might be in compro-
419 mise due to this reason: the available dataset does not include information
420 about the kilometer, so the similarity is not as precise as it would be needed.

421 Nevertheless, the whole solution allows the user to sort the candidates,
422 producing labeled candidates together with the available risk index. If better
423 data is provided, the procedure described in this research would allow to
424 obtain finest models and a more precise risk index. A careful selection of the
425 most suitable features should be performed in order to obtain a risk index
426 for each road and crossing types.

427 As future work we plan to perform a deeper study about risk factors and
428 study if there is any relation between it and the accident severity that it is
429 predicted by a classifier. In addition, it is also necessary to introduce more
430 information to the system for representing forks in a more accurate way.
431 On the other hand, feature selection techniques as well as other parameter
432 settings (minimum number of examples per leaf or pruning level) can be
433 deeper studied in order to improve the efficiency of the method.

434 **Acknowledgements**

435 Funding: This work was supported by the Innterconecta call of the Eu-
436 ropean Union Structural Funds (FEDER) INTERCONNECTA CDTI project
437 ABECATIM (SOL-00082271 / ITC-20151039) and by the Spanish Govern-
438 ment through the MINECO project TIN2017-87600-P.

439 **References**

- 440 [1] M. Cornelissen, P. Salmon, N. A. Stanton, R. McClure, Assessing the
441 'system' in safe systems-based road designs: using cognitive work anal-
442 ysis to evaluate intersection designs, *Accident Analysis & Prevention* 74
443 (2015) 324–338.
- 444 [2] D. J. Gabauer, X. Li, Influence of horizontally curved roadway sec-
445 tion characteristics on motorcycle-to-barrier crash frequency, *Accident*
446 *Analysis & Prevention* 77 (2015) 105–112.

- 447 [3] J. Strandroth, Validation of a method to evaluate future impact of road
448 safety interventions, a comparison between fatal passenger car crashes
449 in sweden 2000 and 2010, *Accident Analysis & Prevention* 76 (2015)
450 133–140.
- 451 [4] Q. Yu, Z. Guo, Z. Zhang, J. Wang, Assistant decision-making system for
452 road safety strategy, *Procedia-Social and Behavioral Sciences* 96 (2013)
453 320–328.
- 454 [5] O. Basile, L. Persia, Tools for assessing the safety impact of interventions
455 on road safety, *Procedia-Social and Behavioral Sciences* 53 (2012) 682–
456 691.
- 457 [6] C. Roque, F. Moura, J. L. Cardoso, Detecting unforgiving roadside con-
458 tributors through the severity analysis of ran-off-road crashes, *Accident*
459 *Analysis & Prevention* 80 (2015) 262–273.
- 460 [7] D. J. Gabauer, H. C. Gabler, Comparison of roadside crash injury
461 metrics using event data recorders, *Accident Analysis & Prevention* 40
462 (2008) 548 – 558.
- 463 [8] P. K. Agarwal, P. K. Patil, R. Mehar, A methodology for ranking road
464 safety hazardous locations using analytical hierarchy process, *Procedia-*
465 *Social and Behavioral Sciences* 104 (2013) 1030–1037.
- 466 [9] B. Burdett, Improving Accuracy of KABCO Injury Severity Assessment
467 by Law Enforcement Officers, Master’s thesis, Civil and Environmental
468 Engineering School, University of Wisconsin-Madison, 2014.
- 469 [10] F. H. A. (FHWA), Kabco injury classification scale and defi-
470 nitions, [https://safety.fhwa.dot.gov/hsip/spm/conversion_tbl/
471 pdfs/kabco_ctable_by_state.pdf](https://safety.fhwa.dot.gov/hsip/spm/conversion_tbl/pdfs/kabco_ctable_by_state.pdf), 2017.
- 472 [11] F. Chen, J. Wang, Y. Deng, Road safety risk evaluation by means of
473 improved entropy topsis–rsr, *Safety science* 79 (2015) 39–54.
- 474 [12] S. Kumar, D. Toshniwal, Analysing road accident data using association
475 rule mining, in: *Computing, Communication and Security (ICCCS)*,
476 2015 International Conference on, IEEE, pp. 1–6.

- 477 [13] L. Martín, L. Baena, L. Garach, G. López, J. de Oña, Using data mining
478 techniques to road safety improvement in spanish roads, *Procedia-Social
479 and Behavioral Sciences* 160 (2014) 607–614.
- 480 [14] S. Kumar, D. Toshniwal, A data mining framework to analyze road
481 accident data, *Journal of Big Data* 2 (2015) 26.
- 482 [15] N. S. S. Al-Bdairi, S. Hernandez, An empirical analysis of run-off-road
483 injury severity crashes involving large trucks, *Accident Analysis & Pre-
484 vention* 102 (2017) 93–100.
- 485 [16] L. Mussone, M. Bassani, P. Masci, Analysis of factors affecting the
486 severity of crashes in urban road intersections, *Accident Analysis &
487 Prevention* 103 (2017) 112–122.
- 488 [17] M. Imprialou, M. Quddus, Crash data quality for road safety research:
489 current state and future directions, *Accident Analysis & Prevention*
490 (2017).
- 491 [18] D. G. de Tráfico, Dirección general de tráfico, portal estadístico,
492 accidentes 2015, https://sedeapl.dgt.gob.es/WEB_IEST_CONSULTA/,
493 2017.
- 494 [19] P. Michalaki, M. A. Quddus, D. Pitfield, A. Huetson, Exploring the
495 factors affecting motorway accident severity in england using the gener-
496 alised ordered logistic regression model, *Journal of Safety Research* 55
497 (2015) 89 – 97.
- 498 [20] R. Agrawal, R. Srikant, Fast algorithms for mining association rules in
499 large databases, in: *Proceedings of the 20th International Conference
500 on Very Large Data Bases VLDB '94*, pp. 487–499.
- 501 [21] I. H. Witten, E. Frank, M. A. Hall, *Data Mining: Practical Machine
502 Learning Tools and Techniques*, Morgan Kaufmann Publishers Inc., San
503 Francisco, CA, USA, 3rd edition, 2011.
- 504 [22] R. J. Quinlan, *Data Mining Tools See5 and C5.0*, 2000.
- 505 [23] L. Breiman, Random forests, *Machine Learning* 45 (2001) 5–32.

- 506 [24] I. Díaz, J. Ranilla, E. Montañés, J. Fernández, E. F. Combarro, Im-
507 proving performance of text categorisation by combining filtering and
508 support vector machines, *Journal of the American Society for Informa-
509 tion Science and Technology* 55 (2004) 579–592.
- 510 [25] V. Rosolino, I. Teresa, A. Vittorio, F. D. Carmine, T. Antonio,
511 R. Daniele, Z. Claudio, Road safety performance assessment: a new
512 road network risk index for info mobility, *Procedia-Social and Behav-
513 ioral Sciences* 111 (2014) 624–633.
- 514 [26] C. Roque, J. L. Cardoso, Investigating the relationship between run-
515 off-the-road crash frequency and traffic flow through different functional
516 forms, *Accident Analysis & Prevention* 63 (2014) 121–132.
- 517 [27] C. Roque, J. L. Cardoso, Safeside: a computer-aided procedure for
518 integrating benefits and costs in roadside safety intervention decision
519 making, *Safety science* 74 (2015) 195–205.
- 520 [28] F. Russo, M. Busiello, G. Dell’Acqua, Safety performance functions for
521 crash severity on undivided rural roads, *Accident Analysis & Prevention*
522 93 (2016) 75–91.
- 523 [29] K. Jamroz, M. Budzyński, W. Kustra, L. Michalski, S. Gaca, Tools for
524 road infrastructure safety management–polish experiences, *Transporta-
525 tion Research Procedia* 3 (2014) 730–739.
- 526 [30] J. Kolodner, *Case-based reasoning*, Morgan Kaufmann, 2014.
- 527 [31] G. Menardi, N. Torelli, Training and assessing classification rules with
528 imbalanced data, *Data Mining and Knowledge Discovery* 28 (2014) 92–
529 122.
- 530 [32] M. Kuhn, The caret package, [http://topepo.github.io/caret/
531 index.html](http://topepo.github.io/caret/index.html), 2017.